

1 はじめに

GAN の最適解はナッシュ均衡として知られている。ここでいうナッシュ均衡とは、Discriminator の Loss が Discriminator のパラメータに関して最小であり、かつ Generator の Loss も Generator のパラメータに関して最小な点の事である。従来の GAN は、アルゴリズムが収束するという保証が無い Gradient descent を使用していたためナッシュ均衡を得ることができなかった。しかし、GAN のようなコスト関数が非凸であり、パラメータが連続値かつ非常に高次元な場合で、ナッシュ均衡を求めるアルゴリズムは未だ知られていない。

Feature matching, Minibatch Discrimination はそんな GAN に対して、収束を促進させるようにヒューリスティックに動機付けられた手法を提案している。つまり、近似的に解を得ようとする発想である。

2 Feature Matching

feature matching は現在の Discriminator に関しての overtraining を防ぐ新しい目的関数を作ることで、GAN の不安定さに対処する。具体的には、Generator が Discriminator の中間層の特徴量の期待値にマッチするように学習させる手法であり、以下のように Generator の目的関数を定義する。

$$\min_G \|\mathbb{E}_{\mathbf{x} \sim p_{data}} \mathbf{f}(\mathbf{x}) - \mathbb{E}_{\mathbf{z} \sim p_z} \mathbf{f}(G(\mathbf{z}))\|_2^2 \quad (1)$$

ただし、 $\mathbf{f}(\mathbf{x})$ は Discriminator の中間層の出力を表す。

Discriminator はトレーニングによって、本物か生成されたものかを最も判別可能な特徴量を見つけようとするので、この手法は自然な選択である。オリジナルの GAN の Loss ではデータ分布と全く同じになる最適解が存在したが、feature matching は Generator の損失関数に手を加えるため、この最適解に達するかの保証はない。しかし、オリジナルでは不安定だったシチュエーションにおいて実験的に feature matching は効果的であったと言っている。

疑問点

- feature matching を使うときの activation とは活性化関数をかけた後のことを言っているのか。
- feature matching を適用するときにはどの層の出力後に適用すべきか。出力層前の Conv Layer に適用するのが一般的?
- feature matching の Loss に係数 λ をかけ、通常の Generator の Loss に加えるとどうなるのか?

$$\min_G \mathbb{E}_{\mathbf{z} \sim p_z} \log D(G(\mathbf{z}_i)) + \lambda \|\mathbb{E}_{\mathbf{x} \sim p_{data}} \mathbf{f}(\mathbf{x}) - \mathbb{E}_{\mathbf{z} \sim p_z} \mathbf{f}(G(\mathbf{z}))\|_2^2$$

Algorithm 1 に、Feature matching の Generator の最適化アルゴリズムを示す。

Algorithm 1 Feature Matching

Require: $\mathbf{x}_n, n \in [N]$: ミニバッチデータセット, $\mathbf{z}_n, n \in [N]$: コード

Ensure: $\mathbf{x}_n \in \mathbb{R}^D, \mathbf{z}_n \in \mathbb{R}^{100}$

- 1: $f(\mathbf{x})$: 本物データを Discriminator に入れ, 適当な中間層 (出力層の一つ手前など) の出力を得る
- 2: $f(G(\mathbf{z}))$: コードを Discriminator に入れ, 適当な中間層 (出力層の一つ手前など) の出力を得る
- 3: 次の関数を計算し, 最小化するように誤差逆伝播を行い最適化する.

$$\left\| \frac{1}{N} \sum_n \mathbf{f}(\mathbf{x}_n) - \frac{1}{N} \sum_n \mathbf{f}(G(\mathbf{z}_n)) \right\|_2^2$$

3 Minibatch Discrimination

GAN の主な失敗の 1 つに, Generator が常に同じ出力をだす mode collapse がある. これは Discriminator が各データを独立に処理し, Generator の出力が各データとより異なるようになれというメカニズムがないために引き起こされる問題である.

このタイプの失敗を避けるのに適した方法は, Discriminator に複数のデータを組み合わせて見せることを可能にさせ, 以下の minibatch discrimination を実行することである.

複数のデータを個別ではなく組み合わせてみる Discriminator は, Generator の mode collapse を避けるのを助けることができるのかもしれない. 実際, DCGAN の論文 (by Radford) であった, batch normalization の応用の成功はこの観点から, うまく説明することができる. ミニバッチ内のデータの近さをモデリングする方法の 1 つは以下の通りである.

ミニバッチを $\mathbf{X} \in \mathbb{R}^{N \times D}$ としたとき, ミニバッチの n 番目のサンプル \mathbf{x}_n を入れた時の, Discriminator のある中間層の特徴ベクトルを $\mathbf{f}(\mathbf{x}_n) \in \mathbb{R}^A$ とする. そして, $\mathbf{f}(\mathbf{x}_n)$ にテンソル $T \in \mathbb{R}^{A \times B \times C}$ をかけ, $M_n \in \mathbb{R}^{B \times C}$ を得る. M_n を c 次元行ベクトルが B 個並んだものと考え, 行ごとに $n \in [N]$ との L_1 -distance をとる. さらに, 負の指数をとり, 以下を得る.

$$c_b(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\|M_{i,b} - M_{j,b}\|_{L_1}) \quad (2)$$

ここで, $b \in [B], i, j \in [N]$ である. minibatch discrimination layer の \mathbf{x}_n についての出力 $o(\mathbf{x}_n)$ は他のミニバッチ内のデータ全てとの和で定義される.

