

# АПРИОРНЫЕ РАСПРЕДЕЛЕНИЯ ДЖЕФФРИСА

Сергей Николенко

СПбГУ — Санкт-Петербург

28 ноября 2024 г.

*Random facts:*



- 28 ноября 1443 г. Георгий Кастриоти, он же Скандербег, торжественно вступил в Крую и был провозглашён старшинами главой княжества Кастриоти и вождём всех албанцев
- 28 ноября 1660 г. двенадцать человек, в том числе Роберт Бойль, Кристофер Рен и сэры Роберт Морэй собрались в Грешем-колледже, первом высшем учебном заведении Лондона, и основали «Colledge for the Promoting of Physico-Mathematicall Experimentall Learning», который потом превратится в the Royal Society
- 28 ноября 1814 г. лондонская The Times стала первой газетой, производившейся на паровом печатном прессе, созданном немецкими инженерами Кёнигом и Бауэром
- 28 ноября 1918 г. в Эстонию вошли российские войска, части РККА заняли Нарву, и началась Освободительная Война (1918—1920)
- 28 ноября 1905 г. Артур Гриффит основал политическую партию под названием Sinn Féin
- 28 ноября 1990 г. Маргарет Тэтчер ушла в отставку как глава консервативной партии и, как следствие, как премьер-министр; её преемником на обеих позициях стал Джон Мейджор

# АПРИОРНЫЕ РАСПРЕДЕЛЕНИЯ ДЖЕФФРИСА

---

## АПРИОРНЫЕ РАСПРЕДЕЛЕНИЯ ДЖЕФФРИСА

- Мы уже говорили о неинформативных априорных распределениях
- Равномерное распределение для шаров, которые достают из урны, оказалось не таким уж неинформативным: оно на самом деле соответствует данным  $D = (2, 1)$ , в которых мы уже достали два шара и увидели, что ровно один из них красный
- А равномерное распределение на параметр монетки  $\theta$  можно рассматривать как результат «виртуального эксперимента» с двумя подбрасываниями; как обобщить эти рассуждения?
- Хотелось бы получить *объективные* априорные распределения, т.е. добиться того, чтобы два субъекта с одной и той же априорной информацией задали одни и те же априорные распределения
- Для этого нужно придумать целевую функцию, условие, которому должны априорные распределения удовлетворять

## АПРИОРНЫЕ РАСПРЕДЕЛЕНИЯ ДЖЕФФРИСА

- Вернёмся для начала к монетке; равномерное распределение  $p(\theta) = \text{Unif}([0, 1])$  мы мотивировали тем, что «не предпочитаем никакого значения  $\theta$  другому»
- Но это не единственный способ параметризовать монетку! Давайте, к примеру, перейдём к логарифму шансов (log-odds)

$$\eta = \log \frac{\theta}{1 - \theta}, \quad \text{то есть} \quad \theta = \frac{1}{1 + e^{-\eta}}$$

- Логарифм шансов параметризует монетку величиной, изменяющейся от  $-\infty$  до  $\infty$ , и мы, по идее, не предпочитаем и никакого значения  $\eta$  другому
- Но если сделать в  $p_\theta(\theta) = \text{Unif}([0, 1])$  замену переменных  $\theta = \frac{1}{1 + e^{-\eta}}$ , получится вовсе не равномерное распределение:

$$p_\eta(\eta) = p_\theta \left( \frac{1}{1 + e^{-\eta}} \right) \left( \frac{1}{1 + e^{-\eta}} \right)' = \frac{e^\eta}{(1 + e^\eta)^2}$$

- Получается, что если мы решаем «не предпочитать» разные значения  $\theta$ , то мы тем самым очень даже предпочитаем разные значения  $\eta$ ; и наоборот
- Равномерное распределение «ничего не предпочитает» только для конкретной параметризации — а это ведь одна и та же монетка, и у неё один и тот же параметр, как его ни преобразуй; по крайней мере, монотонные преобразования вроде перехода от  $\theta$  к  $\eta$  и наоборот уж точно должны быть разрешены
- Долгое время это считалось серьёзным аргументом против байесовского подхода к статистике в целом; сам сэръ Рональд Фишер критиковал байесовский подход с этих позиций

