

Exercise Briefing

Ki Hyun Kim

nlp.with.deep.learning@gmail.com

Objective:

- 모델 두 개를 동시에 학습할 수 있는 환경 갖추기
 - 학습 시작점 (e.g. dual_train.py, continue_dual_train.py)
 - Trainer (e.g. dual_trainer.py)
 - 입출력 tensor를 바꿔가며 반대쪽 모델에게 넣어줌
 - 예전 모델(e.g. Seq2seq, Transformer)을 그대로 활용 가능해야 함
- Loss 구현
 - 언어 모델 및 학습 환경 구현 (e.g. rnnlm.py, lm_trainer.py)
 - 각 모델의 backward() 호출을 위한 적절한 detach() 활용

Review: AutoGrad in PyTorch

- Computation Graph 동적 생성에 따른 AutoGrad
 - 필요에 따라 detach() 기능을 활용하여 back-prop을 통제할 수 있음

Equations

- LM을 미리 학습하여, 추후 freeze 한 상태로 활용

$$\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$$

$$\hat{\phi} = \operatorname{argmax}_{\phi \in \Phi} \sum_{i=1}^N \log P(x_i; \phi)$$

$$\hat{\psi} = \operatorname{argmax}_{\psi \in \Psi} \sum_{i=1}^N \log P(y_i; \psi)$$

Equations

- detach()를 잘 활용하여, 각 loss를 적절히 계산해줘야 함

$$\theta_{x \rightarrow y} \leftarrow \theta_{x \rightarrow y} - \eta \nabla_{\theta_{x \rightarrow y}} \mathcal{L}(\theta_{x \rightarrow y})$$

$$\mathcal{L}(\theta_{x \rightarrow y}) = \sum_{i=1}^N \left(-\log P(y_i | x_i; \theta_{x \rightarrow y}) + \lambda \times \left((\log P(x_i; \phi) + \log P(y_i | x_i; \theta_{x \rightarrow y})) - (\log P(y_i; \psi) + \log P(x_i | y_i; \theta_{y \rightarrow x})) \right)^2 \right)$$

$$\theta_{y \rightarrow x} \leftarrow \theta_{y \rightarrow x} - \eta \nabla_{\theta_{y \rightarrow x}} \mathcal{L}(\theta_{y \rightarrow x})$$

$$\mathcal{L}(\theta_{y \rightarrow x}) = \sum_{i=1}^N \left(-\log P(x_i | y_i; \theta_{y \rightarrow x}) + \lambda \times \left((\log P(x_i; \phi) + \log P(y_i | x_i; \theta_{x \rightarrow y})) - (\log P(y_i; \psi) + \log P(x_i | y_i; \theta_{y \rightarrow x})) \right)^2 \right)$$

