

2-27-2019

## Kaggle and Click-Through Rate Prediction

Todd W. Neller  
*Gettysburg College*

Follow this and additional works at: <https://cupola.gettysburg.edu/csfac>



Part of the [Databases and Information Systems Commons](#)

**Share feedback** about the accessibility of this item.

---

### Recommended Citation

Neller, Todd W. "Kaggle and Click-Through Rate Prediction." Presented at the Gettysburg College Computer Science Colloquium, Gettysburg, PA, February 2019.

This is the author's version of the work. This publication appears in Gettysburg College's institutional repository by permission of the copyright owner for personal use, not for redistribution. Cupola permanent link: <https://cupola.gettysburg.edu/csfac/56>

This open access presentation is brought to you by The Cupola: Scholarship at Gettysburg College. It has been accepted for inclusion by an authorized administrator of The Cupola. For more information, please contact [cupola@gettysburg.edu](mailto:cupola@gettysburg.edu).

---

## Kaggle and Click-Through Rate Prediction

### Abstract

Neller presented a look at Kaggle.com, an online Data Science and Machine Learning learning community, as a place to seek rapid, experiential peer education for most any Data Science topic. Using the specific challenge of Click-Through Rate Prediction (CTRP), he focused on lessons learned from relevant Kaggle competitions on how to perform CTRP.

### Keywords

Data science, machine learning, Kaggle, click-through rate prediction

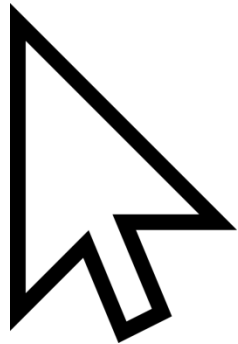
### Disciplines

Computer Sciences | Databases and Information Systems

### Comments

This presentation was given at the Gettysburg College Computer Science Colloquium on February 27th, 2019.

# Kaggle and Click-Through Rate Prediction



Dr. Todd W. Neller  
Professor of Computer Science

Gettysburg  
COLLEGE





**TruWood**

Sponsored ·



Start 2019 with a win! Access our exclusive New Year Sale today! 🌲 ❄️

Use the code NEWYEAR15 at checkout for an additional 15% OFF!

Shop Now -> [www.mytruwood.com](http://www.mytruwood.com)



**vimeo**

**Tons and tons of video storage**



**Get up to 7TB storage**

[vimeo.com](http://vimeo.com)

Your videos deserve more space. Keep rough cuts, final edits, and everything in between.



**Flintstones Vitamins**

[flintstonesvitamins.com](http://flintstonesvitamins.com)

Flintstones Complete Gummies are the #1 pediatrician recommended children's vitamin brand,...

English (US) · Español · Português (Brasil) · Français (France) · Deutsch



[Privacy](#) · [Terms](#) · [Advertising](#) · [Ad Choices](#)

[Cookies](#) · [More](#) ·

Facebook © 2019

# Click-Through Rate (CTR) Prediction

$$\text{CTR} = \frac{\text{Number of click-throughs}}{\text{Number of impressions}} \times 100(\%)$$

- Number of impressions = number of times an advertisement was served/offered
- Given: much data on past link offerings and whether or not users clicked on those links
- Predict: the probability that a current user will click on a given link



# Example Data on Past Link Offerings

- User data:
  - User ID from site login, cookie
  - User IP address, IP address location
- Link context data:
  - Site ID, page ID, prior page(s)
  - Time, date
- Link data:
  - Link ID, keywords
  - Position offered on page











# Example: Facebook Information

Access Your Information

Here is a list of your Facebook information that you can access at any time. We've categorized this information by type so you can easily find what you're looking for. Our [Data Policy](#) has more information about how we collect and use your information, how it's shared and how long we retain it. It also outlines your rights and how you can exercise them, and how we operate and transfer your information as part of our global services.

You can choose to [download your information](#) if you'd like a copy of it.

Your Information ⓘ Expand All

 <b>Posts</b> Posts you've shared on Facebook and posts you've been tagged in >	 <b>Photos and Videos</b> Photos and videos you've shared or been tagged in >
 <b>Comments</b> Comments you've posted on your own posts, on other people's posts or in groups you belong to >	 <b>Likes and Reactions</b> Posts, comments and Pages you've liked or reacted to >
 <b>Friends</b> The people you are connected to on Facebook >	 <b>Following and Followers</b> People, organizations or business you choose to see content from, and people who follow you >
 <b>Messages</b> Messages you've exchanged with other people on Messenger >	 <b>Groups</b> Groups you belong to, groups you manage and your posts and comments within the groups y... >
 <b>Events</b> Your responses to events and a list of the events you've created >	 <b>Profile Information</b> Your contact information, information you've written in your About You section in your profil... >

# Why is CTR Prediction Important?

- Advertising Industry View:
  - Much of online advertising is billed using a **pay-per-click model**.





# New Idea?



## ATTENTION IS THE CURRENCY OF THE INTERNET

Traditional Economy → Scarce Resources

Information Economy → Unlimited Resources → Limited Time

# Benefits Beyond Advertising

- [Herbert Simon](#), 1971:
  - “In an information-rich world, the wealth of information means a dearth of something else: a scarcity of whatever it is that information consumes. What information consumes is rather obvious: the **attention** of its recipients.”
- Better CTR prediction → more *relevance* → better use of scarce time



# Outline

- Click-Through Rate Prediction (CTRP) Introduction
- Kaggle
  - Learning community offerings incentives
  - CTRP Competitions
- Feature Engineering
  - Numbers, Categories, and Missing Values
- Favored regression techniques for CTRP
  - Logistic Regression
  - Gradient Boosted Decision Trees (e.g. xgBoost)
  - Field-aware Factorization Machines (FFMs)
- Future Recommendations

# What is [Kaggle.com](https://www.kaggle.com)?

- Data Science and Machine Learning Community featuring
  - [Competitions](#) → \$\$\$, peer learning, experience, portfolio
  - [Datasets](#)
  - [Kernels](#)
  - [Discussions](#)
  - [Tutorials](#) (“Courses”)
  - Etc.
- Status incentives

## 15 Active Competitions



TWO SIGMA

### Two Sigma: Using News to Predict Stock Movements

Use news analytics to predict stock price performance

**Featured** · Kernels Competition · 5 months to go · news agencies, time series, finance, money

**\$100,000**  
2,897 teams



### LANL Earthquake Prediction

Can you predict upcoming laboratory earthquakes?

**Research** · 4 months to go · earth sciences, physics, signal processing

**\$50,000**  
938 teams



### Elo Merchant Category Recommendation

Help understand customer loyalty

**Featured** · 20 days to go · tabular data, banking, regression

**\$50,000**  
3,487 teams



### Google Analytics Customer Revenue Prediction

Predict how much GStore customers will spend

**Featured** · 9 days to go · regression, tabular data

**\$45,000**  
1,104 teams



### Gendered Pronoun Resolution

Pair pronouns to their correct entities

**Research** · 3 months to go · nlp, text data

**\$25,000**  
30 teams



### PetFinder.my Adoption Prediction

How cute is that doggy in the shelter?

**\$25,000**  
972 teams

# Datasets

[Documentation](#)[New Dataset](#)

Public

Your Datasets

Favorites

Sort by **Hotness**

14,368 Datasets

Sizes



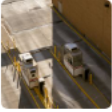
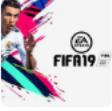

File types

Licenses

Tags

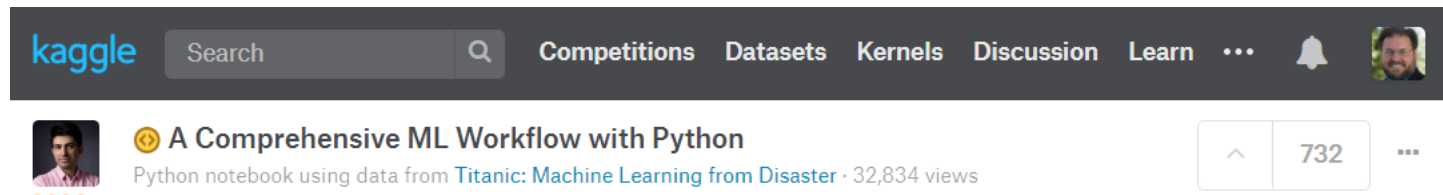
Search datasets



- |     |   |   |  |                           |                      |
|-----|---|---|--|---------------------------|----------------------|
| 150 |    | <b>Heart Disease UCI</b><br><a href="https://archive.ics.uci.edu/ml/datasets/Heart+Disease">https://archive.ics.uci.edu/ml/datasets/Heart+Disease</a><br>ronit updated 7 months ago (Version 1) | biology<br>health<br>classification<br>binary clas...          | CSV<br>3.4 KB<br>Other    | <> 55<br>2<br>24k    |
| 441 |    | <b>Graduate Admissions</b><br>Predicting admission from important parameters<br>Mohan S Acharya updated a month ago (Version 2)   | regression ...<br>model com...<br>random for...<br>+ 2 more... | CSV<br>9.4 KB<br>CC0      | <> 147<br>12<br>59k  |
| 48  |    | <b>Los Angeles Parking Citations</b><br>From Los Angeles Open Data<br>City of Los Angeles <span>Maintained by Kaggle</span> updated a day ago (Version 129)                                     | socrata  | Other<br>245.1 MB<br>ODbL | <> 20<br>0<br>13k    |
| 351 |  | <b>FIFA 19 complete player dataset</b><br>18k+ FIFA 19 players, -90 attributes extracted from the latest FIFA database<br>Karan Gadiya updated 2 months ago                                     | sports<br>data visuali...<br>regression ...<br>+ 2 more...     | CSV<br>2.1 MB<br>CC4      | <> 53<br>4<br>40k    |
| 907 |  | <b>Google Play Store Apps</b><br>Web scraped data of 10k Play Store apps for analysing the Android market.<br>Lavanya Gupta updated 3 days ago (Version 6)                                      | video games<br>computer s...<br>internet<br>mobile web         | CSV<br>1.9 MB<br>Other    | <> 191<br>22<br>208k |

