



Автоэнкодеры

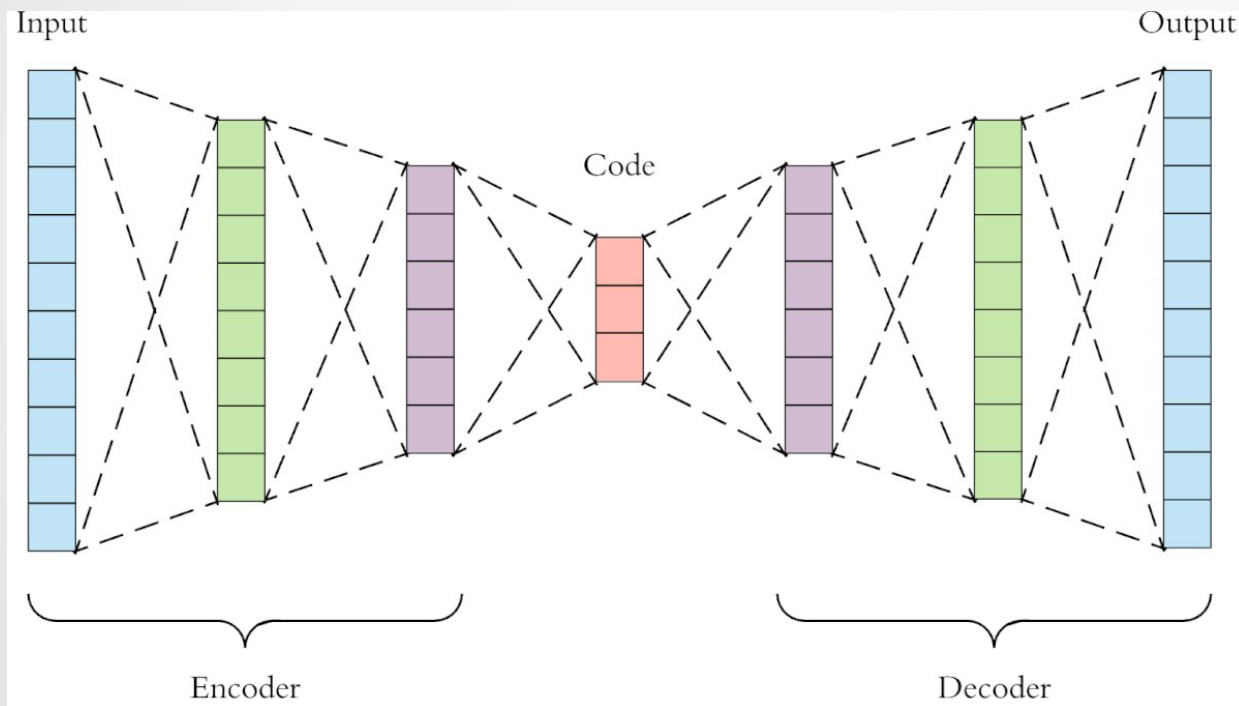
Евгений Борисов

Нейросети

Автоэнкодер

нейронная сеть прямого распространения

вторая половина сети зеркально повторяет первую



при обучении требуем от сети
восстановить исходный образ

цель — получить внутренне
представление входного образа,
(скрытый слой)

<https://neurohive.io/ru/osnovy-data-science/avtojenkoder-tipy-arhitektur-i-primeneniye/>

David Charte et al. A pratial tutorial on autoenoders for nonlinear feature fusion:taxonomy, models, software and guidelines. 2018.

Rumelhart, Hinton, Williams. Learning Internal Representations by Error Propagation.1986.

Нейросети

Задачи для автоэнкодеров

Rumelhart, Hinton, Williams. Learning Internal Representations by Error Propagation. 1986.

David Charpe et al. A practical tutorial on autoencoders for nonlinear feature fusion: taxonomy, models, software and guidelines. 2018.