Equations

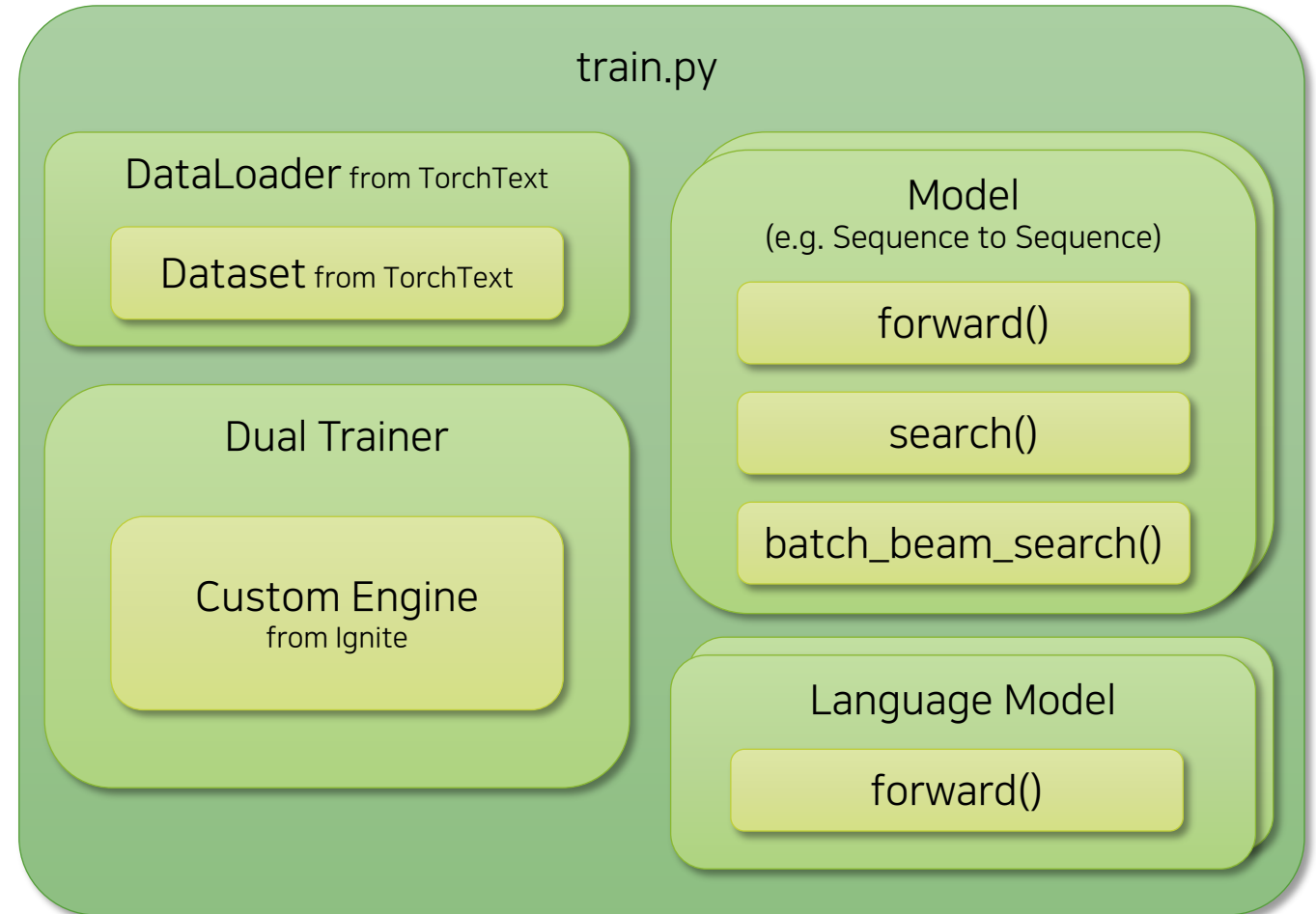
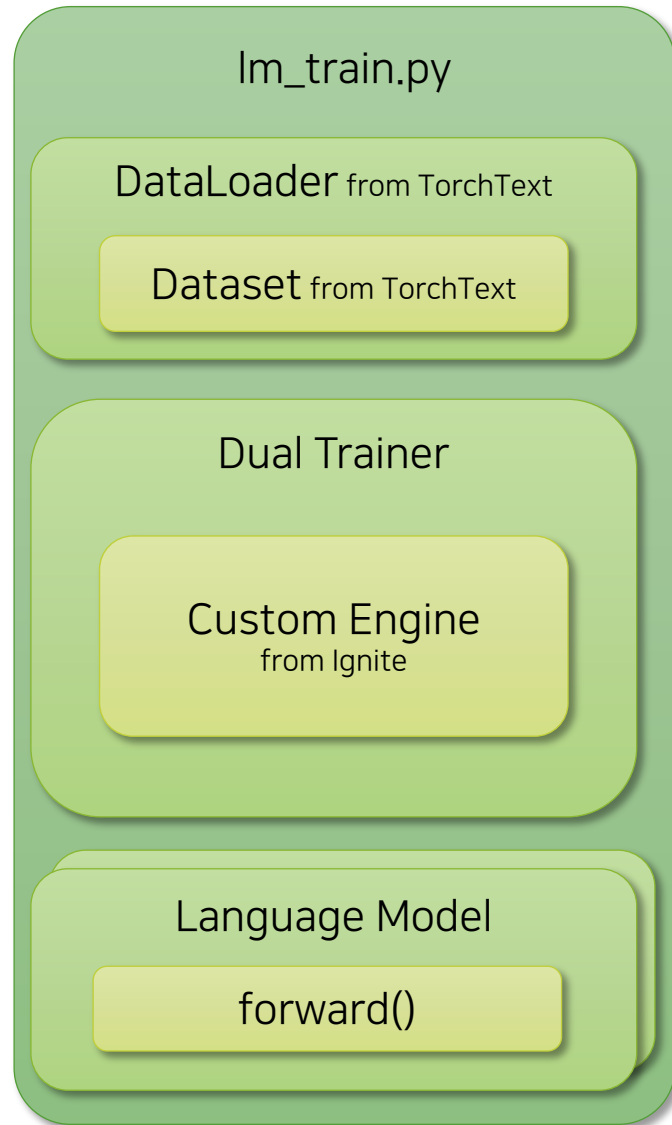
- 입출력 tensor쌍을 반대쪽 모델에 맞게 잘 변환해야 함

$$\begin{aligned}\log P(y_i|x_i; \theta_{x \rightarrow y}) &= \sum_{t=1}^m \log P(y_{i,t}|x_i, y_{i,<t}; \theta_{x \rightarrow y}) \\ &= \sum_{t=1}^m y_{i,t}^\top \cdot \log \hat{y}_{i,t}\end{aligned}$$

$$\begin{aligned}\log P(x_i|y_i; \theta_{y \rightarrow x}) &= \sum_{t=1}^n \log P(x_{i,t}|y_i, x_{i,<t}; \theta_{y \rightarrow x}) \\ &= \sum_{t=1}^n x_{i,t}^\top \cdot \log \hat{x}_{i,t}\end{aligned}$$

where $x_{i,0} = y_{i,0} = \langle \text{BOS} \rangle$ and $x_{i,n} = y_{i,m} = \langle \text{EOS} \rangle$.

Project Implementation



Training Procedure

