

# Zastosowanie metryki Wassersteina w problemie uczenia ograniczonych maszyn Boltzmannna

dr inż. Maksymilian Bujok



Zakład Algebry i Kombinatoryki  
Wydział Matematyki i Nauk Informatycznych

27 marca 2019

# Plan prezentacji

- ▶ Ograniczone maszyny Boltzmannna
- ▶ Problem transportu optymalnego
- ▶ Odległość Wassersteina
- ▶ Uczenie maszyn Boltzmannna z zastosowaniem odległości Wassersteina
- ▶ Podsumowanie

# Ograniczone maszyny Boltzmanna

# Maszyny Boltzmann

- ▶ Maszyny Boltzmann możemy rozpatrywać jako sieć losowych węzłów (wierzchołków).
- ▶ Każdy węzeł (wierzchołek) możemy traktować jak dwustanowy neuron.
- ▶ Do opisu całego układu możemy posługiwać się pojęciem energii układu.
- ▶ Stan każdego neuronu zmienia się w czasie o kroku dyskretnym w sposób losowy.
- ▶ Układ dąży do stanu o niższej energii, który możemy utożsamiać z optimum globalnym.

# Maszyny Boltzmann

- ▶ Rozważmy sieć złożoną z  $N$  połączonych parami wierzchołków (neuronów) i  $I \subseteq \binom{N}{2}$ .
- ▶ Każdy z wierzchołków  $i \in N$  zwany dalej *neuronem* może przyjmować stany  $i \in \{0, 1\}$ .
- ▶ Cała sieć możemy opisać przy pomocy wektora  $x = (x_i)_{i \in N} \in \{0, 1\}^N$ .
- ▶ Każdy neuron ma funkcję  $\theta_i \in \mathbb{R}$ .
- ▶ Każda para neuronów jest połączona krawędzią o wadze  $\theta_{\{i,j\}} \in \mathbb{R}$ .

# Maszyny Boltzmana

- ▶ Dla danych wartości  $\theta_i \in \mathbb{R}$  oraz  $\theta_{\{i,j\}} \in \mathbb{R}$  energia układu, którą możemy utożsamiać z funkcją oceny, przyjmuje postać:

$$\mathbf{E}(x, \theta) = - \sum_{i \in N} \theta_i x_i - \sum_{\{i,j\} \in I} \theta_{\{i,j\}} x_i x_j, \quad x \in \{0, 1\}^N. \quad (1)$$

- ▶ Maszyna Boltzmannna uaktualnia swój stan w dyskretnym czasie.
- ▶ Ten proces nazywany jest próbkowaniem Gibbsa (ang. Gibbs sampling).
- ▶ Przejście ze stanu  $x^{(t)} \in \{0, 1\}^N$  w czasie  $t$  do stanu w czasie  $t + 1$ ,  $x^{(t+1)}$  odbywa się poprzez wybór jednego z neuronów  $i \in N$  i jego przejściem do stanu  $x_i^{(t+1)} = 1$  z prawdopodobieństwem:

$$Pr(x_i^{(t+1)} = 1 | x^{(t)}) = \sigma\left(\sum_{\{i,j\} \in I} \theta_{\{i,j\}} x_j^{(t)} + \theta_i\right), \quad (2)$$

lub  $x_i^{(t+1)} = 0$  z prawdopodobieństwem  $1 - Pr(x_i^{(t+1)} = 1)$ ,  
 a  $\sigma$  jest funkcją logistyczną.

- ▶ Prawdopodobieństwo znalezienia się układu, Maszyny Boltzmana, w określonym stanie opisuje rozkład Gibbsa-Boltzmana

$$p(x; \theta) = \frac{\exp(-E(x; \theta))}{Z(\theta)}, \quad x \in \{0, 1\}^N, \quad (3)$$

gdzie  $Z(\theta) = \sum_x \exp(-E(x; \theta))$  nazywamy funkcją rozkładu (sumą statystyczną).

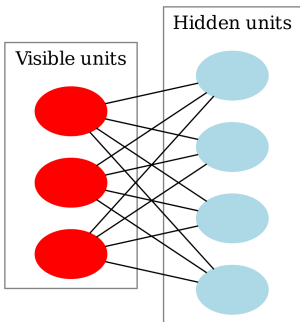


# Ograniczone maszyny Boltzmannna

- ▶ Będą nas interesowały tzw. **ograniczone maszyny Boltzmannna** (ang. Restricted Boltzmann Machines).
- ▶ Zaczniemy od tego, że jest możliwa obserwacja tylko  $V \subseteq N$  spośród wszystkich neuronów.
- ▶ Pozostałe neurony, czyli  $H = N \setminus V$  nazywamy ukrytymi.
- ▶ Z rozkładu prawdopodobieństwa  $p(x; \theta)$  dla wszystkich stanów układu  $x = (x_V, x_H) \in 0, 1^H \times 0, 1^V$  prawdopodobieństwo brzegowe widocznych stanów przyjmuje postać:

$$p(x_V; \theta) = \sum_{x_H \in 0, 1^H} p(x; \theta), \quad x_V \in \{0, 1\}^V. \quad (4)$$

