



**AGH UNIVERSITY OF SCIENCE
AND TECHNOLOGY**

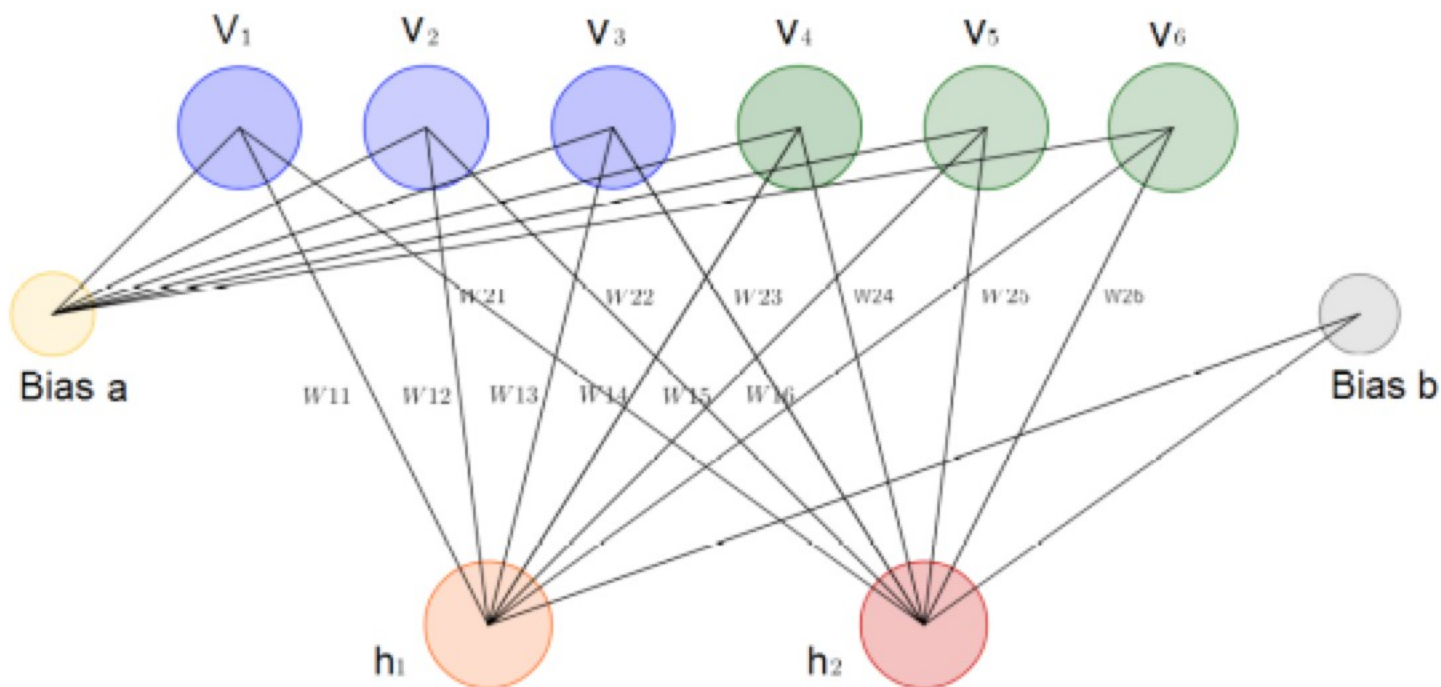
Autoencoders

Dariusz Kucharski
Katedra Automatyki i Robotyki

Kraków, 21.11.2019

Ograniczona maszyna Boltzmanna

•Ogólna architektura sieci

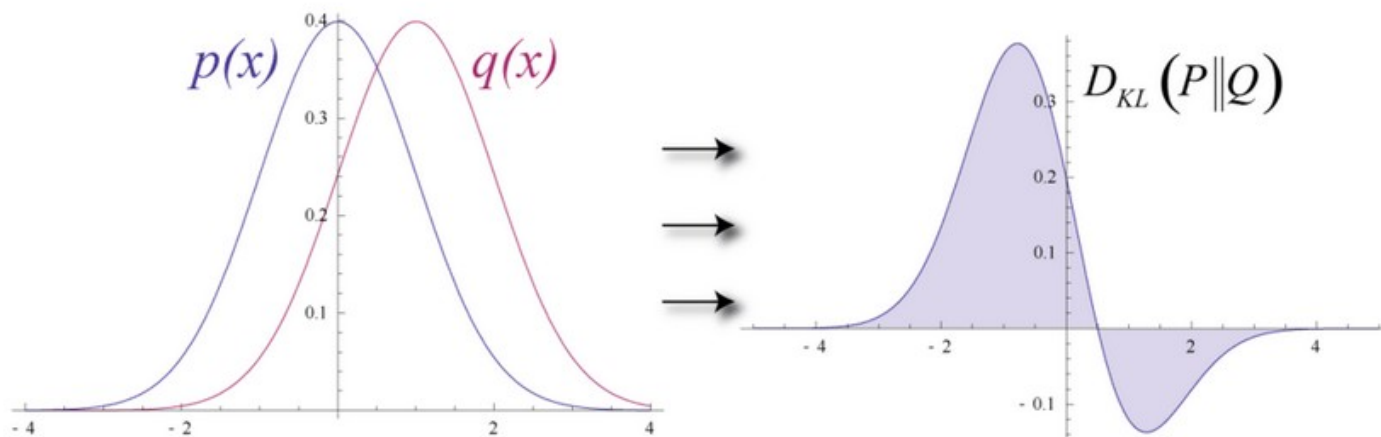


Ograniczona maszyna Boltzmannna

- Odpowiedzi warstwy ukrytej interpretujemy jako prawdopodobieństwo $p(a|x; w)$
- Przy obliczeniu rekonstrukcji możemy wynik rekonstrukcji interpretować jako $p(x|a; w)$
- Suma tych prawdopodobieństw to wspólny rozkład prawdopodobieństwa $p(a, x)$ (joint probability)

Ograniczona maszyna Boltzmanna

• Jeśli klasyfikacja to odgadywanie klasy (discriminative learning), regresja to szacowanie liczby, to rekonstrukcja to szukanie rozkładu prawdopodobieństwa wejścia (generative learning)

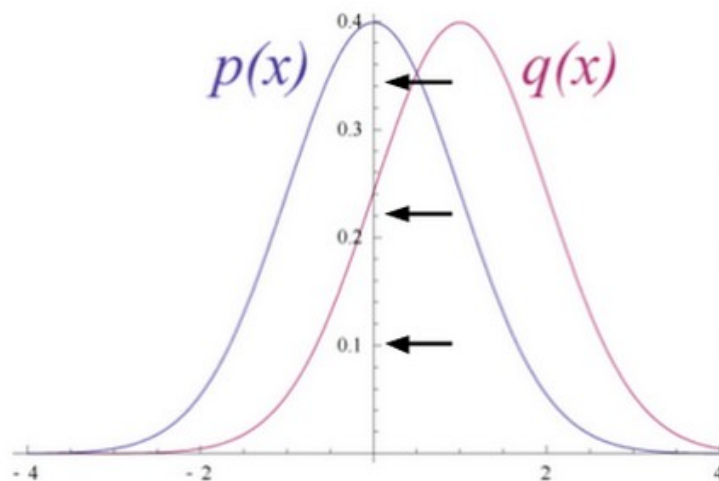


Różnica w rozkładzie prawdopodobieństwa sygnału

oryginalnego i rekonstrukcji [1]

