

# ЭКСПОНЕНЦИАЛЬНОЕ СЕМЕЙСТВО

---

Сергей Николенко

СПбГУ — Санкт-Петербург

22 ноября 2022 г.

---

*Random facts:*

- 22 ноября 1307 г. Климент V издал буллу «Pastoralis Praeeminentiae», в которой утверждал, что обвинения против ордена тамплиеров доказаны, и призывал государей Европы последовать примеру Филиппа Красивого
- 22 ноября 1604 г. король Испании Филипп III отдыхал у камина, и пламя сильно разгорелось; пока специального дворянина, имевшего право затворить заслонку или отодвинуть кресло короля, искали по всему дворцу, Филипп ждал не двигаясь с места и получил серьёзные ожоги (это, конечно, анекдот)
- 22 ноября 1621 г. настоятелем лондонского собора св. Павла был назначен Джон Донн; с тех пор он мог сам выбирать, по ком звонит колокол
- 22 ноября 1641 г. английский парламент принял «Великую ремонстрацию» (The Grand Remonstrance) — перечень из 204 статей, исчисляющих злоупотребления короля
- 22 ноября 1977 г. начались полёты «Конкорда» через Атлантику: из Лондона в Нью-Йорк можно было добраться менее чем за три с половиной часа, билет «туда-обратно» стоил \$10500

# ЭКСПОНЕНЦИАЛЬНОЕ СЕМЕЙСТВО

- Мы видели общий паттерн: найти правдоподобие, посмотреть на его форму и догадаться, как должно выглядеть семейство сопряжённых априорных распределений.
- Это выглядит как достаточно несложная процедура, которая должна обобщаться.
- *Экспоненциальное семейство* распределений (exponential family): параметрическое семейство распределений принадлежит экспоненциальному семейству, если оно имеет вид

$$p(\mathbf{x}|\theta) = h(\mathbf{x})e^{\eta(\theta)^\top \mathbf{t}(\mathbf{x}) - a(\theta)} = h(\mathbf{x})g(\theta)e^{\eta(\theta)^\top \mathbf{t}(\mathbf{x})}$$

для некоторого параметра  $\theta$ ; здесь  $g(\theta) = e^{-a(\theta)}$ .

- Векторная функция  $\mathbf{t}(\mathbf{x})$  выделяет *достаточные статистики* (sufficient statistics), и она играет роль извлечения признаков из  $\mathbf{x}$ .

# ЭКСПОНЕНЦИАЛЬНОЕ СЕМЕЙСТВО

- Если  $\eta(\theta) = \theta$ , то такая параметризация называется *естественной*, а  $\theta$  в таком случае называется *естественным параметром* (natural parameter):

$$p(\mathbf{x}|\theta) = h(\mathbf{x})e^{\theta^\top \mathbf{t}(\mathbf{x}) - a(\theta)} = h(\mathbf{x})g(\theta)e^{\theta^\top \mathbf{t}(\mathbf{x})}.$$

- Определение выглядит очень общим; главное предположение здесь в том, как  $\theta$  и  $\mathbf{x}$  разделяются в этом определении: в экспоненте они связаны друг с другом линейно, а вне экспоненты полностью разнесены по функциям  $h(\mathbf{x})$  и  $g(\theta)$ , то есть единственная зависимость между  $\mathbf{x}$  и  $\theta$  — это скалярное произведение в экспоненте.
- Вообще говоря, почти всё, о чём мы говорили – частные случаи *экспоненциального семейства* распределений.

# ЭКСПОНЕНЦИАЛЬНОЕ СЕМЕЙСТВО

- Например, биномиальное распределение

$$\begin{aligned}\text{Binom}(k|n, p) &= \binom{n}{k} p^k (1-p)^{n-k} = \\ &= \binom{n}{k} e^{k \log p + (n-k) \log(1-p)} = \binom{n}{k} e^{k \log \frac{p}{1-p} + n \log(1-p)}.\end{aligned}$$

- В итоге получается, что биномиальное распределение принадлежит экспоненциальному семейству, и его естественный параметр — это

$$\theta = \log \frac{p}{1-p}, \quad p = \frac{e^\theta}{1+e^\theta},$$

то есть в точности те самые log-odds;  $t(k) = k$ ,  $h(k) = \binom{n}{k}$ ,

$$a(\theta) = -n \log(1-p) = n \log(1+e^\theta), \quad g(\theta) = e^{n \log(1-p)} = (1+e^\theta)^{-n}.$$

