B-jet Momentum Regression Based on Transformer

Zhongqi Zhao, Hengdong Lu, Zhenzhong Tang

UNIVERSITY OF COPENHAGEN





3



Optimization & Techniques





Data Background

B-jet regression

- A b-jet is a beam of particles produced by a collision event between high-energy particles.
- Compare our results with ATLAS prediction.
- Evaluation: average and standard deviation of prediction/truth, mean absolute relative error



Scalar Format Data

• 780206 data

• XGBoost already better than ALTAS!





Scalar Format Data

• 780206 data with 721 characters

• XGBoost already better than ALTAS!

How to make that even better?



Vector format data contain more information!

Vector Format Data

- 3 groups: 'charged_pf', 'neutral_pf', and 'scalar', each containing 780,206 data tensors(highly irregular shape) corresponding to that in scalar format data.
- Scalars in a single vector represent different characters.



Single 'scalar' tensor shape=(1,25) Single 'charged_pf' tensor shape=(70,12) Single 'neutral_pf' tensor shape=(56,46)

Vector Format Data

...,

Lots of invalid data exist in 'charged_pf' and 'neutral_pf', which would cause • huge data noise.

{'char	ged	_pf': a	rray([[(2.2459	9798	, 1.11323	07	, -2.860483	4,	3.79	0048	6, 13	39.57,	-1., 0.,	8208.,	2.245979	7e+03	, 28.127182	, 0.143	0246 ,	True),
	(1.7881	.23 ,	1.2620555	5,	-2.649424	,	3.4143665,	139	9.57,	1.,	0.,	8208.,	1.788123	0e+03,	28.1339	94,	0.01648689,	True),		
	(1.6688	688,	1.1936611	L,	-3.1198075	,	3.0090773,	139	9.57,	1.,	0.,	8208.,	1.668868	8e+03,	-12.7994	23,	0.00056229,	True),		
		,																			
	(0.	,	0.	,	0.	,	0.,	e	ð.,	0.,	0.,	0.,	0.00000	0e+00,	0.	,	0.,	False),		
	(0.	,	0.	,	0.	,	0.,	6	ə.,	0.,	0.,	0.,	0.00000	0e+00,	0.	,	0.,	False),		
	(0.	,	0.	,	0.	,	0. ,	6	ə.,	0.,	0.,	0.,	0.00000	0e+00,	0.	,	0. ,	False)],	_	
	[(2.6027	467,	2.2941675	5,	-1.4546521	. ,	13.036634 ,	139	9.57,	-1.,	0.,	8208.,	2.602746	8e+03,	-2.8193	183,	0.06687441,	True),	1	
	(1.1454	937,	2.1662593	З,	-1.4887234	. ,	5.0651135,	139	9.57,	1.,	0.,	8208.,	1.145493	8e+03,	-10.8400	13,	-0.01959403,	True),		
	(1.0948	3113,	2.4059186	5,	-1.3677806	,	6.1209264,	139	9.57,	1.,	0.,	8208.,	1.094811	3e+03,	-3.3708	525,	-0.04104954,	True),	1.6	onsor in 'cha
	- 22	,																		10	
	(0.	,	0.	,	0.	,	ø. ,	e	ə.,	0.,	0.,	0.,	0.00000	0e+00,	0.	,	0. ,	False),	(to	tally (80000)
	(0.	,	0.	,	0.	ر	0. ,	e	ə.,	0.,	0.,	0.,	0.00000	0e+00,	0.	,	0. ,	False),		
	(0.		0.	,	0.	,	0. ,	e	ə.,	0.,	0.,	0.,	0.00000	0e+00,	0.	,	0. ,	False)],		
	[(4.5994	387,	-1.7090869	Э,	1.790061	,	13.120517 ,	139	9.57,	-1.,	0.,	8208.,	4.599438	5e+03,	-21.3337	2,	0.03466429,	True),		
	(3.5137	36,	-1.572401	,	1.8378757	,	8.830682 ,	139	9.57,	-1.,	0.,	8208.,	3.513736	1e+03,	-19.4075	58,	-0.02446163,	True),		
	(3.0004	964,	-1.9127998	з,	1.7642748	,	10.382189 ,	139	9.57,	1.,	0.,	8208.,	3.000496	3e+03,	-21.3562	28,	-0.05878485,	True),		
		,																			
	(0.	,	0.	,	0.	,	0.,	6	ð.,	0.,	0.,	0.,	0.00000	0e+00,	0.	,	0.,	False),		
noice	(0.	,	0.	,	0.	,	ø. ,	6	ə.,	0.,	0.,	0.,	0.00000	0e+00,	0.	,	0. ,	False),		
	(0.	,	0.	,	0.	,	ø.,	e	э.,	0.,	0.,	0.,	0.00000	0e+00,	0.	,	0.,	False)],		
		• ,															63				
	[(68.2474	з,	-0.2006337	7,	2.8248081	,	69.62579 ,	139	9.57,	1.,	0.,	8208.,	6.824743	0e-16,	53.3183	44,	-0.01279807,	True),		
	(49.4440	96,	-0.2111132	26,	2.8046072	,	50.550217 ,	139	9.57,	-1.,	0.,	8208.,	4.871319	8e+03,	53.2339	36,	0.03639252,	True),		
	(-	47.1472	.85 ,	-0.2091238	32,	2.828629	,	48.18219 ,	139	9.57,	1.,	0.,	8208.,	5.655724	1e+03,	53.3407	в,	0.01992808,	True),		

8

12/06/2024

Vector Format Data

 Lots of invalid data exist in 'charged_pf' and 'neutral_pf', which would cause huge data noise.

