

Методы машинного обучения. Обучение без учителя: векторизация данных

Воронцов Константин Вячеславович
www.MachineLearning.ru/wiki?title=User:Vokov
вопросы к лектору: voron@forecsys.ru

материалы курса:
github.com/MSU-ML-COURSE/ML-COURSE-22-23
орг.вопросы по курсу: ml.cmc@mail.ru

Выявление структуры данных на основе сходства:

- кластеризация (clustering) и квантизация (quantization)
- оценивание плотности распределения (density estimation)
- одноклассовая классификация (anomaly detection)

Преобразование признакового пространства:

- метод главных компонент (principal components analysis)
- автокодировщики (autoencoders)
- многомерное шкалирование (multidimensional scaling)
- матричные разложения (matrix factorization)

Поиск взаимосвязей в данных или синтез учителя:

- поиск ассоциативных правил (association rule learning)
- частичное обучение (semi-supervised learning)
- самостоятельное обучение (self-supervised learning)

1 Сети Кохонена для кластеризации и визуализации

- Задача кластеризации
- Модели конкурентного обучения
- Карты Кохонена

2 Автокодировщики

- Постановка задачи понижения размерности
- Методы регуляризации
- Применение автокодировщиков

3 Векторные представления графов и текстов

- Многомерное шкалирование
- Векторные представления графов
- Векторные представления текстов

Постановка задачи кластеризации и квантизации

Дано:

$X^\ell = \{x_i\}_{i=1}^\ell$ — обучающая выборка объектов, $x_i \in \mathbb{R}^n$
 $\rho^2(x, w) = \|x - w\|^2$ — евклидова метрика в \mathbb{R}^n

Найти:

центры кластеров $w_y \in \mathbb{R}^n$, $y \in Y$; алгоритм кластеризации
«правило жёсткой конкуренции» (WTA, Winner Takes All):

$$a(x) = \arg \min_{y \in Y} \rho(x, w_y)$$

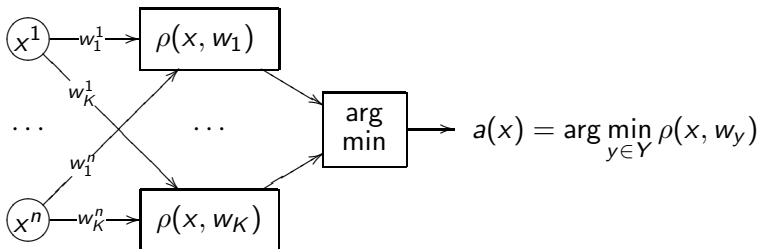
Критерий: среднее внутрикластерное расстояние

$$Q(w; X^\ell) = \sum_{i=1}^{\ell} \rho^2(x_i, w_{a(x_i)}) \rightarrow \min_{w_y: y \in Y}$$

Квантизация данных — замена x_i на ближайший центр $w_{a(x_i)}$

Сеть Кохонена (сеть с конкурентным обучением)

Структура алгоритма похожа на двухслойную нейросеть:



Градиентный шаг в методе SG: для выбранного $x_i \in X^\ell$

$$w_y := w_y + \eta(x_i - w_y)[a(x_i) = y]$$

Если x_i относится к кластеру y , то w_y сдвигается в сторону x_i

Алгоритм SG (Stochastic Gradient)

Вход: выборка X^ℓ ; темп обучения η ; параметр λ ;

Выход: центры кластеров $w_1, \dots, w_K \in \mathbb{R}^n$;

инициализировать центры $w_y, y \in Y$;

инициализировать текущую оценку функционала:

$$Q := \sum_{i=1}^{\ell} \rho^2(x_i, w_{a(x_i)});$$

повторять

выбрать объект x_i из X^ℓ (например, случайно);

определить ближайший центр: $y := \arg \min_{y \in Y} \rho(x_i, w_y)$;

градиентный шаг: $w_y := w_y + \eta(x_i - w_y)$;

оценить значение функционала:

$$Q := (1 - \lambda)Q + \lambda \rho^2(x_i, w_y);$$

пока значение Q и/или веса w не стабилизируются;

Жёсткая и мягкая конкуренция

Правило жёсткой конкуренции WTA (winner takes all):

$$w_y := w_y + \eta(x_i - w_y) [a(x_i) = y], \quad y \in Y$$

Недостатки правила WTA:

- медленная скорость сходимости
- некоторые w_y могут никогда не выбираться

Правило мягкой конкуренции WTM (winner takes most):

$$w_y := w_y + \eta(x_i - w_y) K(\rho(x_i, w_y)), \quad y \in Y$$

где ядро $K(\rho)$ — неотрицательная невозрастающая функция

Теперь центры всех кластеров смещаются в сторону x_i ,
но чем дальше от x_i , тем меньше величина смещения

Карта Кохонена (Self Organizing Map, SOM)

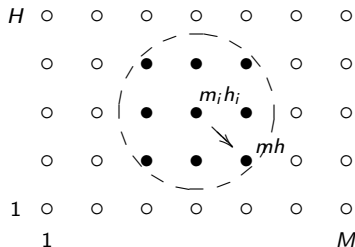
$Y = \{1, \dots, M\} \times \{1, \dots, H\}$ — прямоугольная сетка кластеров

Каждому узлу (m, h) приписан нейрон Кохонена $w_{mh} \in \mathbb{R}^n$

Наряду с метрикой $\rho(x_i, x)$ на X вводится метрика на сетке Y :

$$r((m_i, h_i), (m, h)) = \sqrt{(m - m_i)^2 + (h - h_i)^2}$$

Окрестность (m_i, h_i) :



Обучение карты Кохонена

Вход: X^ℓ — обучающая выборка; η — темп обучения;

Выход: $w_{mh} \in \mathbb{R}^n$ — векторы весов, $m = 1..M$, $h = 1..H$;

$w_{mh} := \text{random} \left(-\frac{1}{2MH}, \frac{1}{2MH} \right)$ — инициализация весов;

повторять

выбрать объект x_i из X^ℓ случайным образом;

WTA: вычислить координаты кластера:

$(m_i, h_i) := a(x_i) \equiv \arg \min_{(m,h) \in Y} \rho(x_i, w_{mh});$

для всех $(m, h) \in \text{Окрестность}(m_i, h_i)$

WTM: сделать шаг градиентного спуска:

$w_{mh} := w_{mh} + \eta(x_i - w_{mh}) K(r((m_i, h_i), (m, h)));$

пока кластеризация не стабилизируется;

Интерпретация карт Кохонена

Два типа графиков — цветных карт $M \times H$:

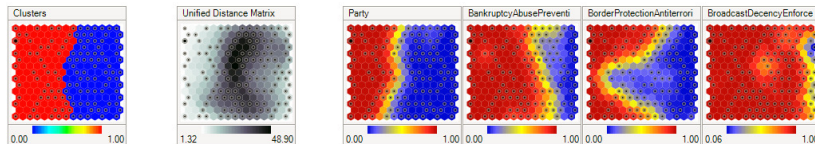
- Цвет узла (m, h) — локальная плотность в точке (m, h) — среднее расстояние до k ближайших точек выборки
- По одной карте на каждый признак:
цвет узла (m, h) — значение j -й компоненты вектора $w_{m,h}$

Пример: задача UCI house-votes (US Congress voting patterns)

Объекты — конгрессмены

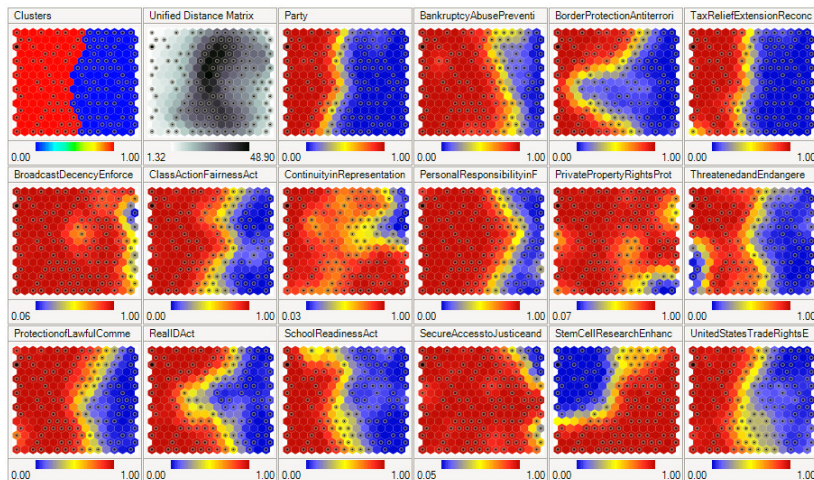
Признаки — результаты голосования по различным вопросам