- Генерация признаков (feature generation)
- Снижение размерности (dimensionality reduction)
- Сжатие данных с минимальными потерями точности
- Более эффективное решение задач обучения с учителем в новом признаковом пространстве
- Обучаемая векторизация объектов, встраиваемая в более глубокие нейросетевые архитектуры
- Послойное предобучение многослойных сетей
- Генерация синтетических объектов, похожих на реальные

Нейросети

Постановка задачи автоэнкодера

$X^\ell = \{x_1, \dots, x_\ell\}$ — обучающая выборка

$f: X \rightarrow Z$ — кодировщик (encoder), кодовый вектор $z = f(x, \alpha)$

$g: Z \rightarrow X$ — декодировщик (decoder), реконструкция $\hat{x} = g(z, \beta)$

Суперпозиция $\hat{x} = g(f(x))$ должна восстанавливать исходные x_i :

$$\mathcal{L}_{\text{AE}}(\alpha, \beta) = \sum_{i=1}^{\ell} \mathcal{L}(g(f(x_i, \alpha), \beta), x_i) \rightarrow \min_{\alpha, \beta}$$

Квадратичная функция потерь: $\mathcal{L}(\hat{x}, x) = \|\hat{x} - x\|^2$

Пример 1. Линейный автокодировщик: $x \in \mathbb{R}^n$, $z \in \mathbb{R}^m$

$$f(x, A) = A x, \quad g(z, B) = B z$$

$m \times n \qquad n \times m$

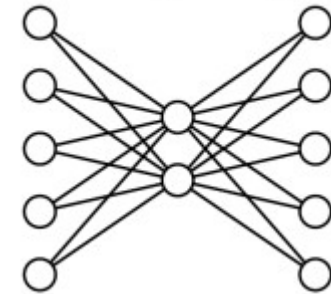
Пример 2. Двухслойная сеть с функциями активации σ_f, σ_g :

$$f(x, A) = \sigma_f(Ax + a), \quad g(z, B) = \sigma_g(Bz + b)$$

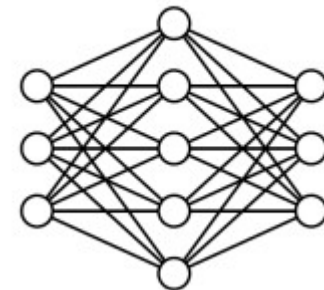
изменение размерности

извлечение признаков

снижение размерности



повышение размерности



Нейросети

Шумоподавляющий автоэнкодер (Denoising AE)

P. Vincent, H. Larochelle, Y. Bengio, P.-A. Manzagol. Extracting and composing robust features with denoising autoencoders. ICML-2008.

при обучении на входной образ накладываем шум

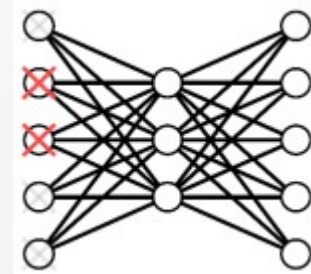
требуем от сети восстановить исходный образ

Устойчивость кодовых векторов z_i относительно шума в x_i :

$$\mathcal{L}_{\text{DAE}}(\alpha, \beta) = \sum_{i=1}^{\ell} \mathbb{E}_{\tilde{x} \sim q(\tilde{x}|x_i)} \mathcal{L}(g(f(\tilde{x}, \alpha), \beta), x_i) \rightarrow \min_{\alpha, \beta}$$

Вместо вычисления $\mathbb{E}_{\tilde{x}}$ в методе SGD объекты x_i сэмплируются и зашумляются по одному: $\tilde{x} \sim q(\tilde{x}|x_i)$. Варианты зашумления:

- $\tilde{x} \sim \mathcal{N}(x_i, \sigma^2 I)$ — гауссовский шум
- обнуление компонент вектора x_i с вероятностью p_0 :



Нейросети

Автоэнкодер для обучения на размеченных данных

Dor Bank, Noam Koenigstein, Raja Giryes. Autoenoders. 202

Данные: неразмеченные $(x_i)_{i=1}^{\ell}$, размеченные $(x_i, y_i)_{i=\ell+1}^{\ell+k}$

