

Power Norm Based Lifelong Learning for Paraphrase Generations

Anonymous ACL submission

Abstract

Seq2seq language generation models are trained with multiple domains in a continue learning manner, where data from each domain being observed in an online fashion. However, continual learning studies often suffer from catastrophic forgetting, a persistent challenge for lifelong learning. To handle this problem, existing work has leveraged experience replay or dynamic architecture to consolidate the past knowledge, which however results in incremental memory space or high computational cost.

In this work, we propose an innovative framework PNLLL that remedies catastrophic forgetting with a power normalization on NLP transformer models. Specifically, PNLLL leverages power norm to achieve a better balance between past experience rehearsal and new knowledge acquisition. These designs enable the knowledge transfer to new tasks while memorizing the experience of past ones. Our experiments on, paraphrase generation, show that PNLLL outperforms SOTA models by a considerable margin and remedy the forgetting greatly.

1 Introduction

Seq2seq language generation is the essential framework for many tasks such as machine translation, summarization, paraphrase, question answering, dialog response generation. In these applications, models are typically trained offline using annotated data from a fixed set of domains. However, in real-world applications, it is desirable for the system to expand its knowledge to new domains and functionalities, i.e., continuously inquiring new knowledges without forgetting the previously learned skills, which is called lifelong learning (LLL) (Ring et al., 1994; Chaudhry et al., 2019).

Neural networks struggle to learn continuously and experience catastrophic forgetting (CF) when optimized on a sequence of learning problems (McCloskey and Cohen, 1989; French, 1999). Some past works in LLL demonstrated that discriminative

models can be incrementally learnt for a sequence of tasks (Chen et al., 2020; Kirkpatrick et al., 2017). In contrast, under generative settings such as language generation, there has been limited research. Recent works in this area include Mi et al. (2020) and Madotto et al. (2020).

Existing work in LLL adopts the *replay based methods* (Pellegrini et al., 2019), such as Latent Replay, or *regularization based methods* (Huszár, 2018), such as Elastic Weight Consolidation (EWC) (Kirkpatrick et al., 2017). Although they can rectify CF in several scenarios, they have some limitations. The replay-based methods require storing samples from previous tasks, and regularization methods often view all the model parameters as equally important and regularize them to the same extent. In addition, those approaches do not explicitly address the data distribution shift that causes the CF problem. The semantic gap between the embedding spaces of two domains is a leading reason of CF (Wang et al., 2021).

In this work, we propose a novel method, power norm based lifelong learning (PNLLL) to alleviate CF in continuous seq2seq language generation. Essentially, power norm, proposed by Shen et al. (2020) is a variant of layer norm (Ba et al., 2016) or batch normalization (Ioffe, 2017). It is proposed to overcome problems of batch normalization, where large distances between batch statistics leads to large fluctuations among batches and thus poor performances in inferences and layer normalization, where running statistics is calculated at batch level, leading large number of outliers being weighted long sentence. In contrast, power normalization overcomes problems of both batch and layer normalization by enforcing unit quadratic mean for the activations and incorporating running statistics for the quadratic mean of the signal in the process of continual learning. Such designing and incorporation enables our lifelong learning improve generalization performances, maintaining a better balance

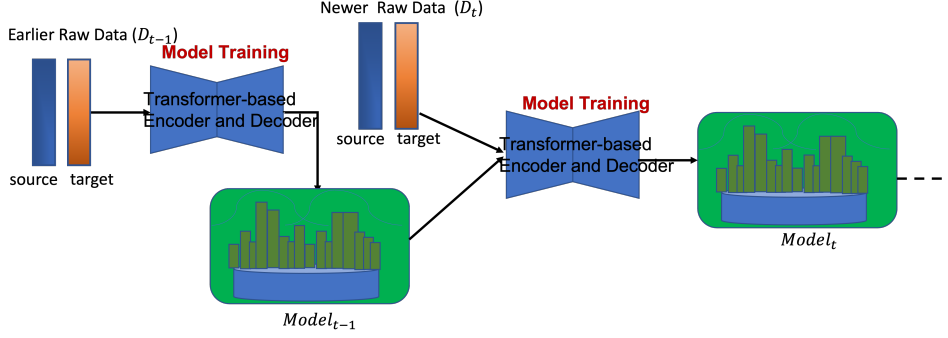


Figure 1: Overview of PNLLL for LLL Seq2seq Language Generation. Figure best viewed in color.

between stability and plasticity in our experiments, showing the effectiveness of PNLLL.