Есть целевой признак «партия» $\in \{\text{демократ, республиканец}\}$



Интерпретация карт Кохонена (продолжение примера)

Пример: задача UCI house-votes (US Congress voting patterns)



Достоинства и недостатки карт Кохонена

Достоинства:

- Возможность визуального анализа многомерных данных

Недостатки:

- **Субъективность.** Карта зависит не только от кластерной структуры данных, но и от...
 - свойств сглаживающего ядра;
 - (случайной) инициализации;
 - (случайного) выбора x_i в ходе итераций.
- **Искажения.** Близкие объекты исходного пространства могут переходить в далёкие точки на карте, и наоборот.

Рекомендуется только для разведочного анализа данных.

Построение автокодировщика — задача обучения без учителя

$X^\ell = \{x_1, \dots, x_\ell\}$ — 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: $x \in \mathbb{R}^n$, $z \in \mathbb{R}^m$

$$f(x, A) = \underset{m \times n}{A} x, \quad g(z, B) = \underset{n \times m}{B} z$$

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

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

Обучение и использование автокодировщиков

Метод обучения:

- Стохастический градиент (SG) по параметрам (α, β)

Способы использования:

- Векторизация данных:
 - понижение размерности (dimensionality reduction)
 - синтез более удачных признаков (feature generation)
 - сжатие данных с минимальной потерей информации
- Обучение с учителем в новом пространстве признаков
- Векторизация объектов для нейронных сетей
- Генерация синтетических объектов, похожих на реальные

Rumelhart, Hinton, Williams. Learning internal representations by error propagation. 1986.

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

Линейный автокодировщик и метод главных компонент

Линейный автокодировщик: $f(x, A) = Ax$, $g(z, B) = Bz$,

$$\mathcal{L}_{AE}(A, B) = \sum_{i=1}^{\ell} \|BAx_i - x_i\|^2 \rightarrow \min_{A, B}$$

Метод главных компонент: $f(x, U) = U^T x$, $g(z, U) = Uz$,
в матричных обозначениях $F = (x_1 \dots x_\ell)^T$, $U^T U = I_m$, $G = FU$,

$$\|F - GU^T\|^2 = \sum_{i=1}^{\ell} \|UU^T x_i - x_i\|^2 \rightarrow \min_U$$

Автокодировщик обобщает метод главных компонент:

- не обязательно $B = A^T$ (хотя часто именно так и делают)
- произвольные A, B вместо ортогональных
- нелинейные модели $f(x, \alpha)$, $g(z, \beta)$ вместо Ax, Bz
- произвольная функция потерь \mathcal{L} вместо квадратичной
- SG оптимизация вместо сингулярного разложения SVD

Разреживающие автокодировщики (Sparse AE)

Применение L_1 или L_2 -регуляризации к векторам весов α, β :

$$\mathcal{L}_{\text{AE}}(\alpha, \beta) + \lambda \|\alpha\| + \lambda \|\beta\| \rightarrow \min_{\alpha, \beta}$$

Применение L_1 -регуляризации к кодовым векторам z_i :

$$\mathcal{L}_{\text{AE}}(\alpha, \beta) + \lambda \sum_{i=1}^{\ell} \sum_{j=1}^m |f_j(x_i, \alpha)| \rightarrow \min_{\alpha, \beta}$$

Энтропийная регуляризация для случая $f_j \in [0, 1]$:

$$\mathcal{L}_{\text{AE}}(\alpha, \beta) + \lambda \sum_{j=1}^m \text{KL}(\varepsilon \| \bar{f}_j) \rightarrow \min_{\alpha, \beta},$$

где $\bar{f}_j = \frac{1}{\ell} \sum_{i=1}^{\ell} f_j(x_i, \alpha)$; $\varepsilon \in (0, 1)$ — близкий к нулю параметр,

$\text{KL}(\varepsilon \| \rho) = \varepsilon \log \frac{\varepsilon}{\rho} + (1 - \varepsilon) \log \frac{1 - \varepsilon}{1 - \rho}$ — KL-дивергенция.