# Kernels

- Jupyter notebooks of mixed text and Python/R
  - Interleaved explanations and free runnable code
- E.g. <https://www.kaggle.com/mjbahmani/a-comprehensive-ml-workflow-with-python>



The screenshot shows the Kaggle website interface. At the top, there is a navigation bar with the Kaggle logo, a search bar, and links for Competitions, Datasets, Kernels, Discussion, and Learn. Below the navigation bar, a kernel titled "A Comprehensive ML Workflow with Python" is displayed. The kernel is authored by a user with a profile picture and is marked with a gold coin icon. The description indicates it is a Python notebook using data from the Titanic dataset, with 32,834 views. A box on the right shows the kernel has 732 upvotes.

## A Comprehensive Machine Learning Workflow with Python

There are plenty of **courses and tutorials** that can help you learn machine learning from scratch but here in **Kaggle**, I want to solve **Titanic competition** a popular machine learning Dataset as a comprehensive workflow with python packages. After reading, you can use this workflow to solve other real problems and use it as a template to deal with **machine learning** problems.

last update: 06/02/2019

*You may be interested have a look at 10 Steps to Become a Data Scientist:*

# Discussions



## Display Advertising Challenge

Predict click-through rates on display ads  
\$16,000 · 718 teams · 4 years ago

Overview Data Kernels Discussion Leaderboard Rules Team

My Submissions

New Topic

81 topics

Follow

Sort by

Hotness

All

Mine

Upvoted

Search topics



13



**Document and code for the 3rd place finish**

[Guocong Song](#) 4 years ago

last comment by

[Vishal Gupta](#) 3mo ago

4

59



**3 Idiots' Solution & LIBFFM**

[Yuchin Juan](#) 4 years ago

last comment by

6mo ago

54

7



**Data Release After Competition Ends**

[joycenv](#) 4 years ago

last comment by

[Olivier Chapelle](#) 3y ago

17

9



**Congratulations to the winners!**

[Dylan Friedmann](#) 4 years ago

last comment by

[marbel](#) 3y ago

64



# Tutorials



## Machine Learning

Machine learning is the hottest field in data science, and this track will get you started quickly.



### Overview

 Free Course

 4 hrs.

 19 Lessons

Prerequisite Skills:  
[Python](#)

Prepares you for these Learn Courses:  
[Deep Learning](#),  
[Machine Learning](#)  
[Explainability](#)

### Instructor



**Dan Becker**  
Data Scientist

### Level 1   Level 2

#### 1. How Models Work

The first step if you're new to machine learning



#### 2. Explore Your Data

Load data and set up your environment for your hands-on project



#### 3. Exercise: Explore Your Data



#### 4. Your First Machine Learning Model

Building your first model. Hurray!



#### 5. Exercise: Your First Machine Learning Model



# Status Incentives



**Michael Jahrer**

Graz, Austria

Joined 9 years ago · last seen in the past day

<http://www.operasolutions.com/>

Followers 1369

Following 4



**Competitions  
Grandmaster**

[Home](#) [Competitions \(64\)](#) [Kernels \(6\)](#) [Discussion \(200\)](#) [Followers \(1,369\)](#)

[Contact User](#)

[Follow User](#)

Competitions Grandmaster



Current Rank

**18**

of 96,011

Highest Rank

**6**

**15**

**11**

**2**

Home Credit Default Risk

🏆 · 5 months ago · Top 1%

**1<sup>st</sup>**  
of 7198

Porto Seguro's Safe Driver ...

🏆 · a year ago · Top 1%

**1<sup>st</sup>**  
of 5169

Cervical Cancer Screening

🏆 · 3 years ago · Top 3%

**1<sup>st</sup>**  
of 40

Kernels Contributor



**Unranked**

**0**

**0**

**0**

XGB Feature Importance (P...

3 years ago

**5**  
votes

naive XGB

2 years ago

**4**  
votes

Simple LightGBM 4!

2 years ago

**2**  
votes

Discussion Expert



Current Rank

**32**

of 82,748

Highest Rank

**9**

**15**

**30**

**99**

1st place with representatio...

🏆 · a year ago

**704**  
votes

Is 0.288 magical and optim...

🏆 · a year ago

**61**  
votes

you were only supposed to ...

🏆 · 2 years ago

**47**  
votes

# Kaggle CTRP Competitions



## Display Advertising Challenge

Predict click-through rates on display ads

Research · 4 years ago



## Avito Context Ad Clicks

Predict if context ads will earn a user's click

Featured · 4 years ago · marketing, tabular data, click prediction



## Outbrain Click Prediction

Can you predict which recommended content each user will click?

Featured · 2 years ago · internet, tabular data, click prediction




## Click-Through Rate Prediction

Predict whether a mobile ad will be clicked

Featured · 4 years ago

# Criteo Display Advertising Challenge




## Display Advertising Challenge

Predict click-through rates on display ads  
\$16,000 · 718 teams · 4 years ago

[Overview](#) [Data](#) [Discussion](#) [Leaderboard](#) [Rules](#)

### Overview

<b>Description</b>	Display advertising is a billion dollar effort and one of the central uses of machine learning on the Internet. However, its data and methods are usually kept under lock and key. In this research competition, CriteoLabs is sharing a week's worth of data for you to develop models predicting ad click-through rate (CTR). Given a user and the page he is visiting, what is the probability that he will click on a given ad?
<b>Evaluation</b>	
<b>Prizes</b>	
<b>About Criteo</b>	
<b>Timeline</b>	
<b>Winners</b>	



The goal of this challenge is to benchmark the most accurate ML algorithms for CTR estimation. All winning models will be released under an open source license. As a participant, you are given a chance to access the traffic logs from Criteo that include various undisclosed features along with the click labels.

<https://www.kaggle.com/c/criteo-display-ad-challenge>

# Criteo Display Advertising Challenge

- Criteo Display Advertising Challenge Data:
  - Features (inputs):
    - 13 numeric: unknown meanings, mostly counts, power laws evident
    - 26 categorical: unknown meanings, hashed (encoding without decoding), few dominant, many unique
  - Target (output): 0 / 1 (didn't / did click through)

# Mysterious Data

Label	I1	I2	...	I13	C1	C2	...	C26
1	3	20	...	2741	68fd1e64	80e26c9b	...	4cf72387
0	7	91	...	1157	3516f6e6	cfc86806	...	796a1a2e
0	12	73	...	1844	05db9164	38a947a1	...	5d93f8ab
					⋮			
?	9	62	...	1457	68fd1e64	cfc86806	...	cf59444f

#Train:  $\approx 45\text{M}$

#Test:  $\approx 6\text{M}$

Source: <https://www.csie.ntu.edu.tw/~r01922136/kaggle-2014-criteo.pdf>

# Mysterious Data

Unknown Labels: meanings of numbers and categories not given

Label	I1	I2	...	I13	C1	C2	...	C26
1	3	20	...	2741	68fd1e64	80e26c9b	...	4cf72387
0	7	91	...	1157	3516f6e6	cfc86806	...	796a1a2e
0	12	73	...	1844	05db9164	38a947a1	...	5d93f8ab
					⋮			
?	9	62	...	1457	68fd1e64	cfc86806	...	cf59444f

#Train:  $\approx$  45M

#Test:  $\approx$  6M

Source: <https://www.csie.ntu.edu.tw/~r01922136/kaggle-2014-criteo.pdf>

# Mysterious Data

Label	I1	I2	...	I13	C1	C2	...	C26
1	3	20	...	2741	68fd1e64	80e26c9b	...	4cf72387
0	7	91	...	1157	3516f6e6	cfc86806	...	796a1a2e
0	12	73	...	1844	05db9164	38a947a1	...	5d93f8ab
					⋮			
?	9	62	...	1457	68fd1e64	cfc86806	...	cf59444f

Categorical data is *hashed*.