Ograniczona maszyna Boltzmannna

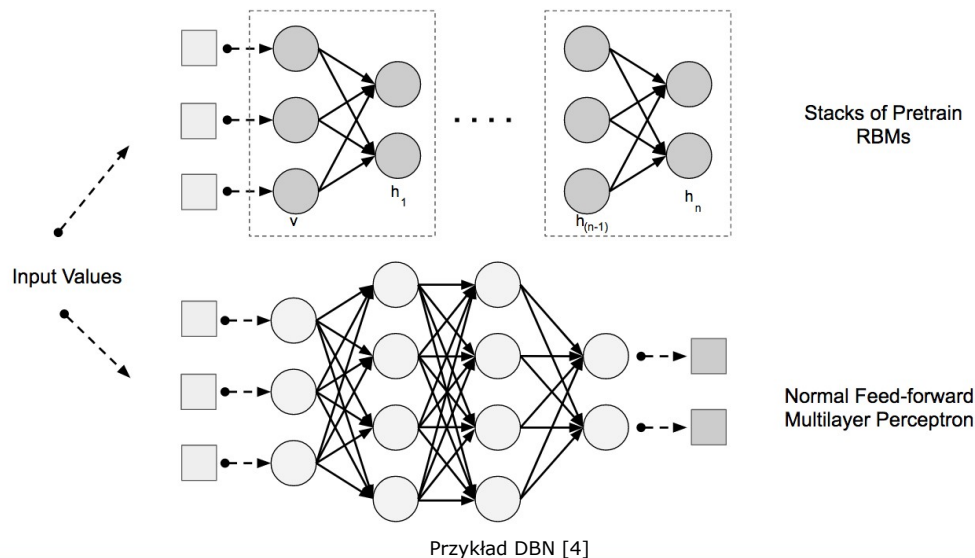
•Jeśli klasyfikacja to odgadywanie klasy (discriminative learning), regresja to szacowanie liczby, to rekonstrukcja to szukanie rozkładu prawdopodobieństwa wejścia (generative learning)



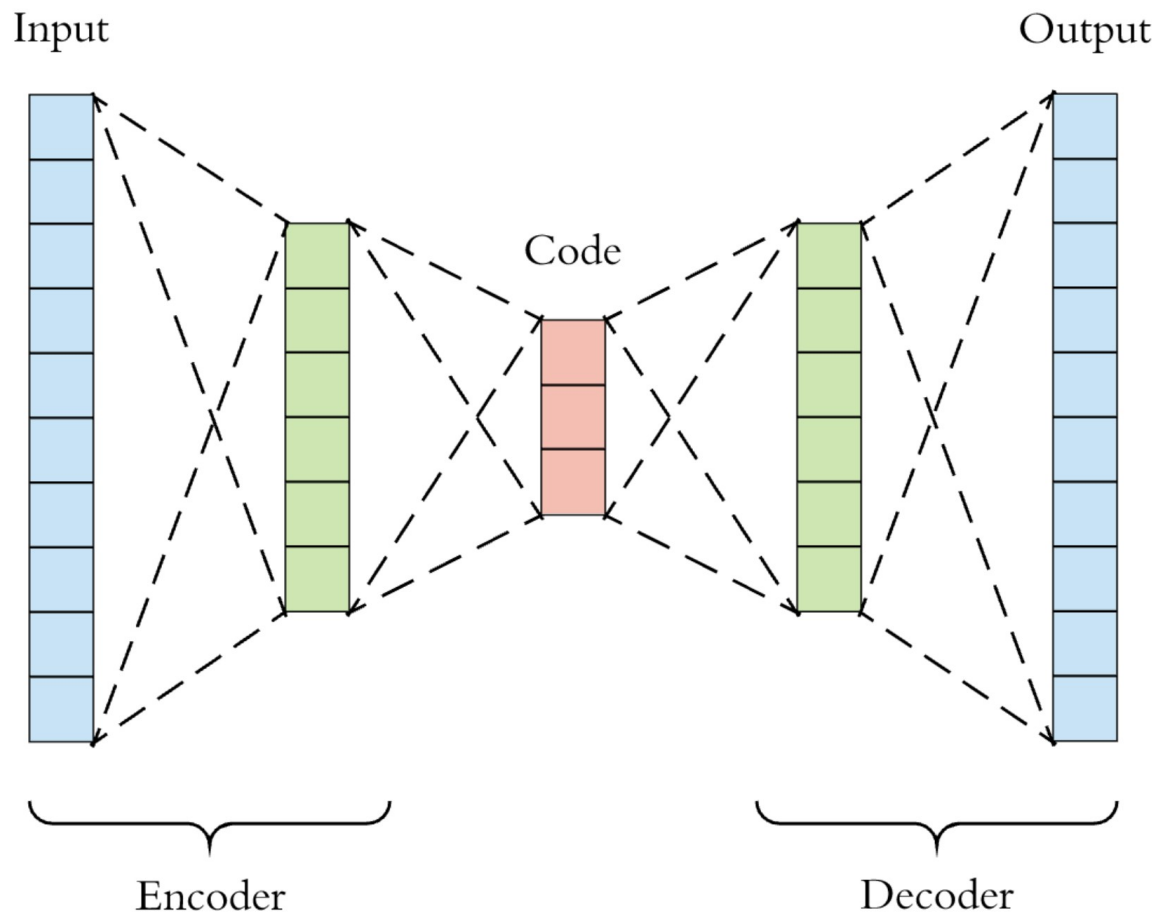
Dążymy to minimalizacji tej rekonstrukcji [1]

Deep Belief Network

- Powstaje poprzez wstępnej uczenie Ograniczonych Maszyn Boltzmann
- Wejście na kolejną maszynę boltzmana jest wyjściem z poprzedniej
- Po wyuczeniu odpowiedniej ilości warstw, ograniczone maszyny boltzmana łączone są w sieć głęboką (często dodaje się też klasyfikator) i następuję fine tuning – douczenie sieci, aby otrzymać jeszcze lepszą skuteczność



Autoenkodery



Przykład autoenkodera[1]

Sieć neuronowa, używana do efektywnego kodowania danych wejściowych (odnalezienia reprezentacji charakteryzującej dany zbiór, często o mniejszym wymiarze niż wejście) w sposób nienadzorowany.

Możemy wyróżnić enkoder (część sieci odpowiedzialna za kodowanie), który dokonuje transformacji wejścia do pewnej ukrytej reprezentacji z (stanowiącej pewnego rodzaju kod):

$$z = f_{\theta}(x) = \sigma(Wx + b)$$

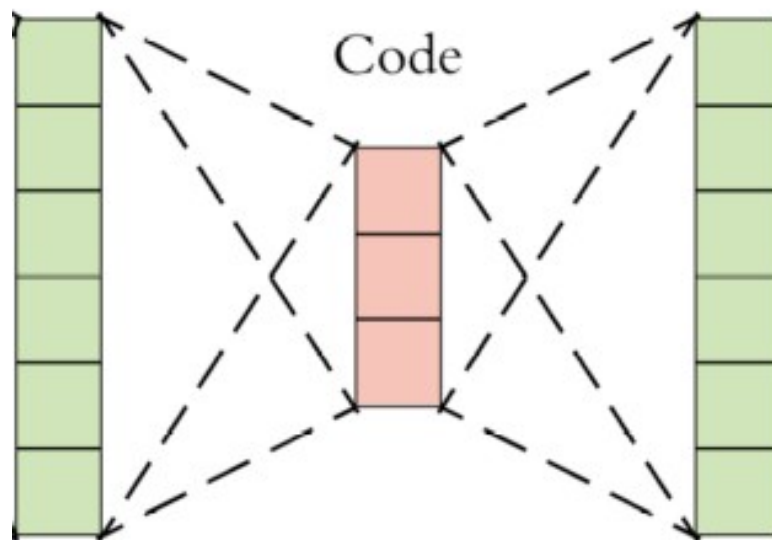
Oraz dekodery, który próbuje otworzyć te dane wejściowe na podstawie ich ukrytej reprezentacji z

$$x' = f_{\theta'}(h) = \sigma'(W'z + b')$$

Minimalizując funkcję błędu między wejściem i rekonstrukcją.

