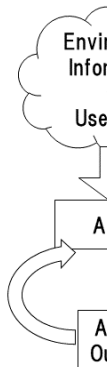
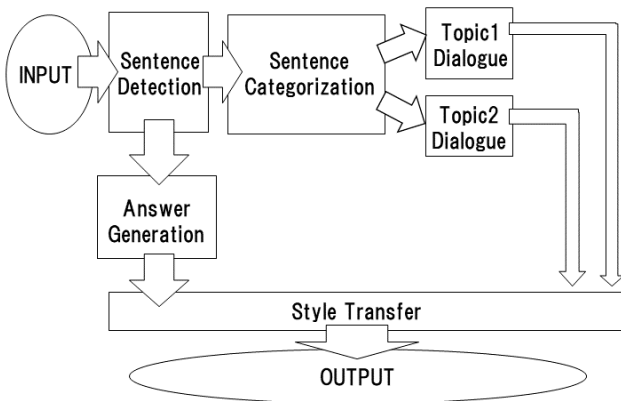




# 全体像



## 今回報告する部分

- スタイル変換 (Style Transfer)
- 1 : 1 対話 (Topic Dialogue)
- 不均衡データ分散なデータセットと不均衡データ数なデータセットの関係 (Question Detection)
- 日本語言語理解タスクのための Subword
- 今後の実験計画

# スタイル変換

- センテンスにおける口調を変換するというタスク。
- 今回主に取り扱うのははです・ます調 $\leftrightarrow$ 話し言葉への変換。
- データセットは自作の **並行な** ものを使用 (データ数が少ないので Loss や Acc への言及はしない)
- 3つの手法で実験し、その結果を観察する。

## Sequence to better sequence

- 対象とするデータ: **非並行な** 2つのスタイルを持つセンテンスの集合

### 学習内容

- 1 入力文と出力文を一つの **VAE** に通し潜在表現を得る
- 2 入力文のスコアを **0**、出力文のスコアを **1** となるように潜在表現を評価する FF を学習する
- 3 何らかの入力  $\mathbf{x}$  に対してそのスコアを 1 になるように潜在表現を調整することで、スタイルが付与された文  $\mathbf{x}'$  を得る。

実験ではこれと、更に多くの汎化性能を得るため、VAE の入力にノイズ (一部を未知語に変換する) をかけたものの実験を行った。

# CopyNet

- 対象とするデータ: **並行な** 2つのスタイルを持つセンテンスの集合

## 学習内容

- 1 入力文を通す RNN から最後の出力 (意味ベクトル) ・それぞれの単語を表すベクトルを得る
- 2 出力文を生成する RNN と演算を行う
- 3 出力文を生成できるようにそれぞれのベクトルに対する重み付けなどを学習する

本来は **機械翻訳** のタスクで使用されている手法だが、今回対象とする問題の、 **入力文・出力文が非常に似ている** という性質からこの手法を用いた。

# 実験結果

## 実験結果

実装	入力	出力
S2BS	<p>応援する。</p> <p>今日は寒かった。</p> <p>夕飯は？</p> <p>早く寝たい。</p> <p>おはようございます。</p>	<p>応援してる。</p> <p>今日は寒かった。</p> <p>夕飯はどうでしょうか？</p> <p>早く寝た方がよいね。</p> <p>おはよう。</p>
S2BS with DAE	S2BS と同じ	
CopyNet	<p>おはようございます。</p> <p>今日は良い天気ですね。</p> <p>こんにちは。</p> <p>頑張るぞい！</p> <p>進捗どうですか？</p>	<p>おはよう。</p> <p>今日は良い天気。</p> <p>こんにちは。</p> <p>頑張るぞい！</p> <p>進捗どう？</p>

## 考察

- S2BS と S2BS with DAE にはほとんど差異が見られなかった。

これが **データ数** によるものなのか、本当になんの意味もないのかは不明である。

- S2BS、CopyNet のいずれでも **入力が正規化されていない** ても問題なく変換ができることがわかった。
- S2BS は CopyNet に比べて **表現力が大きい** ように感じられる。  
→ 学習内容から用意に予測できる。



# 1 : 1 対話

入力文 1 文に対して、前後の文脈は考慮せずに出力文を作成する問題。

## Seq2Seq with attention

以前紹介した Sequence to Sequence という一般的に機械翻訳で用いられる手法に attention という機能を追加したもの。入力文の単語としての性質を出力文により強く影響させることができる。

# Transformer

attention 機能により注目した手法。RNN や CNN は用いていないという点が革新的。現役で優秀な成績を収めている機械翻訳の手法。

# Beam search

今回紹介する Beam search とは、主に出力文を生成する際に貪欲にその場の最大確率を選択せずに、いくつかの出力文を生成し、それらの中で最も良いものを選択するためのアルゴリズムであり、これを用いることで同じモデルでもより良い結果を得ることができる。

# 実験

[間に合えば書きます。]

