# AutoSeM

Automatic Task Selection and Mixing in Multi-Task Learning Han Guo, Ramakanth Pasunuru, Mohit Bansal



### Overview

## 1.Introduction 2.Methods 3. Experiments





2



### Overview

# 1.Introduction 2.Methods

3. Experiments



З



### Introduction

- Multi-task Learning (MTL) is an <u>inductive</u> transfer mechanism which leverages information from related tasks to improve the primary model's generalization performance.
- It achieves this goal by training **multiple** tasks in **parallel** while **sharing representations**, where the training signals from the auxiliary tasks can help improve the performance of the primary task.















<u>Auxiliary Task Selection</u>

 $\mathcal{O}(2^{\kappa})$ 



### **TWO ISSUES**











### Mixing Ratio Learning









### <u>Auxiliary Task Selection</u>



### AutoSeM

### Mixing Ratio Learning







### AutoSeM

### Multi-Armed 777 Bandit

### <u>Auxiliary Task Selection</u>





### Mixing Ratio Learning





## **Related Works**

### Multi-Task Learning







Collobert and Weston, 2008; Girshick, 2015; Luong et al., 2015; Pasunuru and Bansal, 2017; Pasunuru et al., 2017

Misra et al.,2016; Kendall et al., 2017; Dai et al., 2016

Teh et al., 2017; Parisotto et al., 2015; Jaderberg et al., 2016

### Shared Parameter Selection

Ruder et al. (2017); Xiao et al. (2018)

### Identifying Task Relatedness

Ben-David and Schuller (2003), Bingel and Søgaard (2017)



### Data Selection/Reordering

van der Wees et al., 2017; Moore and Lewis, 2010; Duh et al., 2013; Søgaard, 2011; Ruder and Plank (2017); Tsvetkov et al. (2016)

### Multi-Armed Bandit

Graves et al. (2017); Sharma and Ravindran(2017)

### Exploitation/Exploration

Kaelbling et al., 1996; Auer et al., 2002b; Auer et al., 2002a; Houthooft et al., 2016; Russo et al., 2018; Chapelle and Li, 2011

### Gaussian Process

Rasmussen, 2004; Shahriari et al., 2016; Schulz et al.,2018; Snoek et al., 2012; Brochu et al., 2010; Knudde et al., 2017; Cully et al., 2018; Swersky et al., 2013; Golovin et al., 2017)



### Overview

## 1.Introduction 2.Methods 3. Experiments





9











### Base Model

Peters et al., (2018); Conneau et al., (2017)



### Multi-Task Model





### Peters et al., (2018); Conneau et al., (2017)

















## **Auxiliary Task Selection** Beta Distribution



Prior Knowledge





## **Auxiliary Task Selection** Beta Distribution









## **Auxiliary Task Selection** Beta Distribution



















Algorithm 1 BernThompson $(N, \alpha, \beta, \gamma, \alpha_0, \beta_0)$ 

1: for  $t_b = 1, 2, ...$  do # sample model: 2: 3: for k = 1, ..., N do Sample  $\hat{\theta}_k \sim \text{Beta}(\alpha_k, \beta_k)$ 4: 5: end for 6: # select and apply action: 7:  $x_{t_h}^s \leftarrow \arg \max_k \hat{\theta}_k$ 8: Apply  $x_{t_b}$  and observe  $r_{t_b}$ 9: # non-stationarity for k = 1, ..., N do 10: 11:  $\hat{\alpha}_k = (1 - \gamma)\alpha_k + \gamma\alpha_0$  $\hat{\beta}_k = (1 - \gamma)\beta_k + \gamma\beta_0$ 12: 13: if  $k \neq x_{t_h}^s$  then  $(\alpha_k, \beta_k) \leftarrow (\hat{\alpha}_k, \hat{\beta}_k)$ 14: 15: else  $(\alpha_k, \beta_k) \leftarrow (\hat{\alpha}_k, \hat{\beta}_k) + (r_{t_b}, 1 - r_{t_b})$ 16: 17: end if 18: end for 19: **end for** 









Algorithm 1 BernThompson $(N, \alpha, \beta, \gamma, \alpha_0, \beta_0)$ 

1: for  $t_b = 1, 2, ...$  do # sample model: 2: 3: for k = 1, ..., N do Sample  $\hat{\theta}_k \sim \text{Beta}(\alpha_k, \beta_k)$ 4: 5: end for 6: # select and apply action:  $x_{t}^{s} \leftarrow \arg \max_{k} \hat{\theta}_{k}$ 7: Apply  $x_{t_b}^s$  and observe  $r_{t_b}$ 8: 9: # non-stationarity 10: for k = 1, ..., N do  $\hat{\alpha}_k = (1 - \gamma)\alpha_k + \gamma\alpha_0$ 11:  $\hat{\beta}_k = (1 - \gamma)\beta_k + \gamma\beta_0$ 12: 13: if  $k \neq x_{t_h}^s$  then  $(\alpha_k, \beta_k) \leftarrow (\hat{\alpha}_k, \hat{\beta}_k)$ 14: 15: else  $(\alpha_k, \beta_k) \leftarrow (\hat{\alpha}_k, \hat{\beta}_k) + (r_{t_b}, 1 - r_{t_b})$ 16: 17: end if 18: end for 19: **end for** 









