

РЕКОМЕНДАТЕЛЬНЫЕ СИСТЕМЫ II

Сергей Николенко

СПбГУ — Санкт-Петербург

19 ноября 2020 г.

Random facts:

- 19 ноября — Международный мужской день, впервые отмеченный в 1999 году в Тринидаде и Тобаго; кроме того, 19 ноября — Всемирный день туалета, ведь именно 19 ноября 2001 года в славящемся чистотой Сингапуре была учреждена Всемирная туалетная организация (World Toilet Organization, WTO)
- 19 ноября 1824 г. уровень воды в Неве поднялся на 4м 14-21см выше ординара
- 19 ноября 1969 г. на стадионе «Маракана» Эдсон Арантес ду Насименту забил свой 1000-й гол в 909-м официальном матче
- 19 ноября 1979 г. прошёл первый рейс из Ленинграда в Москву скоростного поезда ЭР-200, построенного рижскими вагоностроителями; сначала время в пути составляло 4 часа 50 минут, а потом сократилось до 4ч 30мин

МАТРИЧНЫЕ РАЗЛОЖЕНИЯ И ВЕРОЯТНОСТНЫЕ МОДЕЛИ

- Из чего складывается рейтинг пользователя i , который он выдал продукту a ?
- Вполне может быть, что пользователь добрый и всем подряд выдаёт хорошие рейтинги; или, наоборот, злой и рейтинг зажимает.
- С другой стороны, некоторые продукты попросту лучше других.
- Поэтому мы вводим так называемые *базовые предикторы* (baseline predictors) $b_{i,a}$, которые складываются из базовых предикторов отдельных пользователей b_i и базовых предикторов отдельных продуктов b_a , а также просто общего среднего рейтинга по базе μ :

$$b_{i,a} = \mu + b_i + b_a.$$

- Чтобы найти предикторы, уже нужен байесовский подход: надо добавить нормально распределённый шум и получить модель линейной регрессии

$$r_{i,a} \sim N(\mu + b_i + b_a, \sigma^2).$$

- Можно ввести априорные распределения и оптимизировать; или просто найти среднеквадратическое отклонение с регуляризатором:

$$b_* = \arg \min_b \sum_{(i,a)} (r_{i,a} - \mu - b_i - b_a)^2 + \lambda_1 \left(\sum_i b_i^2 + \sum_a b_a^2 \right).$$

- Как обучить такую модель?

- Правильно, это просто линейная регрессия!
- Можно обучать любым способом, которым мы умеем обучать линейную регрессию, например решить линейную систему:

$$b_* = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T r,$$

где r – вектор рейтингов, \mathbf{X} – матрица из нулей и единиц (в каждой строке по три единицы).

- Часто у нас рейтинги бинарные: лайк/дизлайк, например.
- Тогда разумнее не пытаться приближать 0-1 данные линейной функцией, а рассмотреть их как распределение вероятностей.
- Давайте моделировать через логистический сигмоид:

$$b_{i,a} = \sigma(\mu + b_i + b_a), \quad \sigma(x) = \frac{1}{1 + e^{-x}}.$$

- Теперь это не линейная регрессия, а логистическая; её мы тоже умеем обучать.
- Часто это хороший метод, даже если рейтингов несколько.

- С тем, чтобы напрямую обучать оставшуюся матрицу предпочтений вероятностными методами, опять же есть проблема – матрица X , выражающая рейтинги, содержит $N \times M$ параметров, гигантское число, которое, конечно, никак толком не обучить.
- Более того, обучать их и не надо – как мы уже говорили, данные очень разреженные, и «на самом деле» свободных параметров гораздо меньше, проблема только с тем, как их выделить.
- Поэтому обычно число независимых параметров модели необходимо уменьшать.

