

A Appendix

A.1 Word Dropout as a Special Case

Here, we derive word dropout as an instance of our framework. First, let us introduce a new token, $\langle \text{null} \rangle$, into both the source vocabulary and the target vocabulary. $\langle \text{null} \rangle$ has the embedding of a all-0 vector and is never trained. For a sequence x of words in a vocabulary with $\langle \text{null} \rangle$, we define the *neighborhood* $N(x)$ to be:

$$N(x) = \left\{ x' : |x'| = |x| \text{ and } x'_i \in \{x_i, \langle \text{null} \rangle\} \right\}$$

In other words, $N(x)$ consists of x and all the sentences obtained by replacing a few words in x by $\langle \text{null} \rangle$. Clearly, all augmented sentences \hat{x} that are sampled from x using word dropout fall into $N(x)$.

In (4), the augmentation policy $\mathbf{q}^*(\hat{x}, \hat{y}|x, y)$ was decomposed into two independent terms, one of which samples the augmented source sentence \hat{x} and the other samples the augmented target sentence \hat{y}

$$\mathbf{q}^*(\hat{x}, \hat{y}|x, y) = \underbrace{\frac{\exp\{r_x(\hat{x}, x)/\tau_x\}}{\sum_{\hat{x}'} \exp\{r_x(\hat{x}', x)/\tau_x\}}}_{\mathbf{q}(\hat{x}|x)} \times \underbrace{\frac{\exp\{r_y(\hat{y}, y)/\tau_y\}}{\sum_{\hat{y}'} \exp\{r_y(\hat{y}', y)/\tau_y\}}}_{\mathbf{q}(\hat{y}|y)}$$

Word dropout is an instance of this decomposition, where r_y takes the same form with r_x , given by:

$$r_x(\hat{x}, x) = \begin{cases} -\text{HammingDistance}(\hat{x}, x) & \text{if } \hat{x} \in N(x) \\ -\infty & \text{otherwise} \end{cases}, \quad (5)$$

where $\text{HammingDistance}(\hat{x}, x) = \sum_{i=1}^{|\hat{x}|} \mathbf{1}[\hat{x}_i \neq x_i]$. To see this is indeed the case, let h be the Hamming distance for $\hat{x} \in N(x)$ and set $\lambda_{\text{word}} = \exp\{-1/\tau_x\}$, then we have:

$$\exp\{r_x(\hat{x}, x)/\tau_x\} = \exp\{-h/\tau_x\} = \exp\left\{-h \cdot \log \frac{1}{\lambda_{\text{word}}}\right\} = \exp\{h \cdot \log \lambda_{\text{word}}\} = \lambda_{\text{word}}^h, \quad (6)$$

which is precisely the probability of dropping out h words in x , where each word is dropped with the distribution $\text{Bernoulli}(\lambda_{\text{word}})$.

The difference between word dropout and SwitchOut comes in the fact that $N(x)$ is much smaller than the support of \hat{x} that SwitchOut can sample from, which is $V^{|\hat{x}|}$ where V is the vocabulary. Word dropout concentrates all augmentation probability mass into $N(x)$ while SwitchOut spreads the mass into a larger support, leading to a larger entropy. Meanwhile, both word dropout and SwitchOut are exponentially less likely to diverge a way from x , ensuring the smoothness desiderata of a good data augmentation policy, as we discussed in Section 2.3.

A.2 RAML as a Special Case

Here, we present a detailed description of how RAML is a special case of our proposed framework. For each empirical observation $(x, y) \sim \hat{\mathbf{p}}$, RAML defines a reward aware target distribution $\mathbf{p}_{\text{RAML}}(Y|x, y)$ for the model distribution $\mathbf{p}_\theta(Y|x)$ to match. Concretely, the target distribution in RAML has the form

$$\mathbf{p}_{\text{RAML}}(\hat{y}|x, y) = \frac{\exp\{r(\hat{y}; y)/\tau\}}{\sum_{\hat{y}'} \exp\{r(\hat{y}'; y)/\tau\}},$$

where r is the task reward function. With this definition, RAML amounts to minimizing the expected KL divergence between \mathbf{p}_{RAML} and \mathbf{p}_θ , *i.e.*