Совместное обучение кодировщика, декодировщика и предсказательной модели (классификации, регрессии или др.):

$$\sum_{i=1}^{\ell} \mathcal{L}(g(\mathbf{f}(x_i, \alpha), \beta), x_i) + \lambda \sum_{i=\ell+1}^{\ell+k} \tilde{\mathcal{L}}(\hat{y}(\mathbf{f}(x_i, \alpha), \gamma), y_i) \rightarrow \min_{\alpha, \beta, \gamma}$$

$z_i = \mathbf{f}(x_i, \alpha)$ — кодировщик

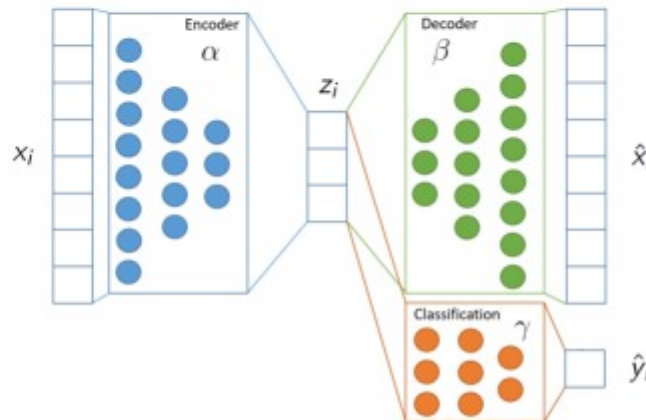
$\hat{x}_i = g(z_i, \beta)$ — декодировщик

$\hat{y}_i = \hat{y}(z_i, \gamma)$ — предиктор

Функции потерь:

$\mathcal{L}(\hat{x}_i, x_i)$ — реконструкция

$\tilde{\mathcal{L}}(\hat{y}_i, y_i)$ — предсказание



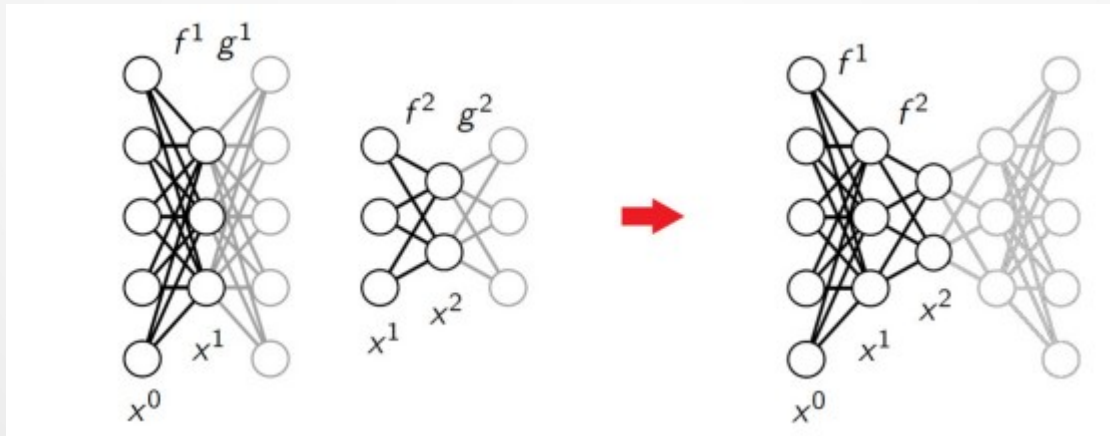
Нейросети

Многослойный автоэнкодер (Staked AE)

Y. Bengio et al. Greedy layer-wise training of deep networks. NIPS 2007.

Послойное обучение: $x^h = f^h(x^{h-1}, \alpha^h)$, $x \equiv x^0$, $z \equiv x^H$

- каждая пара f^h, g^h обучается по выборке $\{x_1^{h-1}, \dots, x_\ell^{h-1}\}$
- декодировщик g^h отбрасывается
- однослойные f^1, \dots, f^H соединяются в H -слойный



