

ResNeXT and Res2Net structures for speaker verification

背景

ResNet ベースの手法は 話者ごとの特徴を広く用いられてきた.

テスト非依存の話者照合システムに必要.

ResNet は

- CNN の残差接続を用いる.

- 残差ブロックを正規化する.

これにより、照合精度を向上.

が、

入力特徴量空間が複雑になる.

ResNet ネットワークの深さが増えるだけで、潜在的なパフォーマンスを發揮できていない.

この論文の時点 (2020) では、バックボーンは CNN (ResNet-based) が主流だった.

→ ResNet の拡張をいふ.

論文の主旨

話者照合タスク向けに、ResNet の拡張として 'ResNeXT' と 'Res2Net' をいふ.

- ResNeXT と Res2Net は、ともに image recognition で使われていたもの.

- モデルの表現力を上げるために、深さを増やす、cardinality と scale という 2 つの次元を導く.

- Scale のみを上げることで、与えられた深さで特徴表現をより深い層まで照合し続けることができる.

Res2Net は、

3 つの方向に拡張.

- Variety (多様性): ResNeXT, Res2Net は、従来の ResNet よりかなり良い性能.

- Res2Net の SER は 18.5% 減少.

- 他 2 つ (内装と、不均衡条件)

- サイズとセグメント (長さ) の変化は、主に初期化 Up.

ResNeXT について.

ResNeXT ブロックは 1 つの層に多くの層を並列に

グループを導出し、ResNet の残差ブロックを

マルチグループに変換して replace.

入力チャンネルと複雑数のグループに分ける 2 つの変更がポイント.

ResNet について.

width, depth の 2 つで学習を調整し、モデルの表現力を上げる.

入力特徴量空間がフラットになる.

学習時とテスト時の条件の差による mismatch.

depth と width を増やすことで overfitting を抑制する.

(一般化能力 down) する.

実験結果

データセット

VoxCelebとMS内の2つのテストセット

- VoxCeleb

- テキスト非依存, 1130名と50環境下の会話を.
 - 多くの音響環境, 短時間発話を収集.
 - 検証に用いたモデルは, VoxCeleb2-dev のデータを用いる.
 - テキスト拡張 (babble, music, noise, reverb) で学習データを拡張.
 - 詳細は
 - ・VoxCeleb1 test.
 - ・VoxCeleb1-E (extended)
 - ・VoxCeleb1-H (hard, includes 同話者, 同環境)
- の3つ.

- MS-SpeechSet

- MS内製.
- 近接マイクで, 短文の発話を収集.
- MSの会議室での会話をマイクで収録したもの.
- 約150人.
- 2-15 sec 発話 (avg. 4sec)

- Carina テストセット

- MS.
- text-dependent SV.
- 13 speakers
- source with Carina or Hey Carina
- 13人3名/42環境下, 13人100% 3名.
- 4人の20%発話者, 10%環境下でマイク5/100%.
- CTC トレーニングで抽出可能, 発話を生成可能.

実験

1. on VoxCeleb2

1.1. Regular test set

→ table 2.

- ・i-スライスのモデルは test, E, H 間の EER は

1.78, 1.76, 5.07 である.

- ・ResNet は test セットで 1.64

ResNet " 1.60 で, i-スライスを越える.

- ・モデルの複雑さを V_p とする.

- ・i-スライズ (deeper, wider) がそれより性能は劣る (+50% complexity)

(ResNet-200 は 100 倍, (+50% complexity))

(+10% complexity, ResNet-200 の i-スライズ).

→ ResNet を用いた cardinality のスケールを学習する効果がある.

特に, スケールを i-スライズで定める.

1. 2. テストセットの~~学習~~に~~関係~~なく、~~学習~~によって~~学習~~される

VoxCeleb2 test 内の~~学習~~は平均 8s.

→ 2, 3, 4 sec ほど、~~学習~~される。

→ Table 3.

「~~学習~~時間が短い、性能が低下する。

が、この~~学習~~の~~長さ~~が EPR が低い。(特に ResNet)

→ ベースラインより低い。

2. MS-SV.

モデルの~~学習~~性能を見る。

「~~学習~~は 20 sec 程度で、1-2s, 2-4s, >4s の~~学習~~の~~長さ~~が異なる。

→ Table 4.

この~~学習~~はベースラインより EPR が低い。

3. Coram Set.

- VoxCeleb2 トリニティ (トモナド) を使う。

- トリニティ ↔ テスト での~~学習~~内容不一致

と、~~学習~~ ↔ テスト " 一致する。

→ Table 5.

全体として、1s が最も低い、2s が最も高い、精度は下がる。

ただし ResNet は ResNet の 4.5% 向上。

4.20 ← 4.40

quiet and TV など、向上。

→ ResNet は ResNet の精度が良い。

が、ResNet はベースラインより精度は低い。

→ ResNet は「ロバスト」ではない。

結論

ResNet と ResNet の~~学習~~者~~学習~~タスクにおける有用性を示す。

- VoxCeleb テストセットに~~学習~~する~~長さ~~は、2-5 sec ベースラインより向上。

特に ResNet が優れている。

- 短い~~学習~~や、異なる条件下、~~学習~~は good。