- Решение предложил сэръ Гарольд Джеффрис: давайте считать это свойство требованием к «неинформативному априорному распределению»!
- Рассмотрим модель машинного обучения с правдоподобием  $p_{\theta}(\mathbf{x}|\theta)$  (пусть пока один параметр  $\theta \in \Theta$ ) и сделаем репараметризацию: применим к параметру  $\theta$  произвольное гладкое монотонное преобразование  $\eta = h(\theta)$
- Тогда заменой переменных мы получим новую модель

$$p_{\eta}(\mathbf{x}|\eta) = p_{\theta}(D|h^{-1}(\eta))$$

## АПРИОРНЫЕ РАСПРЕДЕЛЕНИЯ ДЖЕФФРИСА

- Предположим, что мы так или иначе придумали априорное распределение  $p_\theta(\theta)$  для исходной модели. Тогда у нас есть два способа задать априорное распределение после репараметризации:
  - напрямую применить тот же принцип и взять то же распределение  $p_\eta(\eta)$ , но с параметром  $\eta$ ;
  - сделать замену переменных  $\eta = h(\theta)$  в априорном распределении  $p_\theta(\theta)$  и получить распределение

$$\tilde{p}_\eta(\eta) = \frac{p_\theta(h^{-1}(\eta))}{|h'(h^{-1}(\eta))|}.$$

- Джеффрис предложил потребовать, чтобы для любой гладкой монотонной функции  $h$  эти два распределения совпадали:  $p_\eta(\eta) = \tilde{p}_\eta(\eta)$ , или, что то же самое,

$$\tilde{p}_\eta(h(\theta)) = \frac{p_\theta(\theta)}{|h'(\theta)|} \quad \text{для всех } \theta \in \Theta$$

- Джеффрис дал и конструктивный ответ на вопрос, как построить такое априорное распределение
- В случае  $\theta \in \mathbb{R}$  априорное распределение Джеффриса (Jeffreys prior) для модели  $p_{\theta}(\mathbf{x}|\theta)$  определяется как

$$p^J(\theta) = \text{const} \cdot \sqrt{\mathcal{I}(\theta)}$$

## АПРИОРНЫЕ РАСПРЕДЕЛЕНИЯ ДЖЕФФРИСА

- Проверим нужное свойство, записав определение  $\mathcal{J}(\theta)$  и продифференцировав в этом определении  $\log p_\eta(\mathbf{x}|h(\theta))$  как сложную функцию:

$$\begin{aligned}\mathcal{J}_\theta(\theta) &= \int \left( \frac{\partial}{\partial \theta} \log p_\theta(\mathbf{x}|\theta) \right)^2 p_\theta(\mathbf{x}|\theta) d\mathbf{x} \\ &= \int \left( \frac{\partial}{\partial \theta} \log p_\eta(\mathbf{x}|h(\theta)) \right)^2 p_\eta(\mathbf{x}|h(\theta)) d\mathbf{x} \\ &= \int \left( h'(\theta) \frac{\partial}{\partial \eta} \log p_\eta(\mathbf{x}|\eta) \right)^2 p_\eta(\mathbf{x}|h(\theta)) d\mathbf{x} \\ &= h'(\theta)^2 \mathcal{J}_\eta(h(\theta)), \quad \text{а значит,}\end{aligned}$$

$$p_\eta^J(h(\theta)) = \text{const} \cdot \sqrt{\mathcal{J}_\eta(h(\theta))} = \text{const} \cdot \frac{\mathcal{J}_\theta(\theta)}{|h'(\theta)|} = \frac{p_\theta^J(\theta)}{|h'(\theta)|},$$

что и требовалось

## АПРИОРНЫЕ РАСПРЕДЕЛЕНИЯ ДЖЕФФРИСА

- Если интеграл  $\int_{\Theta} \sqrt{\mathcal{J}(\theta)}$  конечен, то он даст нормировочную константу для плотности, и всё будет хорошо
- Но во многих практических случаях этот интеграл будет расходиться, и распределение Джеффриса будет задавать так называемое *некорректное априорное распределение* (improper prior), то есть такое априорное распределение, которое само по себе не является распределением вероятностей, но при подстановке в формулу Байеса даёт для любого непустого набора данных полноценное апостериорное распределение
- Грубо говоря, в таких случаях мы не будем обращать внимания на то, что интеграл расходится, а будем просто умножать правдоподобие данных на функцию  $\sqrt{\mathcal{J}(\theta)}$  и нормировать результат — и вот у этого результата уже действительно интеграл обязательно должен будет сходиться

- Например, для бросания монетки модель в терминах  $\theta$  выглядит как  $p(x|\theta) = \theta^{[x=1]}(1-\theta)^{[x=0]}$  для данных, где  $x = 1$  означает выпадение орла, а  $x = 0$  — решки
- Априорное распределение Джеффриса можно найти как

$$\begin{aligned} p^J(\theta) \propto \sqrt{\mathcal{J}(\theta)} &= \sqrt{\mathbb{E} \left[ \left( \frac{\partial}{\partial \theta} \log p_\theta(D|\theta) \right)^2 \right]} = \\ &= \sqrt{\mathbb{E} \left[ \left( \frac{[x=1]}{\theta} - \frac{[x=0]}{1-\theta} \right)^2 \right]} = \\ &= \sqrt{\theta \left( \frac{1}{\theta} \right)^2 + (1-\theta) \left( \frac{1}{1-\theta} \right)^2} = \frac{1}{\sqrt{\theta(1-\theta)}} \end{aligned}$$

