

19. 最小誤り率訓練

内山将夫@NICT
mutiyama@nict.go.jp

SMTの構成要素

$$\hat{\mathbf{e}} = \arg \max_{\mathbf{e}} \sum_i \lambda_i h_i(\mathbf{e}, \mathbf{f})$$

- 探索： $\arg \max_{\mathbf{e}}$ なる $\hat{\mathbf{e}}$ の探索
- モデリング： 良い素性 $h_i(\mathbf{e}, \mathbf{f})$ の設計
- パラメタ調整： λ_i の学習

最小誤り率訓練 (MERT, Minimum Error Rate Training)
は、パラメタ調整に利用される。

パラメタ調整の枠組

- 訓練データ： $h_i(e, f)$ を獲得する
- 開発データ： λ_i を獲得する
- テストデータ： 翻訳性能を測定する

パラメタ調整の原則

- 翻訳性能を最大化するパラメタが欲しい
翻訳性能を BLEU で測定するとすると , BLEU を最大化するようなパラメタが欲しい .

開発データにおける入力文を

$$F = \{f_1, f_2, \dots\}$$

参照用の翻訳文を

$$R = \{r_1, r_2, \dots\}$$

F を機械翻訳した結果を

$$E = \{e_1, e_2, \dots\}$$

としたとき ,

$$\hat{\lambda} = \arg \max_{\lambda} \text{BLEU}(R, E)$$

なるパラメタ $\hat{\lambda}$ が欲しい .

最適化としての $\hat{\lambda}$ の探索

1. $\lambda_m =$ 適当な初期値, $C_i = \phi$
2. for $\mathbf{f}_i \in F$
 - (a) $C'_i = \{\mathbf{e}_{i,s} \mid \text{スコア } \sum \lambda_m h_m(\mathbf{e}, \mathbf{f}_i) \text{ が大きい } n \text{ 翻訳文}\}$
 - (b) $C_i = C_i \cup C'_i$. (これまでの翻訳候補に加えて, 今の λ を利用して得られた翻訳候補を追加する)
3. λ を更新する
 - (a) 今の λ を利用して, 拡張された C_i の中から一番スコアが高い $\mathbf{e}_i = \arg \max_{\mathbf{e}} \sum \lambda_m h_m(\mathbf{e}, \mathbf{f}_i)$ なる \mathbf{e}_i を得ることにより, \mathbf{f}_i に対する, 今のパラメタでの翻訳文とする.
 - (b) $E = \{\mathbf{e}_i \mid \text{上記で選ばれた翻訳文}\}$ を利用して, λ に対応する $\text{BLEU}(E, R; \lambda)$ を得る.
 - (c) これにより, $\lambda \rightarrow \text{BLEU}(E, R; \lambda)$ の関係が計算できるので, λ を少しずつ変えながら, 現在の翻訳文集合 C_i から, なるべくBLEUが大きくなるように, \mathbf{e}_i を選択できるような λ を探す
4. goto 2 or exit

多変量最適化の方法

- Simplex 法 , Powell 法等のノンパラメトリック法 (関数勾配が不要な方法)
cf. Numerical Recipes in C
- 対数線形モデルに特有な方法

対数線形モデルに特有な方法

- ある方向 d について，1次元最適化をする
- 上記を，たくさんの方向に繰り返して，少しずつ解を改善して，最適解を求める．

1次元最適化を高速化する．

最小誤り率訓練

BLEU 最大化の代りに，より簡単な，誤り個数最小化の問題を考える．これを，あとで，BLEU 最大化に拡張する

誤り個数 $E(\mathbf{r}_1^s, \mathbf{e}_1^s)$ の定義

$$E(\mathbf{r}_1^s, \mathbf{e}_1^s) = \sum_{s=1}^S E(\mathbf{r}_s, \mathbf{e}_s)$$

$$\text{参照文のリスト} = \{\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_S\}$$

$$\text{翻訳文のリスト} = \{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_S\}$$

$$E(\mathbf{r}_s, \mathbf{e}_s) = \begin{cases} 1 & (\mathbf{r}_s \neq \mathbf{e}_s) \\ 0 & (\mathbf{r}_s = \mathbf{e}_s) \end{cases}$$