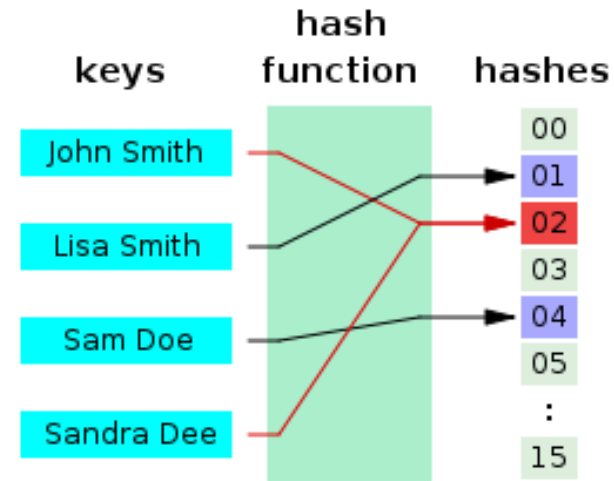
#Train:  $\approx$  45M  
#Test:  $\approx$  6M

Source: <https://www.csie.ntu.edu.tw/~r01922136/kaggle-2014-criteo.pdf>



# Hashing

- A **hash** function takes some data and maps it to a number.
- Example: URL (web address)
  - Representation: string of characters
  - Character representation: a number (Unicode value)
  - Start with value 0.
  - Repeat for each character:
    - Multiply value by 31
    - Add next character Unicode to value
  - Don't worry about overflow – it's just a consistent “mathematical blender”.



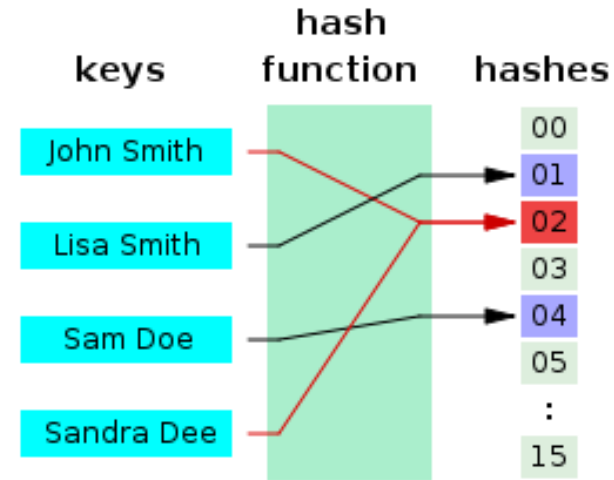
“Hi”

$H = 72, i = 105$

$31 * 72 + 105 = 2337$

# Hash Function Characteristics

- **Mapping:** same input → same output
- **Uniform:** outputs have similar probabilities
  - **Collision:** two different inputs → same output
  - Collisions are allowable (inevitable if  $\#data > \#hash$  values) but not desirable.
- **Non-invertible:** can't get back input from output (e.g. cryptographic hashing, anonymization)



# Missing Data

- The first 10 lines of the training data:

```

0  1  1  5  0  1302  4  15  2  181  1  2  2  68fd1e64  00e26c9b  fb936136  7b4723c4  25c83c90  7e0cccf  de7995b0  1f89b562  a73ee510  a8cd5504  b2cb9c98  37c9c164  2024a5f6  1adcceef  0ba8b39a  891b62e7  e5ba7672  f54016b9  21ddcd9  b1252a9d  07b5194c  3a171ecb  c5c50484  e0b83407  9727dd16
0  2  0  44  1  102  8  2  2  4  1  1  4  68fd1e64  f0cf0014  6f677e5  41274cd7  25c83c90  feeb92e5  922afcc0  0b153074  a73ee510  1b53a5fb  4f1b46f3  6.23E-11  d7020509  b20479f6  e6c5b5cd  c92f3b61  07c540ca  b04e4670  21ddcd9  5840adea  60f6121e  3a171ecb  43f13e8b  e0b83407  731c3655
0  2  0  1  14  767  89  4  2  245  1  3  3  45  207e604f  0a519c5c  02cf9076  c18be101  25c83c90  7e0cccf  c70204a1  0b153074  a73ee510  3b00e40b  9f5e6091  0fe001f4  aa655a2f  07d13a0f  6dc710ed  36103458  0efede7f  3412110d  e507c466  ad3062eb  3a171ecb  3b183c5c
0  893  4392  0  0  0  0  0  0  0  0  0  68fd1e64  2c16a946  a9a87e68  2e17d6f6  25c83c90  feeb92e5  2e0a609b  0b153074  a73ee510  efea433b  e51dd994  a30567ca  3516f6e6  07d13a0f  10231224  52b0600f  1e80c74f  74ef3502  6b3a5ca6  3a171ecb  9117a34a
0  3  -1  0  2  0  3  0  0  0  1  1  0  8cf07265  ae46a29d  c81680bb  f922efad  25c83c90  13718bbd  ad9fa255  0b153074  a73ee510  5282c137  e5d0af57  66a76a26  f0c653ac  1adcceef  8ff4b403  01adbab4  1e80c74f  26b3c7a7  21c9516a  32c7478e  b34f3128
0  -1  12824  0  0  6  0  05db9164  6c9c9cf3  2730ec9c  5400db0b  43b19349  0f6d9be0  53b5f978  0b153074  a73ee510  3b00e40b  91e0fc27  be45b077  9ff13f22  07d13a0f  06969a20  9bc7ff5  776ce399  92555263  242bb710  8ec974f4  be7c41b4  72c78f11
0  1  2  3168  0  1  2  0  439a44a4  ad4527a2  c02372d0  d34ebbaa  43b19349  feeb92e5  4bc6ffea  0b153074  a73ee510  3b00e40b  a4609aab  14d63530  772a00d7  07d13a0f  fd1302e  b00d3dc9  776ce399  cdfa0259  20062612  93bad2c0  1b256e61
1  1  4  2  0  0  1  0  0  1  1  0  68fd1e64  2c16a946  503b9dbc  e4d8ea90  f3474129  13718bbd  30eb9cf4  1f89b562  a73ee510  547c0ffe  bc8c9f21  60ab2f07  46f42a63  07d13a0f  10231224  e0ebbdcc7  e5ba7672  74ef3502  5316a17f  32c7478e  9117a34a
0  44  4  8  19010  249  28  31  141  1  8  05db9164  d833535f  d032c263  c18be101  25c83c90  7e0cccf  d5b6acf2  0b153074  a73ee510  2acdcf4e  008ac2d2  dfbb09fb  41a6ae00  b28479f6  e2502ec9  04090b2a  e5ba7672  42a2ed09  0014c32a  32c7478e  3b183c5c
0  35  1  33737  21  1  2  3  1  1  05db9164  510a40a5  d03e7c24  eb1fd928  25c83c90  52283d1c  0b153074  a73ee510  015ac093  e51dd994  951fe4a9  3516f6e6  07d13a0f  2ae4121c  0ec71479  d4bb7bd8  70d0f5f9  0e63fca0  32c7478e  0e8fe315

```

- Missing numeric and categorical features:

0	1	1	5	0	1382		07b5194c	3a171ecb	c5c50484	e8b83407	9727dd16
0	2	0	44	1	102		60f6221e	3a171ecb	43f13e8b	e8b83407	731c3655
0	2	0	1	14	767		e587c466	ad3062eb	3a171ecb	3b183c5c	
0		893			4392		6b3a5ca6		3a171ecb	9117a34a	
0	3	-1		0	2	• • •	21c9516a		32c7478e	b34f3128	
0		-1			12824		242bb710	8ec974f4	be7c41b4	72c78f11	
0		1	2		3168		20062612		93bad2c0	1b256e61	
1	1	4	2	0	0		5316a17f		32c7478e	9117a34a	
0		44	4	8	19010		0014c32a		32c7478e	3b183c5c	
0		35		1	33737		0e63fca0		32c7478e	0e8fe315	

# Missing Data: Imputation

- One approach to dealing with missing data is to *impute* values, i.e. replace with reasonable values inferred from surrounding data.
- In other words, create predictors for each value based on other known/unknown values.
- Cons:
  - Difficult to validate.
  - In Criteo data, missing values are correlated.
  - So ... we're writing predictors to impute data we're learning predictors from?

# Missing Data: Embrace the “Unknown”

- Retain “unknown” as data that contains valuable information.
- Does the lack of CTR context data caused by incognito browsing mode provide information on what a person is more likely to click?
- Categorical data: For each category C# with missing data, create a new category value “C#:unknown”.

# Missing Data: Embrace the “Unknown”

- Numeric data:
  - Create an additional feature that indicates whether the value for a feature is (un)known.
    - Additionally could impute mean, median, etc., for unknown value.
  - Convert to categorical and add “C#:unknown” category...

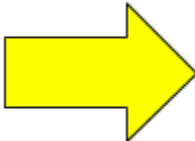
# Numeric to Categorical: Binning

- Histogram-based
  - Uniform ranges: (+) simple (-) uneven distribution, poor for non-uniform data
  - Uniform ranges on transformation (e.g. log): (+) somewhat simple (-) transformation requires understanding of data distribution
- Quantiles
  - E.g. quartiles = 4-quantiles, quintiles = 5-quantiles
  - (+) simple, even distribution by definition, (-) preponderance of few values → duplicate bins (eliminate)

# Categorical to Numeric: One-Hot Encoding

- For each categorical input variable:
  - For each possible category value, create a new numeric input variable that can be assigned numeric value 1 (“belongs to this category”) or 0 (“does not belong to this category”).
  - For each input, replace the categorical value variable with these new numeric inputs.