- ▶ Ograniczonymi maszynami Boltzmann nazywamy maszyny Boltzmann, gdzie oddziaływanie pomiędzy neuronami zachodzi tylko pomiędzy neuronami z warstwy widocznej a neuronami z warstwy niewidocznej. To znaczy:  $I = \{\{i,j\} : i \in V, j \in H\}$ .



Rysunek: [https://commons.wikimedia.org/wiki/File:Restricted\\_Boltzmann\\_machine.svg](https://commons.wikimedia.org/wiki/File:Restricted_Boltzmann_machine.svg)

- ▶ Ten model sieci neuronowej jest przykładem modelu bazującego na pojęciu funkcji energetycznej. Oddziaływanie pomiędzy widocznymi  $x_V^i$  i niewidocznymi  $x_H^j$  neuronami opisuje funkcja energii:

$$\mathbf{E}(x_V^i, x_H^j) = - \sum_j b_j x_H^j - \sum_{ij} x_V^i w_{ij} x_H^j - \sum_i c_i x_V^i, \quad (5)$$

gdzie  $(b_j, W_{ij}, c_i)$  to wagi.

- ▶ Lub w zapisie macierzowym:

$$\mathbf{E}(x_V, x_H) = y^T W x + c^T y + b^T x, \quad (6)$$

gdzie  $x$  to stan widocznych neuronów,  $y$  niewidocznych, a  $Z$  to funkcja rozkładu.

- ▶ Zatem ostatecznie rozkład prawdopodobieństwa dla ograniczonych maszyn Boltzmanna przyjmuje postać:

$$p(x; \theta) = \frac{\sum_{y \in \{0,1\}^H} \exp(y^T Wx + c^T y + b^T x)}{Z(W, c, b)}, x \in \{0, 1\}^V, \quad (7)$$

gdzie, tak jak na poprzednim slajdzie,  $x$  to stan widocznych neuronów,  $y$  niewidocznych, a  $Z$  to funkcja rozkładu.

## Ograniczone maszyny Boltzmanna

- ▶ Do określania rozbieżności pomiędzy rozkładem prawdopodobieństwa dla danych  $P(\{v_i\})$  a rozkładem dla modelu teoretycznego używamy *dywergencji Kullbacka-Leiblera* w postaci:

$$D_{KL}(P(x_v^i) \parallel p_\lambda(x_v^i)) = \sum_{x_v^i} P(x_v^i) \log \left( \frac{P(x_v^i)}{p_\lambda(x_v^i)} \right). \quad (8)$$

- ▶ Zatem jeśli ograniczona maszyna Boltzmanna dokładnie odtwarza rozkład prawdopodobieństwa danych to:

$$D_{KL}(P(x_v^i) \parallel p_\lambda(x_v^i)) = 0. \quad (9)$$

# Uczenie ograniczonych maszyn Boltzmann

- ▶ Proces uczenia, w tym ograniczonych maszyn Boltzmann, możemy również opisać jako

$$\min_{a,W,b} D_{KL}(P(x_v^i) \parallel P_\lambda(x_v^i)). \quad (10)$$

- ▶ Dywergencja KL nie jest metryką. Nie jest bowiem symetryczna:

$$D_{KL}(P \parallel Q) \neq D_{KL}(Q \parallel P), \quad (11)$$

ani nie spełnia nierówności trójkąta:

$$D_{KL}(R \parallel P) \leq D_{KL}(P \parallel Q) + D_{KL}(R \parallel Q). \quad (12)$$

- ▶ Pytanie: czy istnieją lepsze alternatywy?

# Problem transportu optymalnego

- ▶ Problem **transportu optymalnego** (ang. optimal transport) po raz pierwszy sformułował Gaspard Monge w 1781 roku.
- ▶ Następnie latach dwudziestych (dwudziestego wieku) problem został sformalizowany przez A.N. Tolstoj.
- ▶ Najczęściej w literaturze problem występuje pod hasłem **zagadnienia transportowego**.



