### ВАРИАНТЫ И ПРИЛОЖЕНИЯ ЕМ-АЛГОРИТМА

Сергей Николенко СПбГУ— Санкт-Петербург 04 марта 2025 г.





#### Random facts:

- 4 марта 1628 г. королевскую хартию получила колония Массачусетского залива, 4 марта 1681 г. Уильям Пенн получил хартию на будущую Пенсильванию, а уже 4 марта 1789 г. первое заседание Конгресса США в Нью-Йорке ввело в действие конституцию США
- 4 марта 1733 г. был издан указ Анны Иоанновны «Об учреждении полиции в городах», согласно которому в 23 городах России были созданы полицмейстерские конторы
- · 4 марта 1762 г. император Пётр III подписал указ «О свободной для всех торговле»
- 4 марта 1803 г. император Александр I издал «Указ о вольных хлебопашцах», по которому землевладельцы могли освобождать крестьян, наделяя их землёй
- · 4 марта 1837 г. Михаил Лермонтов был арестован за стихотворение «Смерть поэта»
- 4 марта 1990 г. на 5-летний срок был избран Съезд народных депутатов РСФСР, высший законодательный орган РСФСР; он был распущен указом Бориса Ельцина 21 сентября 1993 г. и разогнан вооружённой силой 4 октября 1993 г.
- 4 марта 2012 г. прошли выборы президента Российской Федерации; Владимир Путин был в третий раз избран президентом России

Свойства	И	РАСШИРЕНИЯ	EM

### $\mathsf{T}$ РЕБОВАНИЯ КX И Y

- Что требуется, чтобы ЕМ-алгоритм работал?
- Неформально нужно, чтобы  $p(X\mid\theta)$  было легко максимизировать.
- А формально нужно, чтобы  $p(\mathbf{y} \mid \mathbf{x}, \theta) = p(\mathbf{y} \mid \mathbf{x})$ , т.е. чтобы выполнялось марковское свойство для  $\theta \to \mathbf{x} \to \mathbf{y}$ .
- На самом деле обычно ЕМ применяется тогда, когда  $\mathbf{y} = f(\mathbf{x})$  для детерминированной функции f, и это свойство тривиально выполняется.
- Более того, обычно ЕМ применяется, когда f это просто проекция, т.е. когда  $\mathbf{x}=(\mathbf{y},\mathbf{z})$ , как мы изначально и рассматривали.

# Разложение Q-функции по точкам

• Важное простое свойство — если данные состоят из независимо порождённых  $X=\{{f x}_1,\dots,{f x}_N\}$  (как всегда и бывает), то

$$\begin{split} &Q(\theta \mid \boldsymbol{\theta}^{(m)}) = \mathbb{E}_{X \mid Y, \boldsymbol{\theta}^{(m)}} \left[ \log \prod_{n=1}^{N} p(\mathbf{x}_n \mid \boldsymbol{\theta}) \right] = \\ &= \mathbb{E}_{X \mid Y, \boldsymbol{\theta}^{(m)}} \left[ \sum_{n=1}^{N} \log p(\mathbf{x}_n \mid \boldsymbol{\theta}) \right] = \sum_{n=1}^{N} \mathbb{E}_{\mathbf{x}_n \mid \mathbf{y}_n, \boldsymbol{\theta}^{(m)}} \left[ \log p(\mathbf{x}_n \mid \boldsymbol{\theta}) \right], \end{split}$$

потому что  $p(\mathbf{x}_n \mid Y, \theta) = p(\mathbf{x}_n \mid \mathbf{y}_n, \theta)$ 

• Упражнение: докажите это!

/ı

### GENERALIZED EM

- Другие обобщения могут пригодиться, если всё-таки подсчитать  $\mathbb{E}_{X\mid Y, \theta^{(m)}} \left[\log p(X\mid \theta) 
  ight]$  или оптимизировать  $p(X\mid \theta)$  нелегко.
- · Обобщённый EM (Generalized EM, GEM): вместо  $\arg\max_{\theta}Q(\theta,\theta^{(m)})$  нам достаточно просто выбирать такую  $\theta^{(m+1)}$ , чтобы

$$Q\left(\boldsymbol{\theta}^{(m+1)}, \boldsymbol{\theta}^{(m)}\right) > Q\left(\boldsymbol{\theta}^{(m)}, \boldsymbol{\theta}^{(m)}\right).$$

