



Using Target Encoding to Improve Model Predictions

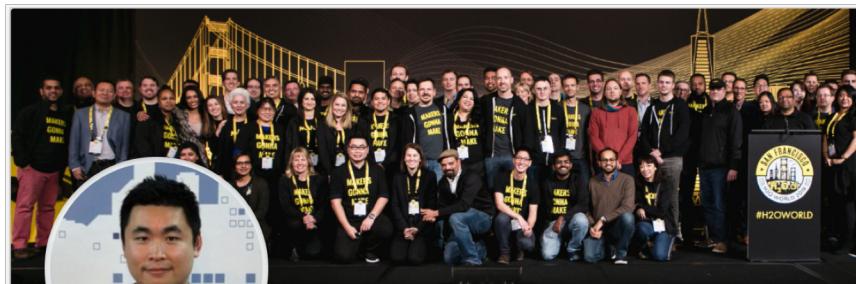
Jo-fai (Joe) Chow

Data Science Evangelist, H2O.ai

joe@h2o.ai / @matlabulous

Download: bit.ly/h2o_meetups

ABOUT ME



Jo-fai Chow

Data Science Evangelist & Community Manager at H2O.ai
United Kingdom

Add profile section ▾

More...

Jo-fai (or Joe) has multiple roles at H2O.ai. He is best known as the #360Selfie g

On LinkedIn, he is the data science evangelist and community manager but every photography skills totally overshadow his data science knowledge these days. Or like a die-hard MATLAB fanboy with the handle @matlabulous (because MATLAB at Uni). Since joining H2O.ai in 2016, Joe has delivered H2O talks/workshops in 4 Europe and US.

06:40 1

Booking.com

Welcome back, Jo-fai!

Search destination/property name

Madrid, Spain

Wed 27 Feb – Fri 1 Mar

1 room · 1 adult · 0 children

Search

H2O. I'm travelling for work

See See Your current trip

Leonardo City Tower Hotel Tel Aviv 24 27
3 nights in Tel Aviv

27 Feb

Search Bookings Profile More

H2O.ai is the open source leader in AI

H2O.ai is a visionary Silicon Valley open source software company that created and reimagined what is possible. We are a company of makers that brought to market new platforms and technologies to drive the AI movement. We are the makers of, H2O, the leading open source data science and machine learning platform used by nearly half of the Fortune 500 and trusted by over 14,000 organizations and hundreds of thousands of data scientists around the world.

Our approach is to be open, transparent and push the bleeding edge. Our philosophy is to create a culture of makers: community, customers, partners, entrepreneurs and our own “makers gonna make”. Our vision is to democratize AI for everyone. Not just a select few. We enable this with our award winning, H2O Driverless AI, the platform that uses AI to do AI to make it easier, faster and cheaper to deliver expert data science as a force multiplier for every enterprise. We want everyone to explore, learn, dream and imagine a new future.

