

A Derivation for Stochastic Pruning

To re-parametrize the discrete binary Bernoulli variable $m_{i,j}^l \sim B(\sigma(g_{i,j}^l))$, denote the approximate differentiable variable as $\tilde{m}_{i,j}^l = \sigma(\frac{g_{i,j}^l + \log U - \log(1-U)}{\tau})$ where τ is a real-valued temperature value, we have the following derivation holds for arbitrary $\epsilon \in (0, 0.5)$:

$$P(m_{i,j}^l = 1) - P(\tilde{m}_{i,j}^l \geq 1 - \epsilon) \leq (\frac{\tau}{4}) \log \frac{1}{\epsilon} \quad (1)$$

Specifically, when temperature τ approaches 0, $\tilde{m}_{i,j}^l = m_{i,j}^l$.

Lemma 1: $\sigma^{-1}(x) = \log \frac{x}{1-x}$.

Lemma 2: $\frac{\sigma(x) - \sigma(y)}{x-y} \leq \frac{1}{4}$.

Proof:

$$\begin{aligned} & P(\tilde{m}_{i,j}^l \geq 1 - \epsilon) \\ &= P(\sigma(\frac{g_{i,j}^l + \log U - \log(1-U)}{\tau}) \geq 1 - \epsilon) \end{aligned} \quad (2)$$

$$= P(\frac{g_{i,j}^l + \log U - \log(1-U)}{\tau} \geq \log(\frac{1}{\epsilon} - 1)) \quad (3)$$

$$= P(g_{i,j}^l - \tau \log(\frac{1}{\epsilon} - 1) \geq \log(\frac{1}{U} - 1)) \quad (4)$$

$$= P(e^{g_{i,j}^l - \tau \log(\frac{1}{\epsilon} - 1)} \geq \frac{1}{U} - 1) \quad (5)$$

$$= P(U \geq \frac{1}{1 + e^{g_{i,j}^l - \tau \log(\frac{1}{\epsilon} - 1)}}) \quad (6)$$

$$= \sigma(g_{i,j}^l - \tau \log(\frac{1}{\epsilon} - 1)) \quad (7)$$

$$= \sigma(g_{i,j}^l) - \sigma(g_{i,j}^l - \tau \log(\frac{1}{\epsilon} - 1)) \quad (8)$$

Then:

$$P(m_{i,j}^l = 1) - P(\tilde{m}_{i,j}^l \geq 1 - \epsilon) \quad (9)$$

$$= \sigma(g_{i,j}^l) - \sigma(g_{i,j}^l - \tau \log \frac{1}{\epsilon} - 1) \quad (10)$$

$$\leq \frac{\tau}{4} \log(\frac{1}{\epsilon} - 1) \quad (11)$$

$$\leq \frac{\tau}{4} \log \frac{1}{\epsilon} \quad (12)$$

The process for deriving $P(m_{i,j}^l = 0) - P(\tilde{m}_{i,j}^l \leq \epsilon) \leq (\frac{\tau}{4}) \log \frac{1}{\epsilon}$ can be analogously obtained. \square

B Knowledge Combination for Fine-tuning, Zero-shot and Triple Classification on Commonsense Reasoning

B.1 Notation for Knowledge Type

HasSubevent: 0

MadeOf: 1

HasPrerequisite: 2

MotivatedByGoal: 3

AtLocation: 4

CausesDesire: 5

IsA: 6

NotDesires: 7

Desires: 8

CapableOf: 9

PartOf: 10

HasA: 11

UsedFor: 12

ReceivesAction: 13

Causes: 14

HasProperty: 15

In the remainder of this section, we use \cup to indicate mask union operation upon multiple commonsense knowledge types.

B.2 Fine-tuning

For fine-tuning on commonsense reasoning tasks, we only experiments with BERT-BASE due and perform hyper-parameter search only in terms of batch size in the range of $\{8, 16, 32\}$ and learning rate in the range of $\{3e^{-5}, 4e^{-5}, 5e^{-5}\}$ due to computational budget. We also adopt early stopping based on accuracy on the development set. The combination achieving highest accuracy is shown in Table 1.

B.3 Zero-shot

In contrast with fine-tuning, zero-shot evaluation is deterministic as long as the model does not involve any stochastic module, thereby averting extensive hyperparameter tuning. Instead we perform exhaustive search over knowledge combinations for each pretrained language model with number of knowledge types in $\{3, 4, 5\}$. The ConceptNet-grounded knowledge type combination achieving highest accuracy is listed in Table 2.

B.4 Triple Classification

In analogy with zero-shot evaluation, here we show the optimal knowledge type combination of each

Task	RTE	COPA	CSQA	SWAG	HellaSWAG	aNLI	CosmosQA
BERT	006014	508014	30408012014	106010011	0030508014	0030508014	0030508014

Table 1: Optimal fine-tuning knowledge type combination for BERT-BASE on commonsense reasoning tasks.

Task	COPA (Dev.)	CSQA	CA	WSC	SM	ARCT1	ARCT2
DISTILBERT	106014	203013	0010709	607010	208013	203014	10207
BERT	4011015	102015	608012	209014	6012015	109010	10508
RoBERTa	20308	00205	00108	1020405011	8011012	205011013	008011013
MPNET	10608010	6012013	203010	1030409	6010013015	20506011	50607011

Table 2: Optimal zero-shot knowledge type combination for each PLM on each commonsense reasoning tasks.

Model	P@1	P@2	P@3	Sparsity	$l_b - l_t$	# Param.
BERT-LARGE w/o pruning	15.1	20.9	24.6	0%	-	336M
BERT-LARGE w/ stochastic pruning	22.1	30.1	35.4	~30%	17-24	336M
BERT-LARGE w/ deterministic pruning	69.2	74.1	76.3	~50%	17-24	284M

Table 3: Macro-averaged precision metrics of BERT-LARGE on the ConceptNet subset of LAMA.

PLM for triple classification task on ConceptNet-100K ¹.

DISTILBERT-BASE: 304012

BERT-BASE: 9013014

RoBERTa-BASE: 00409013

MPNET-BASE: 10409

C Additional Pruning Results

C.1 BERT-LARGE

We also apply our pruning procedure upon BERT-LARGE, the rank-based metrics on LAMA ² is shown in Table 3.

D Extracted Commonsense Triples

Applying the pruned DISTILBERT-BASE model to predict missing objects for triples in ConceptNet-100K test set, we obtain commonsense triples deemed to be novel by three human annotators with Flessi’s Kappa score κ of 0.65. We further filtered out triples that are included in the training or development set of ConceptNet-100K. Here we show some representative cases categorized by their relations:

CapableOf:

(computer, crash), (computer, communicate)

IsA:

(sex, relationship), (submarine, weapon),

(submarine, vessel)

AtLocation:

(knife, war), (knife, dinner), (crab, dinner)

UsedFor:

(stage, fun), (stage, performance),
(literature, education), (literature, research)

HasA:

(book, index), (book, information)

HasProperty:

(music, loud)

Future work involves using seed triples beyond ConceptNet-100K dataset, e.g., the whole ConceptNet knowledge graph, and mining more novel and plausible commonsense knowledge.

¹<https://ttic.uchicago.edu/~kgimpel/commonsense.html>

²<https://github.com/facebookresearch/LAMA>