Color			
Red			
Red			
Yellow			
Green			
Yellow			



Red	Yellow	Green
1	0	0
1	0	0
0	1	0
0	0	1

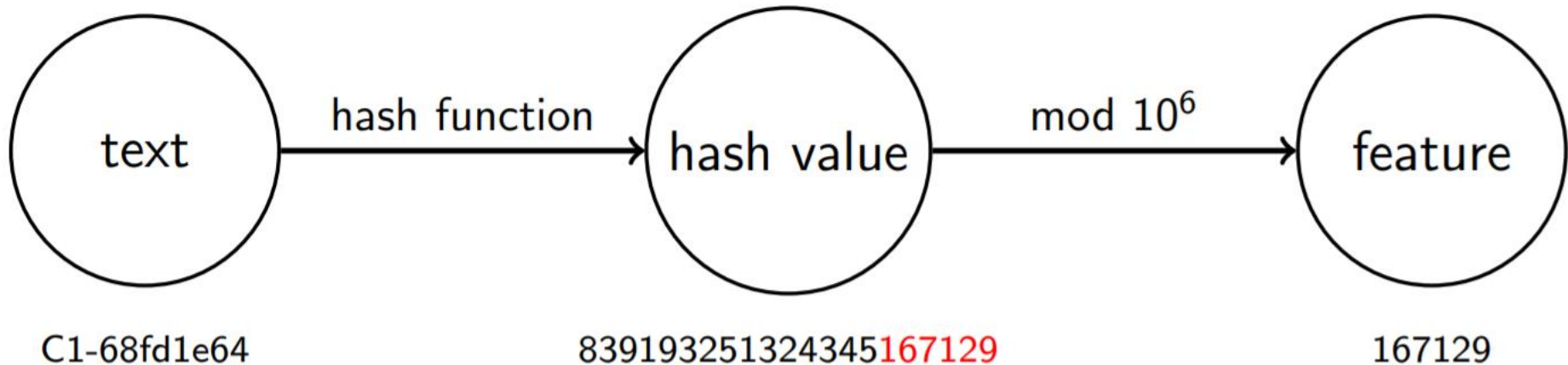


# Categorical to Numeric: Hashing

- When there are a large number of categories, one-hot encoding isn't practical.
  - E.g. Criteo data category C3 in its small sample of CTR data had 10,131,226 distinct categorical values.
  - One approach (e.g. for power law data): one-hot encode few dominant values plus “rare” category.
  - Hashing trick:
    - Append category name and unusual character before category value and *hash* to an integer.
    - Create a one-hot-like category for each integer.

# Hashing Trick Example

- From <https://www.csie.ntu.edu.tw/~r01922136/kaggle-2014-criteo.pdf>:

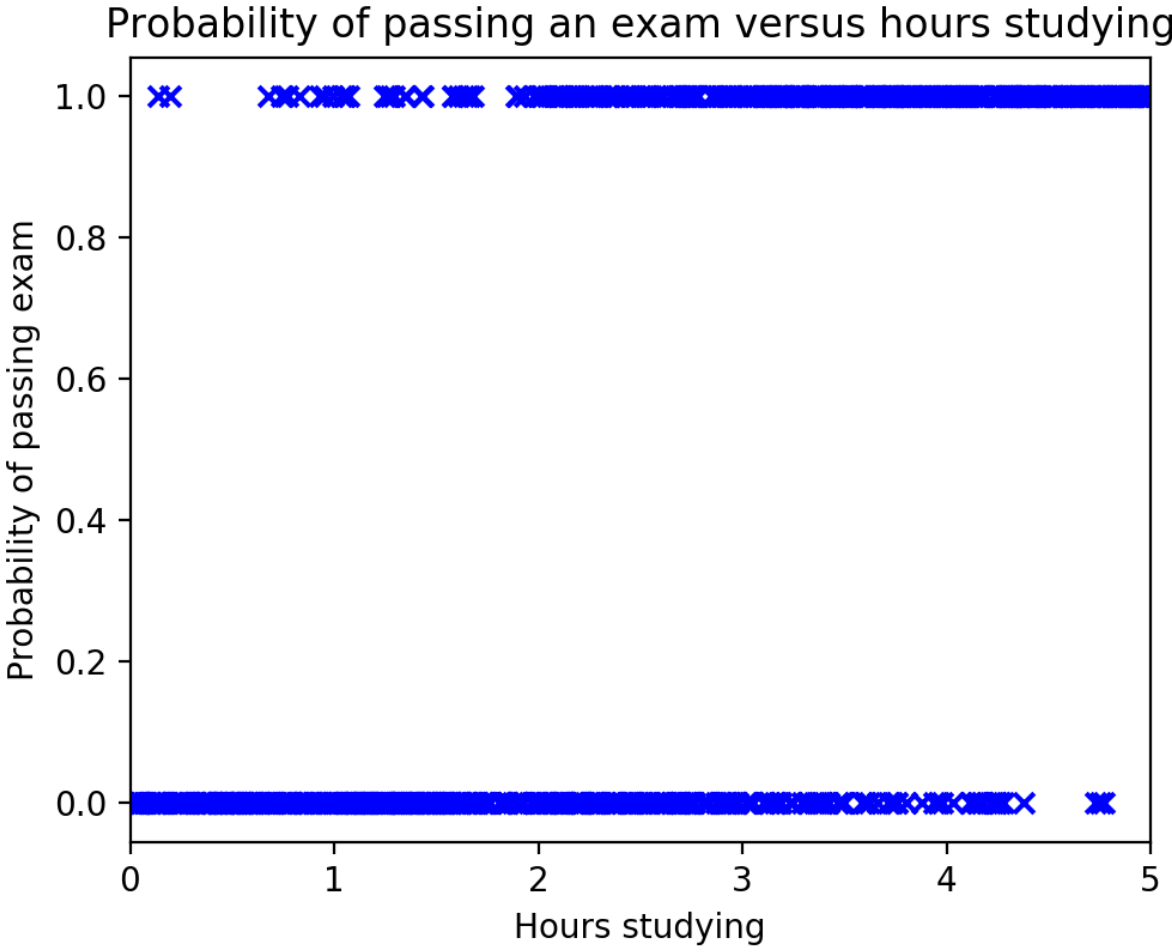


- Fundamental tradeoff: greater/lesser number hashed features results in ...
  - ... more/less expensive computation
  - ... less/more frequent hash collisions (i.e. unlike categories treated as like)

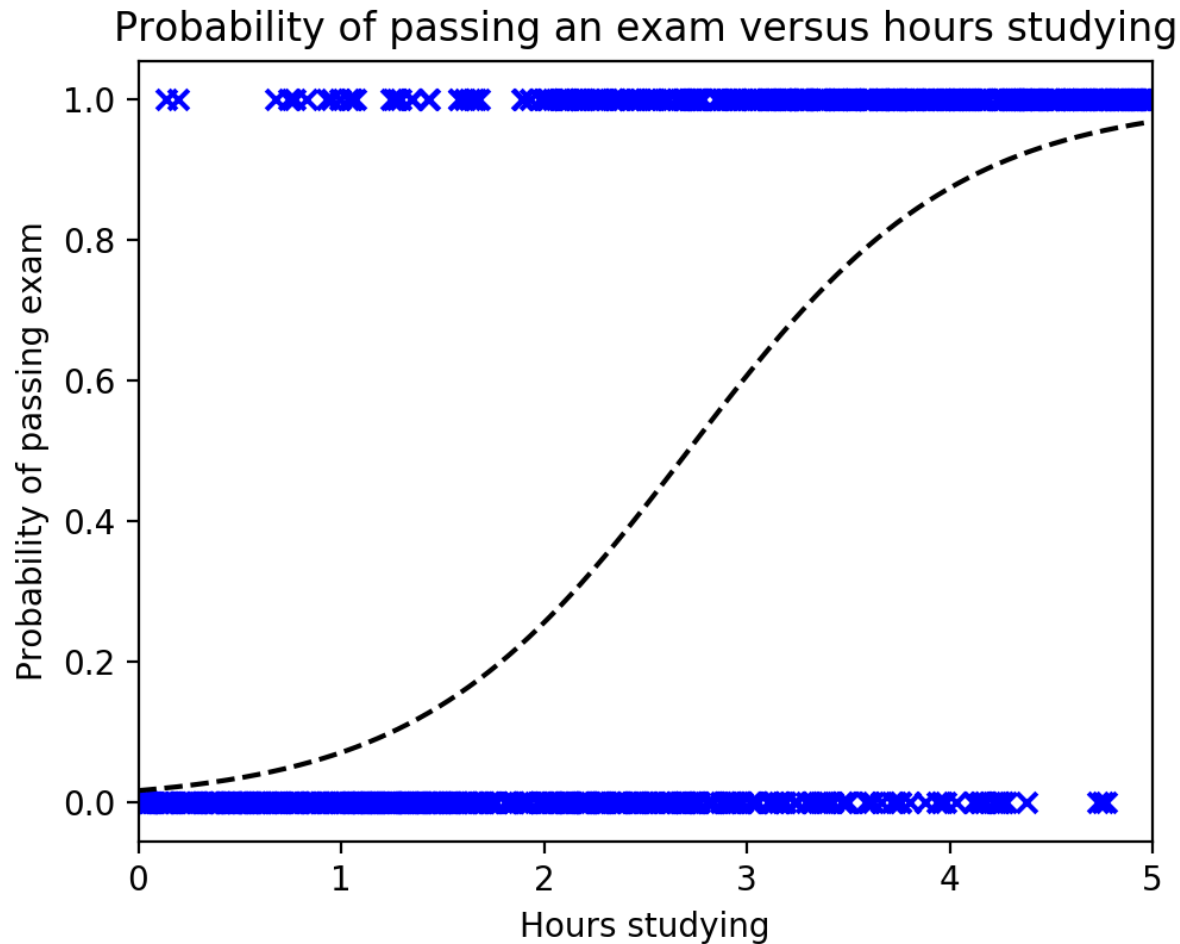
# Logistic Regression Motivation

- Logistic regression is perhaps the simplest technique to beat the Criteo benchmark, scoring  $\sim 42^{\text{nd}}$  percentile on leaderboard:
  - <https://www.kaggle.com/c/criteo-display-ad-challenge/discussion/10322>
  - 100 lines of Python, 200MB RAM, 30 min. training
  - Also: logistic regression recommended for CTRP by researchers of [Criteo](#), [Microsoft](#), [LinkedIn](#), [Google](#), and [Facebook](#) for practical, scalable implementation.

