

## 重みベクトルの適応的正則化手法の発音推定における評価\*

☆久保慶伍, サクティサクリアニ, グラムニュービグ, 戸田智基, 中村哲 (奈良先端大)

### 1 はじめに

文字列の発音推定は、書記素列 (Graphemes) から音素列 (Phonemes) へと変換することから g2p (grapheme-to-phoneme) 変換と呼ばれる (以後、発音推定を g2p 変換と書く)。この技術は未知語の発音を推定ことに使われ、大規模音声認識システムやテキスト音声合成システムにおいて重要な役割を果たす。最近このタスクで用いられている手法として結合系列モデル [1, 2] と Margin Infused Relaxed Algorithm (MIRA) [3] に基づく構造学習手法が挙げられる。結合系列モデルは、書記素列と音素列の断片を合わせて一つの単位とした結合 N-gram を用いる生成モデルである。MIRA は、現在対象としているデータの正解クラスのスコアが誤りのクラスのスコアよりも十分な差で高くなるように特徴量の重みを学習する多値分類のオンライン識別学習手法である。MIRA は g2p 変換のようなクラスの候補数が極端に多い構造学習問題にも拡張されており、先行研究では g2p 変換のタスクにおいて結合系列モデルよりも低い単語誤り率を実現している [4, 5]。しかしながら、MIRA は、もし現在対象としているデータが外れ値または正解ラベルが間違っているデータ (以後、このようなデータをノイズデータと書く) であっても、それを正確に分類できるように特徴量の重みを大きく動かしてしまうため、過学習を引き起こす傾向がある。

このような過学習の問題を解決するために、二値分類において、重みベクトルの適応的正則化手法 (AROW : Adaptive Regularization of Weight Vectors) [6] というオンライン識別学習が提案されている。以後、これを AROW と書く。現在対象としているデータを正しく分類できる特徴量の重みを求める MIRA とは異なり、AROW は現在のデータを正しく分類できることを保証しない代わりに、学習データを正しく分類できる方向へと特徴量の重みを少しずつ動かす。また、他のデータにおいて良く出現する特徴量の重みは、あまり出現しない特徴量の重みよりも動かさない。これにより AROW はノイズデータを分類するために特徴量の重みを大きく動かすことを防ぎ、過学習に対して頑健さを持つ。複数の二値分類タスクにおいて、AROW は、MIRA の二値分類手法と見なすことができる Passive-Aggressive (PA) アルゴリズム [7] を超える性能を示した。そのため、我々は二値分類手法である AROW を構造学習に拡張し、それを構

造学習問題である g2p 変換タスクへと初めて適用した [8]。本報告では様々なデータセットを用いた g2p 変換タスクによる AROW に基づく構造学習の評価実験について報告する。

### 2 線形分類器に基づく g2p 変換

まず最初に線形分類器に基づく g2p 変換について定義する。ある書記素列  $x$  から正しい音素列  $y$  を得るために、以下に定義される線形分類器を用いる。

$$\hat{y} = \arg \max_y w \cdot \Phi(x, y) \quad (1)$$

ここで  $w$  は分類器の特徴量の重みベクトルを意味しており、 $\Phi(x, y)$  は、 $x$  と  $y$  に出現する結合 N-gram の頻度 [5] といった特徴量から構成される特徴量ベクトルを意味している。式 (1) において、 $\hat{y}$  は動的計画法を用いることにより効率的に得ることができる。

### 3 AROW に基づくオンライン構造学習

ここでは線形分類器において用いられる重みベクトル  $w$  を得るための AROW に基づくオンライン構造学習について説明する。AROW は 2 値分類のオンライン識別学習として提案された。AROW は重みベクトルが多次元ガウシアン分布  $\mathcal{N}(\mu, \Sigma)$  に従うと仮定することで、各重みに関する更新量を以下のように制御する。頻繁に出現した特徴量 (他の多くのデータに出てきた特徴量) の重みは現在の位置に信頼性があるので大きく動かさない。逆に、今まであまり出現しなかった特徴量の重みは現在の位置に信頼性がないため大きく動かす。これにより、AROW はノイズデータを学習しても、他のデータに影響を与える重要な特徴量の重みをシステムの性能が落ちる方向へと大きく動かすことを防ぐ。これが、従来手法の MIRA と比べて、AROW が過学習に頑健な理由である。また推定の間、AROW は重みベクトルの期待値  $E[w_t] = \mu_t$  を線形分類器の重みベクトルとして用いる。

