



UMCS
UNIwersytet Marii Curie-Skłodowskiej

UNIwersytet Marii Curie-Skłodowskiej
w Lublinie

Wydział Matematyki, Fizyki i Informatyki

Kierunek: **Informatyka**

Specjalność: **Sztuczna inteligencja**

Filip Ręka

nr albumu: 296595

Variational Autoencoder

Variational Autoencoders

Praca licencjacka
napisana w Katedrze Cyberbezpieczeństwa
pod kierunkiem dr hab. Michała Wydry

Lublin 2021

Spis treści

Wstęp	5
1 Tradycyjny a wariacyjny autoenkoder	7
1.1 Autoencoder	7
1.1.1 Informacje ogólne	7
1.1.2 Zastosowanie	7
1.1.3 Problemy z generacją nowych danych	8
1.2 Wariacyjny autoenkoder	8
1.2.1 Różnice	8
2 Wariacyjny autoenkoder	9
2.1 Teoria za modelem VAE	9
2.1.1 Motywacja statystyczna	9
2.2 Wnioskowanie wariacyjne	9
2.2.1 Dywergencji Kullbacka-Leiblera	9
2.3 Sztuczka reparametryzacyjna	12
3 Implementacja	13
3.1 Tensorflow oraz Keras	13
Spis tabel	15
Spis rysunków	17

Wstęp

Wariacyjne Autoenkodery stają się coraz bardziej popularnymi modelami uczenia maszynowego. Zostały zaproponowane przez Diederika P Kingma i Maxa Wellinga. Najczęściej zostają one skategoryzowane do modeli uczenia częściowo nadzorowanego. Znajdują zastosowanie w generacji obrazów, tekstu, muzyki, odszumianiu obrazków, sygnałów oraz w detekcji anomalii. W przeciwieństwie do tradycyjnych autoenkoderów prezentują pobabilistyczne podejście do generowania zmiennych ukrytych. Swoją popularność zawdzięcza swojej budowie, która jest oparta na sieciach neuronowych oraz możliwości trenowania go przy pomocy metod gradientowych.

Rozdział 1

Tradycyjny a wariacyjny autoenkoder

1.1 Autoencoder

1.1.1 Informacje ogólne

Autoencoder składa się z dwóch części: enkodera, który koduje dane wejściowe oraz dekodera, który na podstawie kodu rekonstruuje wejście. Architektura enkodera wymaga aby jego warstwa wyjściowa generująca reprezentację danych była mniejsza niż warstwa wejściowa. W innym przypadku sieć po prostu przekopiuwałaby wejście. Często to zwężenie nazywa się mianem *bottle neck*. Celem treningu całego autoenkodera jest zminimalizowanie błędu pomiędzy wejściem a wyjściem. W przypadku obrazów funkcją straty może być na przykład błąd średniokwadratowy.

Powiedzmy że mamy dane wejściowe X o wymiarze m oraz chcemy je zakodować do wymiaru n . Formalnie możemy zapisać, że model próbuje nauczyć się funkcji:

$$\text{Enkoder } A : \mathbb{R}^m \rightarrow \mathbb{R}^n$$

$$\text{Dekoder } B : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

Funkcją straty naszego modelu, którą będziemy chcieli zminimalizować, będzie błąd rekonstrukcji danych wejściowych i do tego celu użyjemy błędu średniokwadratowego.

$$\mathcal{L}(x, x') = \frac{1}{m} \sum_{i=0}^m (x_i - x'_i)^2 = \frac{1}{m} \sum_{i=0}^m (x_i - B(A(x_i)))^2$$

1.1.2 Zastosowanie

Dwoma głównymi zastosowaniami tradycyjnych autoenkoderów są:

- Odszumianie

- Redukcja wymiarów

1.1.3 Problemy z generacją nowych danych

Dobrym pytaniem jest czy przy pomocy kodu jesteśmy generować nowe dane bardzo podobne do tych co model otrzymał na wejściu. Aby model mógł generować nowe dane muszą zostać spełnione dwa warunki:

- Nasza przestrzeń kodu (tzw. zmiennych ukrytych) musi być ciągła co znaczy że dwa punkty znajdujące się obok siebie będą dawać podobne dane jak zostaną odkodowane
- Przestrzeń musi być kompletna co znaczy, że punkty wzięte z dystrybucji muszą dawać wyniki mające sens

Korzystając z tradycyjnych autoenkoderów nie mamy zagwarantowanego, że oba te kryteria zostaną spełnione.

Zwykle autoenkodery po prostu nie nadają się do generacji danych ponieważ nigdy nie zostały do tego stworzone. Ich głównym celem jest jak najlepsze odzwierciedlenie danych ze zmiennych ukrytych.

