

# АИС И ЭКСПОНЕНЦИАЛЬНОЕ СЕМЕЙСТВО

---

Сергей Николенко

СПбГУ — Санкт-Петербург

21 октября 2023 г.

---

*Random facts:*

- 21 октября в Великобритании — Apple Day; с 1990 года проводятся ярмарки с конкурсами вроде стрельбы из лука по яблокам или срезания самой длинной кожуры
- 21 октября 63 г. до н.э. избранный консулом Цицерон, получив сведения о намерениях Катилины совершить переворот, произнёс в сенате речь (то самое «Доколе же ты, Катилина, будешь злоупотреблять нашим терпением?») и планы Катилины сорвал
- 21 октября 1727 г. был заключён Кяхтинский договор между Россией и Китаем; «песчаная Венеция» Кяхта сто лет беспошлинно снабжала чаем всю Европу; торговля пришла в упадок после открытия Суэцкого канала, и сейчас в Кяхте живёт 18000 человек
- 21 октября 1805 г. эскадра Горацио Нельсона победила при Трафальгаре превосходящие силы Пьера-Шарля Вильнёва; сам Нельсон погиб, его тело для сохранности поместили в бочку с ромом, и с тех пор выдаваемый на кораблях ром английские моряки называли «адмиральской кровью»
- 21 октября 1824 г. Джозеф Аспдин запатентовал портлендский цемент, 21 октября 1832 г. Павел Шиллинг в своей петербургской квартире впервые продемонстрировал изобретённый им электромагнитный телеграф, а 21 октября 1879 г. Томас Эдисон испытал свою первую лампу накаливания с угольной нитью

# ИНФОРМАЦИОННЫЙ КРИТЕРИЙ АКАИКЕ

---

- Пусть данные  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  были получены из истинного распределения  $p_{\text{data}}(\mathbf{x})$ , а мы пытаемся приблизить их некоторой параметрической моделью  $p(\mathbf{x}|\theta)$ ,  $\theta \in \mathbb{R}^d$ .
- Предположим, что мы обучили модель методом максимального правдоподобия, получив  $p(\mathbf{x}|\theta_{\text{ML}})$ .
- Давайте попробуем оценить, насколько модель  $p(\mathbf{x}|\theta_{\text{ML}})$  отличается от неизвестного истинного распределения  $p_{\text{data}}(\mathbf{x})$ :

$$\begin{aligned} \text{KL}(p_{\text{data}} \| p(\mathbf{x}|\theta_{\text{ML}})) &= \mathbb{E}_{p_{\text{data}}(\mathbf{x})} \left[ \log \frac{p_{\text{data}}(\mathbf{x})}{p(\mathbf{x}|\theta_{\text{ML}})} \right] = \\ &= \mathbb{E}_{p_{\text{data}}(\mathbf{x})} [\log p_{\text{data}}(\mathbf{x})] - \mathbb{E}_{p_{\text{data}}(\mathbf{x})} [\log p(\mathbf{x}|\theta_{\text{ML}})]. \end{aligned}$$

- Модель будет тем лучше, чем больше будет  $\mathbb{E}_{p_{\text{data}}(\mathbf{x})} [\log p(\mathbf{x}|\theta_{\text{ML}})]$ , и для оценки расхождения между  $p(\mathbf{x}|\theta_{\text{ML}})$  и  $p_{\text{data}}(\mathbf{x})$  нужно получить оценку ожидаемого логарифма правдоподобия.
- Во всех критериях важен логарифм правдоподобия в точке его максимума, ведь это как раз выборочная оценка ожидания:

$$\begin{aligned}\mathbb{E}_{p_{\text{data}}(\mathbf{x})} [\log p(\mathbf{x}|\theta_{\text{ML}})] &= \int p_{\text{data}}(\mathbf{x}) \log p(\mathbf{x}|\theta_{\text{ML}}) d\mathbf{x} \approx \\ &\approx \frac{1}{N} \sum_{n=1}^N \log p(\mathbf{x}_n|\theta_{\text{ML}}).\end{aligned}$$

- Но это смещённая оценка как минимум потому, что мы обучаем параметры максимального правдоподобия  $\theta_{\text{ML}}$  на том же датасете  $\mathbf{X}$ , который используется в этой оценке.