Our main contributions are:

- We design an innovative algorithm based on power norm to store distributions of previous tasks while training for the current task for LLL seq2seq generation.
- Our experiments on seq2seq generation benchmark datasets show that our model achieves SOTA in current task learning and reduces forgetting rates for previous tasks.

2 Proposed Method

In this section, we introduce our proposed framework power norm based lifelong learning (PNLLL). In LLL scenario, models are trained for a sequence of domains or tasks. The model of the first task is trained using pretrained models. Starting from the second model, the network is initialized with parameters of its previous model.

2.1 System Architecture

As shown in Figure 1, input data of task1 with source and target pairs are passed into transformer-based encoder and decoder for training (BART is the encoder and decoder in our context). Power normalization is employed to get running statistics of quadratic means rather than the usual batch means and variances. They are updated with a new types of back propagation for better estimate distributions of each layer’s parameters. Trained models’ parameters are deployed as initialization of later models.

2.2 Power Normalization

Power normalization (PN), mentioned in Introduction, enforces unit quadratic mean for the activation to avoid fluctuations brought by using batch normalization in tasks involving small batches (seen often in NLP) (Shen et al., 2018). It has been proven

effective in both machine translation and language modeling. In this work, we make revisions so as to integrate it into our life-long learning framework.

Firstly, we still follow Shen et al. (2018) to enforce quadratic mean for the activations rather than enforce unit variance in order to overcome large variations in the mean. In addition, we pass through running statistics for the quadratic mean during model initialization from past tasks to next ones to facilitate knowledge transfer among related tasks. The above modifications aim to seeking a robust model training process against outlier and noise, meanwhile maintaining stability in parameter updating and consistency of two continuous models.

2.2.1 Replacing batch mean and variance with unit quadratic mean

Technically, for both batch normalization and layer normalization, in their forward inference, a batch norm (BN) layer is added to calculate mean and variances batch by batch as following,

$$\hat{\mathbf{X}} = \frac{\mathbf{X} - \mu_B}{\sigma_B}, \quad \mathbf{Y} = \gamma \odot \hat{\mathbf{X}} + \beta \quad 139$$

$$\text{s.t. } \mu_B = \frac{1}{B} \sum_{i=1}^B \mathbf{x}_i, \quad \sigma_B^2 = \frac{1}{B} (\mathbf{x}_i - \mu_B)^2 \quad 140$$

where B refers to batch, \mathbf{x}_i , \mathbf{X} and \mathbf{y}_i , \mathbf{Y} refer to input and output of BN, respectively. The BN layer enforces zero mean and unit variance and then performs an affine transformation by scaling $\hat{\mathbf{X}}$ with γ and β .

In the PN framework, the feature embedding is scale by quadratic means of the batch and the operation of PN is formally defined as

$$\hat{\mathbf{X}} = \frac{\mathbf{X}}{\psi_B}, \quad \mathbf{Y} = \gamma \odot \hat{\mathbf{X}} + \beta \quad 149$$

$$\text{s.t. } \psi_B^2 = \frac{1}{B} \sum_{i=1}^B \mathbf{x}_i^2 \quad (1) \quad 150$$

where ψ^2 refers to quadratic mean. Compared with BN, there are two modifications in PN: 1) the means of the batch μ_B are removed from the normalization operation; 2) the variance of the batch σ_B is replaced by the quadratic mean of batch ψ_B . This is because enforcing zero-mean and variance in BN may result in instability due to a large variation of the mean in the NLP data (Shen et al., 2020). Thus, PN performs more stable on the NLP tasks.