## STOCHASTIC EM

• Стохастический EM (Stochastic EM): если Q-функцию не получается посчитать в замкнутой форме, но и просто максимум брать не хочется, как в Classification EM, можно попробовать брать  $\mathbf{x}$  случайным образом на E-шаге:

$$\mathbf{x}^{(m)} \sim p(\mathbf{x} \mid \mathbf{y}, \boldsymbol{\theta}^{(m)}),$$

а потом использовать его на М-шаге как обычно:

$$\boldsymbol{\theta}^{(m+1)} = \arg \max_{\boldsymbol{\theta}} p(\mathbf{x}^{(m)} \mid \boldsymbol{\theta}^{(m)}).$$

### Monte Carlo EM

 Монте-Карло EM (Monte Carlo EM): саму Q-функцию тоже можно попытаться приблизить, если подсчитать сложно; можно использовать приближение ожидания в Q-функции через сэмплирование:

$$Q(\theta \mid \boldsymbol{\theta}^{(m)}) \approx \frac{1}{R} \sum_{r=1}^R \log p(\mathbf{x}^{(m,r)} \mid \boldsymbol{\theta}), \text{ где } \mathbf{x}^{(m,r)} \sim p(\mathbf{x} \mid \mathbf{y}, \boldsymbol{\theta}^{(m)}).$$

• Впрочем, в таких случаях часто можно вообще забить на EM и аппроксимировать напрямую апостериорное распределение.

## ЕМ С АПРИОРНЫМ РАСПРЕДЕЛЕНИЕМ

• А можно и априорное распределение в EM добавить, конечно же:

$$\theta_{MAP} = \arg \max_{\theta} \log p(\theta \mid Z) = \arg \max_{\theta} \left( \log p(Y \mid \theta) + \log p(\theta) \right).$$

• При этом базовая схема особо не меняется:

$$\begin{split} Q(\theta \mid \boldsymbol{\theta}^{(m)}) &= \mathbb{E}_{X \mid Y, \boldsymbol{\theta}^{(m)}} \left[ \log p(X \mid \boldsymbol{\theta}) \right], \\ \boldsymbol{\theta}^{(m+1)} &= \arg \max_{\boldsymbol{\theta}} \left( Q(\boldsymbol{\theta} \mid \boldsymbol{\theta}^{(m)}) + \log p(\boldsymbol{\theta})) \right). \end{split}$$

• Например, так можно избежать вырожденных случаев (кластер из одной точки).

### SEMI-SUPERVISED CLUSTERING

- И ЕМ, и k-means хорошо обобщаются на случай частично обученных кластеров.
- То есть про часть точек уже известно, какому кластеру они принадлежат.
- Как это учесть?

### SEMI-SUPERVISED CLUSTERING

- И ЕМ, и k-means хорошо обобщаются на случай частично обученных кластеров.
- То есть про часть точек уже известно, какому кластеру они принадлежат.
- Как это учесть?
- Чтобы учесть информацию о точке  $\mathbf{x}_i$ , достаточно для ЕМ положить скрытую переменную  $g_{nc}$  равной тому кластеру, которому нужно, с вероятностью 1, а остальным с вероятностью 0, и не пересчитывать.
- $\cdot$  Для k-means то же самое, но для  ${
  m clust}_i.$

ПРИМЕР: МОДЕЛИ БРЭДЛИ--ТЕРРИ ——

- Другой подход к рейтинг-системам модели Брэдли-Терри (Bradley-Terry).
- Модель предполагает, что для участников  $1,\dots,n$  можно подобрать такие рейтинги  $\gamma_i,\,i=1..n$ , что вероятность победы участника i над участником j равна

$$p(i \; \text{побеждает} \; j) = \frac{\gamma_i}{\gamma_i + \gamma_j}.$$

• Основная задача заключается в том, чтобы найти  $\gamma = (\gamma_1, \dots, \gamma_m)$  максимального правдоподобия из имеющихся данных D.

• Если принять априорное распределение равномерным, можно просто максимизировать правдоподобие

$$p(D|\gamma) = \prod_{i=1}^{m} \prod_{j=1}^{m} \left(\frac{\gamma_i}{\gamma_i + \gamma_j}\right)^{w_{ij}},$$

где  $w_{ij}$  — то, сколько раз  $x_i$  обыграл  $x_j$  при их попарном сравнении ( $w_{ii}=0$  по определению).