Algorithm 1 BernThompson $(N, \alpha, \beta, \gamma, \alpha_0, \beta_0)$ 

1: for  $t_b = 1, 2, ...$  do # sample model: 2: 3: for k = 1, ..., N do Sample  $\hat{\theta}_k \sim \text{Beta}(\alpha_k, \beta_k)$ 4: 5: end for 6: # select and apply action:  $x_{t_b}^s \leftarrow \arg \max_k \hat{\theta}_k$ 7: Apply  $x_{t}^{s}$ , and observe  $r_{t_{t}}$ 8: 9: # non-stationarity for k = 1, ..., N do 10:  $\hat{\alpha}_k = (1 - \gamma)\alpha_k + \gamma\alpha_0$ 11:  $\hat{\beta}_k = (1 - \gamma)\beta_k + \gamma\beta_0$ 12: 13: if  $k \neq x_{t_h}^s$  then  $(\alpha_k, \beta_k) \leftarrow (\hat{\alpha}_k, \hat{\beta}_k)$ 14: 15: else  $(\alpha_k, \beta_k) \leftarrow (\hat{\alpha}_k, \hat{\beta}_k) + (r_{t_b}, 1 - r_{t_b})$ 16: 17: end if 18: end for 19: **end for** 









# Mixing Ratio Learning Mixing Ratio Multi-Task

Model

Performance



Rasmussen, 2004; Snoek et al., 2012; Shahriari et al., 2016









### Mixing Ratio Learning Mixing Ratio MR-3 MR-2 MR-1 Multi-Task

Model

Performance















### Mixing Ratio Learning Mixing Ratio MR-3 MR-2 MR-1 Multi-Task Next Sample

Model

Performance













### Mixing Ratio Learning Mixing Ratio MR-3 MR-2 MR-1 Multi-Task Next Sample

Model

Performance













### Overview

## 1.Introduction 2.Methods **3. Experiments**



25



### Dataset

Corpus	Train	Test	Task	Metrics	Domain		
	Single-Sentence Tasks						
CoLA SST-2	8.5k 67k	<b>1k</b> 1.8k	acceptability sentiment	Matthews corr. acc.	misc. movie reviews		
Similarity and Paraphrase Tasks							
MRPC QQP	3.7k 364k	1.7k <b>391k</b>	paraphrase paraphrase	acc./F1 acc./F1	news social QA questions		
	Inference Tasks						
MNLI QNLI RTE WNLI	393k 105k 2.5k 634	<b>20k</b> 5.4k 3k <b>146</b>	NLI QA/NLI NLI coreference/NLI	matched acc./mismatched acc. acc. acc. acc.	misc. Wikipedia news, Wikipedia fiction books		



Wang et al., 2018; Warstadt et al., 2018; Socher et al., 2013; Dolan & Brockett, 2005; Cer et al., 2017; Williams et al., 2018; Bowman et al., 2015; Rajpurkar et al. 2016; Dagan et al., 2006; Bar Haim etal., 2006; Giampiccoloetal., 2007; Bentivoglietal., 2009; Levesque et al., 2011



## Dataset Example

### Premise

### Fiction

The Old One always comforted Ca'daan, except t

### Letters

Your gift is appreciated by each and every stude from your generosity.

### Telephone Speech

yes now you know if if everybody like in August v on vacation or something we can dress a little m

### 9/11 Report

At the other end of Pennsylvania Avenue, people a White House tour.



	Label	Hypothesis
t today.	neutral	Ca'daan knew the Old One very well.
ent who will benefit	neutral	Hundreds of students will benefit from your generosity.
when everybody's more casual or	contradiction	August is a black out month for vacations in the company.
le began to line up for	entailment	People formed a line at the end of Pennsylvania Avenue.



Models

BiLSTM+ELMo (Single-Task) (Wang et al., 2 BiLSTM+ELMo (Multi-Task) (Wang et al., 2

Our Baseline Our AUTOSEM



	RTE	MRPC	QNLI	CoLA	SST-2
2018)	50.1	69.0/80.8	69.4	35.0	90.2
2018)	55.7	76.2/83.5	66.7	27.5	89.6
	54.0	75.7/83.7	74.0	30.8	91.3
	58.7	78.5/84.5	79.2	32.9	<b>91.8</b>



Models

BiLSTM+ELMo (Single-Task) (Wang et al., 2 BiLSTM+ELMo (Multi-Task) (Wang et al., 20