In our lifelong learning setting, we address the catastrophic forgetting via balancing the learned parameters on previous tasks and new ones. Besides updating running statistics within current tasks, we update running statistics of model training based on those of previous tasks as well. Formally, we propose an adaptive forward pass for passing through running statistics in the sequential tasks,

$$\hat{\mathbf{X}} = \frac{\mathbf{X}}{\psi^{(t-1)}} \quad \mathbf{Y}^{(t)} = \gamma \odot \hat{\mathbf{X}}^{(t)} + \beta$$

s.t. $(\psi^{(t)})^2 = (\psi^{(t-1)})^2 + (1 - \alpha)(\psi_B^2 - (\psi^{(t-1)})^2)$

where t refers to current task and $t - 1$ refers to previous task, $\alpha \in (0, 1)$ is a moving average coefficient. When $\alpha \approx 0$, the equation reduces to per-batch power normalization, while $\alpha \approx 1$, the PN on current tasks relies much on the previous experiences. Similarly, since forward pass evolves running statistics, the backward propagation cannot be accurately computed. We resort to similar strategies to do the gradient approximation in the backward propagation as following,

$$\nu = \nu^{t-1}(1 - (1 - \alpha)\Gamma^t) + (1 - \alpha)\Lambda^{(t)} \quad (2)$$

where $\Gamma^t = \frac{1}{B} \sum_{i=1}^B \hat{x}_i^{(t)} \hat{x}_i^{(t)}$ and $\Lambda^t = \frac{1}{B} \sum_{i=1}^B \frac{\partial \mathcal{L}}{\partial \hat{x}_i^{(t)}} \hat{x}_i^{(t)}$. Note that the gradient approximation in Eq. (2) is proved to be bounded by a constant (see Theorem 4 in Shen et al. (2020)), which facilitates the robust training process.

3 Experiments on Paraphrase Generation

We apply PNLN to the paraphrase generation task.

3.1 Experimental Setups

For paraphrase generation, we use three existing paraphrase datasets, Quora, Twitter and Wiki_data, in a sequential fashion, that is, the model is first trained on the Quora data, then Twitter, then Wiki_data. We name this experimental setting as QTW. Statistics of the data are provided in Table 1.

	Quora	Twitter	Wiki_Data	total
train	111,947	85,970	78,392	276,309
valid	8,000	1,000	8,154	17,154
test	37,316	3,000	9,324	49,640

Table 1: Dataset stats for QTW

We use a current SOTA generation model, BART, as the seq2seq backbone in our LLL framework, as well as the other methods. We compare our approach with the following baselines.

- **Finetune:** for each task, each model is initialized with the model obtained until the last task, and then fine-tuned with the data of the current task.
- **Full:** we train a model with all the three data sets together.
- **EWC:** the model is trained with the base EWC model on the data from the current task with the initialization of the previous model.

See Appendix for details on the implementation. For evaluation metrics, we use Bleu4, RougeL and Meteor for the generation task. To measure the forgetting rates of different methods, we apply models trained using new data to past data.

3.2 Results

Evaluating on the Current Task

For QTW setting, Table 2 shows results for models evaluated on the data corresponding to the current task.

The first three lines are results from independent models, that is, the BART models are trained on only one of datasets in QTW. As expected, models trained on the matched domain achieve higher performance than otherwise. And there is a large performance drop when using models trained from mismatched domains. This is mostly because of the different writing styles of the three datasets. Wiki is the most formal one, and Twitter is the most informal one.

In the fourth row, the BART model is trained in finetune mode, i.e., in QTW order, the model is initialized with that trained in the previous domain and fine tuned using the subsequent domain. We can see that results on both Twitter and Wiki test data are slightly lower than those when models are trained directly on the corresponding training data. Again, this suggests pretraining the model

Models	Quora Test			Twitter Test			Wiki Test		
	bleu4	rougeL	meteor	bleu4	rougeL	meteor	bleu4	rougeL	meteor
Quora-trained	30.11	55.85	57.17	2.12	6.13	5.49	4.51	11.21	12.13
Twitter-trained	3.18	11.46	9.01	35.47	57.49	54.57	4.60	9.76	7.50
Wiki_data-trained	22.38	43.44	46.23	9.32	17.93	21.03	42.12	73.86	73.10
Finetune	30.11	55.85	57.17	35.79	56.32	54.93	42.12	73.86	73.10
EWC	30.25	56.16	57.98	33.52	54.41	54.21	42.15	73.53	73.59
PNLLL	31.20	58.89	60.33	34.62	57.87	55.12	43.84	74.79	73.72
Full	33.99	59.56	61.67	38.56	58.76	56.	46.86	76.59	75.91