D. Arpit et al. Why regularized auto-encoders learn sparse representation? 2015.

Шумоподавляющий автокодировщик (Denoising AE)

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 SG объекты x_i сэмплируются и зашумляются по одному: $\tilde{x} \sim q(\tilde{x}|x_i)$.

Варианты зашумления $q(\tilde{x}|x_i)$:

- $\tilde{x} \sim \mathcal{N}(x_i, \sigma^2 I)$ — добавление гауссовского шума
- обнуление компонент вектора x_i с вероятностью p_0
- такие искажения $x_i \rightarrow \tilde{x}$, относительно которых реконструкция \hat{x}_i должна быть устойчивой

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

Реляционный автокодировщик (Relational AE)

Наряду с потерями реконструкции объектов минимизируем потери реконструкции отношений между объектами:

$$\mathcal{L}_{\text{AE}}(\alpha, \beta) + \lambda \sum_{i < j} \mathcal{L}(\sigma(\hat{x}_i^\top \hat{x}_j), \sigma(x_i^\top x_j)) \rightarrow \min_{\alpha, \beta}$$

где $\hat{x}_i = g(f(x_i, \alpha), \beta)$ — реконструкция объекта x_i ,
 $x_i^\top x_j$ — скалярное произведение (близость) пары объектов,
 $\sigma(s) = (s - s_0)_+$ — пороговая функция с параметром s_0
(если векторы не близки, то неважно, насколько),
 $\mathcal{L}(\hat{s}, s)$ — функция потерь, например, $(\hat{s} - s)^2$.

Эксперимент: улучшается качество классификации изображений с помощью кодовых векторов на задачах MNIST, CIFAR-10

Автокодировщики для обучения с учителем

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

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

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

$z_i = f(x_i, \alpha)$ — кодировщик

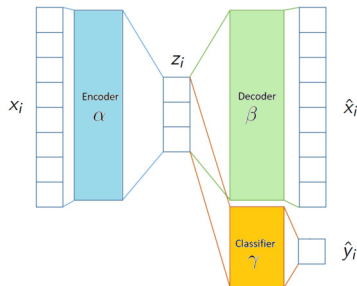
$\hat{x}_i = g(z_i, \beta)$ — декодировщик

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

Функции потерь:

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

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



Dor Bank, Noam Koenigstein, Raja Giryes. Autoencoders. 2020

Многомерное шкалирование (multidimensional scaling, MDS)

Дано: $(i, j) \in E$ — выборка рёбер графа $\langle V, E \rangle$,

R_{ij} — расстояния между вершинами ребра (i, j) .

Например, R_{ij} — длина кратчайшего пути по графу (IsoMAP).

Найти: векторные представления вершин $z_i \in \mathbb{R}^d$, так, чтобы близкие (по графу) вершины имели близкие векторы.

Критерий стресса (stress):

$$\sum_{(i,j) \in E} w_{ij} (\rho(z_i, z_j) - R_{ij})^2 \rightarrow \min, \quad Z \in \mathbb{R}^{V \times d},$$

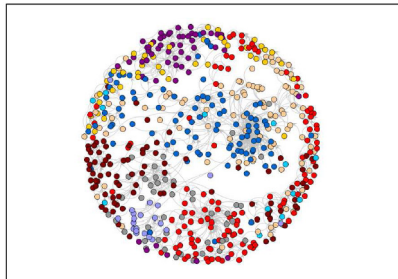
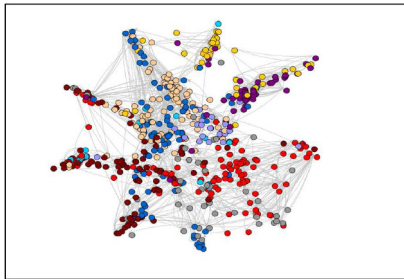
где $\rho(z_i, z_j) = \|z_i - z_j\|$ — обычно евклидово расстояние,
 w_{ij} — веса (какие расстояния важнее, большие или малые).

Обычно решается методом стохастического градиента (SG).