{'chai	rged_pf': arra	ay([[(2.24597	98,	1.113230	97	, -2.86048	34	, 3.79	9004	86,	139.57	, -	-1., 0., 8208.,	2.24	459797e+0	з,	28.127182	, 0.143	30246 ,	True),	
	(1.788123	, 1.	2620555	, -2	2.649424	,	3.4143665	, 1	39.57,	1.	, 0.	, 8208	۰,	1.7881230e+03,	28.	.133904 ,	0.	01648689,	True),			
	(1.6688688	3, 1.	1936611	, -3	3.1198075	,	3.0090773	, 1	39.57,	1.	, 0.	, 8208	.,	1.6688688e+03,	-12.	, 799423	0.	00056229,	True),			
	,																					
	(0.	, 0.		, e	Э.	,	0.	,	0.,	0.	, 0.	, 0	۰,	0.0000000e+00,	0.	. ,	0.	,	False),			
	(0.	, 0.		, e	Э.	,	0.	,	0.,	0.	, 0.	, 0	۰,	0.0000000e+00,	0.	. ,	0.	,	False),			
	(0.	, 0.		, e	э.	,	0.	,	0.,	0.	, 0.	, 0	.,	0.0000000e+00,	0.	. ,	0.	,	False)]			
	[(2.602746]	7, 2.	2941675	, -1	1.4546521	,	13.036634	, 1	39.57,	-1.	, 0.	, 8208	.,	2.6027468e+03,	-2.	.8193183,	0.	06687441,	True),			
	(1.145493	7, 2.	1662593	, -1	1.4887234	,	5.0651135	, 1	39.57,	1.	, 0.	, 8208	۰,	1.1454938e+03,	-10.	.840013 ,	-0.	01959403,	True),			
	(1.0948113	3, 2.	4059186	, -1	1.3677806	,	6.1209264	, 1	39.57,	1.	, 0.	, 8208	.,	1.0948113e+03,	-3.	3708525,	-0.	04104954,	True),	1	tonsor in 'charged	l nf'
	,																				tellsof in charged	<u>_</u> pi
	(0.	, 0.		, e	э.	,	0.	,	0.,	0.	, 0.	, 0	.,	0.0000000e+00,	0.	. ,	0.	,	False),	(10	otally 780000)	
	(0.	, 0.		, e	э.	,	0.	,	Ø.,	0.	, 0.	, 0	۰,	0.0000000e+00,	0.	. ,	0.	,	False),			
	(0.	. 0.		, e	Э.	,	0.	,	0.,	0.	, 0.	, 0	.,	0.0000000e+00,	0.	. ,	0.	,	False)],			
	[(4.599438]	7, -1.	7090869	, 1	1.790061	,	13.120517	, 1	39.57,	-1.	, 0.	, 8208	۰,	4.5994385e+03,	-21.	, 33372	0.	03466429,	True),			
	(3.513736	, -1.	572401	, 1	1.8378757	,	8.830682	, 1	39.57,	-1.	, 0.	, 8208	۰,	3.5137361e+03,	-19.	.407558 ,	-0.	02446163,	True),			
	(3.0004964	4, -1.	9127998	, 1	1.7642748	,	10.382189	, 1	39.57,	1.	, 0.	, 8208	۰,	3.0004963e+03,	-21.	, 356228	-0.	05878485,	True),			
	,																					
	(0.	, 0.		, e	Э.	,	0.	,	0.,	0.	, 0.	, 0	.,	0.0000000e+00,	0.	. ,	0.	,	False),			
noice	(0.	, 0.		, e	э.	,	0.	,	0.,	0.	, 0.	, 0	۰,	0.0000000e+00,	0.	. ,	0.	,	False),			
	(0.	, 0.		, e	Э.	,	0.	,	0.,	0.	, 0.	, 0	.,	0.0000000e+00,	0.	. ,	0.	,	False)],			
	,																		9 <u>85</u> 5950			
	[(68.24743	, -0.	2006337	, 2	2.8248081	,	69.62579	, 1	39.57,	1.	, 0.	, 8208	۰,	6.8247430e-16,	53.	.318344 ,	-0.	01279807,	True),			
	(49.444096	, -0.	21111326	, 2	2.8046072	,	50.550217	, 1	39.57,	-1.	, 0.	, 8208	۰,	4.8713198e+03,	53.	.233936 ,	0.	03639252,	True),			
	(47.147285	, -0.	20912382	, 2	2.828629	,	48.18219	, 1	39.57,	1.	, 0.	, 8208	۰,	5.6557241e+03,	53.	.34078 ,	0.	01992808,	True),			
	,																					

Too difficult for Tree-based Methods !