- Если истинная модель  $p_{\text{data}}$  тоже из семейства  $p(\mathbf{x}|\theta)$  с некоторым истинным параметром  $\theta_0$ , то

$$\theta_0 = \arg \max_{\theta} \mathbb{E}_{p_{\text{data}}(\mathbf{x})} [\log p(\mathbf{x}|\theta)],$$

но это ожидание берётся по всему распределению.

- $\theta_0$  — это «истинная» гипотеза максимального правдоподобия; при некоторых условиях регулярности можно доказать, что:

- $\theta_{\text{ML}}(\mathbf{X}) \rightarrow \theta_0$  при  $N \rightarrow \infty$ ;
- для  $\theta_{\text{ML}}(\mathbf{X})$  верна асимптотическая нормальность, т.е. распределение величины  $\sqrt{N}(\theta_{\text{ML}} - \theta_0)$  сходится по вероятности к распределению  $N(0, I(\theta_0)^{-1})$ , где  $I(\theta)$  — это матрица информации Фишера

$$I(\theta) = \int p(\mathbf{x}|\theta) \frac{\partial \log p(\mathbf{x}|\theta)}{\partial \theta} \frac{\partial \log p(\mathbf{x}|\theta)}{\partial \theta^{\top}} d\mathbf{x}.$$

- Более того, те формулы предполагали, что  $p_{\text{data}}(\mathbf{x}) = p(\mathbf{x}|\theta_0)$ , но аналогичные результаты можно получить и если  $p_{\text{data}}(\mathbf{x})$  не принадлежит параметрическому семейству  $p(\mathbf{x}|\theta)$ .
- Пусть  $\theta_0$  — максимум ожидания логарифма правдоподобия по  $p_{\text{data}}$ , то есть решение системы

$$\int p_{\text{data}}(\mathbf{x}) \frac{\partial \log p(\mathbf{x}|\theta)}{\partial \theta} d\mathbf{x} = 0.$$

- Тогда при тех же условиях можно доказать, что:
  - $\theta_{\text{ML}}(\mathbf{X}) \rightarrow \theta_0$  при  $N \rightarrow \infty$ ;
  - распределение величины  $\sqrt{N}(\theta_{\text{ML}} - \theta_0)$  сходится по вероятности к нормальному распределению

$$\sqrt{N}(\theta_{\text{ML}} - \theta_0) \rightarrow_{N \rightarrow \infty} N(0, J^{-1}(\theta_0)I(\theta_0)J^{-1}(\theta_0)),$$

где  $I(\theta)$  — это та же матрица информации Фишера, только по распределению  $p_{\text{data}}$ :

$$I(\theta) = \int p_{\text{data}}(\mathbf{x}) \frac{\partial \log p(\mathbf{x}|\theta)}{\partial \theta} \frac{\partial \log p(\mathbf{x}|\theta)}{\partial \theta^\top} d\mathbf{x};$$

а  $J(\theta)$  — это ожидание матрицы вторых производных

$$J(\theta) = - \int p_{\text{data}}(\mathbf{x}) \frac{\partial^2 \log p(\mathbf{x}|\theta)}{\partial \theta \partial \theta^\top} d\mathbf{x}.$$

- Иначе говоря, на позиции  $(i, j)$  у матрицы  $I(\theta)$  стоит ожидание произведения  $\frac{\partial \log p(\mathbf{x}|\theta)}{\partial \theta_i}$  и  $\frac{\partial \log p(\mathbf{x}|\theta)}{\partial \theta_j}$ , а у матрицы  $J(\theta)$  — ожидание второй производной  $\frac{\partial^2 \log p(\mathbf{x}|\theta)}{\partial \theta_i \partial \theta_j}$ .
- И если всё-таки  $p_{\text{data}}(\mathbf{x}) = p(\mathbf{x}|\theta_0)$ , то

$$\begin{aligned} \frac{\partial^2 \log p(\mathbf{x}|\theta)}{\partial \theta_i \partial \theta_j} &= \frac{\partial}{\partial \theta_i} \left( \frac{1}{p(\mathbf{x}|\theta)} \frac{\partial p(\mathbf{x}|\theta)}{\partial \theta_j} \right) = \\ &= \frac{1}{p(\mathbf{x}|\theta)} \frac{\partial^2 p(\mathbf{x}|\theta)}{\partial \theta_i \partial \theta_j} - \frac{\partial \log p(\mathbf{x}|\theta)}{\partial \theta_i} \frac{\partial \log p(\mathbf{x}|\theta)}{\partial \theta_j}, \end{aligned}$$