• Взяв логарифм, будем максимизировать

$$l(\gamma) = \sum_{i=1}^m \sum_{j=1}^m \left( w_{ij} \log \gamma_i - w_{ij} \log (\gamma_i + \gamma_j) \right).$$

- Максимизировать будем классическим ММ-алгоритмом (minorization-maximization), фактически вариационным приближением.
- Заметим, что

$$1 + \log \frac{x}{y} - \frac{x}{y} \le 0.$$

**Упражнение.** Докажите это.

• Рассмотрим вспомогательную функцию

$$Q(\gamma, \gamma^{(k)}) = \sum_{i,j} w_{ij} \left[ \log \gamma_i - \frac{\gamma_i + \gamma_j}{\gamma_i^{(k)} + \gamma_j^{(k)}} - \log \left( \gamma_i^{(k)} + \gamma_j^{(k)} \right) + 1 \right].$$

**Упражнение.** Используя предыдущее неравенство, докажите, что  $Q(\gamma, \gamma^{(k)}) \leq l(\gamma).$ 

• Чтобы найти  $\max_{\gamma}Q\left(\gamma,\gamma^{(k)}\right)$ , можно просто взять производные  $\frac{\partial Q}{\partial \gamma}$ :

$$\begin{split} \frac{\partial Q}{\partial \gamma_{l}} &= \sum_{i,j} w_{ij} \left[ \frac{\delta_{il}}{\gamma_{i}} - \frac{\delta_{il} + \delta_{jl}}{\gamma_{i}^{(k)} + \gamma_{j}^{(k)}} \right] = \\ &= \frac{1}{\gamma_{l}} \sum_{j} w_{lj} - \sum_{j} \frac{w_{lj}}{\gamma_{i}^{(k)} + \gamma_{j}^{(k)}} - \sum_{i} \frac{w_{il}}{\gamma_{i}^{(k)} + \gamma_{j}^{(k)}}. \end{split}$$

• Если  $w_l$  — общее количество побед игрока l ( $w_l = \sum_j w_{lj}$ ), и  $N_{ij}$  — количество встреч между игроками i и j ( $N_{ij} = w_{ij} + w_{ji}$ ), получаем

$$\frac{w_l}{\gamma_l} - \sum_j \frac{N_{lj}}{\gamma_i^{(k)} + \gamma_j^{(k)}} = 0.$$

• В результате правило пересчёта на одной итерации выглядит так:

$$\gamma_l^{(k+1)} := w_l \left[ \sum_j \frac{N_{lj}}{\gamma_i^{(k)} + \gamma_j^{(k)}} \right]^{-1}.$$

- Получили алгоритм оценки рейтингов. Правда, он пока работает только для ситуации, когда игроки встречаются один на один и выигрывают или проигрывают.
- Но даже в шахматах бывают ничьи. Как их учесть?

• Если возможны ничьи, их вероятность можно описать дополнительным параметром  $\theta>1$ , и это приводит к модели, в которой

$$p(i \text{ побеждает } j) = \frac{\gamma_i}{\gamma_i + \theta \gamma_j},$$
 
$$p(j \text{ побеждает } i) = \frac{\gamma_j}{\theta \gamma_i + \gamma_j},$$
 
$$p(i \text{ и } j \text{ играют вничью}) = \frac{(\theta^2 - 1)\gamma_i \gamma_j}{\left(\gamma_i + \theta \gamma_j\right) \left(\theta \gamma_i + \gamma_j\right)}.$$

- Аналогично можно вводить другие обобщения.
- Например, если результат может зависеть от порядка элементов в паре (скажем, команды проводят «домашние» матчи и «гостевые»), можно ввести дополнительный параметр θ, характеризующий, насколько большое преимущество дают «родные стены», и рассмотреть модель с

$$p(i \text{ побеждает } j) = \begin{cases} \frac{\theta \gamma_i}{\theta \gamma_i + \gamma_j}, & \text{ если } i \text{ играет дома,} \\ \frac{\gamma_j}{\theta \gamma_i + \gamma_j}, & \text{ если } j \text{ играет дома.} \end{cases}$$

• Можно даже обобщить на случай, когда в одном турнире встречаются несколько игроков: пусть перестановка  $\pi$  подмножества игроков  $A=\{1,\dots,k\}$  (результат турнира) имеет вероятность

$$p_A(\pi) = \prod_{i=1}^k \frac{\gamma_{\pi(i)}}{\gamma_{\pi(i)} + \gamma_{\pi(i+1)} + \ldots + \gamma_{\pi(k)}}.$$