AROW を構造学習へと拡張した我々の提案手法は、 $i$  番目のデータ  $(x_i, y_i)$  と  $n$  番目の仮説  $\hat{y}_n$  が与えられた時、以下の目的関数を最小化する分布  $\mathcal{N}(\mu_t, \Sigma_t)$  を求める。

$$L(\mu_t, \Sigma_t) = \mathbf{D}_{\text{KL}}(\mathcal{N}(\mu_t, \Sigma_t) \| \mathcal{N}(\mu_{t-1}, \Sigma_{t-1})) + \frac{1}{2r} \ell_{n^2}(x_i, y_i, \hat{y}_n, \mu_t) + \frac{1}{2r} \mathbf{u}_{in}^T \Sigma_t \mathbf{u}_{in} \quad (2)$$

\* Evaluation of Adaptive Regularization of Weight Vectors on Grapheme-to-Phoneme Conversion. by Kubo Keigo, Sakriani Sakti, Graham Neubig, Tomoki Toda, Satoshi Nakamura (Nara Institute of Science and Technology)

ここで  $\mathcal{N}(\boldsymbol{\mu}_{t-1}, \boldsymbol{\Sigma}_{t-1})$  は現在の重みベクトルの分布、 $\mathbf{u}_{in}$  は正解と仮説の特徴量ベクトルの差ベクトル  $\Phi(\mathbf{x}_i, \mathbf{y}_i) - \Phi(\mathbf{x}_i, \hat{\mathbf{y}}_n)$ 、 $r > 0$  はパラメータの更新量を制御するためのハイパーパラメータである。 $\ell_{h^2}(\mathbf{x}_i, \mathbf{y}_i, \hat{\mathbf{y}}_n, \boldsymbol{\mu}_t)$  は以下に定義される損失関数である。

$$\ell_{h^2}(\mathbf{x}_i, \mathbf{y}_i, \hat{\mathbf{y}}_n, \boldsymbol{\mu}_t) = (\max\{0, d(\mathbf{y}_i, \hat{\mathbf{y}}_n) - \boldsymbol{\mu}_t \cdot \mathbf{u}_{in}\})^2 \quad (3)$$

ここで  $d(\mathbf{y}_i, \hat{\mathbf{y}}_n)$  は損失値であり、g2p 変換では音素誤り率などが用いられる。

式 (2) を  $\boldsymbol{\mu}_t$  で偏微分し、0 と置くことで、以下に定義される AROW に基づくオンライン構造学習の  $\boldsymbol{\mu}_t$  に関する更新式を得る。

$$\boldsymbol{\mu}_t = \boldsymbol{\mu}_{t-1} + \frac{\max\{0, d(\mathbf{y}_i, \hat{\mathbf{y}}_n) - \boldsymbol{\mu}_t \cdot \mathbf{u}_{in}\}}{\mathbf{u}_{in}^T \boldsymbol{\Sigma}_{t-1} \mathbf{u}_{in} + r} \boldsymbol{\Sigma}_{t-1} \mathbf{u}_{in} \quad (4)$$

g2p 変換における特徴の数は巨大であるため、それらの共分散関係を扱うことは困難である。そのため、我々は  $\boldsymbol{\Sigma}_t$  を対角行列であると仮定する。式 (2) の目的関数を  $\boldsymbol{\Sigma}_t$  の  $p$  番目の対角行列の要素  $(\boldsymbol{\Sigma}_t)_{p,p}$  で偏微分し、0 と置くと、以下のように  $\boldsymbol{\Sigma}_t$  に関する更新式を得る。

$$\frac{\partial}{\partial (\boldsymbol{\Sigma}_t)_{p,p}} L(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t) = \frac{1}{2} \left( \frac{1}{(\boldsymbol{\Sigma}_{t-1})_{p,p}} - \frac{1}{(\boldsymbol{\Sigma}_t)_{p,p}} + \frac{(\mathbf{u}_{in})_p^2}{r} \right) = 0 \quad (5)$$

ここで  $(\mathbf{u}_{in})_p$  は  $\mathbf{u}_{in}$  における  $p$  番目の特徴量を意味する。上記の式を  $(\boldsymbol{\Sigma}_t)_{p,p}$  に関する式に以下のように変形する。

$$(\boldsymbol{\Sigma}_t)_{p,p} = \frac{r(\boldsymbol{\Sigma}_{t-1})_{p,p}}{r + (\mathbf{u}_{in})_p^2 (\boldsymbol{\Sigma}_{t-1})_{p,p}} \quad (6)$$

