変分オートエンコーダ (variational autoencoder)

正田 備也

masada@rikkyo.ac.jp

Contents

オートエンコーダ

変分オートエンコーダ

変分オートエンコーダの実装

オートエンコーダ (AE; autoencoder)

- ▶ dimensionality reduction(次元圧縮、次元削減)の手法の一つ
 - ▶ 高次元ベクトルを、低次元の空間へと写す手法
- ightharpoons 元の空間の次元をd、写す先の空間の次元をkとする
- ightharpoons 元のd次元ベクトルを $oldsymbol{x}_i$ 、写した後のk次元ベクトルを $oldsymbol{z}_i$ と書くことにする
- ightharpoonup AEでは、 z_i の良し悪しを問題にする
- ▶ どういう z_i なら良いと考えられているのか?

復習: 主成分分析 (PCA)

- 1. まず、データ集合を中心化(=重心を原点へ移動)する
- 2. 原点を通るベクトルのうち、データ集合が最も大きく散ら ばっている方向を向いているものを選ぶ
 - ▶ これが第1主成分。
- 3. 次に、その方向に垂直な超平面へ、データ集合を押し潰す
- 4. すると空間の次元が一つ下がるので、次元が下がった空間の中で、2. と 3. を繰り返す
 - ▶ 第2主成分、第3主成分、・・・と続けて、第k主成分まで見つける
- ightharpoonup 全体のばらつきを最もよく表す軸をk本選ぶ、ということ
 - ト これらを座標軸として設定し直し、各データ点 x_i を k 次元空間の点 z_i として表現し直すのが PCA 4 /

オートエンコーダによる次元圧縮

- AE は、データ集合に属する個々のデータ点 x_i について、それ自身をより良く再現できるような低次元表現 z_i を求める AE での低次元表現をコード (code) と呼ぶ
- ▶ オートエンコーダは2つのニューラルネット(NN)から成る
- 1. エンコーダ: 個々のデータ点 x_i を入力とし、コード z_i を出力する NN
- 2. デコーダ: コード \mathbf{z}_i を入力とし、 \mathbf{x}_i と同じ次元のベクトル $\hat{\mathbf{x}}_i$ を出力する NN
- $lacksymbol{ iny}$ 2つの NN は、 $\hat{m{x}}_i$ ができるだけ $m{x}_i$ に近くなるように訓練する

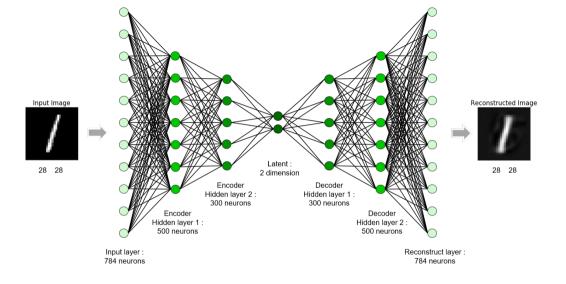
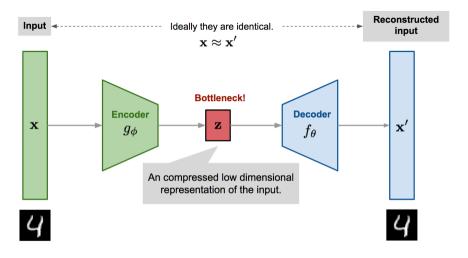


Figure: https://encodebox.medium.com/auto-encoder-in-biology-9264da118b83



 $\label{ligidiscontinuity} \textbf{Figure: https://lilianweng.github.io/lil-log/2018/08/12/from-autoencoder-to-beta-vae.html}$

Contents

オートエンコーダ

変分オートエンコーダ

変分オートエンコーダの実装

変分オートエンコーダ (VAE; variational autoencoder)