Rysunek: Gaspard Monge,  
1746-1818



# Odległość Wassersteina

## Definicja metryki Wassersteina

- ▶ Odległość Wassersteina, niech  $(X, d)$  będzie *przestrzenią metryczną polską* i niech  $p \in [1, \infty)$ . Dla każdych dwóch miar probabilistycznych,  $\mu, \nu$ , na  $X$  **odległością Wassersteina** stopnia  $p$  pomiędzy tymi miarami, definiujemy następująco:

$$\begin{aligned} W_p(\mu, \nu) &= \left( \inf_{\pi \in \Pi(\mu, \nu)} \int_X d(x, y)^p d\pi(x, y) \right)^{\frac{1}{p}} \\ &= \inf \{ [\mathbb{E}d(X, Y)^p]^{\frac{1}{p}} \} \end{aligned}$$

## Definicja $\gamma$ -gładkiej metryki Wassersteina

- ▶ Zakładamy dwa prawdopodobieństwa  $p$  i  $q$  ze zbioru  $X = \{0, 1\}^d$ , mamy funkcję  $D : X \times X \rightarrow \mathbb{R}_+$ , a  $\alpha$  i  $\beta$  to dwie funkcje na  $X$ . Wtedy niech  $\gamma$ -gładka odległość Wassersteina będzie zdefiniowana w następujący sposób:

$$W(p, q)_\gamma = \max_{\alpha, \beta} [\mathbb{E}_X[\alpha(X)] + \mathbb{E}_{X'}[\beta(X')]] - \gamma \sum_{x, x' \in \{0, 1\}^d} \exp\left(\frac{1}{\gamma}(\alpha(x) + \beta(x') - D(x, x')) - 1\right).$$

- ▶ Podobnie jak w przypadku dywergencji Kullbacka-Leiblera proces uczenia można sformułować jako:

$$\min_{\theta} W_\gamma(p_{data}, p_\theta). \quad (13)$$

- ▶ Mając prawdopodobieństwo dla ograniczonych maszyn Boltzmanna::

$$p(x; \theta) = \frac{\sum_{y \in \{0,1\}^H} \exp(y^T Wx + c^T y + b^T x)}{Z(W, c, b)}, x \in \{0, 1\}^V, \quad (14)$$

gdzie, tak jak na poprzednim slajdzie,  $x$  to stan widocznych neuronów,  $y$  niewidocznych, a  $Z$  to funkcja rozkładu.

- ▶ oraz stosując współczynnik regularyzacji:

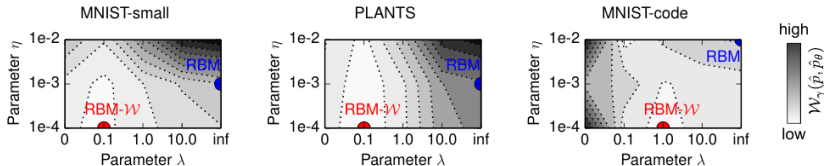
$$\Omega(c, W, b) = KL(p_{dat} || p_{c,W,b}) + \eta \left( \|a\|^2 + \sum_j \|W_j\|^2 \right).$$

- ▶ Zatem ostatecznie problem uczenia można sformułować jako:

$$\min_{c,W,b} W_\gamma(p_{c,W,b}, p_{data}) + \lambda \Omega(c, W, b), \quad (15)$$

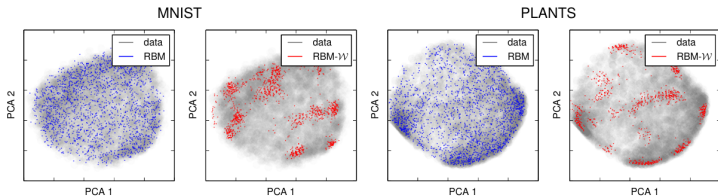
gdzie  $\eta$  i  $\gamma$  to współczynniki regularyzacji.

# Wyniki empiryczne



**Rysunek:** Montavon, Grégoire and Müller, Klaus-Robert and Cuturi, Marco. "Wasserstein Training of Restricted Boltzmann Machines", Advances in Neural Information Processing Systems 29.

# Wyniki empiryczne



**Rysunek:** Montavon, Grégoire and Müller, Klaus-Robert and Cuturi, Marco. "Wasserstein Training of Restricted Boltzmann Machines", Advances in Neural Information Processing Systems 29.

# Wyniki empiryczne






**Rysunek:** Montavon, Grégoire and Müller, Klaus-Robert and Cuturi, Marco. "Wasserstein Training of Restricted Boltzmann Machines", Advances in Neural Information Processing Systems 29.

# Podsumowanie



Dziękuję za uwagę

-  Montúfar, Guido. "Restricted Boltzmann Machines: Introduction and Review." Information Geometry and its Applications IV. Springer, Cham, 2016.
-  Montavon, Grégoire and Müller, Klaus-Robert and Cuturi, Marco. "Wasserstein Training of Restricted Boltzmann Machines", Advances in Neural Information Processing Systems 29.
-  Li Wuchen, and Guido Montúfar. "Natural gradient via optimal transport." Information Geometry 1.2 (2018): 181-214.