$p = 1, \dots, d$  の各対角要素  $(\boldsymbol{\Sigma}_t)_{p,p}$  は式 (6) により更新する。また、 $\ell_{h^2}(\mathbf{x}_i, \mathbf{y}_i, \hat{\mathbf{y}}_n, \boldsymbol{\mu}_{t-1})$  が 0 の時、 $\boldsymbol{\mu}_{t-1}$  と  $\boldsymbol{\Sigma}_{t-1}$  は更新しない。

AROW に基づくオンライン構造学習の手続きを **Algorithm 1** に示す。 $\boldsymbol{\mu}$  と  $\boldsymbol{\Sigma}$  は 0 ベクトルと単位行列により各々初期化される。 $(\boldsymbol{\Sigma}_0)_{p,p} = 1$  と  $r > 0$ 、式 (6) から、 $(\boldsymbol{\Sigma}_{t-1})_{p,p} \geq (\boldsymbol{\Sigma}_t)_{p,p}$  が全ての  $t$  において成り立つ。 $(\boldsymbol{\Sigma}_t)_{p,p} = 0$  の時、 $\boldsymbol{\mu}$  の  $p$  番目の特徴量の重みは固定される。故に **Algorithm 1** の収束は保証される。**Algorithm 1** において、 $N$ -best 仮説  $\hat{\mathbf{y}}_1, \dots, \hat{\mathbf{y}}_N$  は文献 [4] と同様にフレーズ単位デコーダ [9] に基づくビームサーチにより近似的に推定される。

## 4 評価実験

提案手法である AROW に基づくオンライン構造学習を g2p 変換タスクにおいて評価する。表 1 はこの実験において用いたデータセットのデータ名 (Dataset)、出現する書記素と音素の種類数 (g/p: g が書記素, p

---

### Algorithm 1 AROW に基づくオンライン構造学習

---

**Input:** Training dataset  $\mathbf{D} = \{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_{|D|}, \mathbf{y}_{|D|})\}$

**Output:**  $\boldsymbol{\mu}$  as weight vector  $\boldsymbol{w}$

$\boldsymbol{\mu} = \mathbf{0}, \boldsymbol{\Sigma} = \mathbf{I}$

**repeat**

**for**  $i = 1$  to  $|D|$  **do**

    Predict  $N$ -best hypotheses  $\hat{\mathbf{y}}_1, \dots, \hat{\mathbf{y}}_N$  by  $\boldsymbol{\mu} \cdot \Phi(\mathbf{x}_i, \hat{\mathbf{y}})$

**for**  $n = 1$  to  $N$  **do**

**if**  $\ell_{h^2}(\mathbf{x}_i, \mathbf{y}_i, \hat{\mathbf{y}}_n, \boldsymbol{\mu}) > 0$  **then**

        Update  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$  by Eq.(4) and Eq.(6) respectively

**end if**

**end for**

**end for**

**until** Stop condition is met

---

が音素の種類数に対応)、学習データ数 (Train)、開発データ数 (Dev)、テストデータ数 (Test)、交差検定の回数 (K-fold) を示している。データセットの NETtalk (English)、Brulex (French)、Beep (English) は、Pascal Letter-to-Phoneme Conversion Challenge<sup>1</sup> から得た単語の発音辞書である。また、CMUdict (English)<sup>2</sup>、Celex (English, German, Dutch)<sup>3</sup> もまた単語の発音辞書である。文献 [2] の実験で用いられているデータセット (NETtalk, Brulex, Beep, CMUdict) において、我々は、学習データから開発データをランダムに選んだことを除いて、書記素列が 1 文字で構成されるといった例外データの取り除き方、学習データ数 (+ 開発データ数) とテストデータ数の割合に関して、文献 [2] の実験の再現を試みた。また、AROW に基づくオンライン構造学習が過学習に対して頑健であることを確かめるため、我々は学習データの 10% の書記素列に対して辞書内の音素列をランダムに付与することでノイズデータを人工的に作り出し、新しく Noisy NETtalk データセットを作成した。Noisy NETtalk において、過学習に対して頑健性を持たない手法の性能は、ノイズデータを過学習することにより劣化すると考えられる。表 1 の Noisy は人工的に作り出したノイズデータの数を示している。Noisy NETtalk は 17595 個の語彙のうち、1760 個のノイズデータを含んでいる。また、開発データ (Dev) は、ハイパーパラメータなどといった学習により決定できないパラメータを決定するためのデータ数を意味している。

比較手法の g2p 変換ツールとして、Sequitur<sup>4</sup> と Di-

---

<sup>1</sup><http://pascallin.ecs.soton.ac.uk/Challenges/PRONALSYL/Datasets>

<sup>2</sup><http://www.speech.cs.cmu.edu/cgi-bin/cmudict>

<sup>3</sup><http://www ldc.upenn.edu/Catalog/catalogEntry.jsp?catalogId=LDC96L14>

<sup>4</sup><http://sequitur.info/>

Table 1 g2p 変換タスクの評価実験で使用するデータセット.