а в ожидании по  $p_{\text{data}}(\mathbf{x}) = p(\mathbf{x}|\theta_0)$  при подстановке  $\theta = \theta_0$

$$\begin{aligned} \mathbb{E}_{p(\mathbf{x}|\theta_0)} \left[ \frac{1}{p(\mathbf{x}|\theta_0)} \frac{\partial^2 p(\mathbf{x}|\theta_0)}{\partial \theta_i \partial \theta_j} \right] &= \int \frac{p(\mathbf{x}|\theta_0)}{p(\mathbf{x}|\theta_0)} \frac{\partial^2 p(\mathbf{x}|\theta_0)}{\partial \theta_i \partial \theta_j} d\mathbf{x} = \\ &= \frac{\partial^2}{\partial \theta_i \partial \theta_j} \int p(\mathbf{x}|\theta_0) = 0, \quad \text{то есть здесь } I(\theta_0) = J(\theta_0). \end{aligned}$$



- Различные информационные критерии для сравнения моделей оценивают смещение выборочной оценки для величины  $\mathbb{E}_{p_{\text{data}}(\mathbf{x})} [\log p(\mathbf{x}|\theta_{\text{ML}})]$

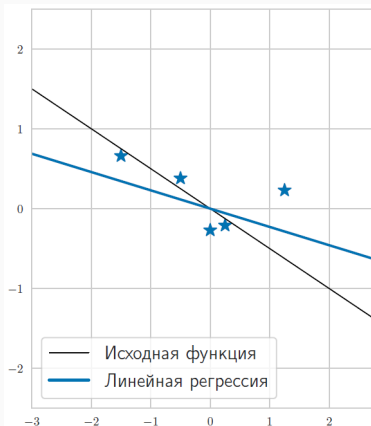
$$b(p_{\text{data}}) = \mathbb{E}_{\mathbf{X} \sim p_{\text{data}}} [\log p(\mathbf{X}|\theta_{\text{ML}}) - N \mathbb{E}_{p_{\text{data}}(\mathbf{z})} [\log p(\mathbf{z}|\theta_{\text{ML}})]] ,$$

где мы взяли ожидание по датасетам  $\mathbf{X} \sim p_{\text{data}}$ .

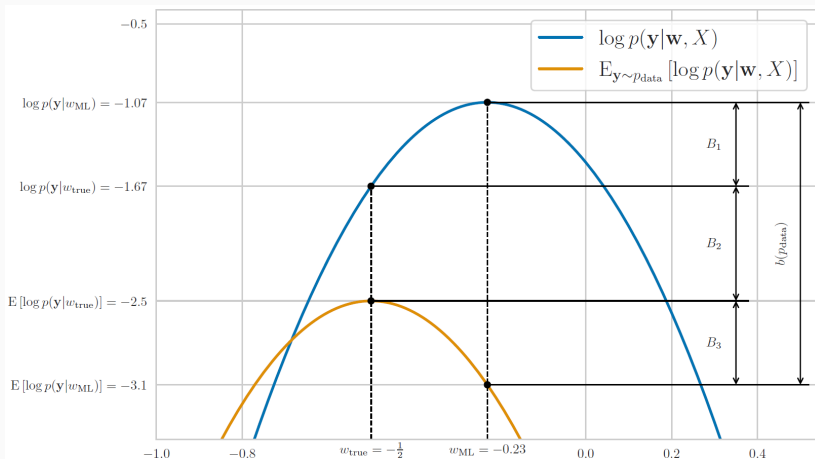
- Если мы сможем оценить смещение  $b(p_{\text{data}})$ , то информационный критерий можно будет построить, умножив на  $-2$  аналогично BIC:

$$\begin{aligned} \text{IC}(\mathbf{X}, \theta) &= \\ &= -2 (\text{логарифм правдоподобия } \mathbf{X} \text{ в } \theta_{\text{ML}} - \text{оценка смещения}) = \\ &= -2 \sum_{n=1}^N \log p(\mathbf{x}_n | \theta_{\text{ML}}) + 2 (\text{оценка } b(p_{\text{data}})) . \end{aligned}$$

- Пример — давайте рассмотрим линейную регрессию с одним параметром  $w$ :