# Example: Passing vs. Studying

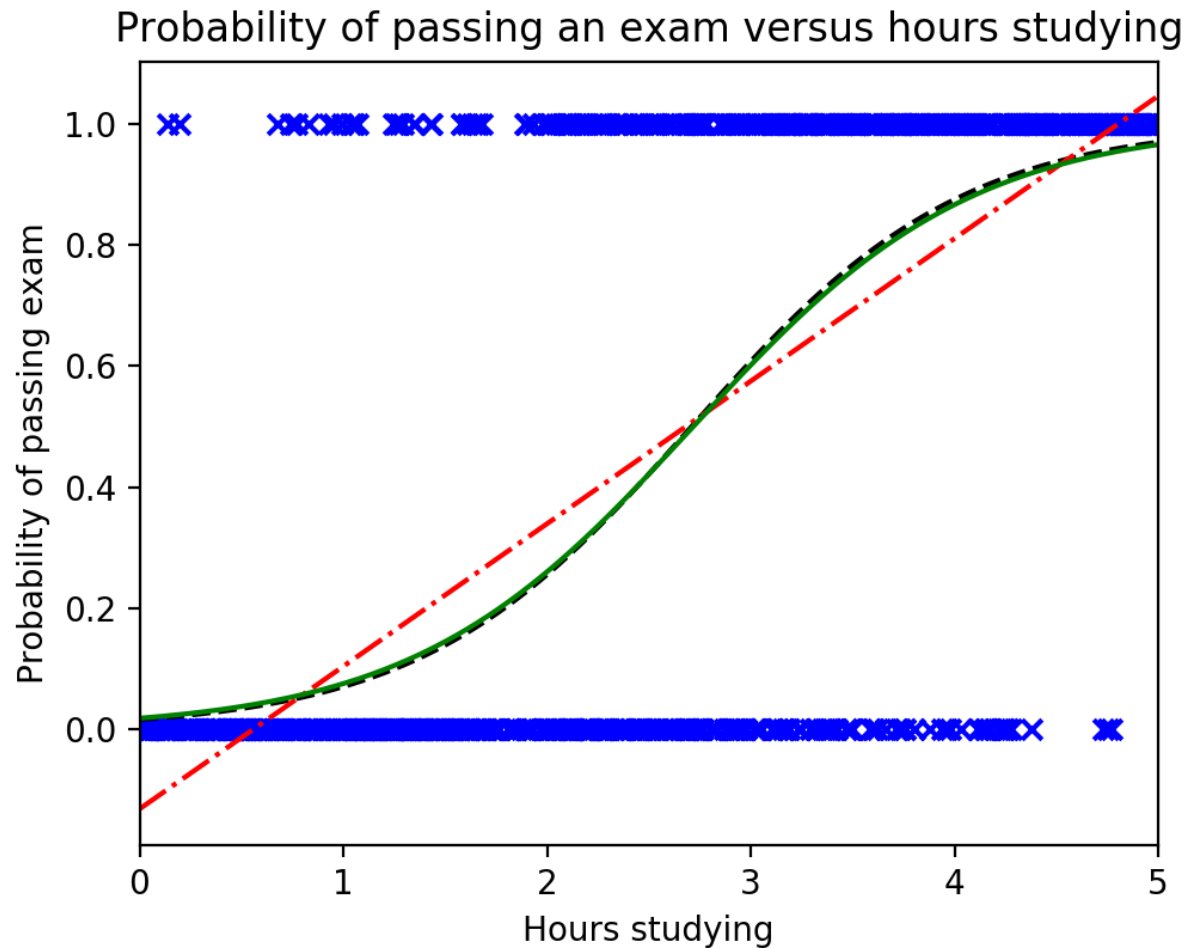


# Unknown Logistic Model





# Logistic Regression Recovering Model



# Logistic Regression with Stochastic Gradient Descent

- Output:  $p(x) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$
- Initially:  $\beta_0 = \beta_1 = 0$
- Repeat:
  - For each input  $x$ ,
    - Adjust intercept  $\beta_0$  by learning rate \* error \*  $p'(x)$
    - Adjust coefficient  $\beta_1$  by learning rate \* error \*  $p'(x) * x$
- Note:
  - Error =  $y - p(x)$
  - $p'(x) = p(x) * (1 - p(x))$  (the slope of  $p$  at  $x$ )
  - This is neural network learning with a single logistic neuron with bias input of 1



# Logistic Regression Takeaways

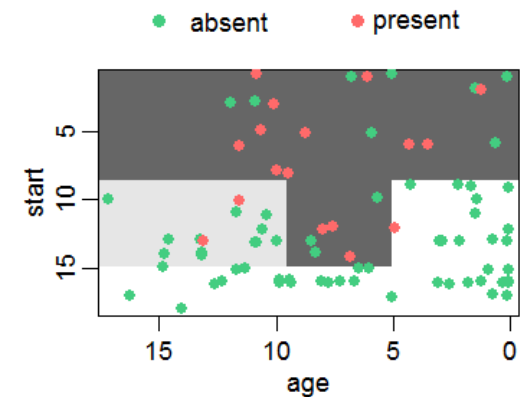
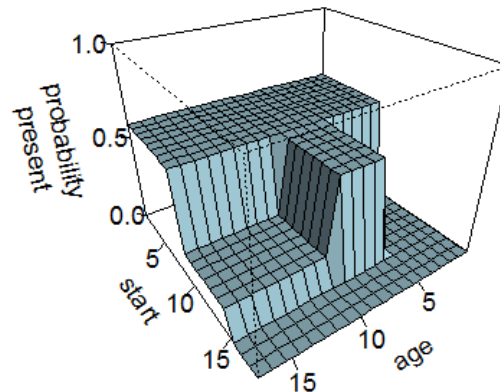
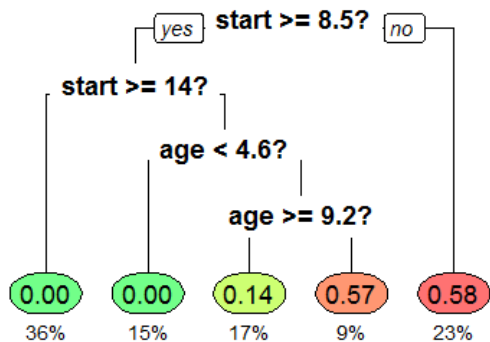
- The previous algorithm doesn't require complex software. (12 lines raw Python code)
- Easy and effective for CTR prediction.
- Key to good performance: *skillful feature engineering of numeric features*
- Foreshadowing: Since logistic regression is a simple special case of neural network learning, I would expect deep learning tools to make future inroads here.

# Maximizing Info with Decisions

- Number Guessing Game example:
  - “I’m thinking of a number from 1 to 100.”
  - Number guess  $\rightarrow$  “Higher.” / “Lower.” / “Correct.”
  - What is the best strategy and why?
- Good play maximizes information according to some measure (e.g. entropy).

# Decision Trees for Regression (Regression Trees)

- Numeric features (missing values permitted)
- At each node in the tree, a branch is decided on according to a features value (or lack thereof)



A regression tree estimating the probability of kyphosis (hunchback) after surgery, given the age of the patient and the vertebra at which surgery was started.