- ▶ 一見、AEと似ている・・・が、かなり違う
- $lackbrack x_i$ 自身を再現するための低次元表現として、AE で言えばエンコーダにあたる NN の出力を、そのまま使うことはない
- なぜなら、NNの出力は、それがそのままコードであるわけではなく、変分事後分布のパラメータだから
- ightharpoonup この変分事後分布から得たサンプルが、コード z_i となる
- ightharpoonup さらに、デコーダの出力は、元のデータ点 $oldsymbol{x}_i$ と直接比較できるようなものであるとは限らない
 - ightharpoons 元のデータ点 $oldsymbol{x}_i$ を生成する分布のパラメータを出力する

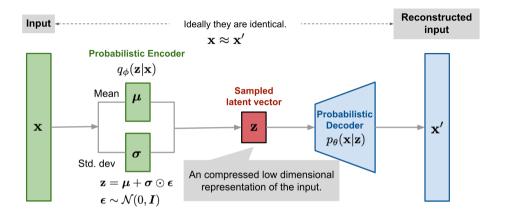


Figure: https://lilianweng.github.io/lil-log/2018/08/12/from-autoencoder-to-beta-vae html

変分ベイズ法の一種としてのVAE

- ▶ VAEを、AEの一種として理解するには、無理がある
- ▶ 変分ベイズ法から理解するのが、吉
- ▶ そうでないと、VAEが、なぜあれでいいのか、分からない
 - ▶ エンコーダからデコーダへ移るときに、なぜ、サンプルを得るというステップが挟まっているのか?
 - エンコーダについて、なぜ、KL 情報量による正則化が考えられているのか?
 - ▶ これらが、AE からの延長で VAE を考えていると、分からない

変分ベイズ法の復習

$$\ln p(\mathcal{X}) \ge \int q(\mathbf{\Theta}) \ln \frac{p(\mathbf{\Theta})p(\mathcal{X}|\mathbf{\Theta})}{q(\mathbf{\Theta})} d\mathbf{\Theta}$$
$$= \int q(\mathbf{\Theta}) \ln \prod_{i=1}^{N} p(\mathbf{x}_{i}|\mathbf{\Theta}) d\mathbf{\Theta} - \int q(\mathbf{\Theta}) \ln \frac{q(\mathbf{\Theta})}{p(\mathbf{\Theta})} d\mathbf{\Theta}$$

$$= \sum_{i=1}^{N} \int q(\mathbf{\Theta}) \ln p(\mathbf{x}_i | \mathbf{\Theta}) d\mathbf{\Theta} - D_{\mathsf{KL}}(q(\mathbf{\Theta}) \parallel p(\mathbf{\Theta})) \quad (1)$$

- $lackbox{m{\Theta}}$ の値が given なら、 $p(\mathcal{X}|m{\Theta}) = \prod_{i=1}^N p(m{x}_i|m{\Theta})$ が成り立つ
- ightharpoonup KL 情報量の項は $q(oldsymbol{\Theta})$ を $p(oldsymbol{\Theta})$ に近づけるはたらきをする

観測データのモデルのパラメータには2種類ある

- 1. 個々のデータ点に対して別々に用意されているパラメータ
 - ▶ 例:潜在的ディリクレ配分法での各文書のトピック確率 $\theta_i = (\theta_{i,1}, \dots, \theta_{i,K})$ for $i = 1, \dots, N$
- 2. データ集合全体に対して用意されているパラメータ
 - ▶ 例:潜在的ディリクレ配分法での各トピックの単語確率 $\phi_k = (\phi_{k,1}, \dots, \phi_{k,W})$ for $k = 1, \dots, K$
- ▶ どちらの種類のパラメータにも事前分布を導入できる
- ▶ VAEでは、個々のデータ点に対して別々に用意されている パラメータのほうに事前分布を導入する