- И нарисуем графики  $\log p(\mathbf{y}|w)$  и  $\mathbb{E}_{\mathbf{y} \sim p_{\text{data}}} [\log p(\mathbf{y}|w)]$ :



- Давайте попробуем оценить  $b(p_{\text{data}})$ , разложив как на картинке на три слагаемых (но теперь в ожидании):

$$\begin{aligned} b(p_{\text{data}}) &= \mathbb{E}_{\mathbf{X} \sim p_{\text{data}}} \left[ \log p(\mathbf{X} | \theta_{\text{ML}}) - N \mathbb{E}_{p_{\text{data}}(\mathbf{z})} [\log p(\mathbf{z} | \theta_{\text{ML}})] \right] = \\ &= \mathbb{E}_{\mathbf{X} \sim p_{\text{data}}} [\log p(\mathbf{X} | \theta_{\text{ML}}) - \log p(\mathbf{X} | \theta_0)] + \\ &+ \mathbb{E}_{\mathbf{X} \sim p_{\text{data}}} \left[ \log p(\mathbf{X} | \theta_0) - N \mathbb{E}_{p_{\text{data}}(\mathbf{z})} [\log p(\mathbf{z} | \theta_0)] \right] + \\ &+ \mathbb{E}_{\mathbf{X} \sim p_{\text{data}}} \left[ N \mathbb{E}_{p_{\text{data}}(\mathbf{z})} [\log p(\mathbf{z} | \theta_0)] - N \mathbb{E}_{p_{\text{data}}(\mathbf{z})} [\log p(\mathbf{z} | \theta_{\text{ML}})] \right] = \\ &= B_1 + B_2 + B_3. \end{aligned}$$

- Будем оценивать слагаемые по отдельности.

- Проще всего оценить  $B_2$ , потому что в нём нет  $\theta_{\text{ML}}$ :

$$\begin{aligned} B_2 &= \mathbb{E}_{\mathbf{X} \sim p_{\text{data}}} \left[ \log p(\mathbf{X} | \theta_0) - N \mathbb{E}_{p_{\text{data}}(\mathbf{z})} [\log p(\mathbf{z} | \theta_0)] \right] = \\ &= \mathbb{E}_{\mathbf{X} \sim p_{\text{data}}} \left[ \sum_{n=1}^N \log p(\mathbf{x}_n | \theta_0) \right] - N \mathbb{E}_{p_{\text{data}}(\mathbf{z})} [\log p(\mathbf{z} | \theta_0)] = 0. \end{aligned}$$

- Это не значит, что  $B_2$  всегда равно нулю; для конкретного датасета  $\mathbf{X}$  значение  $B_2$  будет ненулевым, но в ожидании получится ноль.

- Чтобы оценить  $B_3$ , рассмотрим функцию  $\eta(\theta_{\text{ML}}) = \mathbb{E}_{p_{\text{data}}(\mathbf{z})} [\log p(\mathbf{z}|\theta_{\text{ML}})]$  и разложим её по формуле Тейлора в окрестности точки  $\theta_0$  (её максимума):

$$\eta(\theta_{\text{ML}}) \approx \eta(\theta_0) - \frac{1}{2} (\theta_{\text{ML}} - \theta_0)^\top J(\theta_0) (\theta_{\text{ML}} - \theta_0),$$

где

$$\begin{aligned} J(\theta_0) &= -\mathbb{E}_{p_{\text{data}}(\mathbf{z})} \left[ \frac{\partial^2 \log p(\mathbf{z}|\theta)}{\partial \theta \partial \theta^\top} \Bigg|_{\theta_0} \right] = \\ &= -\int p_{\text{data}}(\mathbf{z}) \frac{\partial^2 \log p(\mathbf{z}|\theta)}{\partial \theta \partial \theta^\top} \Bigg|_{\theta_0} d\mathbf{z}. \end{aligned}$$

