

## 音響モデルの知識蒸留における正解ラベルと温度パラメータの利用

Knowledge Distillation of Acoustic Models with Reference Labels and Temperature Parameter

太刀岡勇気 Yuuki Tachioka

デンソーアイティラボラトリ Denso IT laboratory

## 1 はじめに

小規模(生徒)モデル学習時に、高精度(教師)モデルのソフトラベルを教師ラベルとして使う知識蒸留処理により、書き起こしを元としたハードラベルに基づく学習よりも性能が向上する[1]。本稿では、ハードラベルも付加的に使う方法について、2種のハードラベルの利用法と温度パラメータについて検討した。

## 2 正解ラベルを利用した知識蒸留処理

クロスエントロピー基準の音響モデルの学習では、書き起こしに基づくハードラベル  $w_i$  が使われ、損失関数は  $\mathcal{L}_H = -\sum_i w_i \ln s_i$  となる。 $s_i$  は生徒モデルのHMM状態  $i$  の出力確率である。一方、知識蒸留処理では、ハードラベルではなく、教師ラベル  $t_i$  を用いた損失関数  $\mathcal{L}_S = -\sum_i t_i \ln s_i$  を使う。ハードラベルは正解  $i$  に対してのみ one-hot であるが、 $t_i \geq 0$  ( $for \forall i$ ) である。

**Sequence-level distillation (SD)** では、確率  $(1 - \alpha)$  で  $\mathcal{L}_S$  が、 $\alpha$  で  $\mathcal{L}_H$  が選択される。損失関数は

$$\mathcal{L}_{SD} = \sigma(\alpha - r)\mathcal{L}_H + \sigma(r - \alpha)\mathcal{L}_S \quad (1)$$

となる。 $\sigma$  は階段関数、 $r$  は0から1の一様乱数である。ここでは、それを発話ごとに切り替えることとした。

**Sequence-level interpolation (SI)** は選択ではなく、 $\alpha$  によりソフトラベルとハードラベルを内挿する。

$$\mathcal{L}_{SI} = -\sum_i [\alpha w_i + (1 - \alpha)t_i] \ln s_i \quad (2)$$

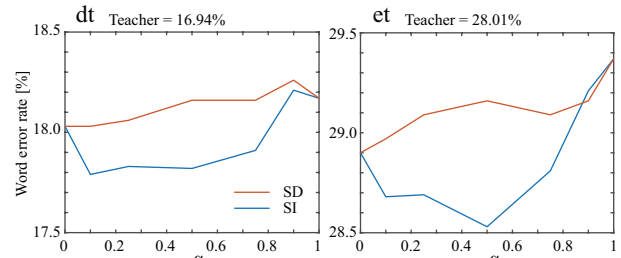
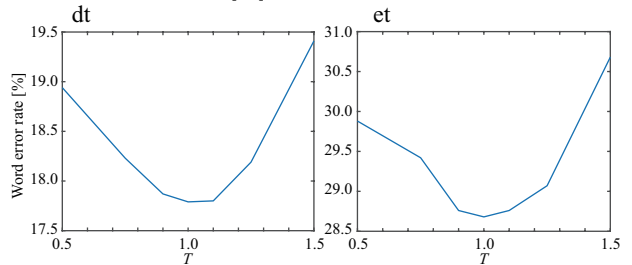
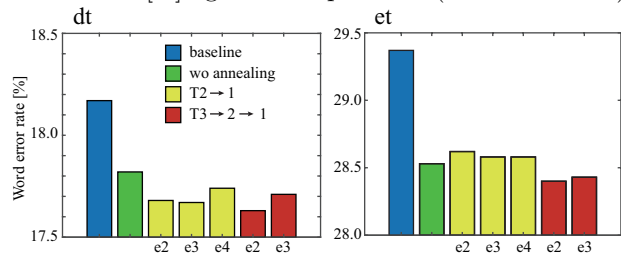
温度パラメータ  $T$  を導入することで、 $t_i$  の分布形状を変化させる。教師モデルの最終層の出力を  $z_i$  として、 $t_{i,T} = \exp(z_i/T) / \sum_j \exp(z_j/T)$  のようにしてラベルを作る。これにより、 $t_{i,1} = t_i$  で、 $t_{i,T>1}$  の分布が平滑に、 $t_{i,T<1}$  の分布が急峻になる。

## 3 実験(第4回 CHiME チャレンジ 1chトラック)

語彙数 5,000 の騒音下音声認識実験により、提案法の有効性を検証した。4種の騒音環境で、開発セット(dt)、評価セット(et)に対する平均単語誤り率(WER)で評価した。音響特徴量は、13次元MFCCとその動的特徴量に、特徴量空間最尤線形回帰を施したものである。

入・出力次元は 440, 1987とした。教師モデルは 30Mパラメータ(2048×7層)である。生徒モデルは 2Mパラメータ(512×4層)である。教師ラベルは  $t_i < 0.01$  となる  $t_i$  は0とし、 $\sum_i t_i = 1$  となるように再度正規化した。

図1は、教師モデルのWERと、SD/SIの各  $\alpha$  における生徒モデルのWERを示す。 $\alpha = 1$  が生徒モデルの原性能、 $\alpha = 0$  が知識蒸留の場合である。SD/SIともに知

図1 WER[%] against  $\alpha$  of SD and SI.図2 WER[%] against temperature (SI and  $\alpha = 0.1$ ).図3 Effect of annealing (SI and  $\alpha = 0.5$ ).

識蒸留は有効だが、SDは  $\alpha = 0$  が最良で内挿により性能が低下した。これに対し、SIは適当な  $\alpha$  で性能が向上した。図2には、温度パラメータ  $T$  とWERの関係を示す。 $T$  を1以外にするとWERが悪化することが分かった。これに対して、図3のように、 $T = 2$  を  $e(= 2, 3, 4)$  エポック分それぞれを繰り返して  $T = 1$  とする  $T \rightarrow 1$ 、もしくは、 $e$  エポックずつ  $T = 3 \rightarrow 2 \rightarrow 1$  と下げていくと、WERが改善することが分かった。また、生徒モデルが学習済の場合、教師ラベルで再学習することで学習時間を短縮できる。再学習でも、同様の設定で 17.77%(dt)、28.42%(et) が得られ、再学習で十分なことが分かった。

## 4 まとめ

知識蒸留において、SDとSIの2手法を検証し、SIの方が性能が一貫してよいことを示した。温度パラメータとともにアニーリングを行うと、さらに性能が向上した。

## 参考文献

- [1] Y. Kim and A.M. Rush, "Sequence-level Knowledge Distillation", Proc. of EMNLP, 1317-1327, 2016