Table 2: Results of model evaluations on QTW setting

Quora test with Model trained with Twitter			
Models	bleu4	rougeL	meteor
Quora-trained	30.11	55.85	57.17
Finetune	15.80	46.59	47.31
EWC	15.63	41.53	46.03
PNLLL	17.58	47.88	49.20

Quora test with Model trained with Wiki_data			
Models	bleu4	rougeL	meteor
Quora-trained	30.11	55.85	57.17
Finetune	19.07	51.76	55.95
EWC	19.63	49.35	53.02
PNLLL	20.34	52.59	56.06

Twitter test with Model trained with Wiki_data			
Models	bleu4	rougeL	meteor
Twitter-based	35.79	56.32	54.93
Finetune	14.09	37.97	45.89
EWC	14.84	38.65	46.33
PNLLL	16.49	39.93	49.28

Table 3: Results of all the methods when testing new models on previous domains (from 2nd row to the last).

with mismatched data is not beneficial. The results from the EWC baseline are not consistently better than the finetune method, showing the limited effectiveness of EWC regularization. In contrast, our proposed approaches obtain better results than Finetune. Even for the first task, Quora, we observe around 1% better results for all three metrics. This demonstrates that even for pretrained models, regularization shows positive effect. For the later tasks, twitter data obtains about 3-4% increase and wiki data even obtains around 7-9% increase. This shows the effectiveness of PNLLL. In addition, six out of nine results from PNLLL win about 1% over MR. This shows that further regularization with quadratic penalty has positive impact on selection of important parameters. The last row is the results of Full. Since the model has seen all the data, it is not surprising that results for both Twitter and

Wiki_data are better than our models, and it may be partly due to similarity in Quora and Wiki data.

Evaluating on Previous Tasks

Table 3 shows the results when models trained on new domains are evaluated on data from past domains. Since we are using the order of QTW, results are presented for evaluating on Quora and Twitter data. For the Quora test set, we show results after training with Twitter data, and then subsequently Wiki_data. The first row of each sub-table is the result of the BART model trained on the only corresponding data. The second row uses the baseline fine tuning fashion.

Each of them yields much better results than the finetune or EWC baselines, with much less drop rates. This shows each module can reduce forget rates. In addition, after the model is trained on Wiki_data, forgetting rates for Quora Test (the first dataset) are even lower than the model trained on Twitter. This again indicates Wiki_data and Quora are more similar in style than Twitter. On the Twitter test set, there is some difference between Bleu, Rouge and Meteor metrics.

4 Conclusion

In this work, we introduce PNLLL, a generic LLL framework for addressing forgetting in seq2seq language generation learning. Our experimental results have shown that it outperformed SOTA in paraphrase generation, a neural seq2seq language generation task. Future work includes applying PNLLL to diverse generation tasks and generation network structures. In addition, improvements of domain shift estimation can be made with the introduction of topic similarity. In order to make the model more discriminative against domain differences, we may add contrastive learning loss function to our current label smoothing cross entropy loss as in [Gunel et al. \(2020\)](#).

292
293
294
295
296
297

298
299
300

301
302
303
304
305

306
307
308
309
310
311

312
313
314
315
316
317
318

319
320
321

322
323
324
325
326

327
328
329

330
331
332
333
334
335

336
337
338
339
340
341
342

343
344
345

References

Rahaf Aljundi, Francesca Babiloni, Mohamed Elhoseiny, Marcus Rohrbach, and Tinne Tuytelaars. 2018. Memory aware synapses: Learning what (not) to forget. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 139–154.

Lei Jimmy Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. 2016. [Layer normalization](#). *CoRR*, abs/1607.06450.