Dataset	g/p	Vocabulary size			
		Train (Noisy)	Dev	Test	K-fold
NETtalk	26/50	17595	1000	1000	10
Noisy NETtalk	26/50	17595 (1760)	1000	1000	10
Brulex	40/39	23353	1373	2747	1
CMUdict	27/39	100886	5941	12000	1
Beep	26/44	169823	8938	19862	1
CELEX English	26/53	39995	15000	5000	1
CELEX German	30/59	206807	25851	77552	1
CELEX Dutch	41/44	196587	24573	73721	1

Table 2 各手法において設定が必要な特徴量とパラメータ.

	Sequitur	DirecTL+	Proposed
joint n-gram	<b>5,6,7,8,9,10</b>	Follow Sequitur	Follow Sequitur
context window	-	<b>4,5,6</b>	Follow DirecTL+
<i>n</i> -best hypotheses	-	<b>1,3,5</b>	Follow DirecTL+
hyperparameter <i>r</i>	-	-	<b>500,1000,1500</b>
beam width	-	<b>150</b>	<b>150</b>

recTL+<sup>5</sup>を用いた。Sequitur は書記素列と音素列の結合 N-gram の生成モデルである結合系列モデルが実装されている。DirecTL+ は MIRA に基づくオンライン構造学習が実装されている。提案手法と DirecTL+ は文献 [5] に従い、文脈特徴量 (Context features), 連鎖特徴量 (Chain features), 結合 N-gram 特徴量 (Joint n-gram features) を用いている。表 2 はそれらの特徴量や設定が必要なパラメータの詳細を示している。文献 [5] の遷移特徴量 (transition features) は NETtalk において性能の劣化が見られたため用いなかった。NETtalk において文脈窓サイズと結合 N-gram サイズ, ハイパーパラメータ *r*, 学習時における *N*-best 仮説, ビームサーチのビーム幅, 学習の繰り返し回数は, 各交差検定において開発データの音素誤り率が最小になるように決定した。太字は 10 回の NETtalk の交差検定中, 一度でも用いられた値を示している。また, NETtalk 以外の他のデータセットに関して, 特徴量とパラメータは NETtalk の実験において多く採

<sup>5</sup><http://code.google.com/p/directl-p/>

用された値を使用し, 学習回数とハイパーパラメータ *r* は NETtalk と同じ方法で決定した。他のデータセットで用いた特徴量とパラメータは文脈窓サイズが 6, 結合 N-gram サイズが 5, 学習時における *N*-best 仮説が 5 である。また, NETtalk の実験においてビーム幅を 50 にしても性能の劣化が見られなかったため, 他のデータセットでは探索のビーム幅を 50 にした。書記素列と音素列の最小単位を決めるアライメントに関して, 我々は mpaligner<sup>6</sup> に実装されている文献 [10] の制約なし多対多アライメント手法を用いた。提案手法と MIRA の損失値は音素誤り率を用いた。

表 3 は評価実験の結果を示している。PER と WER は各々音素誤り率と単語誤り率を意味し, “±” は 90% 信頼区間を示している。NETtalk, CELEX の German と Dutch を除いて, 提案手法は Sequitur と DirecTL+ の音素誤り率と単語誤り率を改善している。NETtalk において, 提案手法は DirecTL+ と同等の性能であるのに対し, Noisy NETtalk では提案手法が DirecTL+ を上回る性能を示している。この結果は AROW に基づくオンライン構造学習が, 二値分類の場合と同様に, MIRA の過学習問題を解決していることを示している。そのため, 他のデータセットにおいても性能の改善が見られたと考えられる。

Sequitur が他の手法と比べて CELEX の German と Dutch において性能を改善したのは, 極端に低い誤り率と Sequitur のバッチ学習によるものだと考えられる。提案手法や MIRA で採用されるオンライン学習では, 個々のデータを使って重みベクトルを更新するたびに, 過去に学習したデータの識別精度が薄れていく。そのため, 極端に低い誤り率を持つ CELEX の German と Dutch において, 提案手法や MIRA は過去に学習したデータの識別精度が薄れる影響により, 全てのデータを同時に学習する Sequitur のバッチ学習よりも高い性能を示すことができなかつたと考えられる。

また, NETtalk 以外のデータセットにおいて, Sequitur よりも DirecTL+ の性能が劣っていた。文献 [4, 5] では Sequitur に実装された結合系列モデルよりも DirecTL+ に実装された MIRA の方が高い性能を示している。これは今回ランダムに選択した学習データやテストデータの影響によるものだと考えられる。より正確な評価のため, NETtalk のようにクロスバリデーションによる評価が必要だと考えられる。

## 5 まとめ

我々は AROW をオンライン構造学習へと拡張し, 様々なデータセットを用いた g2p 変換タスクにおいて評価した。評価実験において提案手法は MIRA に基

<sup>6</sup><http://sourceforge.jp/projects/mpaligner/>

Table 3 g2p 変換タスクにおける評価実験の結果.