VAEにおけるELBO

$$\ln p(\mathcal{X}) \ge \sum_{i=1}^{N} \int q(\mathbf{\Theta}) \ln p(\mathbf{x}_{i}|\mathbf{\Theta}) d\mathbf{\Theta} - D_{\mathsf{KL}}(q(\mathbf{\Theta}) \parallel p(\mathbf{\Theta}))$$

$$= \sum_{i=1}^{N} \int q(\mathbf{z}_{i}) \ln p(\mathbf{x}_{i}|\mathbf{z}_{i}) d\mathbf{z}_{i} - \sum_{i=1}^{N} D_{\mathsf{KL}}(q(\mathbf{z}_{i}) \parallel p(\mathbf{z}_{i})) \quad (2)$$

- ▶ 観測データのモデルのパラメータのうち、各データ点 x_i に対して別々に用意されたパラメータ z_i だけを考慮する
 - ▶ 他のモデルパラメータは自由パラメータのままでもいいし、事前 分布を使ってベイズ化されていてもいいが、VAEには関係しない

VAEにおけるnotation

- ▶ 確率モデリングの世界では、zは、離散値をとる確率変数を表すために使うことが多い
 - ▶ 例:各データ点がどのクラスタに属するかを表す潜在確率変数 *z_i*
- ightharpoonup ところが、VAEの話をするときには、z を、観測データを生成するモデルのパラメータを表すために使う
 - ▶ 気分としては、i番目の観測データ x_i の生成に関与する確率分布のパラメータは θ_i と書きたいが、なぜか z_i と書く
 - ▶ AE の世界でそう書かれていたから?
 - "We assume that the data are generated by some random process, involving an unobserved continuous random variable z." [Kingma+ arXiv:1312.6114v10]

Contents

オートエンコーダ

変分オートエンコーダ

変分オートエンコーダの実装

VAEの実装の目標

- ▶ 目標は ELBO の最大化の計算を実装すること
 - ▶ これによって、事後分布の近似としての変分事後分布が得られる
- ▶ そこで、ELBO の計算全体の計算グラフを作り・・・
- ▶ ELBO 最大化の問題を、通常のニューラルネットの学習と同じように行う(ELBO にマイナスを付けて最小化する)
- ▶ 以下、このようなことを効率的に実現するために、実際には 多くの場合どのように VAE を実装するか、説明する

VAE における変分事後分布 $q(z_i)$

- $\mathbf{p} = q(\mathbf{z}_i)$ は、観測データ \mathbf{x}_i をモデリングする確率分布のパラメータが従う分布である
 - lacktriangle この確率分布の密度関数を使って、尤度 $p(oldsymbol{x}_i|oldsymbol{z}_i)$ が表される
- ightharpoonup VAE の変分事後分布 $q(oldsymbol{z}_i)$ としては、普通、正規分布を使う
- ▶ しかも、共分散行列が対角行列であることを仮定する
- ト よって、 $q(z_i)$ のパラメータは、平均パラメータ μ_i と、分散 パラメータ σ_i^2 で、いずれも K 次元ベクトルとなる
 - ▶ 以下、この場合についてだけ説明する

VAEにおけるELBOの前向き計算

▶ VAE における ELBO は

$$\mathcal{L} = \sum_{i=1}^N \int q(oldsymbol{z}_i) \ln p(oldsymbol{x}_i | oldsymbol{z}_i) doldsymbol{z}_i - \sum_{i=1}^N D_{\mathsf{KL}}(q(oldsymbol{z}_i) \parallel p(oldsymbol{z}_i))$$

- ▶ この ELBO の計算グラフを作り、backpropagation すれば、 様々なパラメータを更新していけるが…
- $ightharpoonup q(z_i)$ についての積分はどう計算すればいい?
- ▶ KL情報量の項はどう計算すればいい?