- Метод SVD (singular value decomposition) – разложим матрицу X в произведение матриц маленького ранга.
- Зафиксируем некоторое число f *скрытых факторов*, которые так или иначе описывают каждый продукт и предпочтения каждого пользователя относительно этих факторов.
- Пользователь – вектор $p_i \in \mathbb{R}^f$, который показывает, насколько пользователь предпочитает те или иные факторы; продукт – вектор $q_a \in \mathbb{R}^f$, который показывает, насколько выражены те или иные факторы в этом продукте.

- Предпочтение в итоге будем подсчитывать просто как скалярное произведение $q_a^\top p_i = \sum_{j=1}^f q_{a,j} p_{i,j}$.
- Таким образом, добавляя теперь сюда baseline-предикторы, получаем следующую модель предсказаний рейтингов:

$$\hat{r}_{i,a} \sim \mu + b_i + b_a + q_a^\top p_i.$$

- Как обучить такую модель?

- SGD – стохастический градиентный спуск.
- Считаём градиент функции правдоподобия, идём по тренировочным примерам, апдейтим на каждом шаге:

$$b_i := b_i + \gamma (e_{i,a} - \lambda_2 b_i),$$

$$b_a := b_a + \gamma (e_{i,a} - \lambda_2 b_a),$$

$$q_{a,j} := q_{a,j} + \gamma (e_{i,a} p_{i,j} - \lambda_2 q_{a,j}) \text{ для всех } j,$$

$$p_{i,j} := p_{i,j} + \gamma (e_{i,a} q_{a,j} - \lambda_2 p_{i,j}) \text{ для всех } j,$$

где γ – скорость обучения.

- ALS – попеременные наименьшие квадраты.
- Заметим, что если в $\hat{r}_{i,a} \sim \mu + b_i + b_a + q_a^\top p_i$ фиксировать p_i , то по q_a будет опять линейная регрессия.
- ALS – повторять до сходимости:
 - фиксируем p_i , обучим q_a ;
 - фиксируем q_a , обучим p_i .
- Обычно более устойчив и быстрее работает, чем SGD.

- То же замечание про логистический вариант – при бинарных рейтингах можно рассмотреть

$$\hat{r}_{i,a} \sim \sigma(\mu + b_i + b_a + q_a^\top p_i).$$

- Тогда в SGD просто добавятся производные сигмоида $\sigma'(x) = \sigma(x)(1 - \sigma(x))$.
- А в ALS вместо линейной регрессии на каждой итерации надо будет обучать логистическую.

- Пусть мы хотим построить разложение матрицы рейтингов на матрицы меньшего ранга:

$$\hat{R} = U^T V.$$

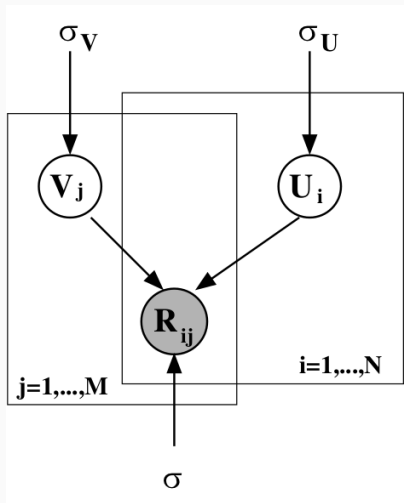
- Вероятностно мы имеем правдоподобие

$$p(R | U, V, \sigma^2) = \prod_i \prod_a (N(r_{i,a} | u_i^T v_j, \sigma^2))^{[i \text{ оценил } a]}.$$

- Добавим гауссовские априорные распределения на U и V :

$$p(U | \sigma_U^2) = \prod_i N(U_i | 0, \sigma_U^2 I), \quad p(V | \sigma_V^2) = \prod_a N(V_a | 0, \sigma_V^2 I).$$