• Можно показать, что такая модель эквивалентна весьма естественной «аксиоме Люса»: для любой модели, в которой вероятности игроков обыграть друг друга в любой паре не равны нулю, для любых подмножеств игроков  $A\subset B$  и любого игрока  $i\in A$ 

```
p_B(i побеждает) = = p_A(i \; {
m nofeexdaet}) p_B({
m nofeexdaet} \; {
m кто-то} \; {
m из} \; {
m множества} \; A).
```

• Но для крупных турниров это перестаёт работать; и совсем трудно что-то осмысленное сделать, если игроки соревнуются не поодиночке, а в командах.

Пример EM: presence-only data

### PRESENCE-ONLY DATA

- Пример из экологии: пусть мы хотим оценить, где водятся те или иные животные.
- Как определить, что суслики тут водятся, понятно: видишь суслика значит, он есть.
- Но как определить, что суслика нет? Может быть, ты не видишь суслика, и я не вижу, а он есть?..



### PRESENCE-ONLY DATA

• Формально говоря, есть переменные  ${f x}$ , определяющие некий регион (квадрат на карте), и мы моделируем вероятность того, что нужный вид тут есть,  $p(y=1\mid {f x})$ , при помощи логит-функции:

$$p(y = 1 \mid \mathbf{x}) = \sigma(\eta(\mathbf{x})) = \frac{1}{1 + e^{-\eta(\mathbf{x})}},$$

где  $\eta(\mathbf{x})$  может быть линейной (тогда получится логистическая регрессия), но может, в принципе, и не быть.

- Заметим, что даже если бы мы знали настоящие y, это было бы ещё не всё: сэмплирование положительных и отрицательных примеров неравномерно, перекошено в пользу положительных.
- Важное замечание: prospective vs. retrospective studies, case-control studies.

### Presence-only data

• Это значит, что ещё есть пропорции сэмплирования (sampling rates)

$$\gamma_0 = p(s=1 \mid y=0), \quad \gamma_1 = p(s=1 \mid y=1),$$

т.е. вероятности взять в выборку положительный/отрицательный пример.

• Их можно оценить как

$$\gamma_0 = \frac{n_0}{(1-\pi)N}, \quad \gamma_1 = \frac{n_1}{\pi N},$$

где  $\pi$  — истинная доля положительных примеров (встречаемость, occurrence).

 $\cdot$  Кстати, эту  $\pi$  было бы очень неплохо оценить в итоге.

### PRESENCE-ONLY DATA

• Тогда, если знать истинные значения всех y, то в принципе можно обучить:

$$\begin{split} p(y=1\mid s=1,\mathbf{x}) &= \\ &= \frac{p(s=1\mid y=1,\mathbf{x})p(y=1\mid \mathbf{x})}{p(s=1\mid y=0,\mathbf{x})p(y=0\mid \mathbf{x}) + p(s=1\mid y=1,\mathbf{x})p(y=1\mid \mathbf{x})} = \\ &= \frac{\gamma_1 e^{\eta(\mathbf{x})}}{\gamma_0 + \gamma_1 e^{\eta(\mathbf{x})}} = \frac{e^{\eta^*(\mathbf{x})}}{1 + e^{\eta^*(\mathbf{x})}}, \end{split}$$
 где  $\eta^*(\mathbf{x}) = \eta(\mathbf{x}) + \log(\gamma_1/\gamma_0)$ , т.е. 
$$\eta^*(\mathbf{x}) = \eta(\mathbf{x}) + \log\left(\frac{n_1}{n_0}\right) - \log\left(\frac{\pi}{1-\pi}\right). \end{split}$$

#### Presence-only data

- Таким образом, если  $\pi$  неизвестно, то  $\eta(\mathbf{x})$  можно найти с точностью до константы.
- $\cdot$  А у нас не  $n_0$  и  $n_1$ , а naive presence  $n_p$  и background  $n_u$ , т.е.

$$\begin{split} &p(y=1\mid s=1) = \frac{n_p + \pi n_u}{n_p + n_u}, \quad p(y=1\mid s=0) = \frac{(1-\pi)n_u}{n_p + n_u}, \\ &\gamma_1 = \frac{p(y=1)s=1)p(s=1)}{p(y=1)} = \frac{n_p + \pi n_u}{\pi(n_p + n_u)} p(s=1), \\ &\gamma_0 = \frac{p(y=0|s=1)p(s=1)}{p(y=0)} = \frac{n_u}{n_p + n_u} p(s=1), \end{split}$$

и в нашей модели  $\log \frac{n_1}{n_0} = \log \frac{n_p+\pi n_u}{\pi n_u}$ , т.е. всё как раньше, но  $n_1=n_p+\pi n_u$ ,  $n_0=(1-\pi)n_u$ .