Transformer

- Self-attention
 - Long distance dependency capture
 - Multi-head attention(Multi self-attention with different weight)
- Inter-layer parallelism

Suitable for such complex and highly irregularly shaped data



Transformer Architecture

Transformer : Vectors input

















How to improve it to beat XGBoost?



• Point 1: Attention WITHIN each group

Individual attention for charged_pf, neutral_f

• Point 2: Is decoder necessary?

Simplify transformer











12/06/2024 23

Optimization & Techniques



- Layer architecture
- Hyperparameter: Dropout rate, Learning rate, Nheads...
- Loss function: MAE, MRE, Log_response, MSE...

Tensor Processing Unit (TPU)

- Specialization
- High-Efficiency Computing
 - Low Power Consumption

Good for Attention Mechanism

[(48.866028	,	-1.0349238	,	-0.8517365	,	77.45633	,	139.57	,	1.,	0.,	8208.,	7.0004043e+03,	12.696693	,	-0.07532337	,	True),
	(39.858517	,	-1.0368809	,	-0.8564458	5,	63.27485	,	139.57	7,	-1.,	0.,	8208.,	1.0329485e+04,	12.742829	,	0.04961777	,	True),
	(31.882801	,	-1.0423672	,	-0.8490515	ŧ,	50.830006	,	139.57	1,	-1.,	0.,	8208.,	1.1534389e+04,	12.609667	,	-0.2534103	,	True),
	,		0		0		0		0		0	0	0	0.000000.000	0		0		
	(0.	,	υ.	,	0.	,	0.	,	υ.	,	v.,	v.,	v.,	0.0000000e+00,	0.	,	0.	,	Faise),
	(0.	,	0.	,	0.	,	0.	,	0.	,	0.,	0.,	0.,	0.0000000e+00,	0.	,	0.	,	False),
	(0.	,	0.	,	0.	,	0.	,	0.	,	0.,	0.,	0.,	0.0000000e+00,	0.	,	0.	,	False)],
Ľ	(81.49732	,	-0.5086934	,	-2.1938472	,	92.27127	,	139.57	7,	-1.,	0.,	8208.,	8.1497317e-16,	20.954634	,	-0.04437601	,	True),
	(51.169197	,	-0.5012035	4,	-2.2005699	,	57.732014	,	139.57	7,	1.,	0.,	8208.,	4.2174072e+03,	20.835516	,	0.0185816	,	True),
	(11.600805	,	-0.510175	,	-2.191241	,	13.144298	,	139.57	Ϊ,	1.,	0.,	8208.,	5.5604448e+03,	21.237577	,	-0.15153563	,	True),
	••••																		
	(0.	,	0.	,	0.	,	0.	,	0.	,	0.,	0.,	0.,	0.0000000e+00,	0.	,	0.	,	False),
	(0.	,	0.	,	0.	,	0.	,	0.	,	0.,	0.,	0.,	0.0000000e+00,	0.	,	0.	,	False),
	(0.	,	0.	,	0.	,	0.	,	0.	,	0.,	0.,	0.,	0.0000000e+00,	0.	,	0.	,	False)]],

Learning Rate

"But, there is no reason to consider a fixed value for the learning rate!" [Troels C. Petersen (NBI)]

```
initial_learning_rate = 0.0001
lr_schedule = tf.keras.optimizers.schedules.ExponentialDecay(
    initial_learning_rate,
    decay_steps=10000,
    decay_rate=0.96,
    staircase=True)
```

Early Stopping and CallBacks

Saves us a lot of time and allows us not to miss the best model.



Results and Comparison

BDT (XGBoost)



Long Short-Term Memory (LSTM)



Transformer

ATLAS_mre = 0.1582 XGB_mre = 0.1400 LSTM_mre = 0.1316 Trans_mre = 0.1247



Resolution Comparison (Transformer VS ATLAS)



Transformer Result in Different Momentum Intervals





Questions & Comments?

e.g.

Dropout(0.1)

- It means that each neuron in this layer has a 10% chance of being dropped out at each training step. This operation only occurs during training.
- And during testing or prediction, the *Dropout layer* adjusts its outputs based on the drop rate to ensure the total output remains consistent.

Bias

It could be related to loss function, layer structure, data volume and embedding dimension.....





Background Kernel Density Estimation of Prediction to True Ratios Transformer: mu = 1.04e + 00std = 1.79e - 013.0 ATLAS: mu = 1.04e+00std=2.35e-01 Embedding dim= XGBoost: mu = 1.04e+00std=1.98e-01 --- Perfect Prediction 2.5 256 2.0 Density 1.5 1.0 0.5 0.0 + 0.5 1.0 0.0 1.5 2.0 2.5 3.0 Prediction / True



Modified Transformer Loss vs. Epoch





Transformer tried



12/06/2024 45

Transformer tried





All participants contributed evenly