ГРАФИЧЕСКАЯ МОДЕЛЬ



ВЕРОЯТНОСТНОЕ РАЗЛОЖЕНИЕ МАТРИЦ

- Если просто зафиксировать σ^2 , σ_V^2 и σ_U^2 , то они будут играть роль регуляризаторов, и нет никакого отличия от «обычного» SVD.
- Разница здесь в том, что теперь мы можем автоматически найти оптимальные $\sigma = (\sigma^2, \sigma_V^2, \sigma_U^2)$, максимизируя общее правдоподобие модели

$$\sigma^* = \arg \max_{\sigma} p(R | \sigma) = \arg \max_{\sigma} \int p(R, U, V | \sigma) dU dV$$

EM-алгоритмом:

- сначала зафиксируем σ и найдём

$$f(\sigma) = E_{U, V | R, \sigma} [\log p(R, U, V | \sigma)];$$

- потом максимизируем

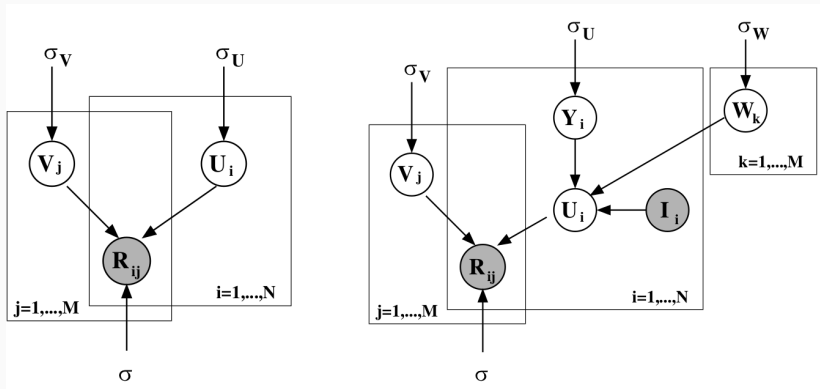
$$\sigma := \arg \max_{\sigma} f(\sigma).$$

- Модификация: пользователи с малым числом оценок в PMF получают апостериорные распределения, очень похожие на «среднего пользователя».
- Чтобы на редких пользователей лучше обобщалось, добавим ещё факторы, которые меняют априорное распределение факторов у пользователя в зависимости от того, сколько и чего он оценил:

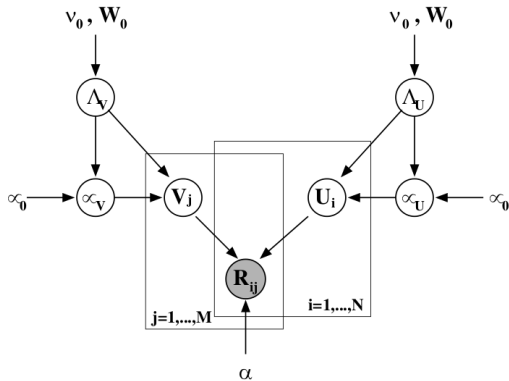
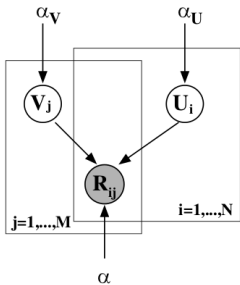
$$U_i = Y_i + \frac{\sum_a [i \text{ оценил } a] W_a}{\sum_a [i \text{ оценил } a]}.$$

- Матрица W показывает, как влияет на априорное распределение
- Тоже в качестве регуляризатора берём априорный гауссиан:
 $p(W | \sigma_W^2) = \prod_i N(W_i | 0, \sigma_W^2 I).$