H2O World in San Francisco

4th & 5th Feb, 2019

Our most-attended event to date

Attendance: 900



222 OF THE 500
H₂O

8 OF TOP 10
BANKS

7 OF TOP 10
INSURANCE COMPANIES

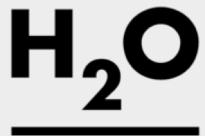
4 OF TOP 10
HEALTHCARE COMPANIES

14,000 Companies using H2O



H2O.ai Product Suite

H₂O.ai



In-memory, distributed
machine learning algorithms
with H2O Flow GUI



H2O AI open source engine
integration with Spark



Lightning fast machine
learning on GPUs

Open Source

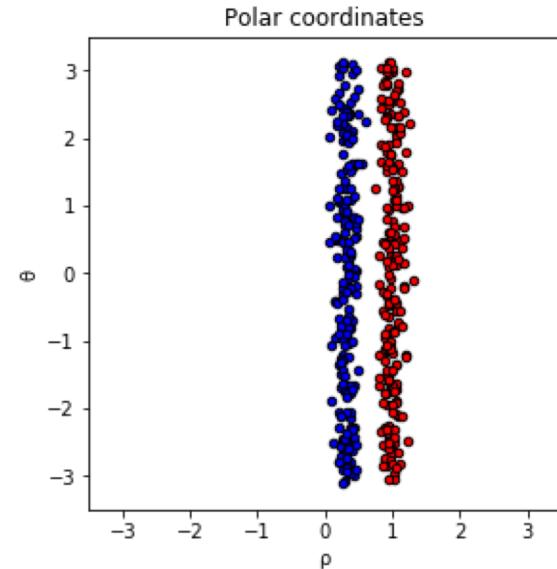
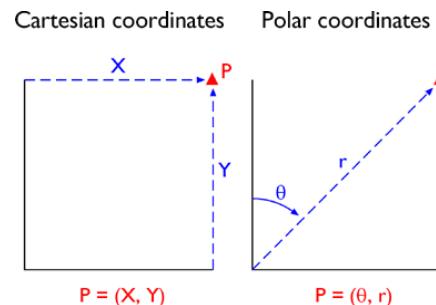
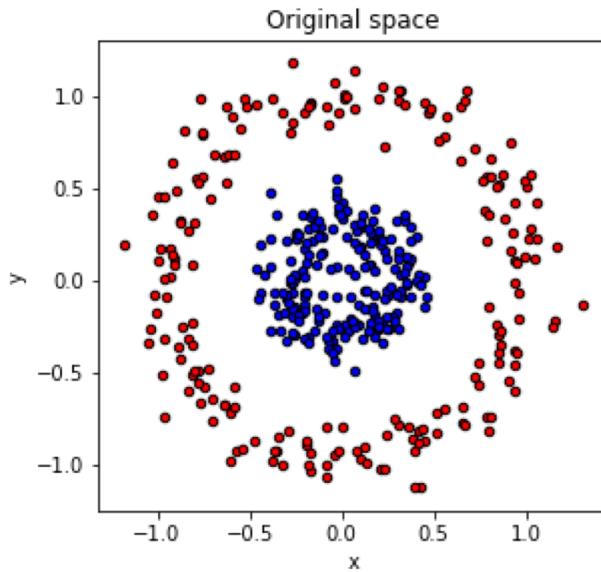
- 100% open source – Apache V2 licensed
- Built for data scientists – interface using R, Python on H2O Flow (interactive notebook interface)
- Enterprise support subscriptions

DRIVERLESSAI

Automatic feature engineering,
machine learning and interpretability

- Enterprise software
- Built for domain users, analysts and data scientists – GUI-based interface for end-to-end data science
- Fully automated machine learning from ingest to deployment
- User licenses on a per seat basis (annual subscription)

Why Feature Engineering?



$$\left[\begin{array}{c} \\ \\ \textbf{X} \\ \\ \end{array} \right]$$



$$\left[\begin{array}{c} \\ \\ \textbf{U} \\ \\ \end{array} \right]$$

Feature Engineering: Target Mean Encoding

What?

- Replace categorical variables with the mean of the response

Why?

- Categorical variables increase the number of features (dummy encoding) and can cause us to overfit

Caution!

- If it is applied without care, it may lead to overfitting

Target Mean Encoding

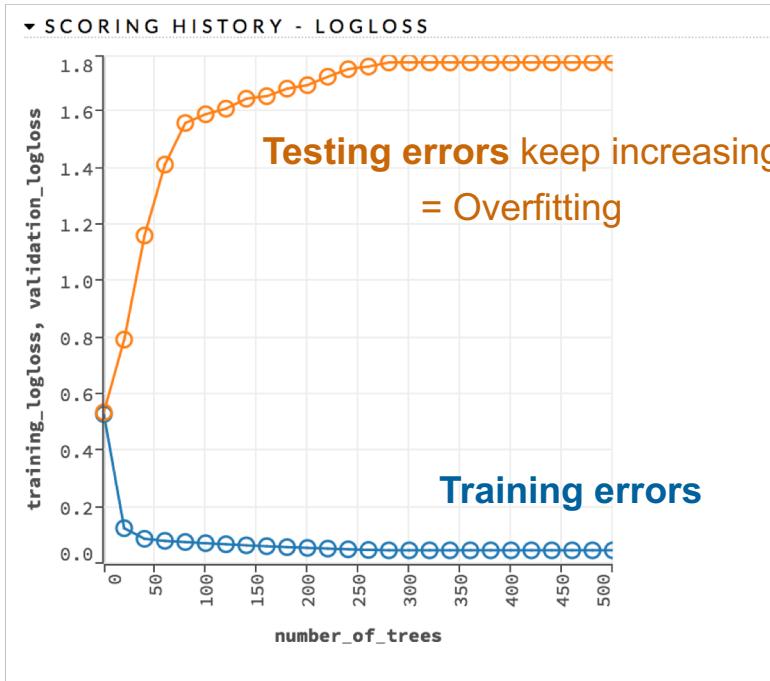
Categorical Feature	Target	Target Mean Encoding
A	0	0.3333333333333333
A	0	
A	0	
B	1	0.6666666666666667
B	0	
B	0	
C	1	0.6666666666666667