この $E(\mathbf{r}_s, \mathbf{e}_s)$ が、文が完全に一致するときに0で、そうでないときに1となっているので、翻訳文の評価としては、ずいぶんと簡略化されている。

誤り個数を最小化するパラメタ $\hat{\lambda}$

$$\hat{\lambda}_1^M = \arg \min_{\lambda_1^M} \sum_{s=1}^S E(\mathbf{r}_s, \mathbf{e}(\mathbf{f}_s; \lambda_1^M))$$

$$\mathbf{e}_s = \mathbf{e}(\mathbf{f}_s; \lambda_1^M) = \arg \max_{\mathbf{e} \in C_s} \sum_{m=1}^M \lambda_m h_m(\mathbf{e}, \mathbf{f}_s)$$

C_s = n 個の翻訳候補

λ_m = 素性 m の重み

$h_m(\mathbf{e}, \mathbf{f}_s)$ = 素性 m の値

\mathbf{f}_s = 入力文

誤り個数の性質

$$E(\mathbf{r}_1^s, \mathbf{e}_1^s) = \sum_{s=1}^S E(\mathbf{r}_s, \mathbf{e}_s)$$

- $E(\mathbf{r}_s, \mathbf{e}_s)$ の和が全体の値となる
- したがって、個々の誤り $E(\mathbf{r}_s, \mathbf{e}_s)$ と λ_1^M の関係がわかれば、それを加算すれば、全体の誤りと λ_1^M の関係がわかる。

さて、一次元の最適化では、ある方向 d に向けての最適化をする。その方向を M 次元ベクトル d_1^M により表現する。すると、ある定数ベクトル g_1^M を利用することにより、素性ベクトル λ_1^M は

$$\lambda_1^M = g_1^M + \gamma d_1^M$$

と表現できる。

したがって、 λ_1^M を、ある与えられた方向 d_1^M に最適化するとは、 $E(\mathbf{r}_1^s, \mathbf{e}_1^s)$ が最小となるような、 g_1^M と γ を求めることである。

ここで、

$$\mathbf{e}_s = \mathbf{e}(\mathbf{f}_s; \lambda_1^M) = \arg \max_{\mathbf{e} \in C_s} \sum_{m=1}^M \lambda_m h_m(\mathbf{e}, \mathbf{f}_s)$$

により、候補 C_s から \mathbf{e}_s を選んで、それにより、 $E(\mathbf{r}_s, \mathbf{e}_s)$ が決まる。この \mathbf{e}_s が、 λ_1^M 、つまり、 g_1^M と γ により異なる。

したがって、 \mathbf{e}_s と g_1^M 、 γ の関係が知りたい。

\mathbf{e}_s と \mathbf{g}_1^M , γ の関係

素性値のベクトルを $\mathbf{h}_1^M = \{h_1(\mathbf{e}, \mathbf{f}), \dots\}$ とする . すると , 翻訳文集合 C_s 中の候補を \mathbf{e}_i とすると , そのスコア s_i は ,