• Обучать по y можно так: обучить модель, а потом вычесть из  $\eta(\mathbf{x})$  константу  $\log \frac{n_p + \pi n_u}{\pi n_u}$ .

### Presence-only data

- Но у нас нет настоящих данных y, чтобы обучить регрессию, а есть только presence-only z: если z=1, то y=1, но если z=0, то неизвестно, чему равен y.
- · (Ward et al., 2009): давайте использовать ЕМ. Правдоподобие:

$$\begin{split} \mathcal{L}(\eta \mid \mathbf{y}, \mathbf{z}, X) &= \prod_i p(y_i, z_i \mid s_i = 1, \mathbf{x}_i) = \\ &= \prod_i p(y_i \mid s_i = 1, \mathbf{x}_i) p(z_i \mid y_i, s_i = 1, \mathbf{x}_i). \end{split}$$

• В нашем случае при  $n_p$  положительных примеров и  $n_u$  фоновых (неизвестных)

$$\begin{array}{lcl} p(z_i=0 \mid y_i=0, s_i=1, \mathbf{x}_i) & = & 1, \\ p(z_i=1 \mid y_i=1, s_i=1, \mathbf{x}_i) & = & \frac{n_p}{n_p+\pi n_u}, \\ p(z_i=0 \mid y_i=1, s_i=1, \mathbf{x}_i) & = & \frac{\pi n_u}{n_p+\pi n_u}. \end{array}$$

#### PRESENCE-ONLY DATA

• А максимизировать нам надо сложное правдоподобие, в котором значения **у** неизвестны:

$$\begin{split} &L(\boldsymbol{\eta} \mid \mathbf{z}, \boldsymbol{X}) = \prod_{i} p(z_{i} \mid s_{i} = 1, \mathbf{x}) = \\ &= \prod_{i} \left( \frac{\frac{n_{p}}{\pi n_{u}} e^{\eta(\mathbf{x}_{i})}}{1 + \left(1 + \frac{n_{p}}{\pi n_{u}}\right) e^{\eta(\mathbf{x}_{i})}} \right)^{z_{i}} \left( \frac{1 + e^{\eta(\mathbf{x}_{i})}}{1 + \left(1 + \frac{n_{p}}{\pi n_{u}}\right) e^{\eta(\mathbf{x}_{i})}} \right)^{1 - z_{i}}. \end{split}$$

• Для этого и нужен ЕМ.

### Presence-only data

 $\cdot$  Е-шаг здесь в том, чтобы заменить  $y_i$  на его оценку

$$\hat{y}_i^{(k)} = \mathbb{E}\left[y_i \mid \eta^{(k)}\right] = \frac{e^{\eta^{(k)}} + 1}{1 + e^{\eta^{(k)}} + 1}.$$

• М-шаг мы уже видели, это обучение параметров логистической модели с целевой переменной  $\mathbf{y}^{(k)}$  на данных X.

- (1) Chose initial estimates:  $\hat{y}_i^{(0)} = \pi$  for  $z_i = 0$ .
- (2) Repeat until convergence:
  - Maximization step:
    - Calculate  $\hat{\eta}^{*(k)}$  by fitting a logistic model of  $\hat{\mathbf{y}}^{(k-1)}$  given X.

- Calculate 
$$\hat{\eta}^{(k)} = \hat{\eta}^{*(k)} - \log\left(\frac{n_p + \pi n_u}{\pi n_u}\right)$$
.