# 不均衡データ分散なデータセットと不均衡データ数なデータセットの関係

**本来の問題** 任意の入力文からいくつかの質問・文を抽出したい。

しかし以下の理由から画像認識の問題として (1) の問題を設定した。

- **データを十分** に用意できない。
- 2つの選択肢が考えられる。
  - 1 シンプルなクラス分類 (文 A、文 B ... その他)
  - 2 **文章類似度** を活用したクラス分類

$$\max_i f(\text{similarity}(x, (Y_i)_j))$$

$x \dots \text{input sentence}$

$Y_i \dots \text{set of sentences in class } i$

## 設定した問題

ネコ画像の集合  $X$  と、イヌの画像の集合  $Y$ 、ランダムな画像の集合  $Z$  を用いる。

$X - Y$   $X - Z$  の 2 値分類問題において、データ数の比率を変化させながらその Loss, Acc を比較する。

# 実験

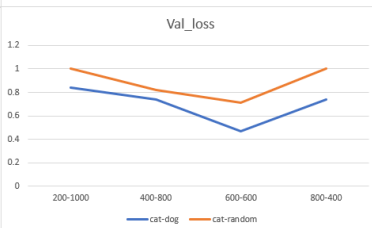
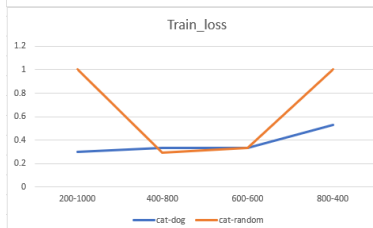
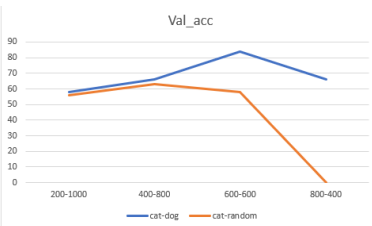
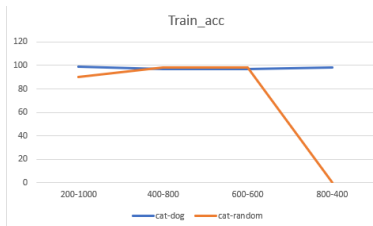
## 6 層 CNN を用いて実験を行った

- 入力画像は  $28 \times 28$  の 3 チャンネル
- 出力は  $2 \times \text{batchsize}$  (多クラス分類への拡張を想定しているため)
- データ数の比率  $x : y$  に対して  $y : x$  となるように loss の重み付けを行った。
- 最適化関数は Adam を用いた
- epoch は十分に学習ができるまでとした
- 検証データのデータ比率は 1:1 とした。



## 実験結果

以下の通り (loss = 1, val=0 は値が荒れて計測が出来なかったことを示している。)





# Subword

- 単語分割を行う手法の一つであり、一般的な単語分割より語彙数を減らすことが出来る。

例えば “subword → sub + word” を挙げることが出来る。

一般に機械翻訳の分野で用いられているが、日本語の場合では上手く分割することが難しい。(ほとんど単語分割に近くになってしまう。)

更に英語などで用いられているプログラムでは漢字かな入り混じり文のせいで上手く分割できない。

# 問題設定

カテゴリ分類やその他の機械学習を用いた自然言語に関する問題を解く際にどのように文を分割すれば良いのだろうか。

- 機械翻訳では subword は優秀だが、言語理解などではあまり優秀ではない。
- 漢字とかなが入り混じっていることで日本語の学習は無駄に難しくなっているのではないだろうか。  
→ かなのみの方が subword という意味では計算しやすい。

# 実験

fasttext の skipgram を用いて単語の分散表現を得る問題で、漢字かな入り混じり文、かなのみの文の2種類を用いて性能比較を行う。

- データセットは十分なデータを用意するために wikipedia のログを用いた
- fasttext の実装は公式が発表しているものを用いた
- loss はそれぞれの単語の分散が正しいのか(つまり意味的に近い単語が近い位置にあるのか)を計算して求めている。
- 計算される類似単語を比較した

## 実験結果

比較として、“日本 (ニホン)” を用いた

### 漢字かな入り混じり文

---

韓国  
米国  
台湾  
にっぽん  
中国  
日本さくらの会  
海外  
実業  
国内  
日本税理士会連合会

---

### かなのみの文

---

ニホンヤモリ  
ニホンバレ  
ニホンシカ  
ニホンウンソウ  
ニッポンザル  
ニホンズイセン  
ヒトツオボエ  
ゴジセイ  
ニホンカジョシユツパン  
ニホンドケン

---

# 考察

漢字かな入り混じり文は国として類似する単語を取り出していることがわかるのに対して、かなのみの文では **生物名** や日本晴れ、といった **慣用的な表現** を多く抽出している。どちらが良いのかを決めることは難しいと考えられる。ただカテゴリ分類の立場に立つのであれば、おそらく前者のほうがより良い結果を導けるのではないかと考えられる。

## ■ 極性判定

日本語言語理解タスクのための適切なフォーマット の項で議論できなかった極性判定について同様の実験を行いたい。

### ■ CoLA タスクを用いた自然言語判定

ある文 が自然なものであるかを判定する CoLA タスクを解く問題を あるモデルから出力される文 に対して同様に処理できるのかを調べる。

### ■ 文章類似度を用いたクラス分類