I. Chami et al. Machine learning on graphs: a model and comprehensive taxonomy. 2020.

Многомерное шкалирование для визуализации данных

При $d = 2$ осуществляется проекция выборки на плоскость



- Используется для визуализации кластерных структур
- Форму облака точек можно настраивать весами и метрикой
- Недостаток — искажения неизбежны
- Популярная современная разновидность метода — t-SNE

Laurens van der Maaten, Geoffrey Hinton. Visualizing data using t-SNE. 2008

Метод векторного представления соседства (Stochastic Neighbor Embedding, SNE)

Дано: исходные точки $x_i \in \mathbb{R}^n$, $i = 1, \dots, \ell$

Найти: точки на карте-проекции $z_i \in \mathbb{R}^d$, $i = 1, \dots, \ell$, $d \ll n$

Критерий: расстояния $\|z_i - z_j\|$ близки к исходным $\|x_i - x_j\|$

Вероятностная модель события « j является соседом i »

на основе перенормированных гауссовских распределений:

$p(j|i) = \text{norm}_{j \neq i} \exp(-\frac{1}{2\sigma_i^2} \|x_i - x_j\|^2)$ — в исходном пространстве;

$q(j|i) = \text{norm}_{j \neq i} \exp(-\|z_i - z_j\|^2)$ — в пространстве проекции;

где $p_j = \text{norm}_j(v_j) = \frac{v_j}{\sum_k v_k}$ — операция нормировки вектора v .

Максимизация правдоподобия (стохастическим градиентом):

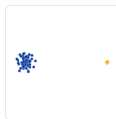
$$\sum_i \sum_{j \neq i} p(j|i) \ln q(j|i) \rightarrow \max_{\{z_i\}}$$

Преимущества метода SNE

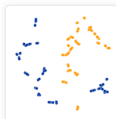
- Преобразование расстояний в вероятности устраняет дисбалансы между большими и малыми расстояниями
- Дисбаланс между точками с большой и малой плотностью соседей выравнивается настройкой σ_i по перплексии

$H(i) = -\sum_j p(j|i) \log_2 p(j|i)$ — энтропия распределения $p(j|i)$;
 $2^{H(i)}$ — перплексия = «эффективное число соседей у x_i »
(если $p(j|i) = \frac{1}{k}$, то $2^{H(i)} = k$); обычно перплексия = 5..50.

Выбор перплексии может существенно влиять на вид проекции:



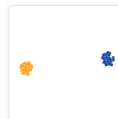
Original



Perplexity: 2



Perplexity: 5



Perplexity: 30

G.E.Hinton, S.T.Roweis. Stochastic Neighbor Embedding. 2002.

Вероятностная модель t-SNE: два усовершенствования SNE

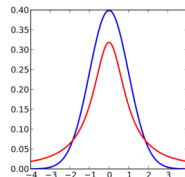
Проблема скученности в SNE: окрестность вмещает гораздо больше точек в n -мерном пространстве, чем в d -мерном

- Использование t -распределения Стьюдента с более тяжёлым хвостом и симметричного совместного распределения $q(i, j)$:

$$q(i, j) = \text{norm}_{(i, j): i \neq j} (1 + \|z_i - z_j\|^2)^{-1}$$

- Использование совместного распределения $p(i, j)$:

$$p(i, j) = \frac{1}{2\ell} (p(j|i) + p(i|j))$$



Максимизация правдоподобия (стохастическим градиентом):

$$\sum_{(i, j): i \neq j} p(i, j) \ln q(i, j) \rightarrow \max_{\{z_i\}}$$

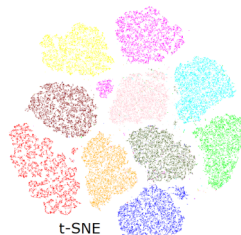
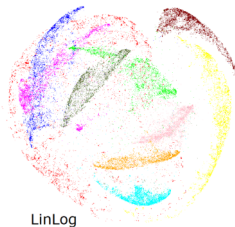
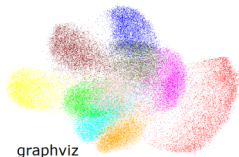
L.J.P. van der Maaten, G.Hinton. Visualizing data using t-SNE. 2008