Expectation step:

$$\hat{y}_i^{(k)} = \frac{e^{\hat{\eta}^{(k)}}}{1 + e^{\hat{\eta}^{(k)}}} \quad \text{for } z_i = 0 \qquad \text{and} \qquad \hat{y}_i^{(k)} = 1 \quad \text{for } z_i = 1$$

### Presence-only data

- У Ward et al. получалось хорошо, но тут вышел любопытный спор.
- Ward et al. писали так: хотелось бы, чтобы можно было оценить  $\pi$ , но " $\pi$  is identifiable only if we make unrealistic assumptions about the structure of  $\eta(\mathbf{x})$  such as in logistic regression where  $\eta(\mathbf{x})$  is linear in  $\mathbf{x}$ :  $\eta(\mathbf{x}) = \mathbf{x}^{\top} \beta$ ".
- Через пару лет вышла статья Royle et al. (2012), тоже очень цитируемая, в которой говорилось: "logistic regression... is hardly unrealistic... such models are the most common approach to modeling binary variables in ecology (and probably all of statistics)... the logistic functions... is customarily adopted and widely used, and even books have been written about it".

### PRESENCE-ONLY DATA

- $\cdot$  Они предложили процедуру для оценки встречаемости  $\pi$ :
  - $\cdot$  для признаков  $\mathbf x$  у нас  $p(y=1\mid \mathbf x)=rac{p(y=1)\pi_1(\mathbf x)}{p(y=1)\pi_1(\mathbf x)+(1-p(y=1))\pi_0(\mathbf x)}$ ;
  - $\cdot$  данные это выборка из  $\pi_1(\mathbf{x})$  и отдельно выборка из  $\pi(\mathbf{x})$ ;
  - $\cdot$  как видно, даже если знать  $\pi$  и  $\pi_1$  полностью, остаётся свобода: надо оценить  $p(y=1\mid \mathbf{x})$ , а  $\pi_0(\mathbf{x})$  мы не знаем;
  - · Royle et al. вводят предположения для  $p(y=1\mid \mathbf{x})$  в виде логистической регрессии:  $p(y=1\mid \mathbf{x})=\sigma(\beta^{\top}\mathbf{x});$
  - · тогда действительно можно записать  $p(y=1)\pi_1(\mathbf{x})=p(y=1\mid\mathbf{x})\pi(\mathbf{x})=p(y=1,\mathbf{x})\text{, и если }\pi(\mathbf{x})$  равномерно (это логично), то

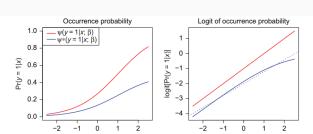
$$\pi_1(\mathbf{x}_i) = \frac{p(y_i = 1 \mid \mathbf{x}_i)}{\sum_{\mathbf{x}} p(y = 1 \mid \mathbf{x})};$$

если подставить сюда логистическую регрессию, то можно оценить  $\beta$  максимального правдоподобия, и не надо знать p(y=1)!

• Что тут не так?

### Presence-only data

- Всё так, но предположение о логистической регрессии здесь выполняет слишком много работы.
- Если рассмотреть две кривые, у которых  $p^*(y=1\mid \mathbf{x},\beta)=\tfrac{1}{2}p(y=1\mid \mathbf{x},\beta)\text{, то у них будет}\\p^*(y=1)=\tfrac{1}{2}p(y=1)\text{, но общее правдоподобие }\pi_1(\mathbf{x}_i)\text{ будет в точности одинаковое, }\tfrac{1}{2}\text{ сократится.}$
- Дело в том, что модель  $p^*$  не будет логистической регрессией, с  $\beta$  произойдёт что-то нелинейное; но откуда у нас настолько сильное предположение? Как отличить синюю кривую справа от пунктирной прямой?



Пример: РЕЙТИНГ СПОРТИВНОГО ЧГК

## Рейтинг спортивного ЧГК

- Пример из практики: в какой-то момент я хотел сделать рейтинг спортивного «Что? Где? Когда?»:
  - · участвуют команды по  $\leq 6$  человек, причём часто встречаются неполные команды;
  - игроки постоянно переходят между командами (поэтому TrueSkill);
  - в одном турнире могут участвовать до тысячи команд (синхронные турниры);
  - командам задаётся фиксированное число вопросов (36, 60, 90), т.е. в крупных турнирах очень много команд делят одно и то же место.

## Рейтинг спортивного ЧГК

- · Первое решение система TrueSkill
- Расскажу о ней потом, когда будем говорить об Expectation Propagation
- У неё есть проблемы в постановке спортивного ЧГК, мы когда-то их отчасти решили, была сложная и интересная модель, она работала лучше базового TrueSkill
- · Ho...