Our Baseline Our AUTOSEM

### Selected Auxiliary Tasks (Stage-1)

RTE: MRPC, QQP, MRPC: QQP, Mul QNLI: MRPC, QQ CoLA: MRPC, QQ SST-2: MRPC, QQ



	RTE	MRPC	QNLI	CoLA	SST-2
2018)	50.1	69.0/80.8	69.4	35.0	90.2
018)	55.7	76.2/83.5	66.7	27.5	89.6
	54.0	75.7/83.7	74.0	30.8	91.3
	58.7	78.5/84.5	79.2	32.9	<b>91.8</b>

- RTE: MRPC, <u>QQP</u>, <u>MultiNLI</u>, QNLI, WNLI
- MRPC: QQP, <u>Multinli</u>, QNLI, <u>RTE</u>, WNLI
- QNLI: MRPC, QQP, <u>MultiNLI</u>, RTE, <u>WNLI</u>
- **Cola**: MRPC, QQP, <u>Multinli</u>, QNLI, RTE, <u>WNLI</u>, SST-2
- SST-2: <u>MRPC</u>, QQP, <u>Multinli</u>, QNLI, RTE, <u>WNLI</u>, CoLA



Models

BiLSTM+ELMo (Single-Task) (Wang et al., 2 BiLSTM+ELMo (Multi-Task) (Wang et al., 20

Our Baseline Our AUTOSEM

### **Selected Auxiliary Tasks (Stage-1)**

**RTE:** MRPC, **<u>QQP</u>**, <u>**MultiNLI**</u>, QNLI, WNLI MRPC: QQP, <u>MultiNLI</u>, QNLI, <u>RTE</u>, WNLI **QNLI**: MRPC, QQP, <u>MultiNLI</u>, RTE, <u>WNLI</u> Cola: MRPC, QQP, <u>Multinli</u>, QNLI, RTE, <u>WNLI</u>, SST-2 SST-2: <u>MRPC</u>, QQP, <u>MultiNLI</u>, QNLI, RTE, <u>WNLI</u>, CoLA



	RTE	MRPC	QNLI	CoLA	SST-2
2018)	50.1	69.0/80.8	69.4	35.0	90.2
.018)	55.7	76.2/83.5	66.7	27.5	89.6
	54.0	75.7/83.7	74.0	30.8	91.3
	58.7	78.5/84.5	79.2	32.9	<b>91.8</b>

**MultiNLI is always** chosen in Stage-1



Models

BiLSTM+ELMo (Single-Task) (Wang et al., 2 BiLSTM+ELMo (Multi-Task) (Wang et al., 20

Our Baseline Our AUTOSEM

### Learned Mixing Ratios (Stage-2)

**RTE:** QQP, MultiNLI = 1:5:5 **MRPC**: RTE, MultiNLI = 9:1:4

**QNLI**: WNLI, MultiNLI = 20:0:5

**CoLA**: MultiNLI, WNLI = 20:5:0



	RTE	MRPC	QNLI	CoLA	SST-2
2018)	50.1	69.0/80.8	69.4	35.0	90.2
.018)	55.7	76.2/83.5	66.7	27.5	89.6
	54.0	75.7/83.7	74.0	30.8	91.3
	58.7	78.5/84.5	79.2	32.9	<b>91.8</b>

- **SST-2:** MultiNLI, MRPC, WNLI = 13:5:0:0



## **Removing Stage-1**

Name	Validation	Test
Baseline	78.3	75.7/83.7
w/o Stage-1	80.3	76.3/83.8
w/o Stage-2	80.3	76.7/83.8
Final MTL	81.2	78.5/84.5

w/o Stage-1: Applying Gaussian process on all candidate auxiliary tasks.

(MRPC Dataset)





## **Removing Stage-2**

Name	Validation	Test
Baseline	78.3	75.7/83.7
w/o Stage-1	80.3	76.3/83.8
w/o Stage-2	80.3	76.7/83.8
Final MTL	81.2	78.5/84.5

w/o Stage-2: Apply manual tuning of mixing ratios on selected auxiliary tasks.

(MRPC Dataset)





## Stability of MTL Models

Name	RTE	MRPC	QNLI	CoLA	SST-2		
	BASELINES						
Mean	58.6	78.3	74.9	74.6	91.4		
Std	0.94	0.31	0.30	0.44	0.36		
MULTI-TASK MODELS							
Mean	62.0	81.1	76.0	75.7	91.8		
Std	0.62	0.20	0.18	0.18	0.29		





### **Educated-Guess Baselines**

Our first educated-guess baseline is to choose other similar paraphrasing-based auxiliary tasks.

MRPC+QQP

Our second educated-guess baseline added MultiNLI as an auxiliary task (in addition to QQP)

MRPC+QQP+MultiNLI





### Visualization of Stage-1





Visualization of task utility estimates from the multiarmed bandit controller on SST-2 (primary task). The xaxis represents the task utility, and the y- axis represents the corresponding probability density. Each curve corresponds to a task and the bar corresponds to their confidence interval.









Acknowledgement: DARPA-YFA, ONR, Google, Facebook, Baidu, Salesforce, and Nvidia.

# Thanks! **IDENTIFICATION OF THE INTERPORT OF THE**

Code: github.com/hanguo97/AutoSeM