



Kaggle Winner Presentation (6th)

Google - Fast or Slow? Predict AI Model Runtime

<https://www.kaggle.com/competitions/predict-ai-model-runtime>

HENG CHER KENG

<https://www.kaggle.com/hengck23>

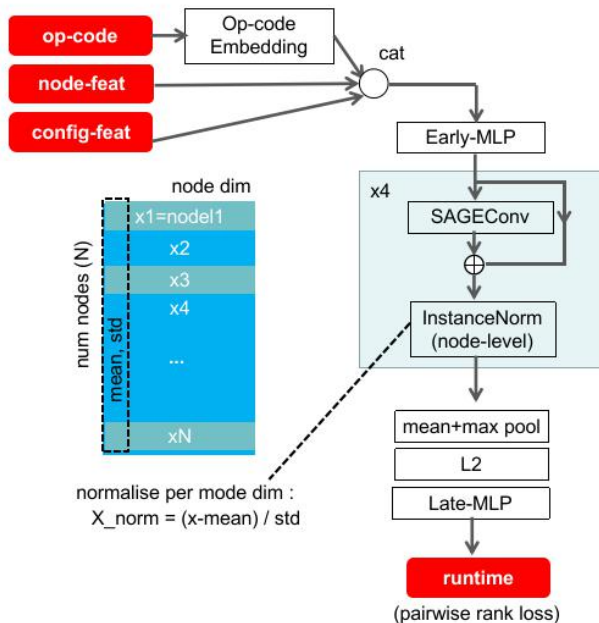
Agenda

1. Background
2. Summary
3. Feature selection & engineering
4. Training methods
5. Important findings
6. Simple model

Background

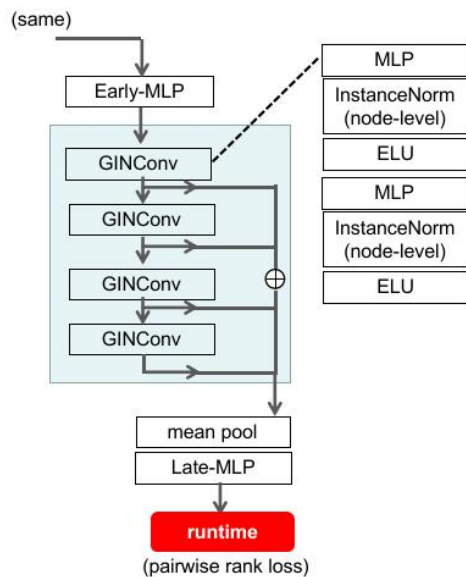
- Contract computer vision and deep learning algorithm engineer.
 - find fracture in x-ray images
 - implement visual slam for robotic navigation.
- Familiar with deep learning and build deep models in my work.
- Experiences in graph neural net GNN from previous Kaggle competitions.

Summary



Layout runtime prediction

Residual SAGEConv Net

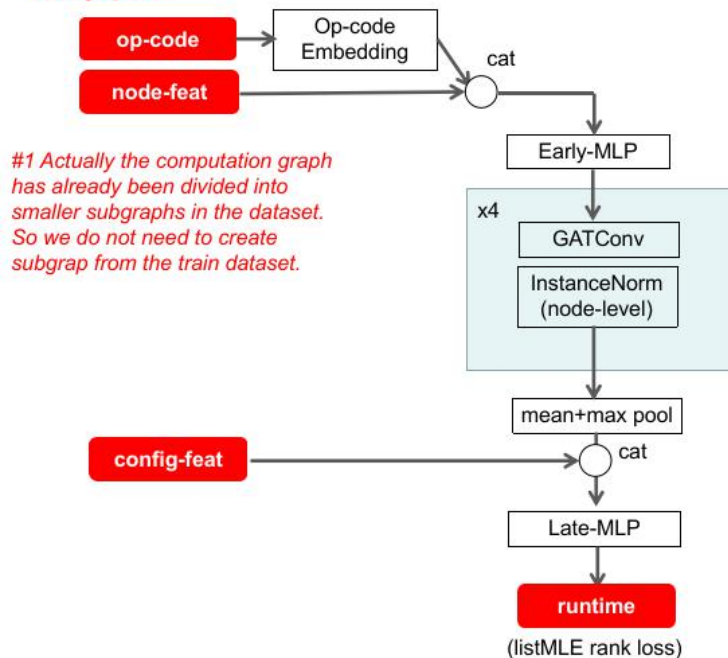


Graph Isomorphism Net (GIN)

For layout runtime prediction:

- Reduce input graph to subgraph by including only the 5-hop neighbours from the configure node
- GNN model with 4-layer SAGE-conv[2] with residual shortcut
- GNN model with 4-layer GIN-conv[3]
- Graph instance normalisation[1] over node

"Full" graph #1



**Tile runtime
prediction**

**GATConv
Net**

For tileruntime prediction:

- GNN model with 4-layer GAT-conv[4]

- Pytorch geometric used for building GNN.
- Training takes abopout 4 to 6 hours for one model in each collection, using Nvidia Quadro RTX 8000 GPU with memory 48 GB.