Target Mean Encoding Done Wrong

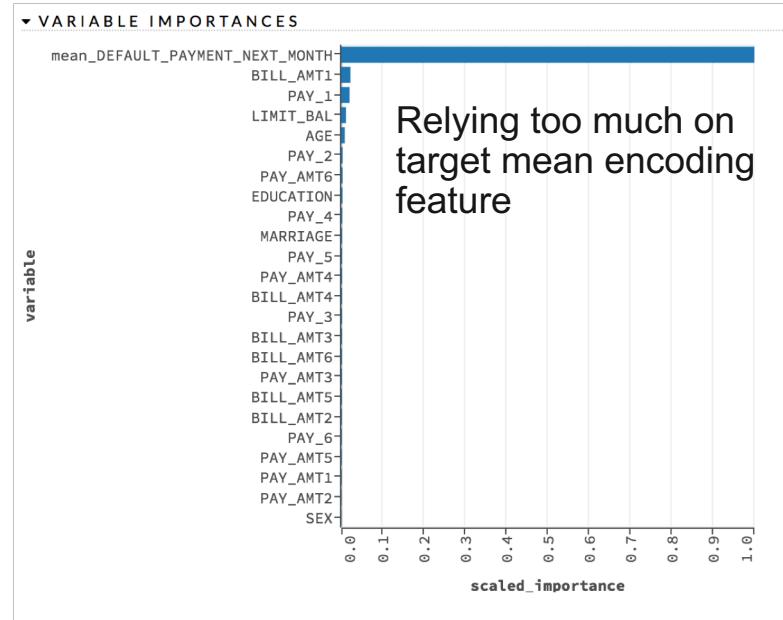
Categorical Feature	Target	Target Mean Encoding
A	0	0
A	0	
A	0	
B	1	0.33
B	0	
B	0	
C	1	1

Worst Case Scenario: Response Column = Mean Target Encoding

Target Mean Encoding Done Wrong



Scoring History: Training vs Testing



Data Leakage Feature is the only important feature

Solution: Cross-Validation Target Encoding

Fold	Categorical Feature	Target
1	A	0
2		0
3		0
2	B	1
1		0
3		0
1	C	1

Solution: Cross-Validation Target Encoding

Fold	Categorical Feature	Target
2	A	0
3	A	0
2	B	1
3	B	0

Fold	Categorical Feature	Target
1	A	0
1	B	0
1	C	1

Solution: Cross-Validation Target Encoding

Fold	Categorical Feature	Target	Target Mean Encoding
2	A	0	0
3	A	0	0
2	B	1	0.5
3	B	0	0.5

Fold	Categorical Feature	Target
1	A	0
1	B	0
1	C	1

Solution: Cross-Validation Target Encoding

Fold	Categorical Feature	Target	Target Mean Encoding
2	A	0	0
3	A	0	0
2	B	1	0.5
3	B	0	0.5

Fold	Categorical Feature	Target	CV Target Encoding
1	A	0	0
1	B	0	
1	C	1	

Solution: Cross-Validation Target Encoding

Fold	Categorical Feature	Target	Target Mean Encoding
2	A	0	0
3	A	0	0
2	B	1	0.5
3	B	0	0.5

Fold	Categorical Feature	Target	CV Target Encoding
1	A	0	0
1	B	0	0.5
1	C	1	

Solution: Cross-Validation Target Encoding

Fold	Categorical Feature	Target	Target Mean Encoding
2	A	0	0
3	A	0	0
2	B	1	0.5
3	B	0	0.5

No info.

Fold	Categorical Feature	Target	CV Target Encoding
1	A	0	0
1	B	0	0.5
1	C	1	NA

Using CV Target Encoding with H2O

H₂O.ai

3.22.1.4

Search docs

Welcome to H2O 3
Quick Start Videos
Cloud Integration
Downloading & Installing H2O
Starting H2O
Getting Data into Your H2O Cluster

Docs • Data Manipulation • Target Encoding [View page source](#)

Target Encoding

Target encoding is the process of replacing a categorical value with the mean of the target variable. In this example, we will be trying to predict `bad_loan` using our cleaned lending club data: <https://raw.githubusercontent.com/h2oai/app-consumer-loan/master/data/loan.csv>.