W, W', b, b' - parametry sieci; σ, σ' - funkcje aktywacji

Autoenkodery

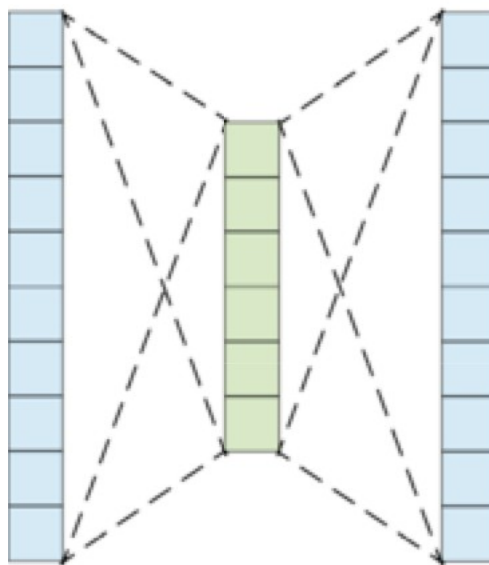


Przykład autoenkodera[1]

Dla $W' = W^T$ w początkowej fazie uczenia – połowa mniej parametrów do optymalizacji

Stacked autoencoders

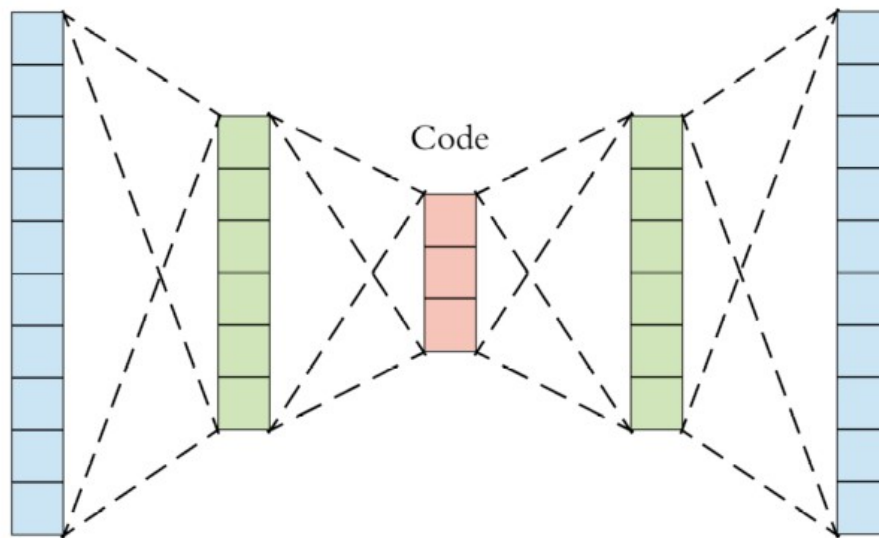
- **Uczenie autoenkoderów wielowarstwowych**
 - **Uczenie sieci z jedną warstwą ukrytą**



Przykład autoenkodera[1]

Stacked autoencoders

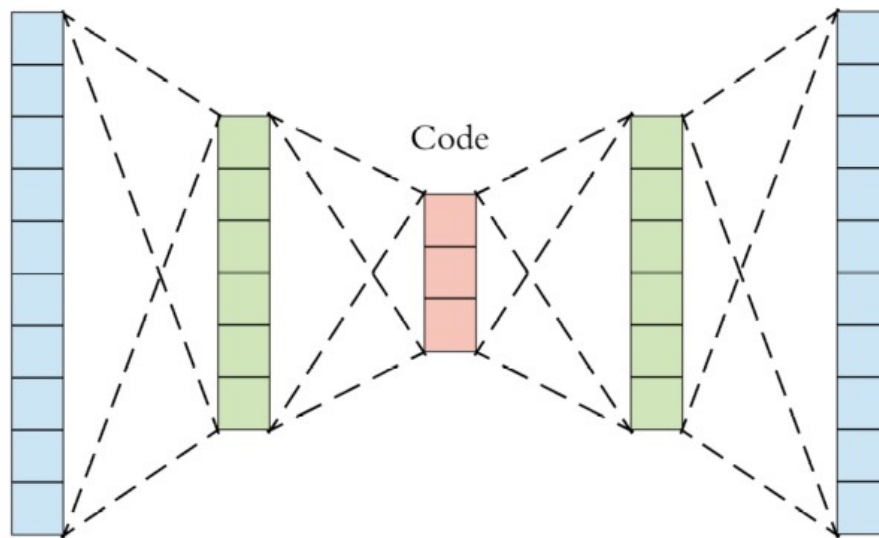
- **Uczenie autoenkoderów wielowarstwowych**
 - **Uczenie sieci z jedną warstwą ukrytą**
 - **Dodanie kolejnej warstwy ukrytej, zablokowanie wag warstwy już wyuczonej, uczenie nowej warstwy ukrytej**



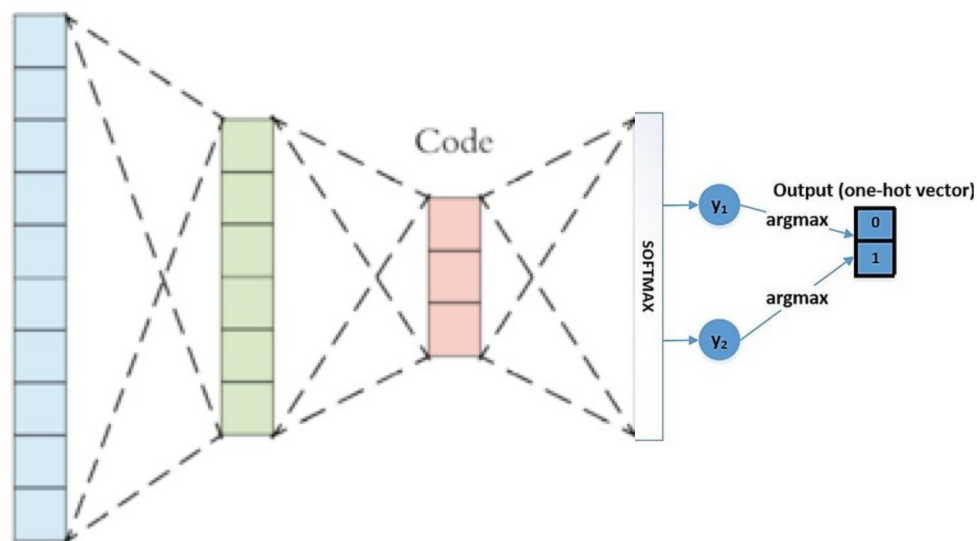
Przykład autoenkodera[1]

Stacked autoencoders

- **Uczenie autoenkoderów wielowarstwowych**
 - **Uczenie sieci z jedną warstwą ukrytą**
 - **Dodanie kolejnej warstwy ukrytej, zablokowanie wag warstwy już wyuczonej, uczenie nowej warstwy ukrytej**
 - **Odblokowanie wag i douczenie całej sieci**



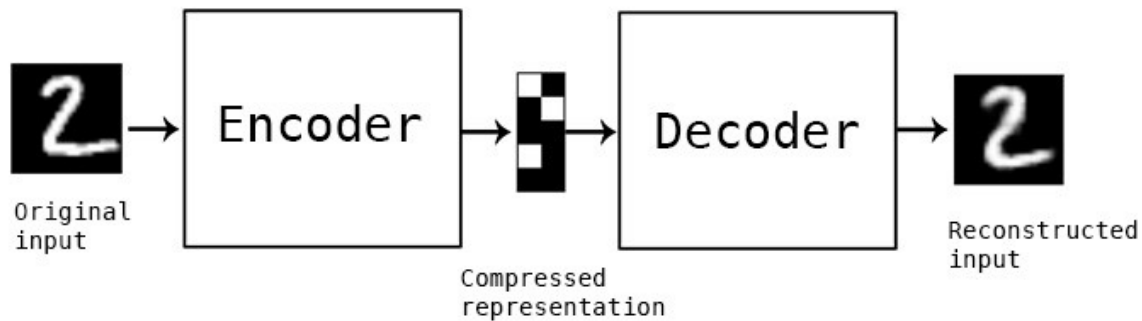
- **Zamiana dekodera na klasyfikator, douczenie go z zablokowanymi wagami enkodera a następnie wykonanie fine tuningu to tzw. uczenie pół nadzorowane**



Uczenie półnadzorowane [1], [https://www.researchgate.net/profile/Babak_Bashari_Rad/publication/326914586/figure/fig1/AS:657671334150145@1533812473734/Proposed-neural-network-classifier-with-softmax-output-function-and-a-bias-unit.png]

- **Należy pamiętać o dobraniu odpowiedniej funkcji kosztu do zadanego problemu**
 - Średni kwadrat
 - Entropia krzyżowa
 - Inne?