$$\begin{aligned} & \min_{\theta} \mathbb{E}_{x, y \sim \hat{\mathbf{p}}} [\text{KL}(\mathbf{p}_{\text{RAML}}(Y|x, y) \parallel \mathbf{p}_\theta(Y|x))] \\ \iff & \max_{\theta} \mathbb{E}_{x, y \sim \hat{\mathbf{p}}} [\mathbb{E}_{\hat{y} \sim \mathbf{p}_{\text{RAML}}(Y|x, y)} [\log \mathbf{p}_\theta(\hat{y}|x)]] \\ \iff & \max_{\theta} \mathbb{E}_{\hat{y} \sim \mathbf{p}_{\text{RAML}}(Y)} [\log \mathbf{p}_\theta(\hat{y}|x)], \end{aligned}$$

where $\mathbf{p}_{\text{RAML}}(Y)$ is the marginalized target distribution, *i.e.* $\mathbf{p}_{\text{RAML}}(Y) = \mathbb{E}_{x,y \sim \hat{\mathbf{p}}}[\mathbf{p}_{\text{RAML}}(Y|x,y)]$. Now, notice that $\mathbf{p}_{\text{RAML}}(Y)$ is a member of the augmentation distribution family in consideration (*c.f.* Section 2.2). Specifically, it is equivalent to a data augmentation distribution where

$$\begin{aligned} \mathbf{q}(\hat{x}, \hat{y} | x, y) &= \mathbf{1}[\hat{x} = x] \cdot \mathbf{p}_{\text{RAML}}(\hat{y}|x, y) \\ \Leftrightarrow \frac{\exp\{s(\hat{x}, \hat{y}; x, y)/\tau\}}{\sum_{\hat{x}', \hat{y}'} \exp\{s(\hat{x}', \hat{y}'; x, y)/\tau\}} &= \mathbf{1}[\hat{x} = x] \cdot \frac{\exp\{r(\hat{y}; y)/\tau\}}{\sum_{\hat{y}'} \exp\{r(\hat{y}'; y)/\tau\}} \\ \Leftrightarrow s(\hat{x}, \hat{y}; x, y) &= \begin{cases} r(\hat{y}; y), & \hat{x} = x \\ -\infty, & \hat{x} \neq x \end{cases}. \end{aligned} \quad (7)$$

The last equality reveals an immediate connection between RAML and our proposed framework. In summary, RAML can be seen as a special case of our data augmentation framework, where the similarity function is defined by (7). Practically, this means RAML only consider pairs with source sentences from the empirical set for data augmentation.

A.3 Datasets Descriptions

	vocab (K)		#sents		
	src	tgt	train	dev	test
en-vi	17.2	7.7	133.3K	1.6K	1.3K
de-en	32.0	22.8	153.3K	7.0K	6.8K
en-de	50.0	50.0	4.5M	2.7K	2.2K

Table 3: Statistics of the datasets.

Table 3 summarizes the statistics of the datasets in our experiments. The WMT 15 en-de dataset is one order of magnitude larger than the IWSLT 16 de-en dataset and the IWSLT 15 en-vi dataset. For the en-vi task, we use the data pre-processed by Luong and Manning (2015). For the en-de task, we use the data pre-processed by Luong et al. (2015), with *newstest2014* for validation and *newstest2015* for testing. For the de-en task, we use the data pre-processed by Ranzato et al. (2016).

A.4 Hyper-parameters

Task	n_{layers}	n_{heads}	d_k, d_v	d_{model}	d_{inner}	init	clip	λ_{drop}	τ_x^{-1}	τ_y^{-1}	λ_{word}
en-de	8	6	64	512	1024	0.04	25.0	0.10	1.00	0.80	0.1
de-en	8	5	64	288	507	0.035	25.0	0.25	0.95	0.90	0.1
en-vi	4	4	64	256	384	0.035	20.0	0.15	1.00	0.90	0.1

Table 4: Hyper-parameters for our experiments.

The hyper-parameters used in our experiments are in Table 4. All models are initialized uniformly at random in the range as reported in Table 4. All models are trained with Adam (Kingma and Ba, 2015). Gradients are clipped at the threshold as specified in Table 4. For the WMT en-de task, we use the legacy learning rate schedule as specified by Vaswani et al. (2017). For the de-en task and the en-vi task, the learning rate is initially 0.001, and is decreased by a factor of 0.97 for every 1000 steps, starting at step 8000. All models are trained for 100,000 steps, during which one checkpoint is saved for each 2500 steps and the final evaluation is performed on the checkpoint with lowest perplexity on the dev set.

Multiple GPUs are used for each experiment. For the de-en and the en-vi experiments, if we use n GPUs, where $n \in \{1, 2, 4\}$, then we only perform $10^5/n$ updates to the models' parameters. We find that this is sufficient to make the models converge.