# Рейтинг спортивного ЧГК

- В какой-то момент база турниров ЧГК стала собирать повопросные результаты.
- Мы теперь знаем, на какие именно вопросы ответила та или иная команда.
- Так что когда я вернулся к задаче построения рейтинга ЧГК, задача стала существенно проще.

## Рейтинг спортивного ЧГК

- Пример аналогичного приложения:
  - есть набор вопросов для теста из большого числа вопросов (например, IQ-тест или экзамен по какому-то предмету);
  - участники отвечают на случайное подмножество вопросов;
  - надо оценить участников, но уровень сложности вопросов нельзя заранее точно сбалансировать.
- «Что? Где? Когда?» это оно и есть, только теперь участники объединяются в команды.

#### BASELINE: ЛОГИСТИЧЕСКАЯ РЕГРЕССИЯ

- Baseline логистическая регрессия:
  - $\cdot$  каждый игрок i моделируется скиллом  $s_i$ ,
  - · каждый вопрос q моделируется сложностью («лёгкостью»)  $c_q$ ,
  - $\cdot$  добавим глобальное среднее  $\mu$ ,
  - обучим логистическую модель

$$p(x_{tq} \mid s_i, c_q) \sim \sigma(\mu + s_i + c_q)$$

для каждого игрока  $i\in t$  команды-участницы  $t\in T^{(d)}$  и каждого вопроса  $q\in Q^{(d)}$ , где  $\sigma(x)=1/(1+e^x)$  — логистический сигмоид,  $x_{tq}$  — ответила ли команда t на вопрос q.

## Модель со скрытыми переменными

- Логистическая модель предполагает фактически, что каждый игрок ответил на каждый вопрос, который взяла команда.
- Это неправда, мы не знаем, кто ответил, только знаем, что кто-то это сделал; а если команда не ответила, то никто.
- Это по идее похоже на presence-only data models (Ward et al., 2009; Royle et al., 2012).

## Модель со скрытыми переменными

- Поэтому давайте сделаем модель со скрытыми переменными.
- Для каждой пары игрок-вопрос, добавим переменную  $z_{iq}$ , которая означает, что «игрок i ответил на вопрос q».
- На эти переменные есть такие ограничения:
  - $\cdot$  если  $x_{ta}=0$ , то  $z_{ia}=0$  для каждого игрока  $i\in t$ ;
  - $\cdot\,\,$  если  $x_{tq}=1$ , то  $z_{iq}=1$  для по крайней мере одного игрока  $i\in t.$

## Модель со скрытыми переменными

• Параметры модели те же — скилл и сложность вопросов:

$$p(z_{iq} \mid s_i, c_q) \sim \sigma(\mu + s_i + c_q).$$

- Обучаем ЕМ-алгоритмом:
  - · Е-шаг: зафиксируем все  $s_i$  и  $c_q$ , вычислим ожидания скрытых переменных  $z_{iq}$  как

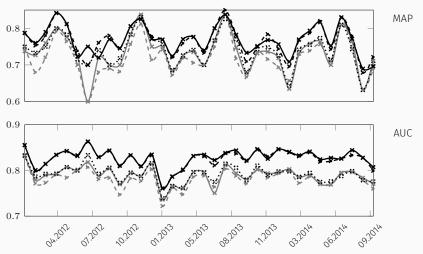
$$\mathbb{E}\left[z_{iq}\right] = \begin{cases} 0, & \text{если } x_{tq} = 0, \\ p(z_{iq} = 1 \mid \exists j \in t \ z_{jq} = 1) = \frac{\sigma(\mu + s_i + c_q)}{1 - \prod_{j \in t} \left(1 - \sigma(\mu + s_j + c_q)\right)}, & \text{если } x_{tq} = 1; \end{cases}$$

· М-шаг: зафиксируем  $\mathbb{E}\left[z_{iq}
ight]$ , обучим логистическую модель

$$\mathbb{E}\left[z_{iq}\right] \sim \sigma(\mu + s_i + c_q).$$

## **РЕЗУЛЬТАТЫ**

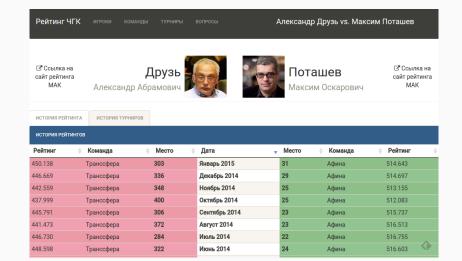
• Ну и, конечно, работает хорошо.