- **Wejście: 28x28x1 → 784**
- **Warstwy ukryte: 32**



Schemat autoenkodera dla zbioru MNIST[2]

- **Wejście: 28x28x1 → 784**
- **Warstwy ukryte: 32**



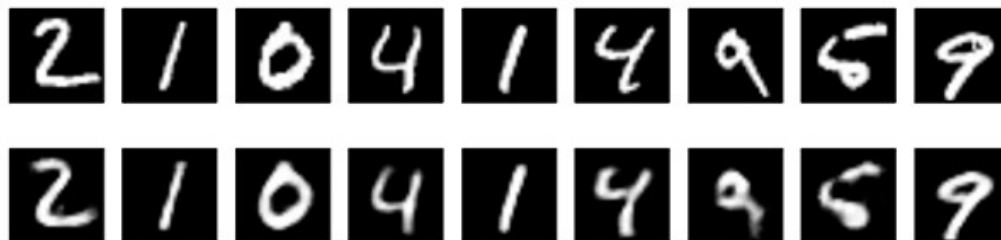
Przykładowy wynik autoenkodera z jedną warstwą ukrytą [2]



AGH

Deep Autoencoder

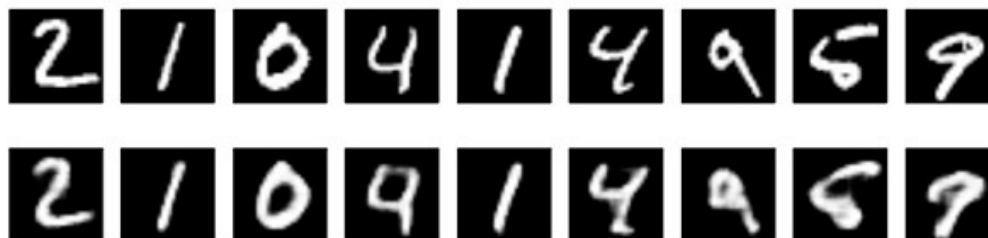
- **Wejście: $28 \times 28 \times 1 \rightarrow 784$**
- **Warstwy ukryte: $128 \rightarrow 64 \rightarrow 32$**



Przykładowy wynik autoenkodera typu deep [2]

Convolutional Autoencoder

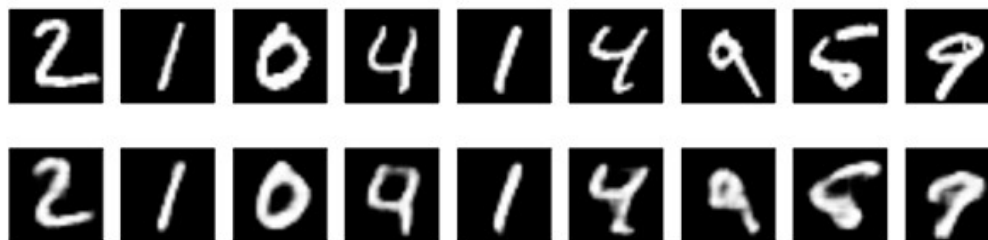
- **Wejście: 28x28x1 → 784**
- **Warstwy ukryte: conv 16 (3, 3) → mp (2, 2) → conv 8 (3, 3) → mp (2, 2) → conv 8 (3, 3) → mp (2, 2)**



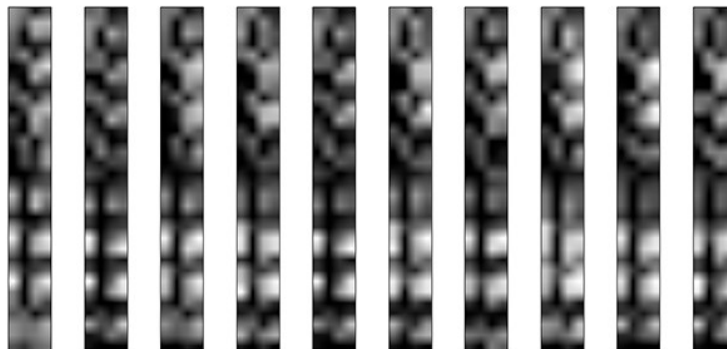
Przykładowy wynik autoenkodera konwolucyjnego [2]

Convolutional Autoencoder

- **Wejście: 28x28x1 → 784**
- **Warstwy ukryte: conv 16 (3, 3) → mp (2, 2) → conv 8 (3, 3) → mp (2, 2) → conv 8 (3, 3) → mp (2, 2)**



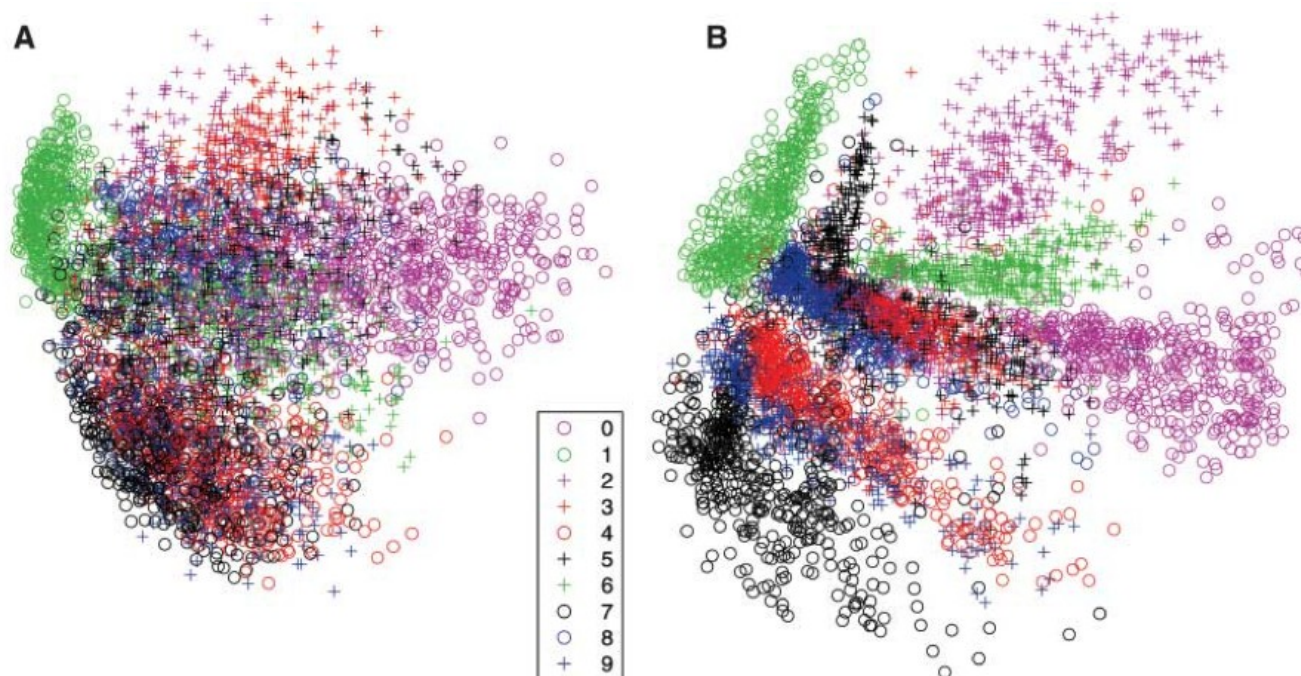
Przykładowy wynik autoenkodera konwolucyjnego [2]