[1] "Learning Graph Normalization for Graph Neural Networks" - Yihao Chen
<https://arxiv.org/abs/2009.11746>

[2] "Inductive Representation Learning on Large Graphs" - William L. Hamilton
<https://arxiv.org/abs/1706.02216>

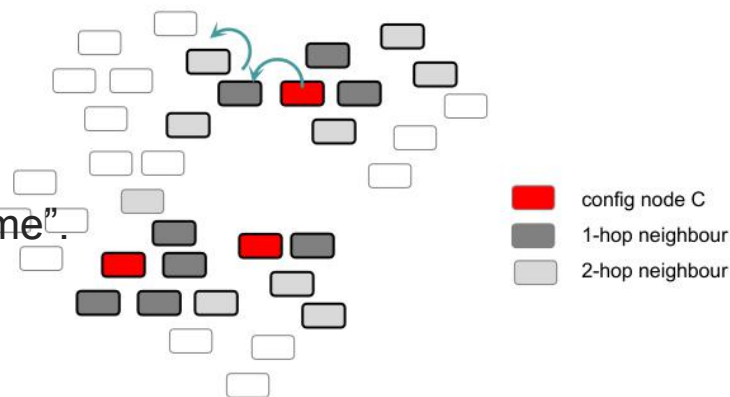
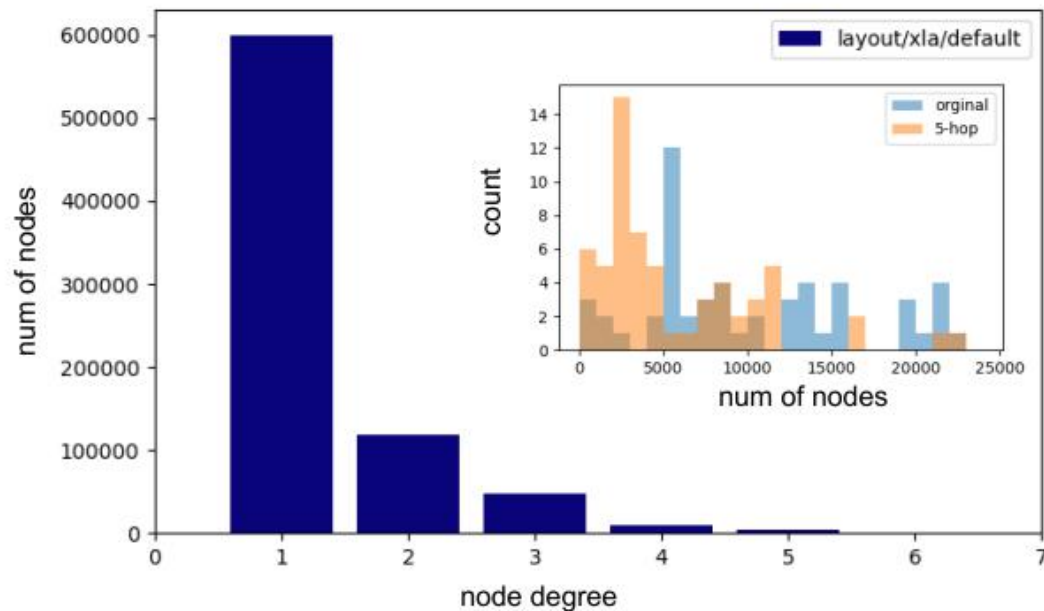
[3] "How Powerful are Graph Neural Networks?" - Keyulu Xu
<https://arxiv.org/abs/1810.00826>

[4] "Graph Attention Networks" - Petar Veličković
<https://arxiv.org/abs/1710.10903>

Features Selection/ Engineering

5-hop neighbourhood subgraph as input:

- config node + their neighbours
- relative runtime ranking, just need to consider the “difference of two graph” to predict the “difference in runtime”.



		XLA/Layout
		Default
graphsage GNN		kendall tau
	3-hop	0.41260
	4-hop	0.42503
	5-hop	0.45508
	6-hop	0.43508

- Ensemble should improve results

		collection				
		NLP/Layout		XLA/Layout		XLA/Tile
		Default	Random	Default	Random	
[a]	4x-gatconv-lstmle					0.97462
[b]	4x-graphsage-pair2	0.53938	0.92654	0.45508	0.67128	
[c]	4x-gin-pair2			0.45958		
[b]+[c]	ensemble			0.46952		
submission		0.53938	0.92654	0.46952	0.67128	0.97462
					avg	0.71627
					public lb	0.69424
					private lb	0.70549

For GIN Net, we don't have time to train for all layout prediction before the competition ends.