Тонкая настройка (fine tuning): результат послойного обучения используется как начальное приближение для BackProp

Нейросети

Пример

снижение размерности, извлечение признаков

размер входного (и выходного слоя) - $784 = 28 \times 28$

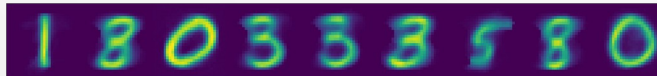
размер представления (скрытого) слоя — 2

датасет MNIST

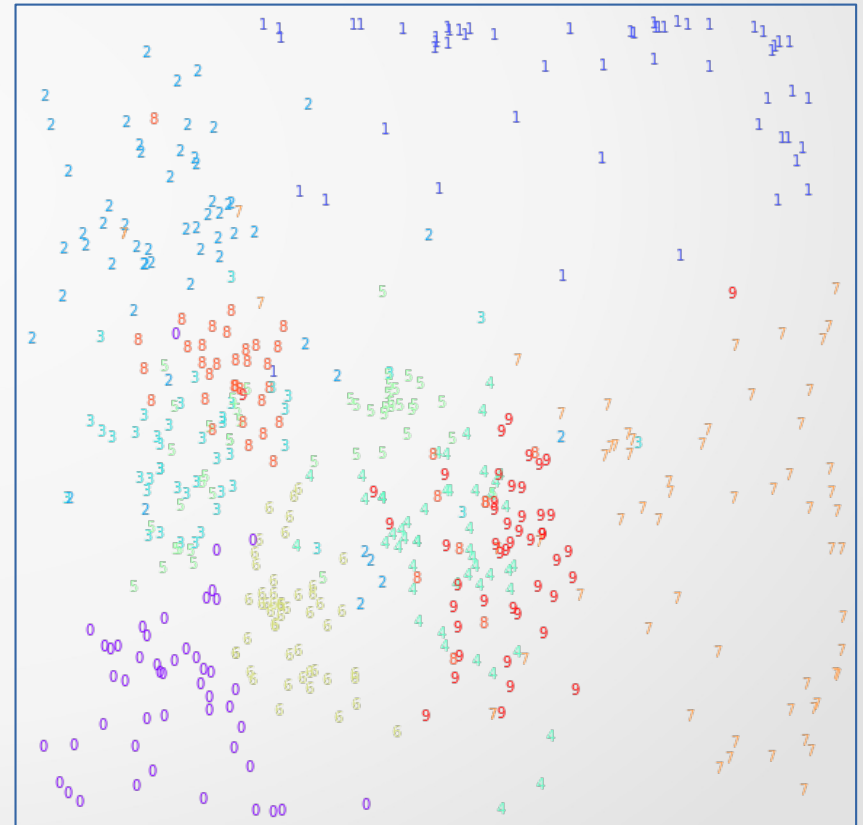
оригинал



восстановленный



карта расположения
объектов в 2D пространстве
признаков



Нейросети

Интерпретация

Изображения цифр mnist можно рассматривать как элементы $28 \times 28 = 784$ -мерного пространства.

Среди всех изображений 28×28 , изображения цифр занимают небольшую часть, остальное это шум.

Для одной выбранной цифры в 784-мерном пространстве можно найти кривую, все точки некоторой области вдоль этой кривой это цифры.

Т.е. в пространстве всех изображений есть подпространство меньшей размерности с цифрами, которое и находит автоэнкодер.

Нейросети: литература

git clone https://github.com/mechanoid5/ml_lectorium.git

Воронцов К. В.

Прикладные модели машинного обучения. 2021.

Лекция 2: Обучение без учителя.

<https://www.youtube.com/watch?v=wfbe2yaXAkl>

Rumelhart, Hinton, Williams. Learning Internal Representations by Error Propagation. 1986.

David Charlet et al. A partial tutorial on autoencoders for nonlinear feature fusion: taxonomy, models, software and guidelines. 2018.

Михаил Сурцуков

Manifold learning и скрытые (latent) переменные

<https://habr.com/ru/post/331500/>