- Получилось бета-распределение с дробными параметрами  $\alpha = \beta = \frac{1}{2}$ :

$$p^J(\theta) = \text{Beta}\left(\theta \middle| \frac{1}{2}, \frac{1}{2}\right) \propto \theta^{-\frac{1}{2}} (1 - \theta)^{-\frac{1}{2}}$$

- Аналогично, априорное распределение Джеффриса для игральной кости с  $k$  гранями – это распределение Дирихле с дробными параметрами  $\alpha = (\frac{1}{2} \dots \frac{1}{2})$

- Другой пример — одномерное нормальное распределение  $N(x|\mu, \sigma^2)$ , где  $\sigma^2$  фиксирована, а вывод ведётся по  $\mu$
- Тогда информация Фишера для одной точки данных  $x$  равна

$$\begin{aligned} \mathcal{J}(\mu) &= \mathbb{E}_{x|\mu} \left[ \left( \frac{\partial}{\partial \mu} \log p(x|\mu) \right)^2 \right] = \mathbb{E}_{x|\mu} \left[ \left( \frac{x - \mu}{\sigma^2} \right)^2 \right] = \\ &= \frac{1}{\sigma^4} \mathbb{E}_{x|\mu} [(x - \mu)^2] = \frac{\sigma^2}{\sigma^4} = \frac{1}{\sigma^2}; \end{aligned}$$

а значит, информация Фишера для  $D = \{x_1, \dots, x_N\}$  равна  $N/\sigma^2$ , и априорное распределение Джеффриса получается пропорциональным  $p^J(\mu) \propto \sqrt{N/\sigma^2}$

- Но  $\sigma^2$  — константа, так что распределение Джеффриса просто равно одной и той же константе для всех  $\mu$
- Это характерный пример некорректного априорного распределения: нельзя взять равномерное распределение на всей прямой  $\theta \in (-\infty, \infty)$ , но для байесовского вывода это не страшно, потому что при умножении на правдоподобие пусть даже одной точки будет получаться суммируемая функция

# АПРИОРНЫЕ РАСПРЕДЕЛЕНИЯ ДЖЕФФРИСА

- В многомерном случае распределение Джеффриса обобщается на векторный параметр  $\theta \in \mathbb{R}^d$  так:

$$p^J(\theta) = \text{const} \cdot \sqrt{\det \mathcal{J}(\theta)}$$

- Базовое свойство проверяется так же, в многомерном случае замена переменных добавляет определитель матрицы частных производных:

$$\begin{aligned} p^J(\eta) &= p^J(\theta) \left| \det \frac{\partial \theta_i}{\partial \eta_j} \right| \propto \sqrt{\det \mathcal{J}(\theta) \left( \det \frac{\partial \theta_i}{\partial \eta_j} \right)^2} \\ &= \sqrt{\det \frac{\partial \theta_k}{\partial \eta_i} \det \mathbb{E} \left[ \frac{\partial \log p_\theta(\mathbf{x}|\theta)}{\partial \theta_k} \frac{\partial \log p_\theta(\mathbf{x}|\theta)}{\partial \theta_l} \right] \det \frac{\partial \theta_l}{\partial \eta_j}} \\ &= \sqrt{\det \mathbb{E} \left[ \sum_{k,l} \frac{\partial \theta_k}{\partial \eta_i} \frac{\partial \log p_\theta(\mathbf{x}|\theta)}{\partial \theta_k} \frac{\partial \log p_\theta(\mathbf{x}|\theta)}{\partial \theta_l} \frac{\partial \theta_l}{\partial \eta_j} \right]} \\ &= \sqrt{\det \mathbb{E} \left[ \frac{\partial \log p_\eta(\mathbf{x}|\eta)}{\partial \eta_i} \frac{\partial \log p_\eta(\mathbf{x}|\eta)}{\partial \eta_j} \right]} = \sqrt{\det \mathcal{J}(\eta)}. \end{aligned}$$

- Для примера давайте найдём распределение Джеффриса  $p^J(\theta) \propto \sqrt{\det \mathcal{J}(\theta)}$  для одномерного гауссиана с двумя параметрами
- Мы находимся в двумерном случае, где  $\theta = (\mu \ \sigma^2)^\top$ , то есть матрица информации Фишера задаётся как

$$\mathcal{J}(\mu, \sigma^2) = -\mathbb{E}_{D|\mu, \sigma^2} \left[ \begin{pmatrix} \frac{\partial^2}{\partial \mu^2} \log p(D|\mu, \sigma^2) & \frac{\partial^2}{\partial \mu \partial \sigma^2} \log p(D|\mu, \sigma^2) \\ \frac{\partial^2}{\partial \sigma^2 \partial \mu} \log p(D|\mu, \sigma^2) & \frac{\partial^2}{\partial (\sigma^2)^2} \log p(D|\mu, \sigma^2) \end{pmatrix} \right].$$