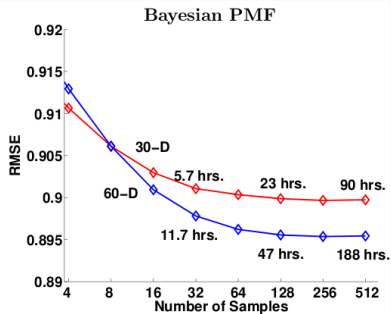
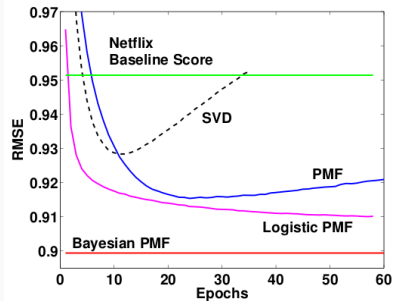
ГРАФИЧЕСКАЯ МОДЕЛЬ



ГРАФИЧЕСКАЯ МОДЕЛЬ



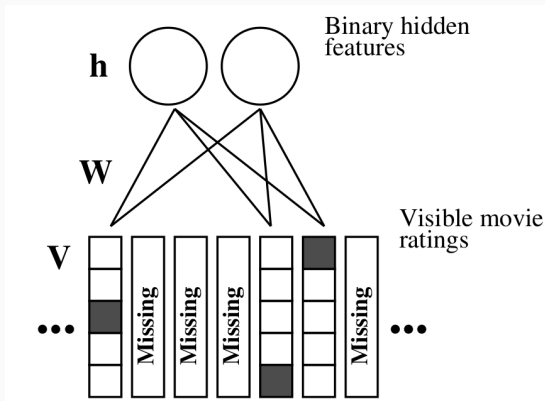
ГРАФИЧЕСКАЯ МОДЕЛЬ



- Ещё один метод вероятностного моделирования – *машины Больцмана* (restricted Boltzmann machine).
- Ненаправленная графическая модель, состоящая из двух уровней, видимого и скрытого.
- Машины Больцмана – один из основных блоков, на которых строят deep learning, так что о них мы поговорим гораздо подробнее позже.

МАШИНЫ БОЛЬЦМАНА

- В коллаборативной фильтрации мы строим машиной Больцмана модель предпочтений пользователя.



- В результате на скрытом слое обучается модель пользователя.
- Метод обучения – contrastive divergence (приближение к максимальному правдоподобию).
- RBM не то чтобы *лучше* SVD, но часто ошибается в *других местах*, поэтому комбинация этих двух моделей даёт значительное улучшение.

- Кстати, что значит «комбинация моделей»?
 1. Просто линейная комбинация (байесовское усреднение или регрессия).
 2. *Бустинг*: метод комбинации простых классификаторов; в качестве простых классификаторов можно брать результаты сложных моделей (оценки вероятности успеха или ожидаемый рейтинг).
- Современные рекомендательные системы – обычно большие *ансамбли* моделей.
- Два уровня: обучение отдельных моделей в ансамбле (*bootstrapping* и т.п.) и обучение комбинации.

- Мы только что описали метод SVD (singular value decomposition) – разложение матрицы X в произведение матриц маленького ранга.
- Это не единственное матричное разложение, и все они интересны по-своему.

- PCA (анализ главных компонент) – один из самых популярных методов сокращения размерности.
- PCA пытается объяснить как можно больше дисперсии исходного датасета.
- Однако часто направления на кластеры в данных не ортогональны, а признаки PCA трудно интерпретировать.

- SVD (сингулярное разложение) мы уже рассматривали.
- Оно делает ровно то, что нам нужно, когда доступны рейтинги:
 - максимизирует правдоподобие имеющихся рейтингов (минимизирует ошибку предсказания рейтингов);
 - умеет работать с разреженными матрицами (минимизирует только по имеющимся рейтингам).
- Но что делать, когда рейтингов никаких нет, а есть только факт использования? SVD тогда ничего хорошего не сделает...

- NMF (неотрицательное разложение): раскладываем по-прежнему в произведение

$$X \approx UV^T,$$

где U размера $n \times f$, V размера $m \times f$, и f гораздо меньше n и m .

- Но теперь требуем, чтобы элементы U и V были неотрицательны.
- Результат – признаки, которые лучше интерпретируются и могут иметь больше смысла.