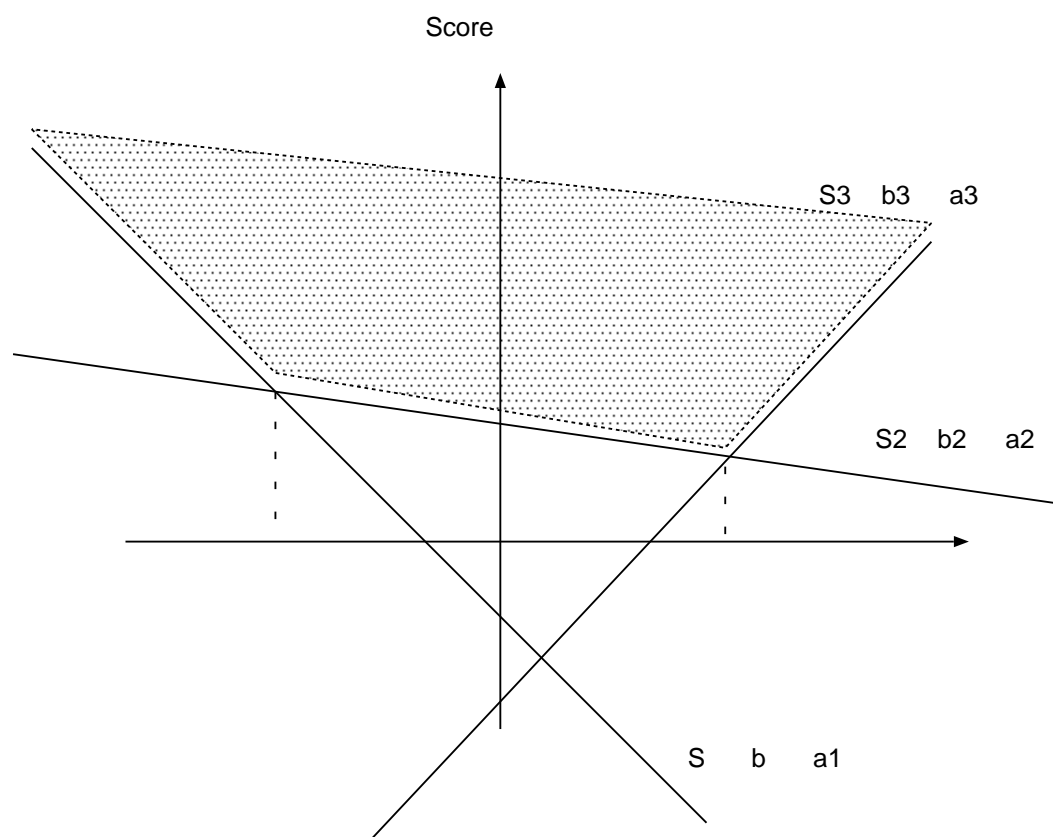
$$\begin{aligned} s_i &= \sum_{m=1}^M \lambda_m h_m(\mathbf{e}_i, \mathbf{f}_s) \\ &= \lambda_1^M \cdot \mathbf{h}_1^M \\ &= (\mathbf{g}_1^M + \gamma \mathbf{d}_1^M) \cdot \mathbf{h}_1^M \\ &= (\mathbf{g}_1^M \cdot \mathbf{h}_1^M + \gamma \mathbf{d}_1^M \cdot \mathbf{h}_1^M) \\ &= (b_i + \gamma a_i) \end{aligned}$$

ただし ,

$$\begin{aligned} b_i &= \mathbf{g}_1^M \cdot \mathbf{h}_1^M \\ a_i &= \mathbf{d}_1^M \cdot \mathbf{h}_1^M \end{aligned} \tag{1}$$

である . つまり , スコア s_i は , ある定数 a_i と b_i から定められる直線 $a_i + \gamma b_i$ の上にある . すなわち , \mathbf{e}_i のスコア s_i は , γ を変えると変わる .

γ の変化による $e_s = \arg \max_{e_i} s_i$ の変化



- $(-\infty, \gamma_1]$ のときは スコア S_1 が最大なので文 e_1 が選ばれる .
- $(-\gamma_1, \gamma_2]$ のときは , e_2 が選ばれる
- $(-\gamma_2, \infty]$ のときは , e_3 が選ばれる

初期値 = $E(\gamma = -\infty) = E(\mathbf{r}, e_1)$

左から右に動いて行って ,

γ_1 になったら $\Delta E = E(\mathbf{r}, e_2) - E(\mathbf{r}, e_1)$

γ_2 になったら $\Delta E = E(\mathbf{r}, e_3) - E(\mathbf{r}, e_2)$

のように , γ の変化と , その時点での ΔE を記録する .