Source: [https://en.wikipedia.org/wiki/Decision\\_tree\\_learning](https://en.wikipedia.org/wiki/Decision_tree_learning)

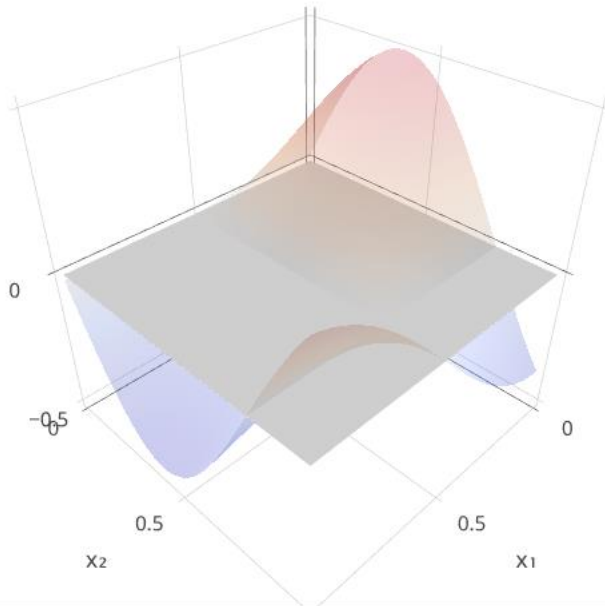
# The Power of Weak Classifiers

- Caveats:
  - Too deep: Single instance leafs → overfitting; similar to nearest neighbor ( $n=1$ )
  - Too shallow: Large hyperrectangular sets → underfitting; poor, blocky generalization
- Many weak classifiers working together can achieve good fit and generalization.
  - “Plans fail for lack of counsel, but with many advisers they succeed.” – Proverbs 15:22
- Ensemble methods: **boosting**, bagging, stacking

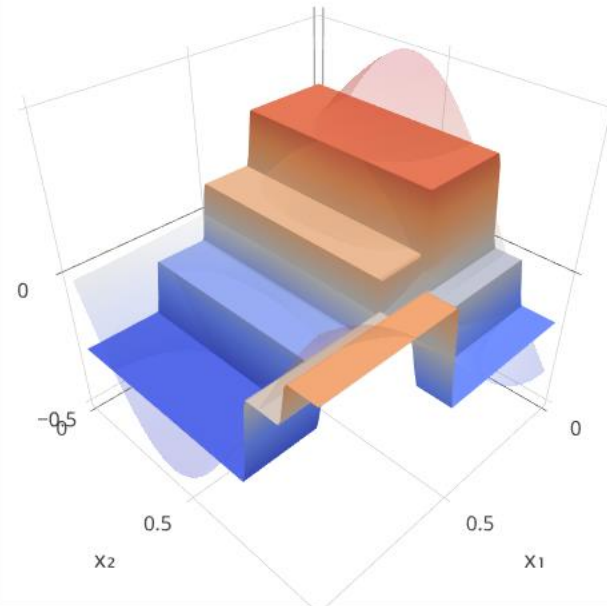
# Gradient Boosting of Regression Trees

- Basic boosting idea:
  - Initially, make a 0 or constant prediction.
  - Repeat:
    - Compute prediction errors from the weighted sum of our weak-learner predictions.
    - Fit a new weak-learner to predict these errors and *add* its weighted error-prediction to our model.
- Alex Rogozhnikov's beautiful demonstration:  
[https://arogozhnikov.github.io/2016/06/24/gradient\\_boosting\\_explained.html](https://arogozhnikov.github.io/2016/06/24/gradient_boosting_explained.html)

target function  $f(\mathbf{x})$  and prediction of previous trees  $D(\mathbf{x})$



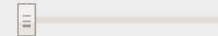
residual  $R(\mathbf{x})$  and prediction of next tree  $d_n(\mathbf{x})$



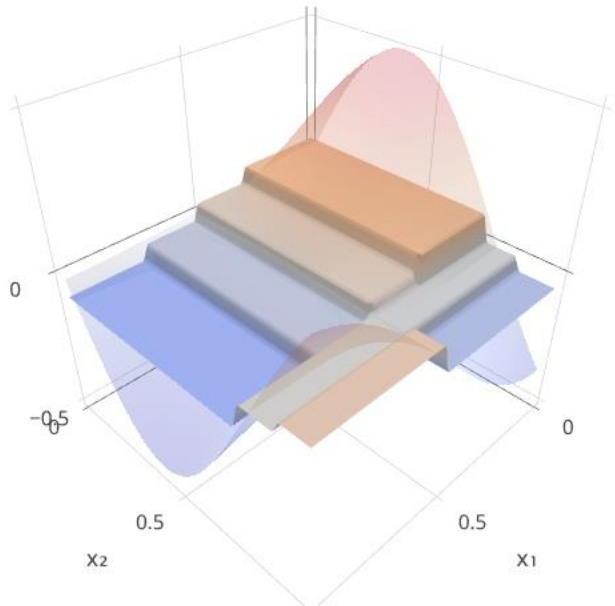
Tree depth: 3



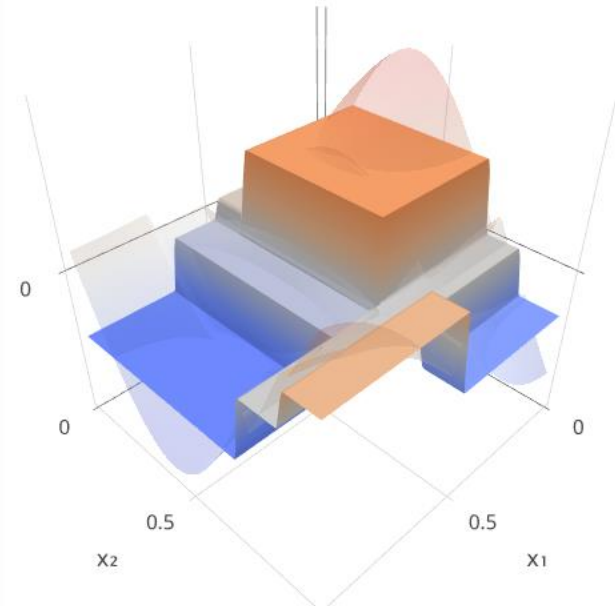
Number of built trees: 0



target function  $f(x)$  and prediction of previous trees  $D(x)$



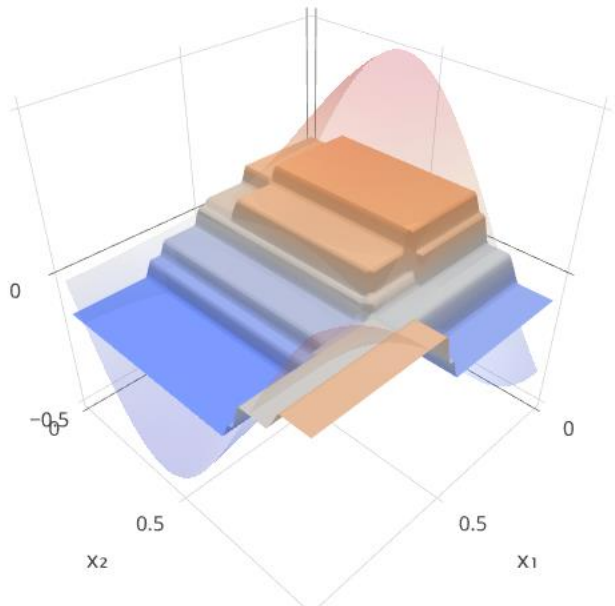
residual  $R(x)$  and prediction of next tree  $d_n(x)$



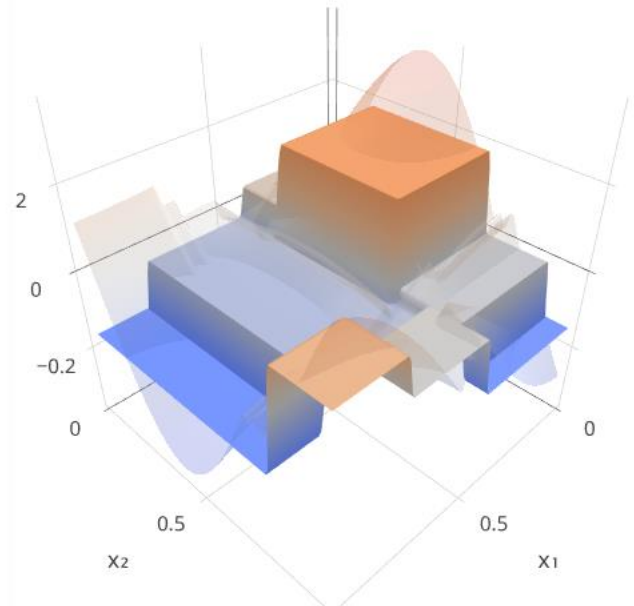
Tree depth: 3

Number of built trees: 1

target function  $f(\mathbf{x})$  and prediction of previous trees  $D(\mathbf{x})$