- NMF можно использовать для рекомендаций, когда рейтингов нет.
- Берём матрицу X , в которой единицы – продукты, использованные пользователем, остальное нули.
- И раскладываем по NMF (методом ALS, нужно итеративно решать системы размерности $n \times f$ и $m \times f$).
- Задача перестаёт быть разреженной, но кое-что сделать можно.

МЕТРИКИ КАЧЕСТВА И РАСШИРЕНИЯ

- Ещё одна важная тема, которую в рекомендательных системах часто замалчивают.
- Как оценить качество рекомендаций? Какая должна быть метрика качества?
- Когда мы обучаем SVD (максимизируем правдоподобие), мы оптимизируем среднеквадратичную ошибку отклонения.
- И Netflix Prize, например, так и был сформулирован: надо было оптимизировать RMSE, среднеквадратичное отклонение предсказаний от истинного рейтинга пользователя.
- Но что нам надо на самом деле? Что у нас в тестовом множестве?

- В тестовом множестве отложены рейтинги некоторых продуктов, которым пользователь дал оценку.
- Наша задача – выдать пользователю топ-рекомендации; не предсказать все рейтинги, а найти, у каких продуктов рейтинг будет самым большим.
- То есть на самом деле это задача *ранжирования*! И метрики качества лучше брать из *information retrieval*; там, когда оценивают качество выдачи, вовсе не пытаются минимизировать отклонение предсказания функции релевантности.
- Дальше для простоты будем рассматривать бинарный случай (лайк-дизлайк).

- Классические метрики:
 - (1) точность (precision) – количество «хороших» (релевантных) запросу в случае поисковой выдачи, отмеченных высокой оценкой в случае рекомендательной системы) документов в выдаче, делённое на общее количество документов в выдаче;
 - (2) полнота (recall) – количество «хороших» документов в выдаче, делённое на общее количество релевантных документов в базе поисковой системы.
- Однако здесь те же проблемы; эти параметры не зависят от ранжирования выдачи, надо знать заранее, сколько потребуется рекомендаций.

- Метрики качества ранжирования:
 - NDCG, Normalized Discounted Commulative Gain; выберем топ- k рекомендаций (k может быть заведомо больше нужного числа) и посчитаем:

$$\text{DCG}_k = \sum_{i=1}^k \frac{2^{\hat{r}_i} - 1}{\log_2(1 + i)},$$
$$\text{NDCG}_k = \frac{\text{DCG}_k}{\text{IDCG}_k},$$

где \hat{r}_i – наша оценка рейтинга продукта на позиции i , а IDCG_k – значение DCG_k при ранжировании по истинным значениям (рейтингам из валидационного набора);

- NDCG от 0 до 1, но ей трудно придумать естественную интерпретацию (как вероятность чего-нибудь, например).

- Метрики качества ранжирования:
 - AUC, Area Under (ROC) Curve; можно считать по всей выдаче сразу;
 - AUC – вероятность того, что случайно выбранная пара продуктов с разными оценками будет отранжирована правильно (понравившийся будет выше в выдаче, чем не понравившийся);
 - в бинарном случае можно посчитать в замкнутом виде:

$$\hat{A} = \frac{S_0 - n_0(n_0 + 1)/2}{n_0 n_1},$$

где n_0, n_1 – число понравившихся и не понравившихся пользователю объектов, $S_0 = \sum p_i$ – сумма номеров позиций понравившихся объектов в выдаче.

- Метрики качества ранжирования:
 - но на самом деле простые метрики тоже важны, потому что обычно пользователь успеет увидеть только несколько самых верхних рекомендаций;
 - WTA (winner takes all) – эта метрика равна 1, если топ-рекомендация (с самым большим предсказанным рейтингом) из просмотренных пользователем получила положительную оценку, и 0 в противном случае;
 - Top k – доля положительных оценок среди топ- k рекомендаций (Top10 часто называют MAP – mean average precision).