Odpowiedź warstwy latent (8x4x4 → 4x32) [2]

Przykładowy autoenkoder dla zbioru mnist

- **Wejście: $28 \times 28 \times 1 \rightarrow 784$**
- **Warstwy ukryte: $100 \rightarrow 500 \rightarrow 250 \rightarrow 2$**



Porównanie dwóch pierwszych składowych PCA (A) oraz zastosowaniu autoenkodera z dwoma neuronami w warstwie latent (B) [3]



AGH

Denoising autoencoders

- » **Autoenkodery przygotowane w celu uzupełnienia pewnych informacji bądź usunięcia tych niechcianych**
- » **Zakłada się, że poprzez odpowiednie uczenie, sieć neuronowa nauczy się jedynie tych cech, które potrzebne są do rekonstrukcji bez zakłóceń**
- » **Czasem stosuje się tę właściwość gdy przy danym zbiorze i pewnej architekturze klasyczny autoenkoder „przepisuje” dane na wyjście (brak wyuczenia ukrytej reprezentacji)**

- » **Uczenie następuje poprzez zakłócenie oryginalnych danych**
 - **Poprzez zastosowanie dropoutu (30%-50% jednostek wyłączonych)**
 - **„ręcznego” zerowania niektórych wartości**
 - **Zastosowanie szumu (np. salt and pepper)**
- » **Błąd rekonstrukcji mierzy się między wyjściem a oryginalnym wejściem (przed dodaniem szumu)**

Denoising autoencoders



Przykład odzsumienia wejścia przez autoenkoder, funkcja kosztu liczona między obrazem przed zaszumieniem a rekonstrukcją [2]



AGH

Denoising autoencoders

A rectangular image showing a paragraph of text that is heavily distorted by a dense pattern of white, diagonal, wavy lines, making it difficult to read. The text is in a serif font and appears to be from a technical document.

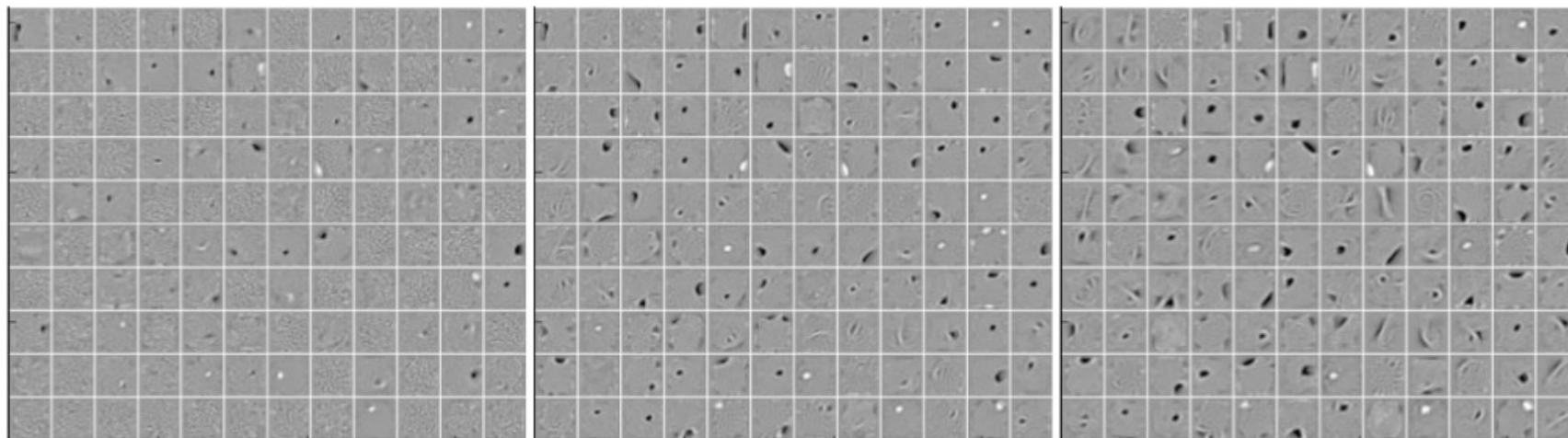
There exist several methods to design for
be filled in. For instance, fields may be surr
ing boxes, by light rectangles or by guiding ru
ods specify where to write and, therefore, n
of skew and overlapping with other parts o
guides can be located on a separate sheet
located below the form or they can be print
form. The use of guides on a separate she
from the point of view of the quality of th
but requires giving more instructions and,
restricts its use to tasks where this type of a
Guiding elements intended for the form designer

A rectangular image showing the same paragraph of text as the left image, but with the noise removed. The text is now clearly legible and matches the original content of the left image.

There are several classic spatial filters for
inating high frequency noise from images.
the median filter and the closing opening fi
used. The mean filter is a lowpass or sm
replaces the pixel values with the neighbor
duces the image noise but blurs the image e
filter calculates the median of the pixel neig
pixel, thereby reducing the blurring effect. F
closing filter is a mathematical morphologic
bines the same number of erosion and dilat
operations in order to eliminate small objec
The main advantage of this approach is

Inny przykład odszumiania [1]

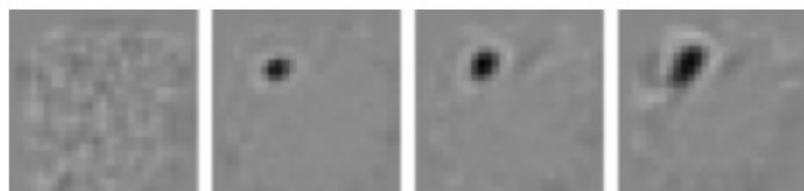
Denoising autoencoders



(a) No destroyed inputs

(b) 25% destruction

(c) 50% destruction



(d) Neuron A (0%, 10%, 20%, 50% destruction)



(e) Neuron B (0%, 10%, 20%, 50% destruction)

Zachowanie neuronów dla zaszumionego wejścia [4]

Sparsity autoencoders

- » **Warstwa latent może być większa od wejścia**
- » **Sieć aktywuje tylko nieliczne neurony dla danych przykładów**
- » **Definiowana przez współczynniki sparsity (ρ)**
 - **Dodanie regularyzatora za odstępstwo od założonego sparsity (ρ) (w uproszczeniu, sieć jest karana za aktywację zbyt dużej lub zbyt małej ilości neuronów)**

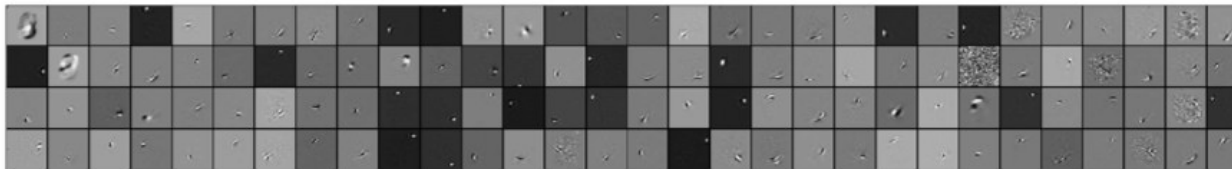


AGH

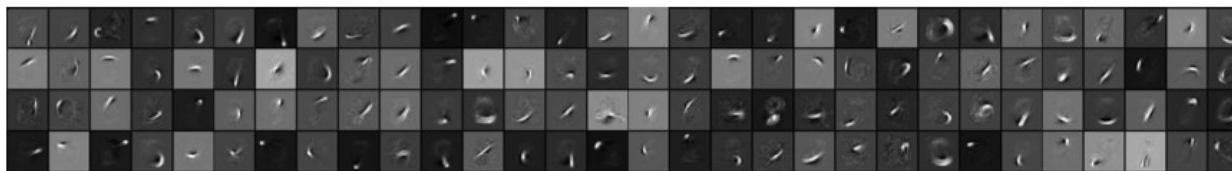
k-Sparse Autoencoders

- » **Obliczenie odpowiedzi warstwy latent**
- » **Wskazanie k największych aktywacji neuronów**
 - **Reszta aktywacji jest zerowana**
- » **Rekonstrukcja z wykorzystaniem k wskazanych neuronów z warstwy latent**
- » **Aktualizacja wag**