- Аналогично, мультиномиальное распределение

$$\text{Mult}(\mathbf{x}|n, p_1, \dots, p_k) = \begin{cases} \frac{n!}{x_1!x_2!\dots x_k!} p_1^{x_1} p_2^{x_2} \dots p_k^{x_k}, & \text{если } \sum_{i=1}^k x_i = n, \\ 0 & \text{в противном случае,} \end{cases}$$

можно переписать как

$$\text{Mult}(\mathbf{x}|n, p_1, p_2, \dots, p_k) = \frac{n!}{x_1!x_2!\dots x_k!} e^{\sum_{i=1}^k x_i \log p_i},$$

то есть на первый взгляд кажется, что в экспоненциальном семействе здесь

$$\mathbf{t}(\mathbf{x}) = \mathbf{x}, \quad \theta = \log \mathbf{p}, \quad a(\theta) = 0, \quad h(\mathbf{x}) = \frac{n!}{x_1!x_2!\dots x_k!}.$$

# ЭКСПОНЕНЦИАЛЬНОЕ СЕМЕЙСТВО

- Но такое представление ведёт к техническим трудностям из-за того, что  $a(\theta) = 0$ , поэтому лучше выразить

$$\begin{aligned} e^{\sum_{i=1}^k x_i \log p_i} &= e^{\sum_{i=1}^{k-1} x_i \log p_i + (n - \sum_{i=1}^{k-1} x_i) \log(1 - \sum_{i=1}^{k-1} p_i)} = \\ &= e^{\sum_{i=1}^{k-1} x_i \log\left(\frac{p_i}{1 - \sum_{i=1}^{k-1} p_i}\right) + n \log(1 - \sum_{i=1}^{k-1} p_i)}. \end{aligned}$$

- Таким образом, в итоге  $\mathbf{t}(\mathbf{x}) = \mathbf{x}$ ,

$$\theta_i = \log\left(\frac{p_i}{1 - \sum_{i=1}^{k-1} p_i}\right) = \log \frac{p_i}{p_k}, \quad p_i = \frac{e^{\theta_i}}{\sum_{j=1}^k e^{\theta_j}},$$

и теперь  $a(\theta) = -n \log\left(1 - \sum_{i=1}^{k-1} p_i\right) = n \log\left(\sum_{j=1}^k e^{\theta_j}\right)$ .

- В обратном выражении для  $p_i$  через  $\theta$  у нас опять получилась как раз та самая softmax-функция.

- С распределением Пуассона совсем нет вопросов:

$$p(x|\lambda) = \frac{1}{x!} \lambda^x e^{-\lambda} = \frac{1}{x!} e^{x \log \lambda - \lambda}$$

сразу же принадлежит экспоненциальному семейству с  $t(x) = x$ ,  $\theta = \log \lambda$ ,  $h(x) = \frac{1}{x!}$ ,  $a(\theta) = \lambda = e^\theta$ .

- Редкий пример распределения, которое *не* принадлежит экспоненциальному семейству — это гипергеометрическое распределение

$$p(x|N, n, K) = \frac{1}{\binom{N}{n}} \binom{K}{x} \binom{N-K}{n-x};$$

его преобразовать к нужной форме никак не получится.

# ЭКСПОНЕНЦИАЛЬНОЕ СЕМЕЙСТВО

- Нормальное распределение:

$$\begin{aligned} N(x|\mu, \sigma^2) &= \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2} = \\ &= \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}x^2 + \frac{\mu}{\sigma^2}x - \frac{\mu^2}{2\sigma^2}} = \frac{1}{\sqrt{2\pi}} e^{\begin{pmatrix} x^2 \\ x \end{pmatrix}^\top \begin{pmatrix} -1/2\sigma^2 \\ \mu/\sigma^2 \end{pmatrix} - \frac{\mu^2}{2\sigma^2} - \log \sigma}. \end{aligned}$$

- Иначе говоря, одномерное нормальное распределение имеет две достаточные статистики,  $\mathbf{t}(x) = \begin{pmatrix} x^2 \\ x \end{pmatrix}$ , и естественный параметр размерности два:

$$\theta = \begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix} = \begin{pmatrix} -1/2\sigma^2 \\ \mu/\sigma^2 \end{pmatrix} = \begin{pmatrix} -\tau/2 \\ \mu\tau \end{pmatrix};$$