One of the predictors is `addr_state`, a categorical column with 50 unique values. To perform target encoding on `addr_state`, we will calculate the average of `bad_loan` per state (since `bad_loan` is binomial, this will translate to the proportion of records with `bad_loan = 1`).

For example, target encoding for `addr_state` could be:

addr_state	average bad_loan
AK	0.1476998
AL	0.2091603
AR	0.1920290
AZ	0.1740675
CA	0.1780015
CO	0.1433022

Instead of using state as a predictor in our model, we could use the target encoding of state.

In this topic, we will walk through the steps for using target encoding to convert categorical columns to numeric. This can help improve machine learning accuracy since algorithms tend to have a hard time dealing with high cardinality columns.

The jupyter notebook, [categorical predictors with tree based model](#), discusses two methods for dealing with high cardinality columns:

- Comparing model performance after removing high cardinality columns
- Parameter tuning (specifically tuning `nbins_cats` and `categorical_encoding`)

H2O-3 Documentation:

<http://docs.h2o.ai/h2o/latest-stable/h2o-docs/data-munging/target-encoding.html>

Comparison

Let's compare all three models:

Evaluation (AUC, Higher = Better):

5-Fold CV (with `addr_state`): 0.7045098 vs. (without `addr_state`): 0.7061583 vs. (with TE): 0.7072099

Test (with `addr_state`): 0.7069701 vs. (without `addr_state`): 0.7076197 vs. (with TE) 0.708911

Higher AUC
with TE

My code example:

https://github.com/h2oai/h2o-meetups/tree/master/2019_02_26_Tel_Aviv

CV Target Encoding + Other Feature Engineering Tricks

Competition Round One (Top 100 to Next Round)

kaggle Search kaggle Competitions Datasets Kernels Discussion Jobs ...

Featured Prediction Competition

Zillow Prize: Zillow's Home Value Prediction (Zestimate)

\$1,200,000 Prize Money

Can you improve the algorithm that changed the world of real estate?

Zillow · 3,779 teams · 2 days ago

40 out of 3779 teams

40	▼ 8	Deal or No Deal		0.0749020	79	3mo
41	▲ 52	SCC		0.0749052	39	3mo
42	▼ 31	KFP		0.0749066	349	3mo