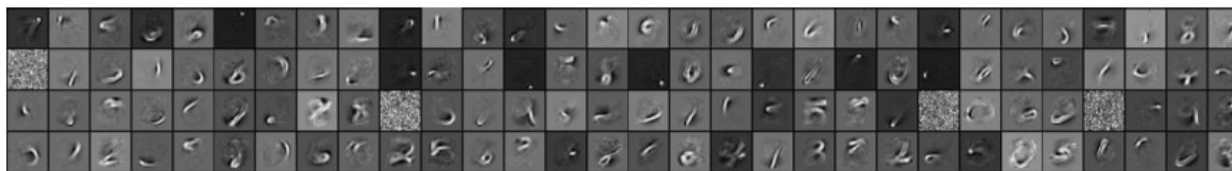
k-Sparse Autoencoders



(a) $k = 70$



(b) $k = 40$



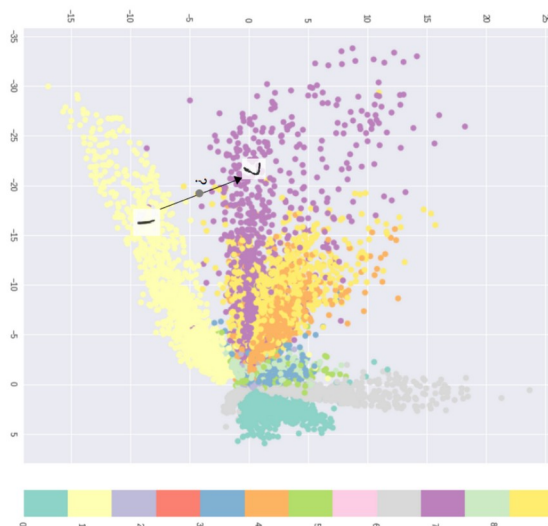
(c) $k = 25$



(d) $k = 10$

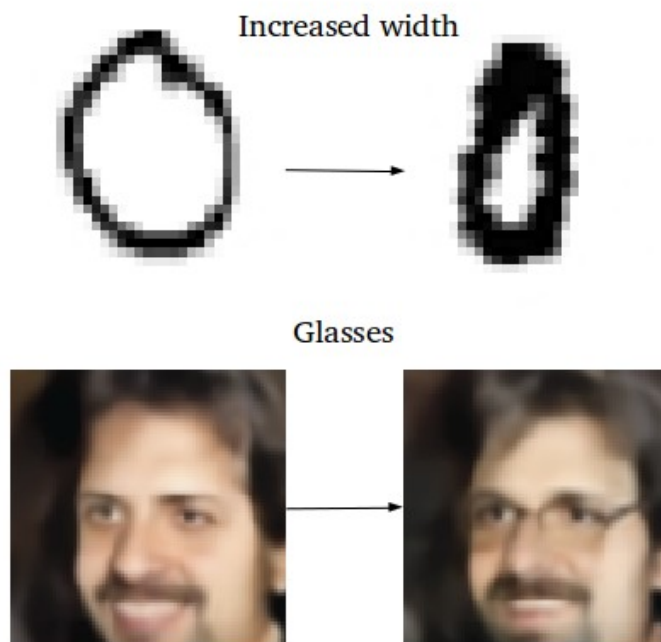
Wyuczona informacja przez sieć w zależności od współczynnika k [5]

- » Problem z klasycznymi autoenkoderami polega na tym, że nie ma reguły co tak naprawdę przedstawia warstwa latent
- » Odległości i zależności przestrzenne między aktywacjami dla danych przykładów często nie mają ze sobą związku – przestrzeń aktywacji warstwy latent nie jest ciągła



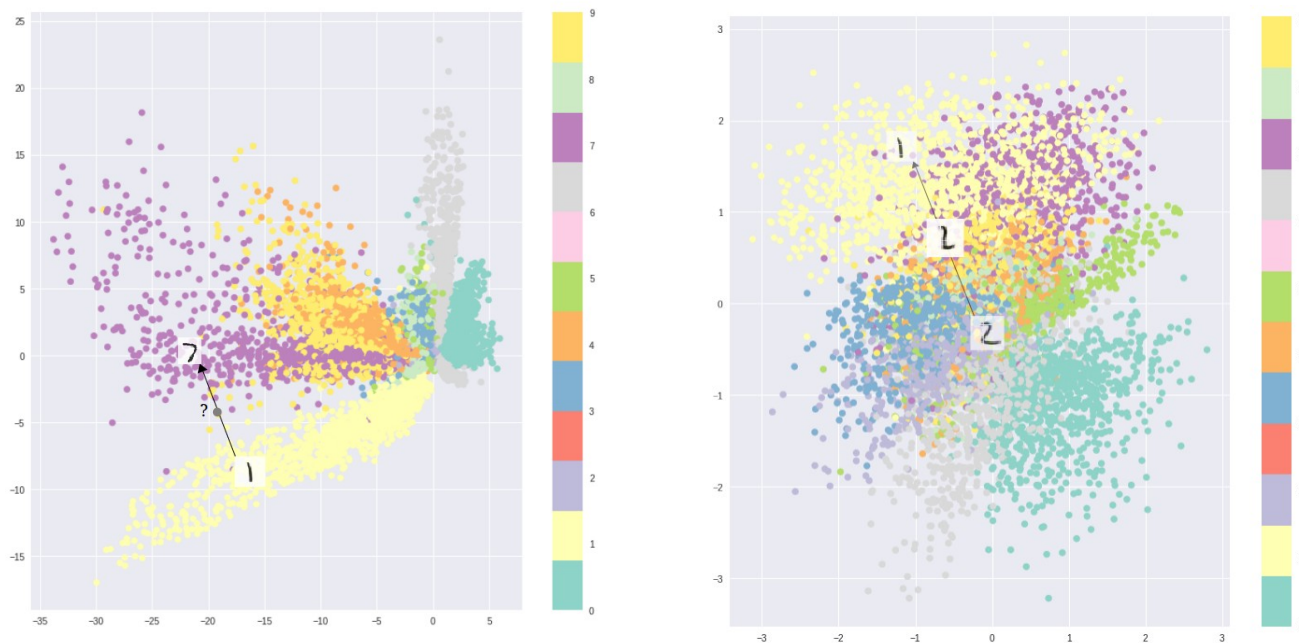
Brak informacji o potencjalnej aktywacji między 7 a 1 [6]

- » **Możliwość generowania (zmieniania) nowych danych podobnych do tych które zostały użyte do uczenia**



Przykład wprowadzonej zmiany w przykładzie ze zbioru [6]

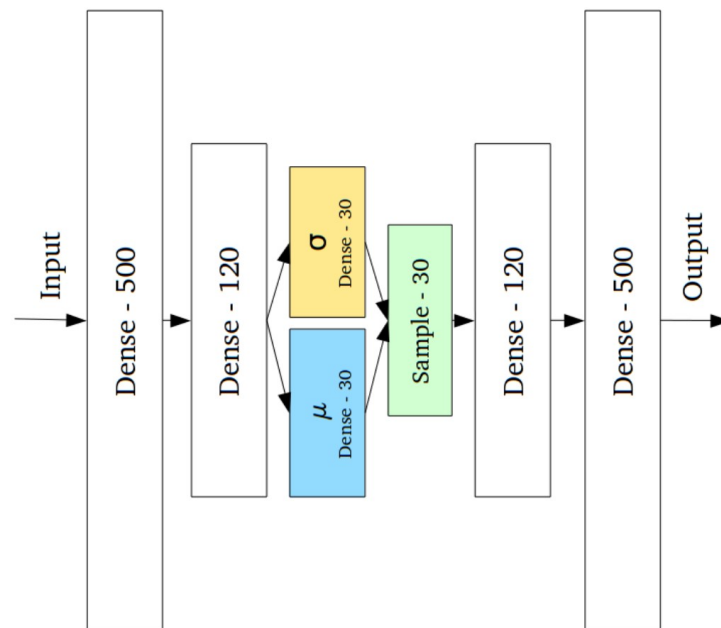
- » Założeniem autoenkoderów wariacyjnych jest ciągła przestrzeń aktywacji warstwy latent
- » Spełnienie założenia umożliwia „zmuszenie” sieci do zakodowania przykładów w postaci rozkładu Gaussa