а остальные функции выглядят как  $h(x) = \frac{1}{\sqrt{2\pi}}$ ,

$$a(\theta) = \frac{\mu^2}{2\sigma^2} + \log \sigma = \frac{\mu^2\tau}{2} - \frac{1}{2} \log \tau = -\frac{\theta_2^2}{4\theta_1^2} - \frac{1}{2} \log(-2\theta_1).$$



# ЭКСПОНЕНЦИАЛЬНОЕ СЕМЕЙСТВО

- Многомерный гауссиан:

$$\begin{aligned} N(\mathbf{x}|\mu, \Sigma) &= \frac{1}{\sqrt{2\pi \det \Sigma}} e^{-\frac{1}{2}(\mathbf{x}-\mu)^\top \Sigma^{-1}(\mathbf{x}-\mu)} = \\ &= e^{-\frac{1}{2}(\mathbf{x}^\top \Sigma^{-1} \mathbf{x} - 2\mathbf{x}^\top \Sigma^{-1} \mu + \mu^\top \Sigma^{-1} \mu + \log(2\pi \det \Sigma))}. \end{aligned}$$

- Нужно представить  $\mathbf{x}^\top \Sigma^{-1} \mathbf{x}$  в виде скалярного произведения; здесь  $\text{vec}(A)$  обозначает разворачивание матрицы в плоский вектор:

$$\mathbf{x}^\top \Sigma^{-1} \mathbf{x} = \sum_{i,j=1}^d (\Sigma^{-1})_{ij} x_i x_j = \text{vec}(\mathbf{x}\mathbf{x}^\top)^\top \text{vec}(\Sigma^{-1}).$$

- В итоге  $h(\mathbf{x}) = 1$ ,  $\mathbf{t}(\mathbf{x}) = \begin{pmatrix} \text{vec}(\mathbf{x}\mathbf{x}^\top) \\ \mathbf{x} \end{pmatrix}$ ,  $\theta = \begin{pmatrix} -\frac{1}{2} \text{vec}(\Sigma^{-1}) \\ \Sigma^{-1} \mu \end{pmatrix}$ ,

$$a(\theta) = \mu^\top \Sigma^{-1} \mu + \log(2\pi \det \Sigma).$$

# ЭКСПОНЕНЦИАЛЬНОЕ СЕМЕЙСТВО

- Теперь интересные результаты. Первый — о среднем и дисперсии распределений из экспоненциального семейства.
- Интеграл от любого распределения равен единице:

$$a(\theta) = \log \int h(\mathbf{x}) e^{\theta^\top \mathbf{t}(\mathbf{x})} d\mathbf{x}.$$

- Возьмём градиент по  $\theta$  слева и справа:

$$\begin{aligned} \nabla_{\theta} a(\theta) &= \nabla_{\theta} \log \int h(\mathbf{x}) e^{\theta^\top \mathbf{t}(\mathbf{x})} d\mathbf{x} = \frac{\int \nabla_{\theta} h(\mathbf{x}) e^{\theta^\top \mathbf{t}(\mathbf{x})} d\mathbf{x}}{\int h(\mathbf{x}) e^{\theta^\top \mathbf{t}(\mathbf{x})} d\mathbf{x}} = \\ &= \frac{\int h(\mathbf{x}) \nabla_{\theta} e^{\theta^\top \mathbf{t}(\mathbf{x})} d\mathbf{x}}{e^{a(\theta)}} = \frac{\int h(\mathbf{x}) e^{\theta^\top \mathbf{t}(\mathbf{x})} \mathbf{t}(\mathbf{x}) d\mathbf{x}}{e^{a(\theta)}} = \\ &= \int \mathbf{t}(\mathbf{x}) \left( h(\mathbf{x}) e^{\theta^\top \mathbf{t}(\mathbf{x}) - a(\theta)} \right) d\mathbf{x}. \end{aligned}$$

- Иначе говоря, мы видим, что  $\nabla_{\theta} a(\theta)$  — это математическое ожидание достаточных статистик исходного распределения:

$$\mathbb{E}[\mathbf{t}(\mathbf{x})] = \nabla_{\theta} a(\theta).$$

- Это очень мощный результат, который нередко приговждается и в машинном обучении.
- Функцию  $a(\theta)$  называют *кумулянтом* (cumulant).