Arslan Chaudhry, Marcus Rohrbach, Mohamed Elhoseiny, Thalaisyasingam Ajanthan, Puneet Kumar Dokania, Philip H. S. Torr, and Marc’Aurelio Ranzato. 2019. [Continual learning with tiny episodic memories](#). *CoRR*, abs/1902.10486.

Sanyuan Chen, Yutai Hou, Yiming Cui, Wanxiang Che, Ting Liu, and Xiangzhan Yu. 2020. Recall and learn: Fine-tuning deep pretrained language models with less forgetting. *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*.

Cyprien de Masson d’Autume, Sebastian Ruder, Lingpeng Kong, and Dani Yogatama. 2019. Episodic memory in lifelong language learning. *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*.

Robert M French. 1999. Catastrophic forgetting in connectionist networks. *Trends in cognitive sciences*, 3(4):128–135.

Beliz Gunel, Jingfei Du, Alexis Conneau, and Ves Stoyanov. 2020. Supervised contrastive learning for pretrained language model fine-tuning. *8th International Conference on Learning Representations, ICLR 2020, Vienna, Austria, May, 2020*.

Ferenc Huszár. 2018. Note on the quadratic penalties in elastic weight consolidation. *Proceedings of the National Academy of Sciences*, page 201717042.

Sergey Ioffe. 2017. [Batch renormalization: Towards reducing minibatch dependence in batch-normalized models](#). In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 1945–1953.

James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. 2017. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526.

Yuncheng Li, Jianchao Yang, Yale Song, Liangliang Cao, Jiebo Luo, and Li-Jia Li. 2017. Learning from noisy labels with distillation. In *ICCV*.

David Lopez-Paz and Marc’Aurelio Ranzato. 2017. Gradient episodic memory for continual learning. *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*. 346
347
348
349
350
351

Andrea Madotto, Zhaojiang Lin, Zhenpeng Zhou, Seungwhan Moon, Paul A. Crook, Bing Liu, Zhou Yu, Eunjoon Cho, and Zhiguang Wang. 2020. [Continual learning in task-oriented dialogue systems](#). *CoRR*, abs/2012.15504. 352
353
354
355
356

Michael McCloskey and Neal J Cohen. 1989. Catastrophic interference in connectionist networks: The sequential learning problem. In *Psychology of learning and motivation*, volume 24, pages 109–165. Elsevier. 357
358
359
360
361

Fei Mi, Liangwei Chen, Mengjie Zhao, Minlie Huang, and Boi Faltings. 2020. Continual learning for natural language generation in task-oriented dialog systems. *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: Findings, EMNLP 2020, Online Event, 16-20 November 2020*. 362
363
364
365
366
367
368

Lorenzo Pellegrini, Gabriele Graffieti, Vincenzo Lomonaco, and Davide Maltoni. 2019. Latent replay for real-time continual learning. *IEEE/RSSJ International Conference on Intelligent Robots and Systems, IROS 2020, Las Vegas, NV, USA, October 24, 2020 - January 24, 2021*. 369
370
371
372
373
374

Mark Bishop Ring et al. 1994. *Continual learning in reinforcement environments*. Ph.D. thesis, University of Texas at Austin Austin, Texas 78712. 375
376
377

David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy P Lillicrap, and Greg Wayne. 2019. Experience replay for continual learning. *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*. 378
379
380
381
382
383
384

Sheng Shen, Zhewei Yao, Amir Gholami, Michael W. Mahoney, and Kurt Keutzer. 2020. [Powernorm: Rethinking batch normalization in transformers](#). In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 8741–8751. PMLR. 385
386
387
388
389
390
391

Yanyao Shen, Hyokun Yun, Zachary C Lipton, Yakov Kronrod, and Animashree Anandkumar. 2018. Deep active learning for named entity recognition. In *ICLR*. 392
393
394
395

Zhuoyi Wang, Yuqiao Chen, Chen Zhao, Yu Lin, and Latifur Khan. 2021. Clear: Contrastive-prototype learning with drift estimation for resource constrained stream mining. In *Proceedings of The Web Conference*. 396
397
398
399
400

Zirui Wang, Sanket Vaibhav Mehta, Barnabás Póczos, and Jaime Carbonell. 2020. Efficient meta lifelong-learning with limited memory. *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*.