- А  $B_3$  — это ожидание  $\eta(\theta_0) - \eta(\theta_{\text{ML}})$  по распределению  $p_{\text{data}}(\mathbf{X})$ :

$$\begin{aligned} B_3 &= \mathbb{E}_{\mathbf{X} \sim p_{\text{data}}} \left[ N \mathbb{E}_{p_{\text{data}}(\mathbf{z})} [\log p(\mathbf{z}|\theta_0)] - N \mathbb{E}_{p_{\text{data}}(\mathbf{z})} [\log p(\mathbf{z}|\theta_{\text{ML}})] \right] = \\ &= \frac{N}{2} \mathbb{E}_{\mathbf{X} \sim p_{\text{data}}} \left[ (\theta_{\text{ML}} - \theta_0)^\top J(\theta_0) (\theta_{\text{ML}} - \theta_0) \right] = \\ &= \frac{N}{2} \mathbb{E}_{\mathbf{X} \sim p_{\text{data}}} \left[ \text{Tr} \left( J(\theta_0) (\theta_{\text{ML}} - \theta_0) (\theta_{\text{ML}} - \theta_0)^\top \right) \right] = \\ &= \frac{N}{2} \text{Tr} \left( J(\theta_0) \mathbb{E}_{\mathbf{X} \sim p_{\text{data}}} \left[ (\theta_{\text{ML}} - \theta_0) (\theta_{\text{ML}} - \theta_0)^\top \right] \right). \end{aligned}$$

- Теперь можно вместо ожидания матрицы ковариаций по датасету  $\mathbf{X}$  подставить асимптотический результат:

$$B_3 = \frac{N}{2} \text{Tr} \left( J(\theta_0) \frac{1}{N} J(\theta_0)^{-1} I(\theta_0) J(\theta_0)^{-1} \right) = \frac{1}{2} \text{Tr} (I(\theta_0) J(\theta_0)^{-1}).$$

- Для оценки  $B_1$  нужно повернуть аналогичный трюк с  $\ell(\theta) = \log p(X|\theta)$ , разложив его вокруг своего максимума  $\theta_{\text{ML}}$ :

$$\ell(\theta) = \ell(\theta_{\text{ML}}) + \frac{1}{2} (\theta - \theta_{\text{ML}})^\top \left. \frac{\partial^2 \ell(\theta)}{\partial \theta \partial \theta^\top} \right|_{\theta_{\text{ML}}} (\theta - \theta_{\text{ML}}).$$

- Мы знаем, что  $\theta_{\text{ML}} \rightarrow \theta_0$  при  $N \rightarrow \infty$ ; а по закону больших чисел можно получить, что

$$-\frac{1}{N} \left. \frac{\partial^2 \ell(\theta)}{\partial \theta \partial \theta^\top} \right|_{\theta} = -\frac{1}{N} \sum_{n=1}^N \left. \frac{\partial^2 \log p(\mathbf{x}_n|\theta)}{\partial \theta \partial \theta^\top} \right|_{\theta_0} \rightarrow J(\theta_0).$$



- Следовательно, в нашу оценку можно подставить

$$\ell(\theta_{\text{ML}}) - \ell(\theta_0) \approx -\frac{N}{2} (\theta - \theta_{\text{ML}})^\top J(\theta_0) (\theta - \theta_{\text{ML}}).$$

- А затем и оценить  $B_1$  так же, как оценивали  $B_3$ :

$$\begin{aligned} B_1 &= \mathbb{E}_{\mathbf{X} \sim p_{\text{data}}} [\log p(\mathbf{X} | \theta_{\text{ML}}) - \log p(\mathbf{X} | \theta_0)] = \\ &= \frac{N}{2} \mathbb{E}_{\mathbf{X} \sim p_{\text{data}}} [(\theta - \theta_{\text{ML}})^\top J(\theta_0) (\theta - \theta_{\text{ML}})] = \\ &= \frac{N}{2} \text{Tr} \left( J(\theta_0) \mathbb{E}_{\mathbf{X} \sim p_{\text{data}}} [(\theta - \theta_{\text{ML}})^\top (\theta - \theta_{\text{ML}})] \right), \end{aligned}$$

ТО ЕСТЬ

$$B_3 = \frac{1}{2} \text{Tr} (I(\theta_0) J(\theta_0)^{-1}). \quad (1)$$

- Осталось только объединить три оценки:

$$b(p_{\text{data}}) = B_1 + B_2 + B_3 = \text{Tr} (I(\theta_0)J(\theta_0)^{-1}).$$

- $I(\theta_0)$  и  $J(\theta_0)$  нам неизвестны, т.к. зависят от  $p_{\text{data}}$ ; если взять оценки  $\hat{I}$  и  $\hat{J}$ , это приведёт нас к *информационному критерию Такеучи* (Takeuchi information criterion, TIC):

$$\text{TIC} = -2 \sum_{n=1}^N \log p(\mathbf{x}_n | \theta_{\text{ML}}) + 2 \text{Tr} (\hat{I} \hat{J}^{-1}).$$