residual  $R(\mathbf{x})$  and prediction of next tree  $d_n(\mathbf{x})$

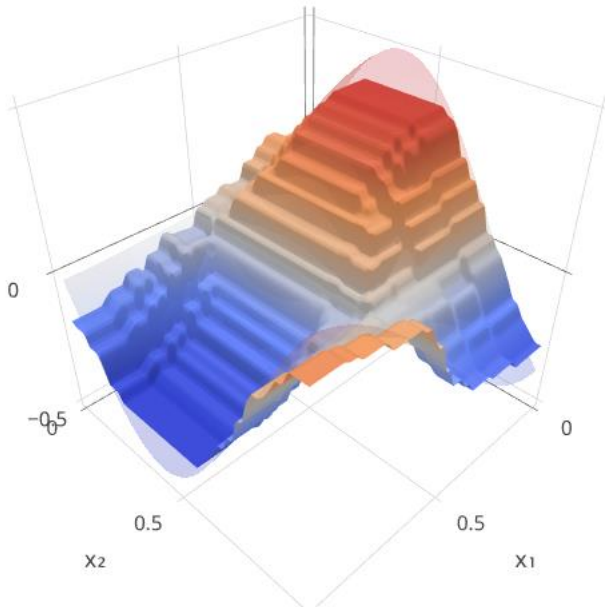


Tree depth: 3

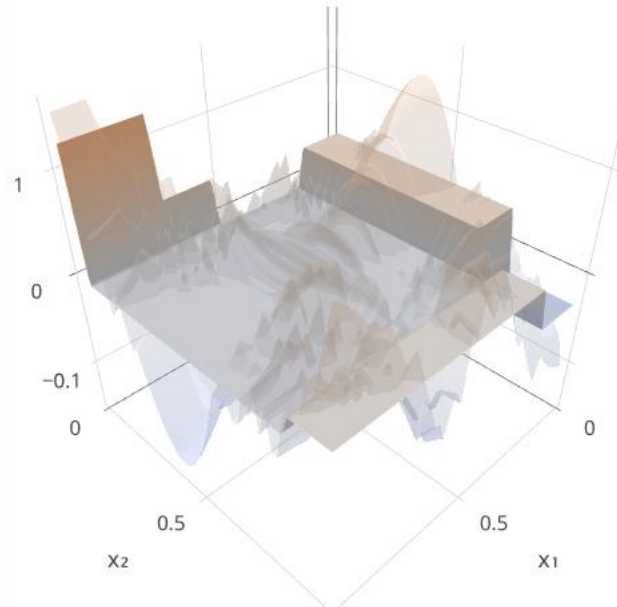
Number of built trees: 2



target function  $f(\mathbf{x})$  and prediction of previous trees  $D(\mathbf{x})$



residual  $R(\mathbf{x})$  and prediction of next tree  $d_n(\mathbf{x})$



Tree depth: 3



Number of built trees: 10



# XGBoost

- *“Among the 29 challenge winning solutions published at Kaggle’s blog during 2015, 17 solutions [~59%] used XGBoost. Among these solutions, eight [~28%] solely used XGBoost to train the model, while most others combined XGBoost with neural nets in ensembles.”* - Tianqi Chen, Carlos Guestrin. [“XGBoost: A Scalable Tree Boosting System”](#)

# XGBoost Features

- **XGBoost** is a specific implementation of **gradient boosted decision trees** that:
  - Supports a command-line interface, C++, Python (scikit-learn), R (caret), Java/JVM languages + Hadoop platform
  - A range of computing environments with parallelization, distributed computing, etc.
  - Handles sparse, missing data
  - Is *fast* and *high-performance* across diverse problem domains
  - <https://xgboost.readthedocs.io>

# Field-aware Factorization Machines (FFMs)

- Top-performing technique in 3 of 4 Kaggle CTR prediction competitions plus RecSys 2015:
  - Criteo: <https://www.kaggle.com/c/criteo-display-ad-challenge>
  - Avazu: <https://www.kaggle.com/c/avazu-ctr-prediction>
  - Outbrain: <https://www.kaggle.com/c/outbrain-click-prediction>
  - RecSys 2015: <http://dl.acm.org/citation.cfm?id=2813511&dl=ACM&coll=DL&CFID=941880276&CFTOKEN=60022934>

# What's Different? Field-Aware Latent Factors

- Latent factor
  - learned weight; tuned variable
  - How much an input contributes to an output
- Many techniques learn “latent factors”:
  - Linear regression: one per feature + 1

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_p X_p$$

- Logistic regression: one per feature + 1

$$p(X) = \frac{e^{\beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p}}{1 + e^{\beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p}}$$

# What's Different? Field-Aware Latent Factors (cont.)

- Many techniques learn “latent factors”:
  - Degree-2 polynomial regression: one per pair of features

$$\phi_{\text{Poly2}}(\mathbf{w}, \mathbf{x}) = \sum_{j_1=1}^n \sum_{j_2=j_1+1}^n w_{h(j_1, j_2)} x_{j_1} x_{j_2}$$

- Factorization machine (FM):

- $k$  per feature
- “latent factor vector”, a.k.a. “latent vector”

$$\phi_{\text{FM}}(\mathbf{w}, \mathbf{x}) = \sum_{j_1=1}^n \sum_{j_2=j_1+1}^n (\mathbf{w}_{j_1} \cdot \mathbf{w}_{j_2}) x_{j_1} x_{j_2}$$

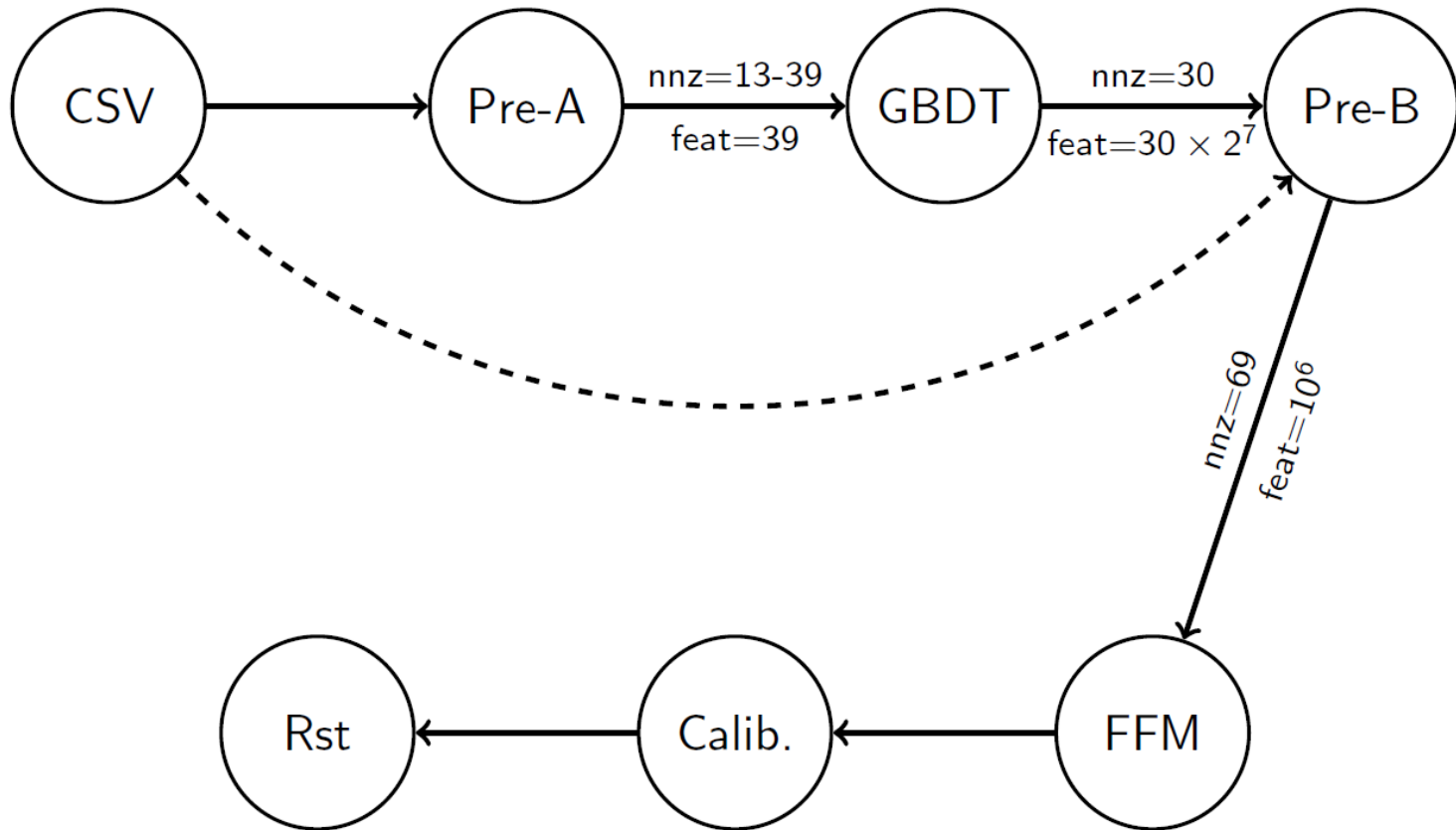
# What's Different? Field-Aware Latent Factors (cont.)