Lu Yu, Bartłomiej Twardowski, Xialei Liu, Luis Heranz, Kai Wang, Yongmei Cheng, Shangling Jui, and Joost van de Weijer. 2020. Semantic drift compensation for class-incremental learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6982–6991.

5 Appendix

5.1 Metrics Details

Throughout the paper, we use those evaluation metrics that have been widely used in the previous work to measure the quality of the paraphrases. In general, BLEU measures how much the words (and/or n-grams) in the machine generated summaries appeared in the human reference summaries. Rouge measures how much the words (and/or n-grams) in the human reference summaries appeared in the machine generated summaries. Specifically, we use the library¹ from HuggingFace to compute BLEU scores and *py-rouge*² to compute ROUGE scores. As BLEU and ROUGE could not measure the diversity between the generated and the original sentences, we follow unsupervised paraphrasing methods and adopt meteor to measure the diversity of expression in the generated paraphrases by penalizing copying words from input sentences. The introduction of Slot error rate, Ω_{all} and Ω_{first} can be seen in the data setting of MultiWoz2.

5.2 Implementation Details

Packages Used for Implementation. The relevant packages that we use in the implementation and their corresponding versions are as following: python==3.6.6, fairseq==1.0, torch==1.4.0, cuda==10.2, tensorboard==1.10.0, numpy==1.14.5, scipy==1.1.0, NLTK==3.4.5 and scikit-learn==0.21.3.

5.3 Related Work

5.3.1 life-long Learning (LLL)

life-long learning has been studied from a few perspectives, including data buffering, regularization and prototype keeping. Replay based methods can be used in data buffering or prototype keeping.

¹<https://huggingface.co/metrics/sacrebleu>

²<https://pypi.org/project/py-rouge/>

It usually keeps a small amount of real samples from old tasks or distills the knowledge from old data and recreates pseudo-data of old tasks for later training. Using these sampled data or pseudo data can prevent weights from deviating from previous status (Rolnick et al., 2019; Wang et al., 2020; Lopez-Paz and Ranzato, 2017). The main idea of this approach is to assign a dedicated capacity inside a model for each task. After a task is completed, the weights are frozen as one prototype (Wang et al., 2021; d’Autume et al., 2019). Both data buffering and prototype keeping need storage of either data samples or model weights, i.e., they require extra memory to memorize important information of previous tasks. Another LLL method is regularization based, which adds a regularization term to weights when learning them for a new task in order to minimize deviation from previously trained weights. Most regularization based methods estimate the importance of each parameter and add them as a constraint to the loss function. Different algorithms have been designed to achieve this goal. For example, elastic weight consolidation (EWC) calculates a Fisher information matrix to estimate the sensitivity of parameters (Kirkpatrick et al., 2017); memory aware synapses (MAS) (Aljundi et al., 2018) uses the gradients of the model outputs; and episodic memory or gradient episodic memory (GEM) (Li et al., 2017; Lopez-Paz and Ranzato, 2017) allows positive backward transfer and prevents the loss on past tasks from increasing. These methods all attempt to slow down the learning of parameters that are important for previous tasks.

5.4 LLL in Seq2seq Language Generation

In Seq2seq language generation, not much work has been done in LLL. The most relevant work is from Mi et al. (2020) where a framework of sequential learning is designed for task-oriented dialogues. Specifically, they replay prioritized exemplars together with an adaptive regularization technique based on EWC. They store representative utterances from previous data (exemplars), and replay them to the Seq2seq language generation model each time it needs to be trained on new data. They achieved good results on the MultiWoZ-2.0 dataset. Nonetheless, their work requires to store data from previous tasks, which leads to poor scalability on large-scale datasets. In addition, their system is specifically designed for the MultiWoz

