深層学習抽出特徴量から生成した擬似特徴量 を用いた不均衡データ多クラス画像分類

Pseudo-feature generation from feature map in deep learning for imbalanced data multi-class image classification

紺野 友彦	藤井 秀明	岩爪 道昭
Tomohiko Konno	Hideaki Fujii	Michiaki Iwazume

国立研究開発法人 情報通信研究機構 知能科学融合研究開発推進センター AI Science R&D Promotion Center in National Institute of Information and Communications Technology

If some classes of the data have an only small number of samples, the accuracies of the classes become too low. It is well known as an imbalanced data problem. We often encounter imbalanced data in reality. In a sense, all the wild data are imbalanced.

In this paper, we make pseudo-feature from feature map in lower layers of deep neural networks, and we augment the data of minor classes to improve the imbalanced-data problem. We compare our proposed method with existing ones in imbalanced data multi-class image classification problems.

1. はじめに

データの中に少量しかないクラスがあると、その クラスの識別率が極端に低下してしまう場合がある。 そのようなデータは不均衡データと呼ばれており、 現実のデータ分析では頻繁に遭遇する。不均衡デー タ問題が起きる理由は直感的には以下の通りである。 例えば、がん検診で負例が 99%、正例が 1%であっ たとしよう。学習機は全データに対して負例である と言うだけで 99%の精度で予測できてしまい、その ような学習をしてしまう傾向がある。不均衡データ 問題の解決方法としては例えば以下のような方法が 知られている 。

- アップサンプリング:マイナークラス(少量し かデータの無いクラス)のデータを重複してサ ンプリングして他のクラスのデータ量に合わせ る方法
- ダウンサンプリング:マイナークラスのデータ 量に合わせて他のクラスのサンプル数も少なく 取る方法
- コスト関数の重みの変更:クラスのサンプル数 に合わせてコスト関数の重みを変更する方法
- 人工的なマイナークラスデータの生成 (参考 SMOTE(Synthetic Minority Over-sampling Technique)[Chawla 02])
- 5. Bagging: bootstrap sampling を行い、アンサ ンブルで決定する方法。

連絡先: 紺野 友彦 tomohiko@nict.go.jp

実際には特定の方法のみを用いるのではなく、ダ ウンサンプリング+Bagging など複数の方法を組み 合わせて不均衡データ問題に取り組むことが多い。 ([Wallace 11] など参照)

ダウンサンプリングを行うとマイナークラスの識 別率が高くなる。しかしながら、マイナークラスに 合わせて他クラスの利用データも減らすことになる ため学習データの総量は少なくなってしまう。結果 として識別問題全体の予想精度が低下してしまう。

我々は以下に示すように、深層学習の途中層で得 られた特徴量を基に人工的な擬似特徴量を生成し、 マイナークラスのデータ数を増強する方法を提案す る。この方法による擬似特徴量生成は不均衡データ 問題のみに限られる手法ではないと期待される。

提案手法

提案手法の概要を図1に示した。手法は次の手続 きからなる。

- Step1: (図1左) 全データを用いて Neural Network 1 で学習を行う。
- Step2: (図1中) 学習済み Neural Network 1 の途中層 から訓練データの特徴量分布を取り出す。
- Step3: (図1中) 取り出した特徴量分布をガウス分布 と仮定し、それに従う擬似特徴量を生成する。 (2.1 で解説)
- Step4: (図1中) 生成された擬似特徴量と実訓練データ(の抽出特徴量)を合わせて Neural Network
 2 で学習を行う。



図 1: 提案手法の概要

- Step5: (図 1 右) 学習済み Neural Network 1 の上部 (特徴量を抽出した層以上の部分)と学習済み Neural Network 2 を結合して Neural Network 3 を作る。
- Step6: (図 1 右) Neural Network 3 を用いて推論を 行う。

今回は不均衡データの分析であるので、擬似特徴量 で増強するのはマイナークラスである。

2.1 擬似特徴量の生成方法

擬似特徴量の作り方を説明する。クラス毎にそれ ぞれ別の擬似特徴量を生成する。特徴量分布がガウ ス分布に従うと仮定して2つの方法で擬似特徴量を 生成した。特徴量分布の次元をDとする。1つは特 徴量がD次元それぞれで独立なガウス分布に従う と仮定して生成した擬似特徴量分布である。つまり、 $\mu_{i,c}, \sigma_{i,c}^2$ をそれぞれクラスcの特徴量分布のi次元 目の平均と分散として、クラスcの擬似特徴量のi次元目の値 $x_{i,c}$ は以下の確率で次元ごとに独立なガ ウス分布に従うと仮定して得たものである。

$$\Pr(x_{i,c}) = \mathcal{N}(\mu_{i,c}, \sigma_{i,c}^2) \text{ for } i = 1, 2, \cdots, D$$
 (1)