- Кроме того, для минимальных представлений распределений из экспоненциального семейства, когда в  $\theta$  нет лишних параметров, это можно обратить, то есть выразить  $\theta$  через  $\mathbb{E}[\mathbf{t}(\mathbf{x})]$ .
- Здесь будет естественно рассмотреть *параметризацию средним* (mean parametrization), в которой параметром будет

$$\mu = \mathbb{E}[\mathbf{t}(\mathbf{x})] = \nabla_{\theta} a(\theta).$$

- Например, привычная нам параметризация гауссиана именно такова.

# ЭКСПОНЕНЦИАЛЬНОЕ СЕМЕЙСТВО

- Результат можно продолжить на другие моменты.
- Рассмотрим матрицу вторых производных  $a(\theta)$ :

$$\begin{aligned}\frac{\partial^2 a}{\partial \theta_i \partial \theta_j} &= \frac{\partial \mathbb{E}[t_i(\mathbf{x})]}{\partial \theta_j} = \frac{\partial}{\partial \theta_j} \int t_i(\mathbf{x}) h(\mathbf{x}) e^{\theta^\top \mathbf{t}(\mathbf{x}) - a(\theta)} d\mathbf{x} = \\ &= \int t_i(\mathbf{x}) h(\mathbf{x}) \frac{\partial e^{\theta^\top \mathbf{t}(\mathbf{x}) - a(\theta)}}{\partial \theta_j} d\mathbf{x} = \int t_i(\mathbf{x}) h(\mathbf{x}) e^{\theta^\top \mathbf{t}(\mathbf{x}) - a(\theta)} \left( t_j(\mathbf{x}) - \frac{\partial a(\theta)}{\partial \theta_j} \right) d\mathbf{x} = \\ &= \int t_i(\mathbf{x}) (t_j(\mathbf{x}) - \mathbb{E}[t_j(\mathbf{x})]) h(\mathbf{x}) e^{\theta^\top \mathbf{t}(\mathbf{x}) - a(\theta)} d\mathbf{x} = \\ &= \mathbb{E}[t_i(\mathbf{x}) t_j(\mathbf{x})] - \mathbb{E}[t_i(\mathbf{x})] \mathbb{E}[t_j(\mathbf{x})].\end{aligned}$$

- Иначе говоря, гессиан функции  $a(\theta)$  — это в точности матрица ковариаций вектора достаточных статистик  $\mathbf{t}(\mathbf{x})$ :

$$\text{Var}[\mathbf{t}(\mathbf{x})] = \mathbf{H}(a(\theta)) = \left( \frac{\partial^2 a}{\partial \theta_i \partial \theta_j} \right)_{i,j=1}^k.$$

- Аналогичные результаты верны и для других моментов.

# ЭКСПОНЕНЦИАЛЬНОЕ СЕМЕЙСТВО

- Второй результат – о правдоподобии:

$$p(\mathbf{x}_1, \dots, \mathbf{x}_N | \theta) = \prod_{n=1}^N h(\mathbf{x}_n) e^{\eta(\theta)^\top \mathbf{t}(\mathbf{x}_n) - a(\theta)},$$

$$\log p(\mathbf{x}_1, \dots, \mathbf{x}_N | \theta) = \sum_{n=1}^N \log h(\mathbf{x}_n) + \eta(\theta)^\top \left( \sum_{n=1}^N \mathbf{t}(\mathbf{x}_n) \right) - Na(\theta).$$

- Всё, что в логарифме правдоподобия зависит от  $\theta$ , содержит  $\mathbf{x}_n$  только в виде  $\sum_{n=1}^N \mathbf{t}(\mathbf{x}_n)$ .
- Т.е. достаточно сохранять суммы  $\sum_{n=1}^N \mathbf{t}(\mathbf{x}_n)$  и число  $N$ , а сами точки  $\mathbf{x}_n$  можно «забывать»; поэтому  $\mathbf{t}(\mathbf{x}_n)$  называются *достаточными статистиками*.
- Например, для одномерного гауссиана достаточно хранить  $\sum_{n=1}^N x_n$  и  $\sum_{n=1}^N x_n^2$ , и из них можно найти гипотезу максимального правдоподобия и для  $\mu$ , и для  $\sigma^2$ .