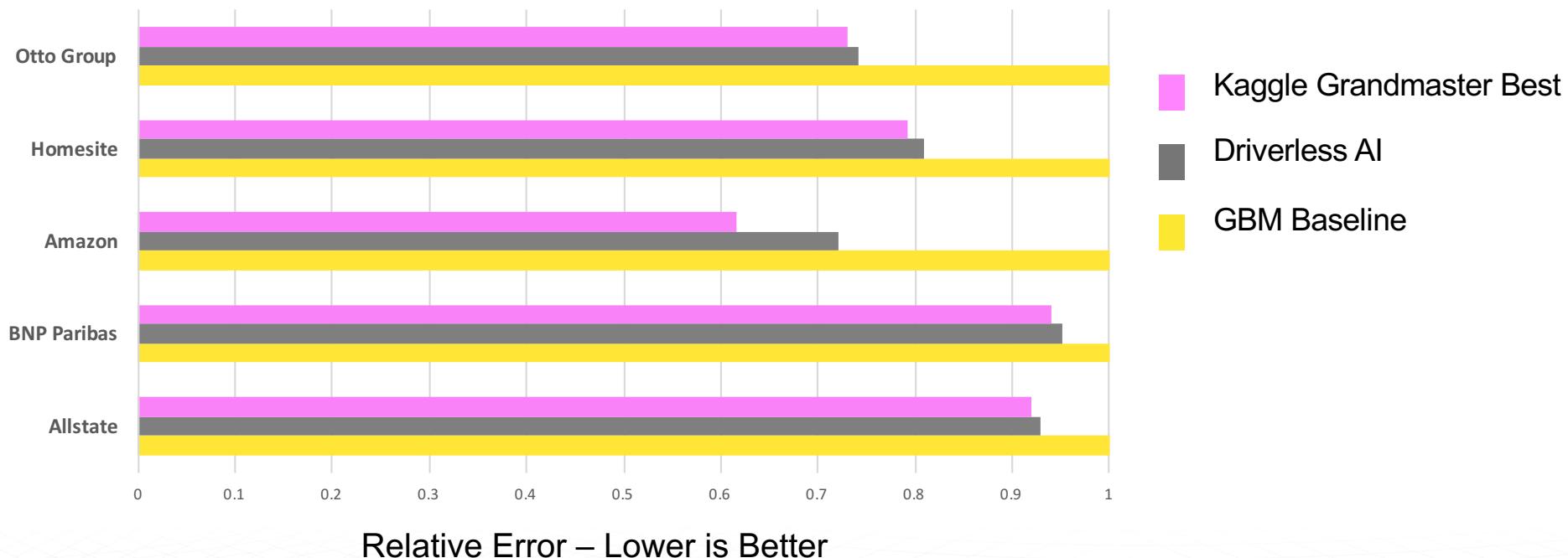


Finished above my H2O Kaggle Grandmasters colleagues 😊

Appendix

How Does Feature Engineering Effect Accuracy?

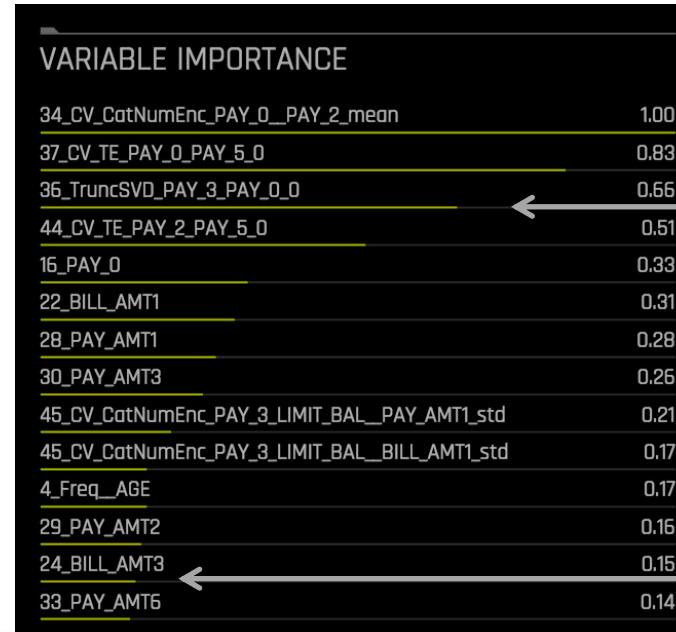
H₂O.ai



Kaggle Grandmaster Out-of-the-Box

Feature Transformations

- Automatic text handling
- Frequency encoding
- Cross validation target encoding
- Truncated singular value decomposition
- Clustering and more



Examples of Generated Features

Examples of Original Features

Accuracy

- Automatic feature engineering to increase accuracy - AlphaGo for AI
- Automatic Kaggle Grandmaster recipes in a box for solving wide variety of use-cases
- Automatic machine learning to find and tune the right ensemble of models

Driverless AI: top 5% in Amazon Kaggle competition

Driverless AI products

Driverless AI: top 1% in BNP Paribas Kaggle competition

single run, **fully automated**; 6h on 3 GPUs

BNP Paribas Cardiff Claims Management

Can you accelerate BNP Paribas Cardiff's claims management process?

\$30,000 · 2,926 teams · 2 years ago

Submission and Description

test_preds.csv

a few seconds ago by [Amao Candel](#)

Driverless AI 1.0.10 10/10/5 on 3 GPUs

Private Score

0.43316

Driverless AI: 18th place in private LB (out of 2926)

Hours for Driverless AI — Weeks for grandmasters

Rank	Team Name	Team Members	Score (0)	Weeks	Last
1	Dexter's Lab	[Icons]	C.42037	198	2w
2	escorted chi	[Icons]	C.42079	162	2w
3	Exploding Kimmie	[Icons]	C.41652	124	2w
4	Brenden Michel utility	[Icons]	C.42219	251	2w
5	the flying bunnies brothers	[Icons]	C.42150	261	2w
6	n_m	[Icons]	C.42358	4	2w
7	PATY	[Icons]	C.42557	310	2w
8	KAME	[Icons]	C.42698	121	2w
9	Jack Upsilon	[Icons]	C.42744	22	2w
10	Dmitry & Bohdan	[Icons]	C.43000	182	2w
11	L1-Der	[Icons]	C.43006	56	2w
12	BNK429RS	[Icons]	C.43059	338	2w
13	x2d4dB	[Icons]	C.43107	65	2w
14	Frenchiez	[Icons]	C.43168	134	2w
15	Altis	[Icons]	C.43174	151	2w
16	SLR 2	[Icons]	C.43213	129	2w
17	no one	[Icons]	C.43217	88	2w

H₂O WORLD 2017