A.5 Source Code for Sampling in TensorFlow

Hamming distance sampling in TensorFlow

```
1 import tensorflow as tf
2 def hamming_distance_sample(sents, tau, bos_id, eos_id, pad_id, vocab_size):
3     """Sample a batch of corrupted examples from sents.
4
5     Args:
6         sents: Tensor [batch_size, n_steps]. The input sentences.
7         tau: temperature.
8         vocab_size: to create valid samples.
9
10    Returns:
11        sents: Tensor [batch_size, n_steps]. The corrupted sentences.
12    """
13
14    # mask
15    mask = [
16        tf.equal(sents, bos_id),
17        tf.equal(sents, eos_id),
18        tf.equal(sents, pad_id),
19    ]
20    mask = tf.stack(mask, axis=0)
21    mask = tf.reduce_any(mask, axis=0)
22
23    # first, sample the number of words to corrupt for each sentence
24    batch_size, n_steps = tf.unstack(tf.shape(sents))
25    logits = -tf.range(tf.to_float(n_steps), dtype=tf.float32) * tau
26    logits = tf.expand_dims(logits, axis=0)
27    logits = tf.tile(logits, [batch_size, 1])
28    logits = tf.where(mask,
29                      x=tf.fill([batch_size, n_steps], -float("inf")), y=logits)
30
31    # sample the number of words to corrupt at each sentence
32    num_words = tf.multinomial(logits, num_samples=1)
33    num_words = tf.reshape(num_words, [batch_size])
34    num_words = tf.to_float(num_words)
35
36    # <bos> and <eos> should never be replaced!
37    lengths = tf.reduce_sum(1.0 - tf.to_float(mask), axis=1)
38
39    # sample corrupted positions
40    probs = num_words / lengths
41    probs = tf.expand_dims(probs, axis=1)
42    probs = tf.tile(probs, [1, n_steps])
43    probs = tf.where(mask, x=tf.zeros_like(probs), y=probs)
44    bernoulli = tf.distributions.Bernoulli(probs=probs, dtype=tf.int32)
45
46    pos = bernoulli.sample()
47    pos = tf.cast(pos, tf.bool)
48
49    # sample the corrupted values
50    val = tf.random_uniform(
51        [batch_size, n_steps], minval=1, maxval=vocab_size, dtype=tf.int32)
52    val = tf.where(pos, x=val, y=tf.zeros_like(val))
53    sents = tf.mod(sents + val, vocab_size)
54
55    return sents
```

A.6 Source Code for Sampling in PyTorch

Hamming distance sampling in Pytorch

```
1 """
2 Sample a batch of corrupted examples from sents.
3
4 Args:
5     sents: Tensor [batch_size, n_steps]. The input sentences.
```

```

6     tau: Temperature.
7     vocab_size: to create valid samples.
8 Returns:
9     sampled_sents: Tensor [batch_size, n_steps]. The corrupted sentences.
10    """
11
12 mask = torch.eq(sents, bos_id) | torch.eq(sents, eos_id) | torch.eq(sents, pad_id)
13 lengths = mask.float().sum(dim=1)
14 batch_size, n_steps = sents.size()
15 # first, sample the number of words to corrupt for each sentence
16 logits = torch.arange(n_steps)
17 logits = logits.mul_(-1).unsqueeze(0).expand_as(
18     sents).contiguous().masked_fill_(mask, -float("inf"))
19 logits = Variable(logits)
20 probs = torch.nn.functional.softmax(logits.mul_(tau), dim=1)
21 num_words = torch.distributions.Categorical(probs).sample()
22
23 # sample the corrupted positions.
24 corrupt_pos = num_words.data.float().div_(lengths).unsqueeze(
25     1).expand_as(sents).contiguous().masked_fill_(mask, 0)
26 corrupt_pos = torch.bernoulli(corrupt_pos, out=corrupt_pos).byte()
27 total_words = int(corrupt_pos.sum())
28 # sample the corrupted values, which will be added to sents
29 corrupt_val = torch.LongTensor(total_words)
30 corrupt_val = corrupt_val.random_(1, vocab_size)
31 corrupts = torch.zeros(batch_size, n_steps).long()
32 corrupts = corrupts.masked_scatter_(corrupt_pos, corrupt_val)
33 sampled_sents = sents.add(Variable(corrupts)).remainder_(vocab_size)
34
35 return sampled_sents

```