Training Methods

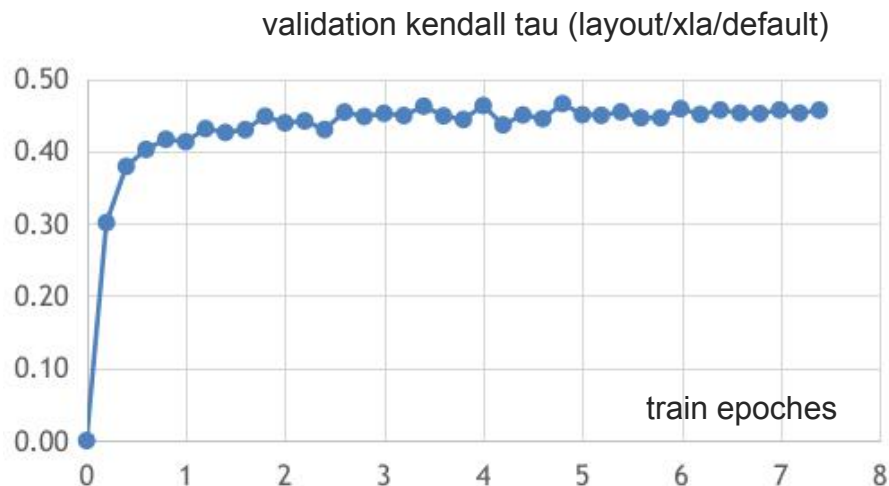
- For layout target kendall tau, use pairwise rank loss.
- Tile top5 slowdown target use listMLE
- ADAM optimizer with fixed lr 0.0005
- Use stochastic weight averaging SWA.
 final weights = average last 10 trained weights
 +0.01 improvement over best model
- batch size = 32, sample 80 to 100 configurations per subgraph

- gradient accumulation

```
optimizer.zero_grad()

for b in range(batch_size):
    r = batch[r]
    loss = net(r) # forward one subgraph
    scaler.scale(loss).backward() #backward accumulate gradient

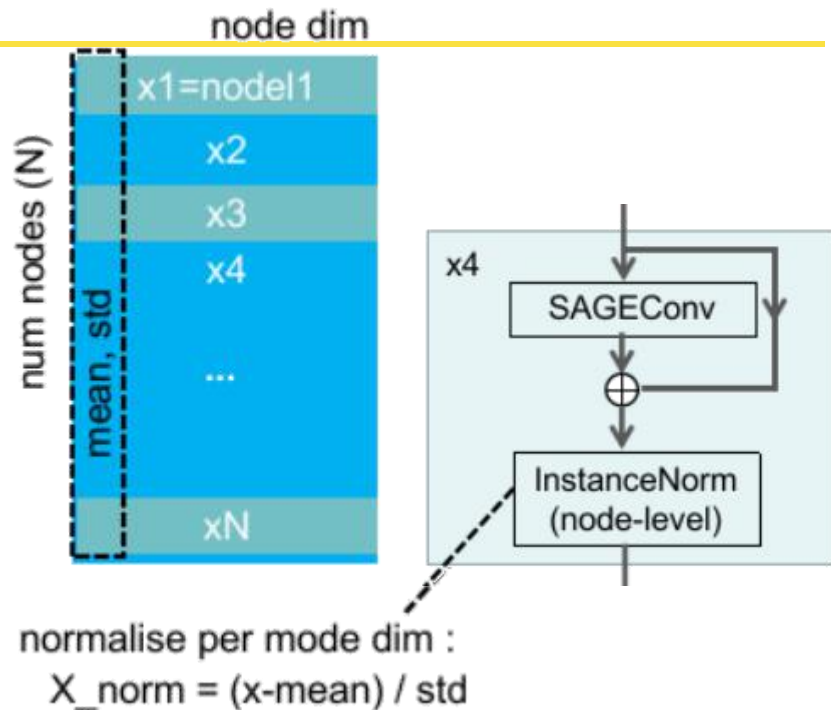
scaler.step(optimizer) #update net parameters
scaler.update()
```



Important and Interesting Findings

- Graph instance normalisation gives faster convergence and improve validation kendall tau by 0.02
- Attention pooling at graph readout for layout runtime prediction improves validation 0.01 for some cases

this solution is not selected because it did not improve public LB. (But private LB is better)



- validation kendall tau has large variance. Careful to select final submission to avoid “shakeup” for private LB.

res-graphsage4-001-xla-default-pair-hop5			
/checkpoint/00003060.pth'			
0	0.433932	tf2_bert_pretrain_dynamic_batch_size.npz	
1	0.506814	bert_pretraining.4x4.fp16.npz	
2	0.348933	mlperf_bert_batch_24_2x2.npz	
3	0.681526	inception_v3_batch_128_train.npz	
4	0.552380	resnet_v1_50_official_batch_128_bf16.npz	
5	0.593138	resnet50.4x4.fp16.npz	
6	0.089296	unet_3d.4x4.bf16.npz	
kendall_tau		0.4580026617688869	
opa_acc		0.7290013308844434	

Simple Model

- Possible to use even smaller 4-hop neighbourhood subgraph for better speed at the expense of 5 to 10% loss in accuracy
- Number of graph conv layers can also be reduced. We find
num of layer = num of hop + 1

		XLA/Layout
		Default
graphsage GNN		kendall tau
	3-hop	0.41260
	4-hop	0.42503
	5-hop	0.45508
	6-hop	0.43508

Question and Answer



kaggle