# ПРИМЕР

Рейтинг	•ЧГК	<b>ИГРОКИ</b> КОМА!	нды турнирь	І ВОПРОСЫ	Рейтинг-лист иг	роков	РЕЛИЗ: ЯІ	НВАРЬ 2015
РЕЙТИНГ-ЛИСТ ИГРОКОВ НА ЯНВАРЬ 201S								
Показывать по 100 ▼ записей					Введите id или начало фамилии:			
Место 🛊	id 💠	Фамилия 💠	ф кмИ	Отчество ф	Команда ф	Сыграно 💠	Взято 🔅	Рейтинг 🔻
1	27177	Ромашова	Вероника	Михайловна	ЛКИ	5278	3784	545.753
2	3083	Белявский	Дмитрий	Михайлович	лки	4831	3396	543.520
3	27403	Руссо	Максим	Михайлович	лки	7180	5205	534.540
4	4270	Брутер	Александра	Владимировна	лки	8244	5974	533.405
5	18332	Либер	Александр	Витальевич	Рабочее название	8658	6210	532.921
6	1585	Архангельская	Юлия	Сергеевна	Ксеп	7542	5286	531.866
7	24384	Пашковский	Евгений	Александрович	ЛКИ	7552	5374	531.327
8	8333	Губанов	Антон	Александрович	Команда Губанова	4475	3153	530.871
9	16332	Крапиль	Николай	Валерьевич	Ксеп	6927	4899	530.559
10	21487	Моносов	Борис	Яковлевич	Команла Губанова	5115	3645	530.175

#### ПРИМЕР





## Марковские цепи

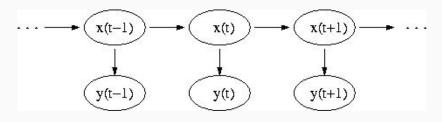
- Марковская цепь задаётся начальным распределением вероятностей  $p^0(x)$  и вероятностями перехода  $T(x^\prime;x)$ .
- · T(x';x) это распределение следующего элемента цепи в зависимости от следующего; распределение на (t+1)-м шаге равно

$$p^{t+1}(x') = \int T(x';x)p^t(x)dx.$$

• В дискретном случае T(x';x) — это матрица вероятностей p(x'=i|x=j).

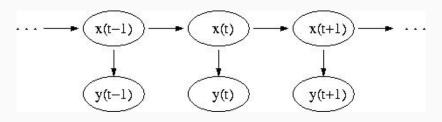
# Дискретные марковские цепи

- Мы будем находиться в дискретном случае.
- Марковская модель это когда мы можем наблюдать какие-то функции от марковского процесса.



# Дискретные марковские цепи

- $\cdot$  Здесь x(t) сам процесс (модель), а y(t) то, что мы наблюдаем.
- Задача определить скрытые параметры процесса.



# Дискретные марковские цепи

• Главное свойство — следующее состояние не зависит от истории, только от предыдущего состояния.

$$\begin{split} p(x(t) = x_j | x(t-1) = x_{j_{t-1}}, \dots, x(1) = x_{j_1}) = \\ &= p(x(t) = x_j | x(t-1) = x_{j_{t-1}}). \end{split}$$

- · Более того, эти вероятности  $a_{ij}=p(x(t)=x_j|x(t-1)=x_i)$  ещё и от времени t не зависят.
- $\cdot$  Эти вероятности и составляют матрицу перехода  $A=(a_{ij}).$

# Вероятности перехода

- Естественные свойства:
- ·  $a_{ij} \geq 0$ .
- $\sum_{j} a_{ij} = 1$ .

# ПРЯМАЯ ЗАДАЧА

- Естественная задача: с какой вероятностью выпадет та или иная последовательность событий?
- Т.е. найти нужно для последовательности  $Q=q_{i_1}\dots q_{i_k}$

$$p(Q|{\rm MOДЕЛЬ}) = p(q_{i_1})p(q_{i_2}|q_{i_1})\dots p(q_{i_k}|q_{i_{k-1}}).$$

- Казалось бы, это тривиально.
- Что же сложного в реальных задачах?

## Скрытые марковские модели

- А сложно то, что никто нам не скажет, что модель должна быть именно такой.
- И, кроме того, мы обычно наблюдаем не x(t), т.е. реальные состояния модели, а y(t), т.е. некоторую функцию от них (данные).
- Пример: распознавание речи.

# Спасибо за внимание!