- Проблема: холодный старт.
- Надо как-то инициализировать; если вообще ничего не знаем, сделать ничего нельзя, конечно.
- Но так не бывает; обычно есть набор признаков – можно пытаться предсказывать значения факторов:
 - просто регрессией по признакам;
 - (обычно для продуктов) выделяя темы при помощи topic modeling.

- Получается, что для признаков пользователя x_i и продукта x_a мы рассматриваем модель

$$r_{i,a} \sim \mu + b_{\text{user}}(x_i) + b_{\text{item}}(x_a) + q_a^\top p_i(t),$$

где

$$b_{\text{user}}(x_i) \sim N(u(x_i), \sigma_u^2),$$

$$b_{\text{item}}(x_i) \sim N(v(x_i), \sigma_v^2),$$

и в качестве u и v может выступать любая регрессия [Agarwal, Chen, 2009].

- Ещё один вариант, через контент:
 - выделить темы из продуктов (LDA), получится распределение $z_{a,k}$ для каждого a ;
 - обучить факторы $s_{i,k}$ того, насколько пользователю “нравятся” эти темы;
 - затем для нового продукта оценить темы $\hat{z}_{a,k}$ по контенту, а потом добавлять в модель слагаемое

$$r_{i,a} \sim \dots + \sum_k s_{i,k} \hat{z}_{a,k},$$

что помогает для холодного старта по продуктам.

- Есть модели, связанные с тем, как обучить не абы какие темы, а хорошо выражающие предпочтения (fLDA).

ВРЕМЯ В КОЛЛАБОРАТИВНОЙ ФИЛЬТРАЦИИ

- Пример: давайте добавим время, т.е. будем рассматривать базовые предикторы и характеристики пользователя как функции от времени:

$$\hat{r}_{i,a} = \mu + b_i(t) + b_a(t) + q_a^\top p_i(t),$$

где

$$b_a(t) = b_a + b_{a, \text{Bin}(t)},$$

$$b_i(t) = b_i + \alpha_i \text{dev}_i(t) + b_{i,t},$$

$$p_{i,f}(t) = p_{i,f} + \alpha_{i,f} \text{dev}_i(t) + p_{i,f,t} + \frac{1}{\sqrt{|V(i)|}} \sum_{b \in V(i)} y_b,$$

$$\text{dev}_i(t) = \text{sign}(t - t_i) |t - t_i|^\beta.$$

- Это называется timeSVD++, и эта модель была одним из основных компонентов модели, взявшей Netflix Prize.

- Предположим, что пользователи приходят из социальной сети.
- Т.е. есть друзья, есть социальный граф (его часть) и т.д. Это тоже можно добавить в рекомендательную модель:
 - фильтр/перевзвешивание в методе ближайших соседей;
 - дополнительные слагаемые в разложение типа SVD;
 - разложение матрицы доверия (из социального графа) вместе с матрицей рейтингов, меняем априорное распределение для PMF и т.д.

- Filter bubble: как вывести человека за его привычный круг.
- Можно до конца жизни рекомендовать одно и то же;
метрики:
 - diversity – разнообразие, мера похожести элементов списка;
 - novelty – новизна для пользователя, распространённость продукта, доля его рейтингов;
 - serendipity – неожиданность, сюрприз, похожесть на историю пользователя.
- Для всего этого нужно уметь распознавать похожесть контента рекомендованных товаров.

- CARS (context-aware recommender systems) – мы рекомендуем в контексте:
 - временном;
 - ситуативном;
 - географическом;
 - предшествующего поведения пользователей и т.д.

- Формально контекст – это новые измерения в матрице предпочтений.
- Получается “гиперкуб” данных, есть методы тензорного разложения, аналогичного SVD.
- Но часто не хуже работают простые решения – отфильтровать контексты и обучить модели только по этим данным, добавить полученные модели и сам контекст как факторы в бленд.

Спасибо за внимание!