すると , ある参照文 \mathbf{r} について , γ を動かしていったときに , どの時点で , 誤りの個数が変化するかがわかる .

各入力文についての結果を統合する

- 入力文 1

$$\gamma_1^1 \rightarrow \Delta E_1^1$$

$$\gamma_2^1 \rightarrow \Delta E_2^1$$

...

- 入力文 2

$$\gamma_1^2 \rightarrow \Delta E_1^2$$

$$\gamma_2^2 \rightarrow \Delta E_2^2$$

...

これらをみんなあわせると

$$\gamma_1 \rightarrow \Delta E_1$$

$$\gamma_2 \rightarrow \Delta E_2$$

...

のように，どの γ において，どの程度，誤りの個数が変化したかがわかる．

これより， $\gamma = -\infty$ のときの誤り個数に対して， $\gamma_1, \gamma_2, \dots$ と γ を変えていったときの誤りの個数の変化がわかるので，そのときの最小誤りのところの γ を利用する．この γ を利用すると，一次元方向での最小化が達成できる．そのため，この結果を利用することにより，多次元の最適化ができる．

BLEUへの拡張

$$\text{BLEU} = BP(\cdot) \exp\left(\sum_{n=1}^N \frac{\log p_n}{N}\right) \quad (2)$$

$BP(\cdot)$ = 長さの短い文へのペナルティ

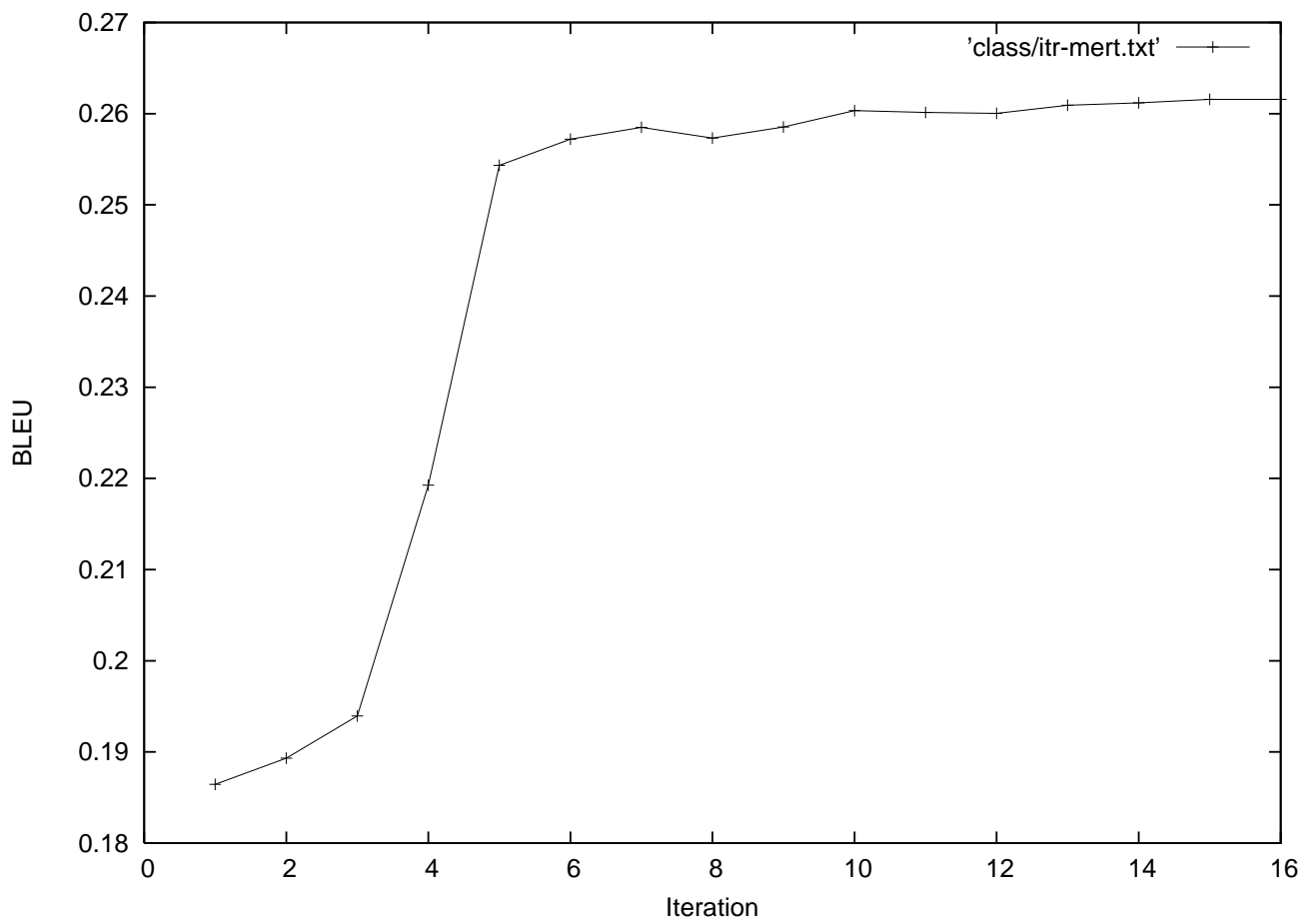
$$N = 4$$

p_n = ngram 精度

$$= \frac{\Sigma_{\text{MT 訳}} \Sigma_{\text{MT 訳の ngram 共有 ngram 数}}}{\Sigma_{\text{MT 訳}} \Sigma_{\text{MT 訳の ngram ngram 数}}}$$

拡張へのポイントは、BLEU においては、 p_n のような部分が、共有する ngram の各文に対する総和として表されることである。したがって、誤りの場合と同様に γ が変化するたびに、共有する ngram の総和が変化するので、そのたびに、共有する ngram の総和等から p_n 等を計算すれば、 γ が変化するたびの BLEU の変化がわかる。したがって、単純な誤り個数の場合と同様に、BLEU が最大となる γ がわかる。