Porównanie przestrzeni latent autoenkodera klasycznego i autoenkodera wariacyjnego[6]

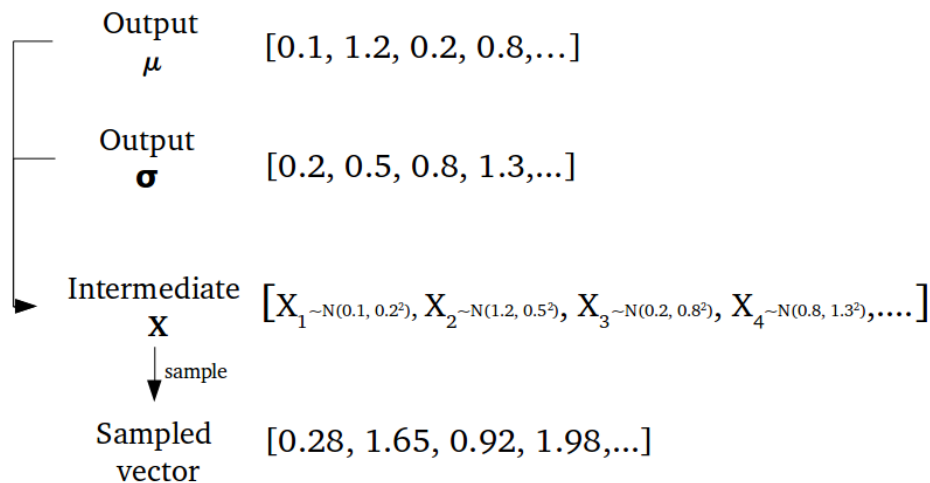
Variational Autoencoders

- » Założeniem autoenkoderów wariacyjnych jest ciągła przestrzeń aktywacji warstwy latent
- » Spełnienie założenia umożliwia „zmuszenie” sieci do zakodowania przykładów w postaci rozkładu Gaussa



Warstwa latent jest reprezentowana przez wektor średnich i odchyień[6]

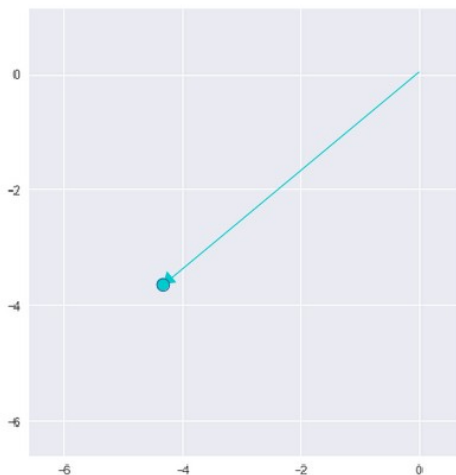
- » Aktywacja jest obliczana poprzez losowanie próbki (przejście stochastyczne)
- » Przez co wartość aktywacji może być różna dla tego samego przykładu



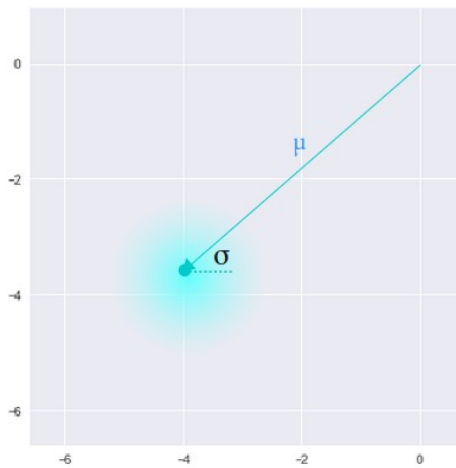
Przykład przejścia przez warstwę latent [6]

Variational Autoencoders

- » Średnia wskazuje gdzie zakodowany przykład ma być umiejscowiony w warstwie latent
- » Odchylenie natomiast wskazuje obszar w jakim powinien znajdować się przykład z danego typu



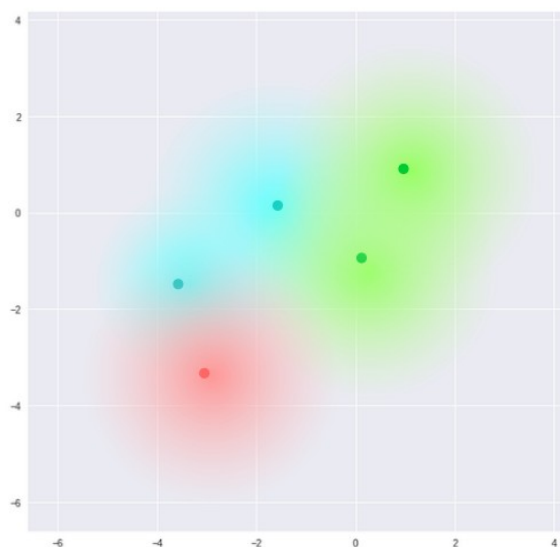
Standard Autoencoder
(direct encoding coordinates)



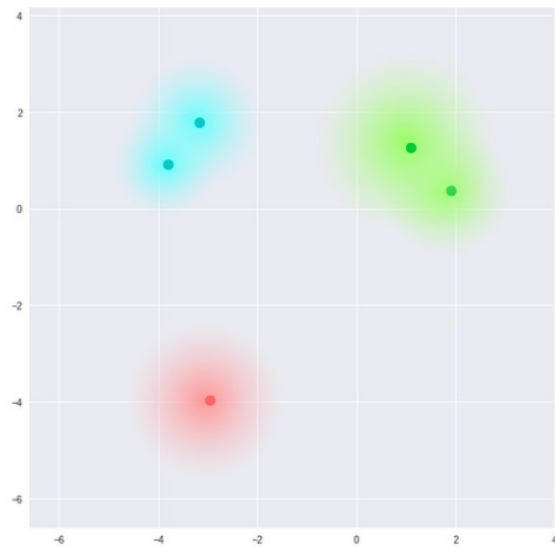
Variational Autoencoder
(μ and σ initialize a probability distribution)

Dany przykład kodowany jest na otoczeniu [6]

- » Aby uniknąć sytuacji w której przejścia między kodowaniami nie są ciągłe, wprowadza się dywergencję KL do funkcji kosztu



What we require



What we may inadvertently end up with

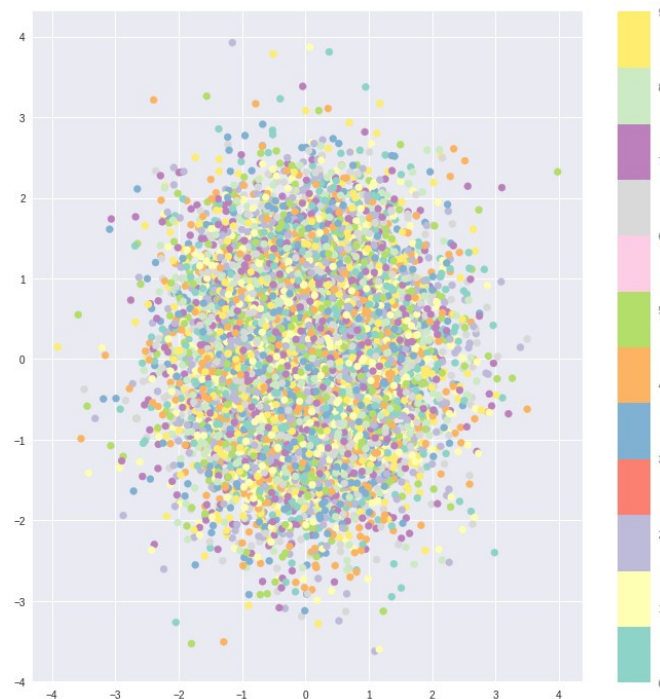
Po lewej stronie rezultat oczekiwany [6]