Dataset	Measure	Sequitur	DirecTL+	Proposed
NETtlak	PER	7.63%±0.24	<b>6.75%±0.22</b>	<b>6.75%±0.20</b>
	WER	31.54%±0.80	<b>28.15%±0.76</b>	<b>28.56%±0.62</b>
Noisy NETtlak	PER	<b>9.78%±0.23</b>	10.33%±0.27	<b>9.79%±0.45</b>
	WER	34.01%±0.85	<b>33.52%±0.46</b>	<b>33.02%±0.95</b>
Brulex	PER	1.30%	1.97%	<b>1.12%</b>
	WER	6.70%	8.26%	<b>5.75%</b>
CMUdict	PER	6.80%	7.25%	<b>6.09%</b>
	WER	28.83%	28.99%	<b>26.38%</b>
Beep	PER	2.85%	4.45%	<b>2.22%</b>
	WER	15.37%	19.58%	<b>12.00%</b>
CELEX English	PER	2.83%	4.23%	<b>2.51%</b>
	WER	13.33%	16.76%	<b>11.83%</b>
CELEX German	PER	<b>0.08%</b>	0.31%	0.13%
	WER	<b>0.67%</b>	1.86%	1.04%
CELEX Dutch	PER	<b>0.08%</b>	1.05%	0.27%
	WER	<b>0.67%</b>	5.28%	1.85%

づくオンライン構造学習よりも過学習に頑健で、g2p 変換の性能を改善することを示した。今後の課題として、NETtalk 以外のデータセットをクロスバリデーションにより評価することや、提案手法の性能をさらに改善するために、メモリの制限内で  $\Sigma$  における 2 つの特徴量間の共分散関係を近似的に扱う手法を考えることが挙げられる。

**謝辞** 本研究の一部は、JSPS 科研費 24240032 および（独）情報通信研究機構の委託研究「知識・言語グリッドに基づくアジア医療交流支援システムの研究開発」の助成を受けたものである。

## 参考文献

- [1] S. Deligne and F. Bimbot, "Inference of variable-length linguistic and acoustic units by multigrams," *Speech Communication*, vol.23, no.3, pp.223–241, 1997.
- [2] M. Bisani and H. Ney, "Joint-sequence models for grapheme-to-phoneme conversion," *Speech Communication*, vol.50, no.5, pp.434–451, 2008.
- [3] K. Crammer and Y. Singer, "Ultraconservative online algorithms for multiclass problems," *Journal of Machine Learning Research*, vol.3, pp.951–991, 2003.
- [4] S. Jiampojamarn and G. Kondrak, "Online discriminative training for grapheme-to-phoneme conversion," *Proc. INTERSPEECH*, pp.1303–1306, 2009.
- [5] S. Jiampojamarn, C. Cherry, and G. Kondrak, "Integrating joint n-gram features into a discriminative training framework," *Proc. NAACL-HLT*, pp.697–700, 2010.
- [6] K. Crammer, A. Kulesza, and M. Dredze, "Adaptive regularization of weight vectors," *Advances In Neural Information Processing Systems*, vol.23, pp.414–422, 2009.
- [7] K. Crammer, O. Dekel, J. Keshet, S. Shalev-Shwartz, and Y. Singer, "Online passive-aggressive algorithms," *Journal of Machine Learning Research*, vol.7, pp.551–585, 2006.
- [8] 久保慶伍, サクティサクリアニ, グラムニュービグ, 戸田智基, 中村哲, "重みベクトルの適応的正則化に基づく発音推定," *信学技報*, 第 113 巻, pp.25–30, 2013.
- [9] R. Zens and H. Ney., "Improvements in phrase-based statistical machine translation," *Proc. NAACL HLT*, pp.257–264, 2004.
- [10] K. Kubo, H. Kawanami, H. Saruwatari, and K. Shikano, "Unconstrained many-to-many alignment for automatic pronunciation annotation," *Proc. AP-SIPA*, pp.1–4, 2011.