MERTの繰り返し回数とBLEUの関係



BLEU の変化が大きいことがわかる．パラメタ値の調整は，質の良い翻訳を達成するために，必要不可欠である．

繰り返し回数と訳文の変化1

Input: thus , the left input data of node nd15 is obtained .

Reference: これにより、ノード N D 1 5 の左入力データが得られたことになる。

itr1: したがって、左入力データがノード N D 1 5 が得られる。

itr2: 以上のようにして構成されているがされているノード N D 1 5 のようにして得られたままに放置しているとされているのが、、、のデータのデータのようにして、入力されているようにされている。

itr3: このようにして、ノード N D 1 5 の左入力データが得られるようになっているようになっている。

itr4: これにより、ノード N D 1 5 の左入力データが得られる。

繰り返し回数と訳文の変化2

Input: the system arrangement shown in fig. 1 will be described in more detail below .

Reference: 図 1 のシステム構成について更に詳細に説明する。

itr1: この構成に詳しく説明する。

itr2: より以上のように構成されているのは、図1の第1の図に示すようにした場合について説明したように構成されているされているシステムの詳細な説明されるようになっているの下方に位置して説明するように構成されているようにされている。

itr3: 図1に示すように構成されて、システムのより詳細に説明する以下のようにになっている。

itr4: 図1に示すように構成され、システムの更に詳細に説明する。

itr5: 図1に示すシステム構成の詳しく説明する。

itr6: 図1に示す構成のシステム詳しく説明する。

itr7: 図1に示すシステム構成の更に詳細に説明する。

まとめ

- BLEU が最大化するようにパラメタを調整することにより，評価値に沿ったパラメタを獲得できる

これより，適正な評価を与えることが重要であることがわかる．

- 自動的に適正な評価を与えることができれば，その評価を最大化するようにパラメタを調整することにより，良いシステムを作成することができる．