積分のモンテカルロ近似

- $ightharpoonup \int q(\boldsymbol{z}_i) \ln p(\boldsymbol{x}_i | \boldsymbol{z}_i) d\boldsymbol{z}_i$ はモンテカルロ近似する
- ▶ つまり、 $q(z_i)$ からサンプル $\{z_i^{(1)},\dots,z_i^{(S)}\}$ を生成し、以下のように近似する

$$\int q(\boldsymbol{z}_i) \ln p(\boldsymbol{x}_i|\boldsymbol{z}_i) d\boldsymbol{z}_i \approx \frac{1}{S} \sum_{s=1}^{S} \ln p(\boldsymbol{x}_i|\boldsymbol{z}_i^{(s)})$$
(3)

▶ 通常、S=1と設定する

VAEのKL情報量

- ト いま、変分事後分布 $q(z_i)$ として、共分散行列が対角行列である正規分布を使っている
- ightharpoonup さらに、事前分布 $p(\mathbf{z}_i)$ として、成分ごとに標準正規分布を使う(大体こうする)とすると…
- **ELBO** に現れている KL 情報量の項 $-D_{\mathsf{KL}}(q(\boldsymbol{z}_i) \parallel p(\boldsymbol{z}_i))$ は、以下のように解析的に計算できてしまう

$$-D_{\mathsf{KL}}(q(\boldsymbol{z}_i) \parallel p(\boldsymbol{z}_i)) = \frac{1}{2} \sum_{k=1}^{K} (1 + \ln((\sigma_{i,k})^2) - (\mu_{i,k})^2 - (\sigma_{i,k})^2)$$

▶ 計算グラフの一部として上の式の計算が入ってくる

VAEにおけるELBOの前向き計算(続)

▶ ELBOは、いまや以下のように計算される

$$\mathcal{L} = \sum_{i=1}^{N} \ln p(\boldsymbol{x}_{i} | \boldsymbol{z}_{i}^{(1)}) + \frac{1}{2} \sum_{k=1}^{K} (1 + \ln((\sigma_{i,k})^{2}) - (\mu_{i,k})^{2} - (\sigma_{i,k})^{2})$$

- ▶ 変分事後分布 $q(z_i)$ のパラメータ μ_i と σ_i^2 はどう準備する?
 - ▶ エンコーダによって準備する
- ▶ 尤度 $p(\boldsymbol{x}_i|\boldsymbol{z}_i^{(1)})$ をどう表現する?
 - ▶ デコーダによって表現する

VAEのエンコーダ

- ▶ VAE では μ_i と σ_i^2 を NN の出力として得る
- ▶ この NN を、VAE ではエンコーダと呼ぶ
- ightharpoonup VAEのエンコーダは x_i を入力とする
- ▶ VAEのエンコーダは μ_i と σ_i^2 とを出力する
 - ▶ 観測データのモデルのパラメータ z_i を出力するのではない!
 - lacktriangle 通常は $m{\sigma}_i^2$ ではなく $\ln m{\sigma}_i^2$ を出力するように実装する
- $oldsymbol{x}_i$ に対応するコード $oldsymbol{z}_i$ は、正規分布 $oldsymbol{\mathcal{N}}(oldsymbol{\mu}_i,oldsymbol{\sigma}_i^2)$ からランダムにサンプルを生成することで得られる
 - ▶ このサンプリングが、モンテカルロ近似のためのサンプリング

VAEのデコーダ

- ightharpoonup VAEのデコーダは $oldsymbol{x}_i$ に対応するコード $oldsymbol{z}_i$ を入力とする
 - ▶ エンコーダの出力をそのまま入力とするのではない
- ▶ VAEのデコーダは、確率分布のパラメータを出力する
- ト そして、デコーダの出力をパラメータとする確率分布を使って x_i の尤度 $p(x_i|z_i)$ を求める
- 例 デコーダの出力を平均パラメータ、単位行列(の定数倍)を 共分散行列とする正規分布を使って、 x_i の尤度を求め、そ れを ELBO 最大化を通して、最大化する
 - lacktriangle これは、デコーダの出力 \hat{x}_i を、データ点 x_i に、ユークリッド距離 $\|\hat{x}_i x_i\|$ の意味で近づけることと、全く同じことになる