## АПРИОРНЫЕ РАСПРЕДЕЛЕНИЯ ДЖЕФФРИСА

- Для гауссиана все эти производные легко подсчитать; для  $D = \{x_1, \dots, x_n\}$  сначала перейдём к достаточным статистикам  $\bar{x} = \frac{1}{n} \sum_i x_i$ ,  $s^2 = \frac{1}{n-1} \sum_i (x_i - \bar{x})^2$ :

$$\log p(D|\mu, \sigma^2) = \text{const} - \frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} ((n-1)s^2 + n(\bar{x} - \mu)^2),$$

а затем возьмём производные:

$$\frac{\partial \log p}{\partial \mu} = \frac{n(\bar{x} - \mu)}{\sigma^2}, \quad \frac{\partial \log p}{\partial \sigma^2} = -\frac{n}{2\sigma^2} - \frac{(n-1)s^2 + n(\bar{x} - \mu)^2}{4\sigma^4},$$

а значит,

$$\begin{aligned} \frac{\partial^2}{\partial \mu^2} \log p(D|\mu, \sigma^2) &= -\frac{n}{\sigma^2}, \\ \frac{\partial^2}{\partial \mu \partial \sigma^2} \log p(D|\mu, \sigma^2) &= \frac{\partial^2}{\partial \sigma^2 \partial \mu} \log p(D|\mu, \sigma^2) = -\frac{n(\bar{x} - \mu)}{\sigma^4}, \\ \frac{\partial^2}{\partial (\sigma^2)^2} \log p(D|\mu, \sigma^2) &= \frac{n}{2\sigma^4} - \frac{(n-1)s^2 + n(\bar{x} - \mu)^2}{\sigma^2}. \end{aligned}$$

## АПРИОРНЫЕ РАСПРЕДЕЛЕНИЯ ДЖЕФФРИСА

- Теперь нужно взять ожидания; поскольку  $\mathbb{E}[\bar{x}] = \mu$ ,  $\mathbb{E}[(\bar{x} - \mu)^2] = \sigma^2/n$  и  $\mathbb{E}[\sum_i (x_i - \bar{x})^2] = (n-1)\sigma^2$ , матрица информации Фишера и её определитель равны

$$\mathcal{J}(\mu, \sigma^2) = \begin{pmatrix} \frac{n}{\sigma^2} & 0 \\ 0 & \frac{n}{2\sigma^4} \end{pmatrix}, \quad \det \mathcal{J}(\mu, \sigma^2) = \frac{n^2}{2\sigma^6},$$

и априорное распределение Джеффриса будет

$$p^J(\mu, \sigma^2) \propto \sqrt{\det \mathcal{J}(\theta)} \propto (\sigma^2)^{-3/2}$$

- Обратите внимание, что это распределение совсем не содержит  $\mu$  — по  $\mu$  опять получается некорректное равномерное априорное распределение на всей прямой
- А вот по  $\sigma^2$  это полноценное распределение, интеграл по  $(0, \infty)$  сходится, и нормировочную константу можно подсчитать

- Из априорного распределения Джеффриса получается такое апостериорное распределение:

$$\begin{aligned} p(\mu, \sigma^2 | x_1, \dots, x_n) &\propto (\sigma^2)^{-3/2} \cdot (\sigma^2)^{-n/2} \cdot e^{-\frac{1}{2\sigma^2}((n-1)s^2 + n(\bar{x} - \mu)^2)} = \\ &= (\sigma^2)^{-(n+3)/2} e^{-\frac{1}{2\sigma^2}((n-1)s^2 + n(\bar{x} - \mu)^2)} = \\ &= \left( \frac{1}{\sqrt{\sigma^2}} e^{-\frac{(\bar{x} - \mu)^2}{2(\sigma/\sqrt{n})^2}} \right) \cdot \left( (\sigma^2)^{-(n+1)/2} e^{-\frac{(n-1)s^2}{2\sigma^2}} \right) \end{aligned}$$

- Первый множитель — это (с точностью до константы) нормальное распределение на  $\mu$  со средним  $\bar{x}$  и дисперсией  $(\sigma/\sqrt{n})^2$
- Получился вполне логичный результат: среднее постепенно уточняется с ростом  $n$ , то есть при получении новых данных

# ЭКСПОНЕНЦИАЛЬНОЕ СЕМЕЙСТВО

---

# 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.

СПАСИБО!

Спасибо за внимание!