- Many techniques learn “latent factors”:
  - Field-aware Factorization machine (FFM):
    - $k$  per feature *and field* pair
    - Field:
      - Features are often one-hot encoded
      - Continuous block of binary features often represent different values for the same underlying “field”
      - E.g. Field: “OS”, features: “Windows”, “MacOS”, “Android”
      - libffm: FFM library (<https://github.com/guestwalk/libffm>)

$$\phi_{\text{FFM}}(\mathbf{w}, \mathbf{x}) = \sum_{j_1=1}^n \sum_{j_2=j_1+1}^n (\mathbf{w}_{j_1, f_2} \cdot \mathbf{w}_{j_2, f_1}) x_{j_1} x_{j_2}$$

# Winning Team Process

- From <https://www.csie.ntu.edu.tw/~r01922136/kaggle-2014-criteo.pdf>:



*“nnz” means the number of non-zero elements of each impression; “feat” represents the size of feature space.*



# Is the Extra Engineering Worth it?

- Kaggle Criteo leaderboard based on **logarithmic loss** (a.k.a. logloss)
  - 0.69315 → 50% correct in binary classification (random guessing baseline)
- Simple logistic regression with hashing trick:
  - 0.46881 (private leaderboard) ~62.6% correct
- FFM with feature engineering using GBDT:
  - 0.44463 (private leaderboard) ~64.1% correct

# Computational Cost

- ~1.5% increase in correct prediction, but greater computational complexity:
  - Logistic regression:  $n$  factors to learn and **relearn** in **dynamic** context
  - FFM:  $kn^2$  factors to learn and **relearn**

# Model Ensemble for Click Prediction in Bing Search Ads

Xiaoliang Ling  
Microsoft Bing  
No. 5 Dan Ling Street  
Beijing, China  
xiaoling@microsoft.com

Weiwei Deng  
Microsoft Bing  
No. 5 Dan Ling Street  
Beijing, China  
dedeng@microsoft.com

Chen Gu  
Microsoft Bing  
No. 5 Dan Ling Street  
Beijing, China  
chengu@microsoft.com

Hucheng Zhou  
Microsoft Research  
No. 5 Dan Ling Street  
Beijing, China  
huzho@microsoft.com

Cui Li<sup>\*</sup>  
Microsoft Research  
No. 5 Dan Ling Street  
Beijing, China  
v-culi@microsoft.com

Feng Sun  
Microsoft Bing  
No. 5 Dan Ling Street  
Beijing, China  
fengsun@microsoft.com

## ABSTRACT

Accurate estimation of the click-through rate (CTR) in sponsored ads significantly impacts the user search experience and business revenue, even 0.1% of accuracy improvement would yield greater earnings in the hundreds of millions of dollars. CTR prediction is generally formulated as a supervised classification problem. In this paper, we share our experience and learning on model ensemble design and our innovation. Specifically, we present 8 ensemble methods and evaluate them on our production data. Boosting neural networks with gradient boosting decision trees turns out to be the best. With larger training data, there is a nearly 0.9% AUC improvement in offline testing and significant click yield gains in online traffic. In addition, we share our experience and learning on improving the quality of training.

## Keywords

click prediction; DNN; GBDT; model ensemble

## 1. INTRODUCTION

Search engine advertising has become a significant element of the web browsing experience. Choosing the right ads for a query and the order in which they are displayed greatly affects the probability that a user will see and click on each ad. Accurately estimating the click-through rate (CTR) of ads [10, 16, 12] has a significant impact on the revenue of search businesses; even a 0.1% accuracy improvement in our production would yield hundreds of millions of dollars in additional earnings. An ad's CTR is usually modeled as a classification problem, and thus can be estimated by machine learning models. The training data is collected from historical impressions and the corresponding clicks. Because of the simplicity, scalability and online learning capability, logistic regression (LR) is the most widely used model that has been studied

<sup>\*</sup>This work was done during her internship in Microsoft Research.

©2017 International World Wide Web Conference Committee (IW3C2) published under Creative Commons CC BY 4.0 License.  
WWW 2017, April 3–7, 2017, Perth, Australia.  
ACM 978-1-4503-4914-7/17/04.  
http://dx.doi.org/10.1145/3041021.3054192



## ABSTRACT

Online advertising allows advertisers to only bid and pay for measurable user responses, such as clicks on ads. As a consequence, click prediction systems are central to most online advertising systems. With over 750 million daily active users and over 1 million active advertisers, predicting clicks on Facebook ads is a challenging machine learning task. In this paper we introduce a model which combines decision trees with logistic regression, outperforming either of these methods on its own by over 3%, an improvement with significant impact to the overall system performance. We then explore how a number of fundamental parameters impact the final prediction performance of our system. Not surprisingly, the most important thing is to have the right features: those capturing historical information about the user or ad dominate other types of features. Once we have the right features and the right model (decisions trees plus logistic regression), other factors play small roles (though even small improvements are important at scale). Picking the optimal handling for data freshness, learning rate schema and data sampling improve the model slightly, though much less than adding a high-value feature, or picking the right model to begin with.

## 1. INTRODUCTION

Digital advertising is a multi-billion dollar industry and is growing dramatically each year. In most online advertising platforms the allocation of ads is dynamic, tailored to user interests based on their observed feedback. Machine learning plays a central role in computing the expected utility of a candidate ad to a user, and in this way increases the

<sup>\*</sup>BL works now at Square, TX and YS work now at Quora, AA works in Twitter and RH works now at Amazon.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

## Simple and scalable response prediction for display advertising

OLIVIER CHAPPELLE, <sup>Criteo</sup>  
EREN MANAVOGLU, <sup>Microsoft</sup>  
ROMER ROSALES, <sup>LinkedIn</sup>

Clickthrough and conversion rates estimation are two core predictions tasks in display advertising. We present in this paper a machine learning framework based on logistic regression that is specifically designed to tackle the specifics of display advertising. The resulting system has the following characteristics: it is simple (we have trained it on terabytes of data); and it provides

Storage And Retrieval); Online Information

g, machine learning, click prediction, hashing, @ature

uary YYYY), 34 pages.  
1145/0000000.0000000

vertising where advertisers pay publishers for ads. The traditional method of selling display advertising contracts between the advertisers and publishers has emerged as a popular alternative for publishers, and increased reach with the advertisers [Muthukrishnan 2009]. A wide range of payment options. If the goal is to reach the target audience (for instance per impression (CPM) with targeting on the advertiser. However, many advertising models unless that impression leads the user to a pre-defined action on their website (such as an email list). An auction that supports short advertiser bids to *Expected price per im-*

ere at Yahoo! Labs.

l of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted, to redistribute to lists, or to use any component of this work without fee. Permissions may be requested from the ACM Press, 400 7th Street, New York, NY 10121-0701 USA, fax +1 (212)

1145/0000000.0000000

ogy, Vol. V, No. N, Article A, Publication date: January YYYY.

efficiency of the marketplace.

The 2007 seminal papers by Varian [11] and by Edelman et al. [4] describe the bid and pay per click auctions pioneered by Google and Yahoo! That same year Microsoft was also building a sponsored search marketplace based on the same auction model [9]. The efficiency of an ads auction depends on the accuracy and calibration of click prediction. The click prediction system needs to be robust and adaptive, and capable of learning from massive volumes of data. The goal of this paper is to share insights derived from experiments performed with these requirements in mind and executed against real world data.

In sponsored search advertising, the user query is used to retrieve candidate ads, which explicitly or implicitly are matched to the query. At Facebook, ads are not associated with a query, but instead specify demographic and interest targeting. As a consequence of this, the volume of ads that are eligible to be displayed when a user visits Facebook can be larger than for sponsored search.

In order to tackle a very large number of candidate ads per query, where a request for ads is triggered whenever a user visits Facebook, we would first build a cascade of classifiers of increasing computational cost. In this paper we focus on the last stage click prediction model of a cascade classifier, that is the model that produces predictions for the final set of candidate ads.

We find that a hybrid model which combines decision trees with logistic regression outperforms either of these methods on its own by over 3%. This improvement has significant impact to the overall system performance. A number of fundamental parameters impact the final prediction performance of our system. As expected the most important thing is to have the right features: those capturing historical information about the user or ad dominate other types of features. Once we have the right features and the right model (decisions trees plus logistic regression), other factors play small roles (though even small improvements are important at scale). Picking the optimal handling for data freshness, learning rate schema and data sampling improve the model slightly, though much less than adding a high-value feature.

# Published Research from the Trenches

- Initial efforts focused on **logistic regression**
- Most big production systems reportedly kept it simple in the final stage of prediction:
  - [Google \(2013\)](#): prob. feature inclusion + Bloom filters → **logistic regression**
  - [Facebook \(2014\)](#): boosted decision trees → **logistic regression**
  - [Yahoo \(2014\)](#): hashing trick → **logistic regression**
- However...

# Towards Neural Network Prediction

- More recently, [Microsoft \(2017\)](#) research
  - reports “factorization machines (FMs), gradient boosting decision trees (GBDTs) and deep neural networks (DNNs) have also been evaluated and gradually adopted in industry.”
  - recommends boosting **neural networks** with gradient boosting decision trees

# Perspective

- The last sigmoid layer of a neural network (deep or otherwise) for binary classification is logistic regression.
- Previous layers of a deep neural network learn an internal representation of inputs, i.e. perform automatic *feature engineering*.
- Thus, most efforts to engineer successful, modern CTR prediction systems focus on layered feature engineering using:
  - Hashing tricks
  - Features engineered with GBDTs, FFMs, and deep neural networks (DNNs), or a layered/ensembled combination thereof.
- Future: Additional automated feature representation learning with deep neural networks

# CTRP Conclusions

- To get prediction performance quickly and easily, hash data to binary features and apply logistic regression.
- For + few % of accuracy, dig into **Kaggle forums** and the **latest industry papers** for a variety of means to engineer features most helpful to CTR prediction. We've surveyed a number here.
- Knowledge is power. ( $\uparrow$  data  $\rightarrow$   $\uparrow$  predictions)
- Priority of effort:  $\uparrow$  data  $>$   $\uparrow$  feature engineering  $>$   $\uparrow$  learning/regression algorithms.

# Next Steps

- Interested in learning more about Data Science and Machine Learning?
  - Create a Kaggle Account
  - Enter a learning competition, e.g. “[Titanic: Machine Learning from Disaster](#)”
  - Take related tutorials, learn from kernels and discussions, steadily work to improve your skills, and share it forward