SOURCE	BART	RMR_DSE	TARGET
Why is German Shepherd/Great Pyrenees mix coveted <i>among breeders?</i>	Why is German Shepherd/Great Pyrenees mix coveted <i>from browns?</i>	Why is German Shepherd/Great Pyrenees mix coveted <i>among breeders?</i>	Why is German Shepherd/Great Pyrenees mix coveted <i>among breeders?</i>
eyeing trump, Obama takes new action to ban <i>arctic drilling</i>	president Obama takes new action to ban <i>arctic drilling</i>	Obama takes new action to ban <i>arctic drilling</i>	please save the earth mr. president . Obama takes new action to ban <i>arctic drilling</i>
death toll in 6.5 - magnitude earthquake in indonesia's <i>aceh province</i> increase to at least 52	a 6.5 earthquake in kills at least 26 people @cnn	death toll in 6.5 - magnitude earthquake in <i>aceh province</i> increase to at least 52	powerfull quake kills dozens at least 25 people were killed in an earthquake that struck indonesia's <i>aceh province</i>
pipeline 150 miles from dakota access protests <i>leaks</i> gallons of oil	the new york times pipeline 150 miles from dakota <i>access pipeline .</i>	pipeline 150 miles from dakota access <i>leaks</i> gallons of oil	of oil, or gallons, have <i>leaked</i> from the pipeline

Table 4: Examples of the generated paraphrases by BART and RMR_DSE on QTW data setting.

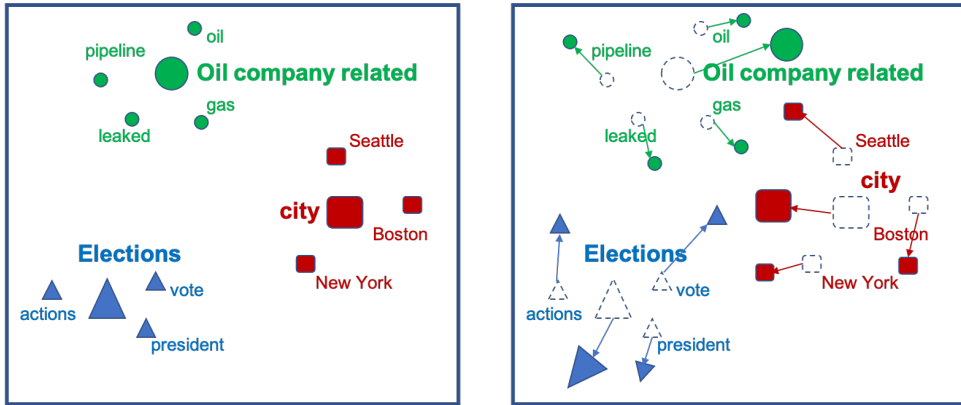


Figure 2: Illustration of Domain Shift: (a) Data with three relevant topic/cluster in the embedding space after model trained on task 1. (b) Data with previous topics in the embedding space after the model trained on task 2, the arrow indicates the domain shift between two tasks.

task and lacks generalization to other tasks. In contrast, our proposed PNLLL method aims to fit different seq2seq language generation applications, therefore it is easy to be integrated to tasks such as summarization, translation, paraphrases, dialog response generation.

5.4.1 Illustrations of Semantic Drift

As illustrated in Figure 2, each data point and their cluster centers trained in Task 1 are shifted after training for Task 2. Yu et al. (2020) proposed to compensate this gap without using any exemplars via domain shift. Nonetheless, these studies mainly focused on classification tasks, which limited their application on language generation model.

5.5 Case Studies

In Table 4, we show some generated samples from QTW setting using the baseline Bart model and our PNLLL model. All examples are results generated by $model_t$ on $data_{t-1}$. Among the five examples, the first one is from Quora, the last one from Wiki data and the other three from Twitter. The reason

that we select more samples from Twitter is that we find Twitter is the most informal in style with quite many fragments. Hence, it is the hardest for the generation task and has lowest metrics and lower forgetting reduction rates. In the four samples, the italicised parts are the key words. From the table, we can observe that compared to *BART*, *PNLLL* has better performances on all of the three datasets. The *BART* model misses all of them except *drilling*. In contrast *PNLLL* succeeds in all cases without forgetting the previously learned patterns.