2つ目は特徴量が多変量ガウス分布に従うと仮定して得られた擬似特徴量である。つまり $\vec{\mu_c}, \Sigma_c$ をそれ

ぞれクラス c の D 次元特徴量分布の平均ベクトルと 分散共分散行列とし、クラス c の擬似特徴量の分布 x_c は以下の多変量ガウス分布に従うと仮定して得ら れたものである。

$$\Pr(\vec{x_c}) = \mathcal{N}_D(\vec{\mu_c}, \Sigma_c) \tag{2}$$

図2で我々の提案した擬似特徴量が識別決定境界を 押し出す直感的な様子をスケッチする。(図2左)不 均衡データの状況ではデータ数が過小なクラスは本 来の境界(破線)よりも内側にサンプルが収まって しまう。識別決定境界(黒実線)を押せれば良い。 (図2右)我々の擬似特徴量生成は実データが作る 境界(赤実線)を超える事が出来る。SMOTEでは この境界を超えることはできない。この効果によっ て識別決定境界をより押し出すことが出来る。

3. 実験

結果は表1と表2にまとめた。現時点で、我々は入 手の容易な Cifar10 [Krizhevsky 09] データセットを 用いて人工的に不均衡データの状況を作り実験を行 なった。Cifar10 は 10 クラスの画像を持ち訓練デー タ5万、テストデータ1万から構成される。クラス 毎にそれぞれ 5000 枚、1000 枚の画像がある。値は 試行の平均値である。

A I. 主体切뻐加平 (Accuracy)						
//Minon Degelin	Deceline	Down Compling	Proposed:	Proposed:		
#1111101	FMINOr Baseline Down-Sampling	Multivariate Normal	Independent Normal			
4	0.711	0.518	0.699	0.697		
5	0.662	0.513	0.660	0.655		
6	0.608	0.525	0.603	0.600		
7	0.585	0.521	0.589	0.583		
8	0.558	0.517	0.549	0.556		
9	0.511	0.522	0.501	0.532		

表 1: 全体の識別率 (Accuracy)

// Ъ.Г .	D I	ne Down-Sampling	Proposed:	Proposed:
#Minor	Baseline		Multivariate Normal	Independent Normal
4	0.344	0.356	0.479	0.444
5	0.188	0.327	0.351	0.314
6	0.227	0.349	0.323	0.296
7	0.207	0.368	0.291	0.260
8	0.250	0.333	0.298	0.280
9	0.219	0.350	0.268	0.296

表 2: 識別率が最低のクラスの識別率 (Min of class-accuracy)



図 2: 直感的なスケッチ 青: メジャークラス 赤: マ イナークラス オレンジ: 擬似的に生成されたデータ

表中の#Minor とはデータの少ないクラスの数で ある。例えば、#Minor が4であるとは学習時には 10種類の画像クラスのうち4クラスが500枚しかな く、残り6クラスは5000枚ある状況である。テスト 時には不均衡データではなく1万枚のテストデータ をそのまま用いた。表中のBaselineとは何もせずに 全てのデータを用いて Neural Network 1を使って 学習したものである。Down-Samplingとは Neural Network 1でダウンサンプリングで学習したものであ る。Proposed: Multivariate Normalとは特徴量分 布から多次元ガウス分布を用いて擬似特徴量を生成し た提案手法の結果である。Proposed: Independent Normalとは特徴量分布のそれぞれの次元で独立な ガウス分布を用いて擬似特徴量を生成した提案手法 の結果である。

3.1 実験結果のまとめ

結果は表3にまとめた。マイナークラス数が4以 上の場合、我々の提案手法は全体の識別率 (Accu-

表 3: 結果のまとめ (Summary)

	全体の識別率	最低識別率
	(Accuracy)	(Min of Acc.)
Baseline	0	×
Down-Sampling	×	0
Proposed	0	≈○

cary) は全データを利用した場合 (Baseline) と同じ 程度に良く、識別率が最低なクラスの識別率 (Min of class-accuracy) もダウンサンプリングと同じ程 度に良い識別率を得ることができた。つまり識別率 (Accucary) も識別率が最低なクラスの識別率 (Min of class accuracy) も共に良い *1 。

4. おわりに

深層学習の途中層から抽出した特徴量の分布から 擬似的な特徴量分布を人工的に生成し、それを用い てマイナークラスのデータを増強した。それによっ て実験では全体の精度も最低精度も共に良い分類を 行うことができた。我々が人工的に生成したデータ は特徴量分布の境界を超えるものである。

^{*1} 現時点の実験ではマイナークラス数が3以下の場合は 我々の提案手法はそれほど良い結果を収めなかった。例え ばマイナークラスが1クラスで残りが9クラスがメジャー の場合、これは転移学習のように残り9クラスで十分に 学習してしまい、マイナーな1クラスのサンプル数が少 なくてもマイナークラスの識別能力が高くなってしまう からでは無いかと推測している。

参考文献

- [Chawla 02] Chawla, N. V., Bowyer, K. W., Hall, L. O., and Kegelmeyer, W. P.: SMOTE: synthetic minority over-sampling technique, *Journal of artificial intelligence research*, Vol. 16, pp. 321–357 (2002)
- [Krizhevsky 09] Krizhevsky, A. and Hinton, G.: Learning multiple layers of features from tiny images (2009)
- [Wallace 11] Wallace, B. C., Small, K., Brodley, C. E., and Trikalinos, T. A.: Class imbalance, redux, in *Data Mining (ICDM), 2011 IEEE 11th International Conference on*, pp. 754–763IEEE (2011)