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

Dywersgencja KL [6]

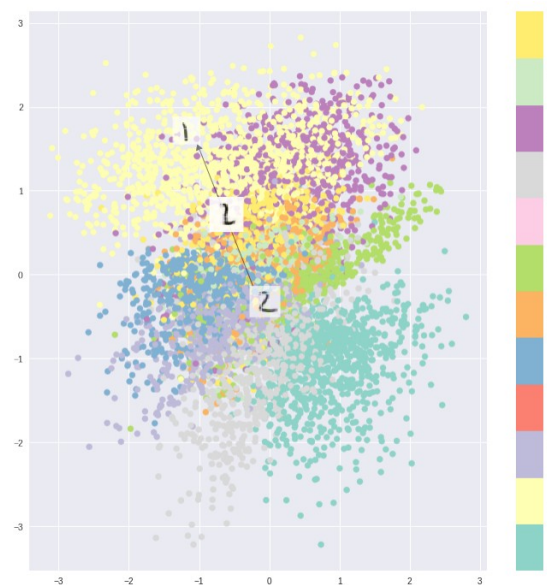
- » Można interpretować jako sumę wszystkich dywersgencji KL między rozkładem przykładu wejściowego (X_i) a wartością rozkładu normalnego dla średniej i odchylenia
- » Intuicyjnie, taka funkcja kosztu „zmusza” vae do rozprzestrzenienia wartości aktywacji równo wokół środka przestrzeni latent

- » **Zastosowanie tylko sumy dywergencji KL jako funkcji kosztu spowodowałoby jednak losowe skupienie wszystkich kodowań wokół środka**



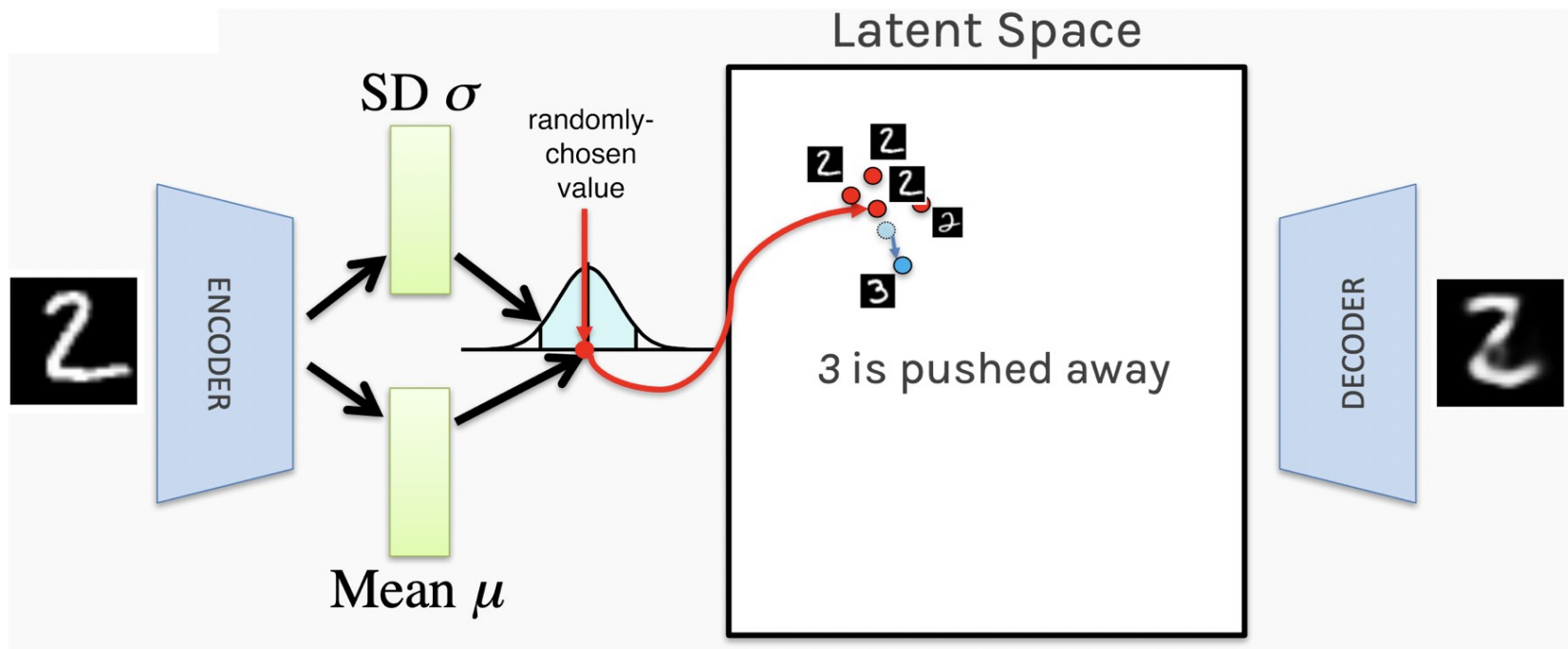
Przykład latent na zbiorze MNIST w przypadku zastosowania samej dywergencji KL jako funkcji kosztu [6]

- » **Dodanie błędu rekonstrukcji do funkcji kosztu powoduje klasteryzowanie kodowań**
- » **Dywergencja KL odpowiada za skupienie wartości wokół centrum, błąd rekonstrukcji odpowiada za klasteryzację**



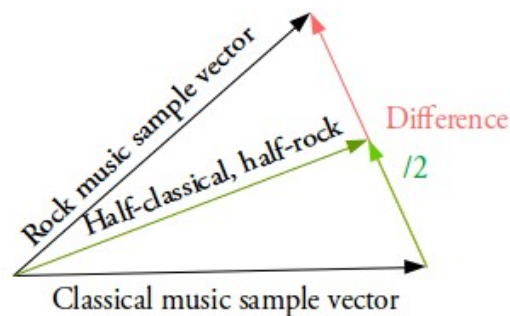
Aktywacje warstwy latent w przypadku zastosowania funkcji kosztu jako sumy dywergencji KL oraz błędu rekonstrukcji[6]

Variational Autoencoders



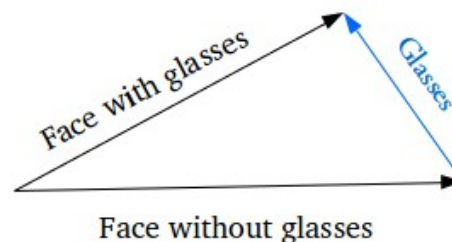
W procesie uczenia, elementy o podobnych cechach są skupiane blisko siebie[1]

- » **Generowanie przykładów poprzez wyznaczenie średniej z dwóch wektorów**



Jeśli dany przykład ma nosić cechy dwóch różnych przykładów, należy obliczyć średnią z wartości ich średnich [6]

- » **Odejmując dwa wektory jesteśmy w stanie wygenerować ich różnicę**



Różnica dwóch kodowań pozwoli natomiast wygenerować różniące ich cechy [6]

- [1] <https://towardsdatascience.com/generating-images-with-autoencoders-77fd3a8dd368>
- [2] <https://blog.keras.io/building-autoencoders-in-keras.html>
- [3] Hinton G. E., Salakhutdinov R. R, Reducing the Dimensionality of Data with Neural Networks, 2006
- [4] Vincent et al., Extracting and Composing Robust Features with Denoising Autoencoders, 2008
- [5] Makhzani A., Frey B., k-Sparse Autoencoders, 2013
- [6] <https://towardsdatascience.com/intuitively-understanding-variational-autoencoders-1bfe67eb5daf>