Преимущества и недостатки t-SNE

Лучшее представление структур сходства по сравнению с другими методами многомерного шкалирования (mnist)

| | | | | | |
|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 1 | 1 | 1 | 1 | 1 |
| 2 | 2 | 2 | 2 | 2 | 2 |
| 3 | 3 | 3 | 3 | 3 | 3 |
| 4 | 4 | 4 | 4 | 4 | 4 |

| | | | | | |
|---|---|---|---|---|---|
| 5 | 5 | 5 | 5 | 5 | 5 |
| 6 | 6 | 6 | 6 | 6 | 6 |
| 7 | 7 | 7 | 7 | 7 | 7 |
| 8 | 8 | 8 | 8 | 8 | 8 |
| 9 | 9 | 9 | 9 | 9 | 9 |



Ложные кластерные структуры при низкой перплексии
Размеры кластеров и расстояния между ними неинформативны
Трудно отличить реальные структуры от артефактов метода
Нет ясного критерия качества для подбора перплексии

M. Wattenberg, F. Viegas, I. Johnson (Google). How to use t-SNE effectively. 2016.
<https://distill.pub/2016/misread-tsne>

Матричные разложения графа (graph factorization)

Дано: $(i, j) \in E$ — выборка рёбер графа $\langle V, E \rangle$,

S_{ij} — близость между вершинами ребра (i, j) .

Например, $S_{ij} = [(i, j) \in E]$ — матрица смежности вершин.

Найти: векторные представления вершин, так, чтобы близкие (по графу) вершины имели близкие векторы.

Критерий для неориентированного графа (S симметрична):

$$\|S - ZZ^T\|_E = \sum_{(i,j) \in E} (\langle z_i, z_j \rangle - S_{ij})^2 \rightarrow \min_Z, \quad Z \in \mathbb{R}^{V \times d}$$

Критерий для ориентированного графа (S несимметрична):

$$\|S - \Phi\Theta^T\|_E = \sum_{(i,j) \in E} (\langle \varphi_i, \theta_j \rangle - S_{ij})^2 \rightarrow \min_{\Phi, \Theta}, \quad \Phi, \Theta \in \mathbb{R}^{V \times d}$$

Обычно решается методом стохастического градиента (SG).

I. Chami et al. Machine learning on graphs: a model and comprehensive taxonomy. 2020.

Векторные представления графов как автокодировщики

Все рассмотренные выше методы векторных представлений графов суть автокодировщики данных о рёбрах:

- многомерное шкалирование: $R_{ij} \rightarrow \|z_i - z_j\|$
- SNE и t-SNE: $p(i, j) \rightarrow q(i, j) \propto K(\|z_i - z_j\|)$
- матричные разложения: $S_{ij} \rightarrow \langle \varphi_i, \theta_j \rangle$

Вход кодировщика:

- W_{ij} — данные о ребре графа (i, j)

Выход кодировщика:

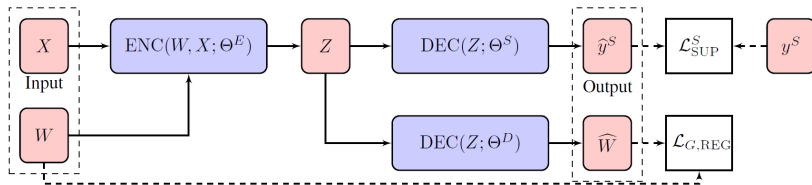
- векторные представления вершин z_i

Выход декодировщика:

- аппроксимация \hat{W}_{ij} , вычисляемая по (z_i, z_j)

GraphEDM: обобщённый автокодировщик на графах

Graph Encoder Decoder Model — обобщает более 30 моделей:



$W \in \mathbb{R}^{V \times V}$ — входные данные о рёбрах

$X \in \mathbb{R}^{V \times n}$ — входные данные о вершинах, признаковые описания

$Z \in \mathbb{R}^{V \times d}$ — векторные представления вершин графа

$DEC(Z; \Theta^D)$ — декодер, реконструирующий данные о рёбрах

$DEC(Z; \Theta^S)$ — декодер, решающий supervised-задачу

y^S — (semi-)supervised данные о вершинах или рёбрах

\mathcal{L} — функции потерь

I. Chami et al. Machine learning on graphs: a model and comprehensive taxonomy. 2020.

Постановка задачи векторизации слов

Дано: текст $(w_1 \dots w_n)$, состоящий из слов словаря W

Найти: векторные представления (embedding) слов $v_w \in \mathbb{R}^d$, так, чтобы близкие по смыслу слова имели близкие векторы

Модель Skip-gram для предсказания вероятности слов контекста $C_i = (w_{i-k} \dots w_{i-1} w_{i+1} \dots w_{i+k})$ по слову w_i :

$$p(w|w_i) = \underset{w \in W}{\text{SoftMax}} \langle u_w, v_{w_i} \rangle \equiv \underset{w \in W}{\text{norm}} (\exp \langle u_w, v_{w_i} \rangle),$$

v_w — вектор предсказывающего слова,

u_w — вектор предсказываемого слова, в общем случае $u_w \neq v_w$.