VAE における ELBO の前向き計算(続々)

- ightharpoonup エンコーダが表す関数を $\operatorname{Enc}(\boldsymbol{x}_i)$ 、デコーダが表す関数を $\operatorname{Dec}(\boldsymbol{z}_i)$ と書くことにする
- ▶ VAEにおける ELBO の前向き計算は、以下の通り
- 1. $\operatorname{Enc}(\boldsymbol{x}_i)$ を計算して $(\boldsymbol{\mu}_i, \boldsymbol{\sigma}_i^2)$ を得る
- 2. $(\boldsymbol{\mu}_i, \boldsymbol{\sigma}_i^2)$ を使って \boldsymbol{z}_i を生成
- 3. $Dec(z_i)$ を計算して尤度 $p(x_i|z_i)$ のパラメータを得る
- 4. 尤度 $p(\boldsymbol{x}_i|\boldsymbol{z}_i)$ とKL項を計算
- ightharpoonup しかし、 z_i を単なる数値として生成すると、ightharpoonup を越えてエンコーダへと遡れない!

reparametrization trick

- ▶ VAEのエンコーダが出力した (μ_i, σ_i^2) によって、そこからサンプルを生成すべき正規分布は確定する
- ▶ しかし、サンプル z_i を単なる数値として生成してしまうと、 BPがデコーダの入り口で止まり、エンコーダへ遡れない!
- ▶ そこで・・・
- 1. $\mathcal{N}(0,\mathbf{I}_k)$ から単なる数値としてサンプル $\boldsymbol{\epsilon}_i^{(1)}$ を生成し…
- 2. その $\epsilon_i^{(1)}$ を $\mu_i+\sigma_i\odot\epsilon_i^{(1)}$ という計算によって、本当に欲しかった $q(z_i)$ からのサンプルへ変換する
 - $m \mu_i + m \sigma_i \odot m \epsilon_i^{(1)}$ という計算は、BP を行う計算グラフの一部になる

VAEにおけるELBOの前向き計算(完)

- ightharpoonup エンコーダが表す関数を $\operatorname{Enc}(oldsymbol{x}_i)$ 、デコーダが表す関数を $\operatorname{Dec}(oldsymbol{z}_i)$ と書くことにする
- ▶ VAEにおける ELBO の前向き計算は、以下の通り
- 1. $\operatorname{Enc}(oldsymbol{x}_i)$ を計算して $(oldsymbol{\mu}_i, oldsymbol{\sigma}_i^2)$ を得る
- 2. $oldsymbol{\epsilon}_i^{(1)}$ を $\mathcal{N}(0,\mathbf{I}_k)$ から単なる数値として生成
- 3. $oldsymbol{\mu}_i + oldsymbol{\sigma}_i \odot oldsymbol{\epsilon}_i^{(1)}$ によって $oldsymbol{z}_i$ を得る
- 4. $Dec(z_i)$ を計算して尤度 $p(x_i|z_i)$ のパラメータを得る
- 5. 尤度 $p(x_i|z_i)$ と KL 項を計算

VAEにおける amortized inference

- \blacktriangleright μ_i と σ_i^2 は、 x_i を入力とする NN の出力として得た
- ▶ だが μ_i と σ_i^2 をデータ点 x_i ごとに単なる未知数として準備し、ELBO最大化によって更新するので良かったのでは?
 - ▶ 変分ベイズ法では普通こうする
- ightharpoonup 全てのデータ点に同じ一つの NN を共有させて、その出力としてデータ点ごとのパラメータ μ_i と σ_i^2 を得る、という考え方は、変分ベイズ法には元々はなかった考え方!
- ightharpoonup 一般に、変分事後分布のパラメータをデータ点 $oldsymbol{x}_i$ の関数として表現してベイズ推論することを、amortized inference と呼ぶ(エンコーダがこの関数を表現している) 28 /