# ЭКСПОНЕНЦИАЛЬНОЕ СЕМЕЙСТВО

- И саму гипотезу максимального правдоподобия можно найти, взяв градиент по  $\theta$  и приравняв нулю:

$$\nabla_{\theta} a(\theta) = \frac{1}{N} \nabla_{\theta} \eta(\theta)^{\top} \left( \sum_{n=1}^N \mathbf{t}(\mathbf{x}_n) \right).$$

- Для естественного параметра  $\eta(\theta) = \theta$  правая часть не зависит от  $\theta$ . А при параметризации средним,  $\eta(\theta) = \theta$ , сразу

$$\mu = \nabla_{\theta} a(\theta) = \frac{1}{N} \sum_{n=1}^N \mathbf{t}(\mathbf{x}_n).$$

- *Теорема Купмана–Питмана–Дармуа* (Pitman–Koornman–Darmois): при некоторых условиях регулярности экспоненциальное семейство — это *единственное* семейство распределений с конечным набором достаточных статистик, т.е. с набором достаточных статистик, размер которого не зависит от  $N$ .

# ЭКСПОНЕНЦИАЛЬНОЕ СЕМЕЙСТВО

- Третий интересный результат – уже о байесовском выводе: оказывается, есть универсальный способ найти семейство сопряжённых априорных распределений для любого распределения из экспоненциального семейства.
- Для семейства распределений с параметром  $\theta$  семейством сопряжённых априорных распределений будет

$$p(\theta|\chi, \nu) = f(\chi, \nu)e^{x^T \eta(\theta) - \nu a(\theta)},$$

где  $\chi$  и  $\nu$  – гиперпараметры, функция  $a(\theta)$  та же самая, что в исходном распределении, а функция  $f(\chi, \nu)$  – это нормировочная константа.



# ЭКСПОНЕНЦИАЛЬНОЕ СЕМЕЙСТВО

- Давайте найдём апостериорное распределение:

$$p(\theta | \mathbf{x}_1, \dots, \mathbf{x}_N, \chi_0, \nu_0) \propto p(\theta | \chi_0, \nu_0) p(\mathbf{x}_1, \dots, \mathbf{x}_N | \theta).$$

- Подставим:

$$\begin{aligned} \log p(\theta | \mathbf{x}_1, \dots, \mathbf{x}_N, \chi_0, \nu_0) &= \sum_{n=1}^N \log h(\mathbf{x}_n) + \eta(\theta)^\top \left( \sum_{n=1}^N \mathbf{t}(\mathbf{x}_n) \right) - Na(\theta) + \\ &+ \log f(\chi_0, \nu_0) + \chi_0^\top \eta(\theta) - \nu_0 a(\theta) + \text{const} = \\ &= \text{const} + \left( \chi_0 + \sum_{n=1}^N \mathbf{t}(\mathbf{x}_n) \right)^\top \eta(\theta) - (N + \nu_0) a(\theta), \end{aligned}$$

и в итоге у нас получилось апостериорное распределение  $p(\theta | \chi_N, \nu_N)$  того же вида, но с параметрами

$$\nu_N = \nu_0 + N, \quad \chi_N = \chi_0 + \sum_{n=1}^N \mathbf{t}(\mathbf{x}_n).$$

## ЭКСПОНЕНЦИАЛЬНОЕ СЕМЕЙСТВО

- Снова видим, что для апостериорного распределения нужно знать только достаточные статистики  $\sum_{n=1}^N \mathbf{t}(\mathbf{x}_n)$ .
- Более того, можно и предсказательное распределение найти:

$$\begin{aligned} p(\mathbf{x} | \mathbf{x}_1, \dots, \mathbf{x}_n, \chi_0, \nu_0) &= \int p(\mathbf{x} | \theta) p(\theta | \mathbf{x}_1, \dots, \mathbf{x}_N, \chi_0, \nu_0) d\theta = \\ &= \int h(\mathbf{x}) e^{\mathbf{t}(\mathbf{x})^\top \eta(\theta) - a(\theta)} f(\chi_N, \nu_N) e^{\chi_N^\top \eta(\theta) - \nu_N a(\theta)} d\theta = \\ &= f(\chi_N, \nu_N) \int h(\mathbf{x}) e^{(\mathbf{t}(\mathbf{x}) + \chi_N)^\top \eta(\theta) - (\nu_N + 1)a(\theta)} d\theta = \frac{f(\chi_N, \nu_N)}{f(\chi_N + \mathbf{t}(\mathbf{x}), \nu_N + 1)}. \end{aligned}$$

- Получилось, что предсказательное распределение — это отношение нормировочных констант для сопряжённого априорного распределения с разными параметрами.
- Сравните это с выводом предсказательного распределения для испытаний Бернулли, который мы делали в начале курса.

Спасибо за внимание!