Критерий максимума log-правдоподобия, $U, V \in \mathbb{R}^{|W| \times d}$:

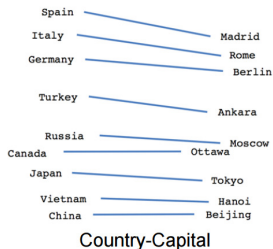
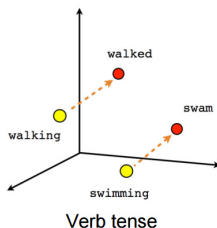
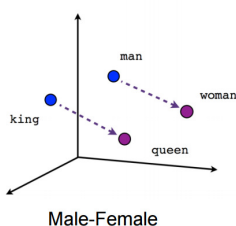
$$\sum_{i=1}^n \sum_{w \in C_i} \log p(w|w_i) \rightarrow \max_{U, V}$$

T. Mikolov et al. Efficient estimation of word representations in vector space, 2013.

Почему эмбединги слов отражают их смыслы

Основная гипотеза дистрибутивной семантики [Harris, 1954]:
«Смысл слова определяется множеством его контекстов»

Задача семантической аналогии слов:
по трём словам угадать четвёртое



Z.Harris. Distributional structure. 1954.

J.R.Firth. A synopsis of linguistic theory 1930-1955. Oxford, 1957.

P.Turney, P.Pantel. From frequency to meaning: vector space models of semantics. 2010.

Подмена задачи: классификация пар слов на два класса

Критерий log-loss для SGNS (Skip-gram Negative Sampling):

$$\sum_{i=1}^n \sum_{w \in C_i} \left(\log p(+1|w, w_i) + \log p(-1|\bar{w}, w_i) \right) \rightarrow \max_{U, V}$$

где $p(y|w, w_i) = \sigma(y \langle u_w, v_{w_i} \rangle)$ — модель классификации, $y = \pm 1$;
 $y = +1$, если пара слов (w, w_i) находится в общем контексте;
 $y = -1$, если пара слов (w, w_i) не находится в общем контексте;
 $\bar{w} \sim p(w)^{3/4}$ сэмплируется из $W \setminus C_i$ в методе SG.

Эвристики и прочие замечания:

- Dynamic window: случайный выбор $k \sim [3..10]$
- Итоговые векторы слов: $\alpha v_w + (1 - \alpha) u_w$
- Приём NS полезен, когда не хватает второго класса

Модель векторных представлений FastText

Идея: векторное представление слова w определяется как сумма векторов всех его буквенных n -грамм $G(w)$:

$$u_w = \sum_{g \in G(w)} u_g$$

В Skip-gram вместо векторов слов u_w обучаются векторы u_g

Пример: $G(\text{дармолюб}) = \{\langle \text{да, арм, рмо, мол, олю, люб, юб} \rangle\}$

Преимущества:

- Это решает проблемы новых слов и слов с опечатками
- Подходит для обработки текстов социальных медиа
- Словарь 2- и 3-грамм обычно меньше словаря W
- Существует много предобученных моделей

Разновидности векторизации данных:

- *Квантизация* — сокращение объёма выборки, замена объектов ближайшими центрами кластеров
- *Автокодировщики* — синтез векторных представлений (эмбедингов) объектов, обычно с понижением размерности
- *MDS, t-SNE, GF, word2vec, graph2vec* и др. — синтез эмбедингов объектов по данным об их взаимодействии

Методы обучения — на основе SG

Тексты — это разновидность графов:

- слова — вершины графа, сочетаемость пары слов — ребро
- слова и документы — вершины двух разных долей графа, вхождение слова в документ — ребро двудольного графа