

Efficient Multi-Task Auxiliary Learning: Selecting Auxiliary Data by Feature Similarity



http://github.com/MiuLab/FastMTL

Po-Nien Kung, Yi-Cheng Chen, Sheng-Siang Yin
Tse-Hsuan Yang, Yun-Nung (Vivian) Chen





- Background
- Two-Stage Multi-Task Auxiliary Learning
- Efficient Data Selection
- Experiments
- Conclusion



Background

- What is Multi-Task Auxiliary Learning
- Two-Stage Multi-Task Auxiliary Learning
- Efficient Data Selection
- Experiments
- Conclusion

Background: Multi-Task Auxiliary Learning

Multi-task learning
Multi-task auxiliary learning



All tasks are important!



One primary task and multiple auxiliary tasks.

Background: Multi-Task Auxiliary Learning

 To achieve better performance on the primary task

5

 More useful when the size of the primary task is small



Background: Multi-Task Learning Chronology

More tasks, more **power**...





More tasks, more **power**... Also, more **computing**!

Compare single-task finetuning and multi-task auxiliary learning when the primary task is RTE:

In MTDNN setting...

In *Muppet* setting...

RTE (2.4k) vs Auxiliary(960k) data

400x Computing cost!

RTE (2.4k) vs Auxiliary (4.8mil) data

2000x Computing cost!



Background

- What is Multi-Task Auxiliary Learning
- Why we need an efficient Multi-Task Auxiliary Learning method
- Two-Stage Multi-Task Auxiliary Learning
- Efficient Data Selection
- Experiments
- Conclusion

Inefficient Multi-Task Auxiliary Learning



9

Less Auxiliary Data is Possible?

The first question: Are all auxiliary data **beneficial**?

- Single-task vs. multi-task
 - GLUE dataset (similar with MTDNN)
 - 10 random seeds
 - Use 1 STDEV as the threshold for performance MPROVED and DROPPED Otherwise, NEUTRAL.



Negative Transfer !

Less Auxiliary Data is Possible?

The first question: Are all auxiliary data **beneficial**?

- Single-task vs. multi-task
 - GLUE dataset (similar with MTDNN)



Some auxiliary dataset might be unhelpful or even harmful!

Negative Transfer !



Background

Two-Stage Multi-Task Auxiliary Learning

- We need a **data-sampling** method to shrink the size of the auxiliary data
- Efficient Data Selection
- Experiments
- Conclusion

— Two-Stage Multi-Task Auxiliary Learning

13

Multi-task Auxiliary Learning Task-Oriented Predictors Sampling Method Select the most beneficial auxiliary data Auxiliary Sub-set Primary Data Goal: reducing the cost of **Auxiliary Data** training auxiliary data

Prior Work: AutoSeM (Guo et al., 2019)

- Idea: automatically select the most beneficial (related) auxiliary tasks
 - Beta-Bernoulli multi-armed bandit with Thompson Sampling
- Decide the mixing ratio of auxiliary tasks
 - Gaussian Process
 - Trial and error

Avoid negative transfer and further improve the performance!

Target 1: Reduce auxiliary dataset size ?

Target 2: Reduce the total computing cost ?



The sampling method itself is **Computationally Expensive**!

Challenges of the Sampling Method







- Background
- Two-Stage Multi-Task Auxiliary Learning
- Efficient Data Selection
 - Select the most beneficial auxiliary data by feature similarity
- Experiments
- Conclusion

Feature Similarity Assumption

Assumption: More **similar** is an auxiliary data to the primary task, more **benefit** it can bring.



Primary task data



Useful auxiliary data





Peature Similarity Assumption

Toy experiment

- MT-DNN setting: multi-task train 500 data for each GLUE task
- T-SNE Visualization
 - Last hidden state features of BERT

MNLI	RTE	MRPC	STS-B	QQP	QNLI	SST-2	CoLA
		\mathbf{C}	\mathbf{C}	_			

SST-2 and CoLA are more separate from others.



Feature Similarity Assumption

- Adding task-discriminator
 - Force the model to encode more task information into features



STS-B, RTE, MRPC and MNLI have more data **overlap** with each other.

The tasks with more similar auxiliary data improve most!



Usefulness of Auxiliary Data

 Feature similarity may indicate the usefulness of auxiliary data to a primary task.

How do we get the rank the feature similarity?

Train a small proxy model with a task discriminator to predict similarity.

21—Data Selection: Similarity Ranking



- Two-Stage Multi-Task Auxiliary Learning

Stage 2: Multi-task Auxiliary

Efficient!

Stage 1: Similarity Ranking

22



Goal: multi-task auxiliary learning on less auxiliary data but comparable (or even improved) performance



- Background
- Two-Stage Multi-Task Auxiliary Learning
- Efficient Data Selection
- Experiments
- Conclusion



Primary tasks

- Similar with MT-DNN
- Data: GLUE (960K)
- Model: Bert-base
- Baselines:
 - No-MTL (Weak)
 - Random Sampling (Surprisingly Strong)
 - Fully-trained (Strong)

MNLIRTEMRPCSTS-BQQPQNLISST-2CoLAImage: Color of the state of the

They are improved by MTL, so there exist useful data in auxiliary tasks













Findings Ours > Random Ours > Fully-Trained Random > Fully Trained (STS-B)

Our method can use **less data** to achieve **better performance**!

28 Efficiency Evaluation

Item With the second state of the second st

Ours: 50%, 60%, 0.05% Random: 100%, 100%, 1%

	Similarity	Sampling	Multi-Task Au		How	
STS-B Runtime(s)	Training small proxy model	Predict similarity	MTL	Finetuning	Total	much faster?
Fully-trained			15801		15991	
Random	-	-	260	190	450	35x
Ours	95	775	200		670	23x



- Background
- Two-Stage Multi-Task Auxiliary Learning
- Efficient Data Selection
- Experiments
- Conclusion



- Address the efficiency importance in multi-task auxiliary learning
- Propose a data sampling method to shrink the size of the auxiliary data → computing cost reduction
- First use **feature similarity** to determine the data usefulness
- Our method outperforms random sampling and further surpass fullytrained model using less data

First work for time-efficiency of multi-task auxiliary learning: http://github.com/MiuLab/FastMTL