$$o(\mathbf{x}_n)_b = \sum_{j=1}^N c_b(\mathbf{x}_i, \mathbf{x}_j) \in \mathbb{R} \quad (3)$$

$$o(\mathbf{x}_n) = [o(\mathbf{x}_n)_1, o(\mathbf{x}_n)_2, \dots, o(\mathbf{x}_n)_B] \in \mathbb{R}^B \quad (4)$$

$$o(\mathbf{X}) = \begin{pmatrix} o(\mathbf{x}_1) \\ o(\mathbf{x}_2) \\ \vdots \\ o(\mathbf{x}_N) \end{pmatrix} \in \mathbb{R}^{N \times B} \quad (5)$$

最後に, minibatch discrimination layer の出力 $o(\mathbf{x}_n)$ と中間層の特徴ベクトル $\mathbf{f}(\mathbf{x}_n)$ とを結合し, 結果を

Algorithm 2 Minibatch Discrimination

Require: $\mathbf{x}_n, n \in [N]$: ミニバッチデータセット, $\mathbf{z}_n, n \in [N]$: コード

Ensure: $\mathbf{X} \in \mathbb{R}^{N \times D}$: ミニバッチ, $\mathbf{x}_n \in \mathbb{R}^D, \mathbf{z}_n \in \mathbb{R}^{100}$

- 1: $\mathbf{f}(\mathbf{X}) \in \mathbb{R}^{N \times A}$: ミニバッチデータを Discriminator に入れ, 適当な中間層 (出力層の一つ手前など) の出力を得る
- 2: $\mathbf{f}(\mathbf{X})$ の個々のデータ $\mathbf{f}(\mathbf{x}_n)$ について以下のように $o(\mathbf{x}_n)$ を計算する

$$\begin{aligned} M &= \mathbf{f}(\mathbf{x}_n) \cdot T \in \mathbb{R}^{B \times C} \\ c_b(\mathbf{x}_n, \mathbf{x}_k) &= \exp(-\|M_{n,b} - M_{k,b}\|) \in \mathbb{R} \\ o(\mathbf{x}_n)_b &= \sum_{k=1}^N c_b(\mathbf{x}_n, \mathbf{x}_k) \in \mathbb{R} \\ o(\mathbf{x}_n) &= [o(\mathbf{x}_n)_1, o(\mathbf{x}_n)_2, \dots, o(\mathbf{x}_n)_B] \in \mathbb{R}^B \end{aligned}$$

- 3: そして, $o(\mathbf{x}_n)$ を行ベクトルとみなし, ミニバッチ全てについての結果 $o(\mathbf{X})$ を以下のようにする.

$$o(\mathbf{X}) = \begin{pmatrix} o(\mathbf{x}_1) \\ o(\mathbf{x}_2) \\ \vdots \\ o(\mathbf{x}_N) \end{pmatrix} \in \mathbb{R}^{N \times B}$$

- 4: 最後に, $o(\mathbf{X})$ を元の活性 $\mathbf{f}(\mathbf{X})$ と結合させ, 次の層に渡す
-

次の層の入力として与える. この計算を Generator からのサンプルのミニバッチ内と学習データからのサンプルのミニバッチ内でそれぞれ別々に行う. 従来のように, Discriminator は各データに対して, そのデータが本物である確率を表す数字を出力するようにする. Discriminator のタスクは実質的に, そのデータが本物かどうかを識別することである. しかし, 今ではサイド情報として, ミニバッチ内の他のサンプルを使うことができる. Minibatch Discrimination は視覚的に魅力的なサンプルを素早く生成させることを可能にさせる. Algorithm2. に Minibatch Discrimination のアルゴリズムを示す.

4 Inception Score

GAN はモデル同士を比べる評価関数が十分でない. 最も直感的で理解しやすい指標は, Annotator にサンプルの視覚的なクオリティを測ってもらうことである Tim Salimans, et al. による論文 [?] では, GAN のモデル比較のために, はじめは Amazon Mechanical Turk を用いて人に生成したデータが本物かどうかを識別してもらい, その正解率を指標としていた. しかし, 人による識別では問題設定, Annotator のモチベーションに左右されてしまう. そこで Inception score を提案した. まず条件付き分布 $p(y|x)$ を得るために全ての生成画像に Inception Model を適用する. y はあるラベルを表す. 論文では TensorFlow の事前学習済みのモデルを使用した^{*1}.

意味のある画像ならば条件付き分布 $p(y|x)$ は, 高い確率を持って識別されるため, 低いエントロピーを持つ

^{*1} <http://download.tensorflow.org/models/image/imagenet/inception-2015-12-05.tgz>.

はずである。また、私たちはモデルが多様な画像を生成することを期待しているため、 $\int_{\mathbf{z} \sim p_{\mathbf{z}}} p(y|\mathbf{x} = G(\mathbf{z})) d\mathbf{z}$ は低いエントロピーを持つはずである。以上2点を考慮すると、Inception score は以下のように定義される。

$$\exp(\mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}} KL(p(y|\mathbf{x})||p(y))) \quad (6)$$

この Inception score は人間の判断と強い相関ををもち、良い指標と言うことがわかった。また、この基準は多様性を測るので十分に大きいデータ数 (50000) で評価することが重要であるということもわかった。

参考文献

- [1] Tim Salimans, et al., Improved Techniques for Training GANs
<https://arxiv.org/abs/1606.03498>