1.2 Wariacyjny autoenkoder

1.2.1 Różnice

Wariacyjny autoenkoder ma inne podejście do generowania zmiennych ukrytych. Zamiast generować jedną zmienną dla każdego wymiaru, generuje dwie liczby, σ oraz μ , które traktujemy jako odchylenie standardowe oraz średnią rozkładu normalnego. Dla przykładu jeśli uznamy że chcemy dane reprezentować jako siedmio-wymiarowy wektor, nasz enkoder wygeneruje dwa wektory siedmio-wymiarowe, z którego jeden będzie przechowywał wartości średniej a drugi odchylenia standardowego dla każdego z siedmiu rozkładu normalnego. Kolejną istotną zmianą jest funkcja straty, która oprócz błędu rekonstrukcji obrazu składa się z dywergencji Kullbacka-Leiblera.

Rozdział 2

Wariacyjny autoenkoder

2.1 Teoria za modelem VAE

2.1.1 Motywacja statystyczna

Powiedzmy że istnieje zmienna ukryta z , która generuje obserwację x . Mamy tylko informację o x i chcemy się dowiedzieć jakiej jest z . Aby to zrobić powinniśmy policzyć $p(z|x)$. Z twierdzenia Bayesa wiemy że:

$$p(z|x) = \frac{p(x|z)p(z)}{p(x)}$$

Aby obliczyć rozkład marginalny $p(x)$ musimy policzyć:

$$p(x) = \int_z p(x, z) dz$$

Obliczenie tej całki jest bardzo trudne ponieważ z jest często wielowymiarowe i przestrzeń przeszukiwań jest zwyczajnie kombinatorycznie za duża aby korzystać z takich metod jak próbkowanie Monte Carlo łańcuchami Markowa.

2.2 Wnioskowanie wariacyjne

Rozwiązaniem tego problemu jest próba policzenia rozkładu $q(z|x)$, które będzie jak najlepiej odzwierciedlać $p(z|x)$ i będzie miał rozkład, który będziemy mogli policzyć.

2.2.1 Dywergencji Kullbacka-Leiblera

Jest to miara określająca rozbieżność między dwoma rozkładami prawdopodobieństwa. Nie można określić jej mianem metryki ponieważ nie jest symetryczna

$$(D_{KL}(P\|Q) \neq D_{KL}(Q\|P)).$$

Naszym celem będzie zminimalizowanie jej.

$$q^*(z|x) = \operatorname{argmin}_{q(z|x) \in Q} (D_{KL}(q(z|x)\|p(z|x)))$$

gdzie Q to rodzina prostych dystrybucji, na przykład rozkładu Gaussa

Policzmy:

$$D_{KL}(q(z|x)\|p(z|x)) = \mathbb{E}_{z \sim q(z|x)} \log \frac{p(z|x)}{q(z|x)} = \int_z q(z|x) \log \frac{q(z|x)}{p(z|x)} dz$$

Natrafiamy na kolejny problem ponieważ nie możemy $p(z|x)$ jednak jesteśmy w stanie to przepisać jako:

$$p(z|x) = \frac{p(z, x)}{p(x)}$$

Tu jest dużo matmy której nie chce mi się na razie pisać ale tu będzie ELBO (dolna granica dowodów).

Wybieramy sobie że nasza funkcja $q(z|x)$ będzie $\mathcal{N}(0, \mathbf{I})$. Dywergencja dla dwóch rozkładów normalnych wygląda w następujący sposób.

$$\frac{1}{2} \left\{ \left(\frac{\sigma_0}{\sigma_1} \right)^2 + \frac{(\mu_1 - \mu_0)^2}{\sigma_1^2} - 1 + 2 \log \frac{\sigma_1}{\sigma_0} \right\}$$

Co w naszym przypadku gdzie $\mu_1 = 0$ oraz $\sigma_1 = 1$ uprości się do:

$$\frac{1}{2} \sum_{i=1}^m \sigma_i^2 + \mu_i^2 - \log(\sigma_i) - 1$$

Jest to pierwsza część naszej funkcji straty.

2.3 Sztuczka reparametryzacyjna

Model VAE po zakodowaniu wejścia dokonuje operacji próbkowania (*sampling*) z dystrybucji na nauczonych parametrach. Przy propagacji do przodu nie jest to problem, jednak podczas propagacji wstecznej jest to nie możliwe. Operacja próbkowania nie jest różniczkowalna co sprawia, że nie możemy policzyć gradientu. Sposobem obejścia tego problemu jest zastosowanie sztuczki (*reparameterization trick*). Próbkowanie z dystrybucji $z \sim \mathcal{N}(\mu, \sigma)$ jesteśmy w stanie zapisać jako:

$$\begin{aligned}\epsilon &\sim \mathcal{N}(0, 1) \\ z &= \mu + \sigma \odot \epsilon\end{aligned}$$

Pozornie nic się nie zmieniło, jednak teraz jesteśmy w stanie poprowadzić gradient przez z , które jest teraz deterministycznie. W poprzednim przypadku było ono losowe wybierane z dystrybucji.

Rozdział 3

Implementacja

3.1 Tensorflow oraz Keras

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.[1]

Spis tabel

Spis rysunków

Bibliografia

- [1] C. Doersch, “Tutorial on variational autoencoders,” 2021.