- В качестве  $\hat{I}$  и  $\hat{J}$  можно подставить просто усреднённые значения по датасету в точке максимума правдоподобия:

$$\hat{I}_{i,j} = \frac{1}{N} \sum_{n=1}^N \left. \frac{\partial \log p(\mathbf{x}_n | \theta)}{\partial \theta_i} \frac{\partial \log p(\mathbf{x}_n | \theta)}{\partial \theta_j} \right|_{\theta_{\text{ML}}}, \quad \hat{J}_{i,j} = \frac{1}{N} \sum_{n=1}^N \left. \frac{\partial^2 \log p(\mathbf{x}_n | \theta)}{\partial \theta_i \partial \theta_j} \right|_{\theta_{\text{ML}}}.$$

- А если можно всё-таки предполагать, что истинное распределение данных  $p_{\text{data}}$  лежит в параметрическом семействе  $p(\mathbf{x}|\theta)$ , то есть  $p_{\text{data}}(\mathbf{x}) = p(\mathbf{x}|\theta_0)$ , то, как мы обсуждали выше,  $I(\theta_0) = J(\theta_0)$ , и информационный критерий Такеучи превращается в *информационный критерий Акаике* (Akaike information criterion, AIC):

$$\begin{aligned} \text{AIC} &= -2 \sum_{n=1}^N \log p(\mathbf{x}_n | \theta_{\text{ML}}) + 2\text{Tr}(\mathbf{I}_d) = \\ &= -2 \sum_{n=1}^N \log p(\mathbf{x}_n | \theta_{\text{ML}}) + 2d. \end{aligned}$$

- Очень простая формула!

- Пример: вернёмся к полиномиальной регрессии с логарифмом правдоподобия

$$\ell(\mathbf{w}) = -\frac{N}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{n=1}^N (t_n - \mathbf{w}^\top \mathbf{x}_n)^2.$$

- Давайте теперь для разнообразия будем дисперсию тоже обучать:

$$\mathbf{w}_{\text{ML}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t}, \quad \sigma_{\text{ML}}^2 = \frac{1}{N} \sum_{n=1}^N (t_n - \mathbf{w}_{\text{ML}}^\top \mathbf{x}_n)^2.$$

- Тогда при подстановке гипотезы максимального правдоподобия получится

$$\ell(\mathbf{w}_{\text{ML}}) = -\frac{N}{2} \log(2\pi\sigma_{\text{ML}}^2) - \frac{N}{2},$$

а параметров в модели будет  $d + 2$ , где  $d$  — степень

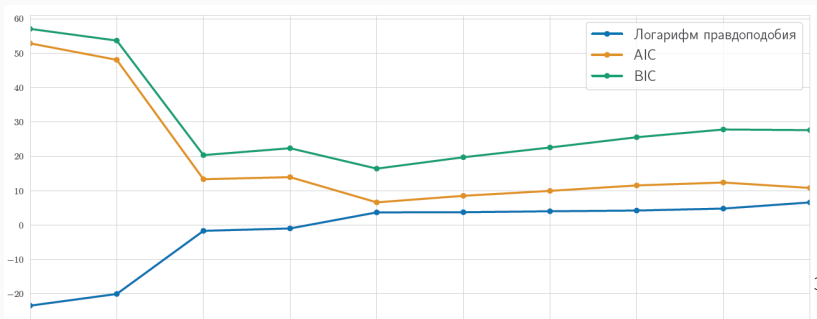
# Идея и вывод AIC

- Тогда AIC и BIC получаются такими:

$$AIC(d) = N (\log (2\pi) + \log \sigma_{ML}^2 + 1) + 2(d + 2),$$

$$BIC(d) = N (\log (2\pi) + \log \sigma_{ML}^2 + 1) + (d + 2) \log N.$$

- AIC и BIC в таком примере будут, скорее всего, выбирать примерно одну и ту же модель, хотя разница между ними всё-таки есть:



# 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  $\mathbf{x}$ .

- 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.



- 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

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

- 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.

- 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  $h(x) = \frac{1}{\sqrt{2\pi}}$ ,

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

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

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

СПАСИБО!

Спасибо за внимание!