

ИНФОРМАЦИОННЫЕ КРИТЕРИИ

Сергей Николенко

СПбГУ — Санкт-Петербург

01 ноября 2022 г.

Random facts:

- 1 ноября в Мексике и других странах — День мёртвых (El Día de Muertos); ещё индейцы майя и ацтеки приносили дары богине Миктлансиуатль и сооружали стены с изображением черепов — цомпантли; а в России 1 ноября — День судебного пристава
- 1 ноября 1478 г. Сикст IV направил Изабелле Кастильской и Фердинанду Арагонскому буллу, разрешившую учреждение Испанской инквизиции, а 1 ноября 1512 г. был впервые показан публике расписанный Микеланджело потолок Сикстинской капеллы
- 1 ноября 1800 г. Джон Адамс первым из президентов переехал в Исполнительный особняк (позднее переименованный в Белый дом)
- 1 ноября 1851 г. Огюст Мариэтт обнаружил в Саккаре под развалинами Серапеума вырубленные в скале грандиозные погребальные катакомбы священных быков — Аписов, вестников бога Птаха
- 1 ноября 1938 г. Seabiscuit победил War Admiral в «матче века»; по ходу гонки он отставал, но на финише опередил на четыре корпуса

БАЙЕСОВСКИЙ ИНФОРМАЦИОННЫЙ КРИТЕРИЙ

- Мы хотим сравнить несколько моделей $\mathcal{M}_1, \dots, \mathcal{M}_K$ с наборами параметров $\theta_1, \dots, \theta_K$ на наборе данных D , т.е. сравнить между собой $p(\mathcal{M}_k|D)$:

$$p(\mathcal{M}_k|D) \propto p(\mathcal{M}_k) p(D|\mathcal{M}_k).$$

- Будем полагать $p(\mathcal{M}_k)$ равномерными. А $p(D|\mathcal{M}_k)$ — это как раз знаменатель теоремы Байеса:

$$p(\theta_k|D, \mathcal{M}_k) = \frac{p(\theta_k|\mathcal{M}_k) p(D|\theta_k, \mathcal{M}_k)}{p(D|\mathcal{M}_k)}.$$

- Нам нужно оценить интеграл

$$p(D) = \int p(\theta) p(\theta|D) d\theta = \int p(\theta) e^{\ell(\theta)} d\theta,$$

где $\ell(\theta) = \log p(\theta|D)$.

- Применим лапласовскую аппроксимацию в окрестности точки максимума правдоподобия θ_{ML} :

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

где

$$J(\theta_{\text{ML}}) = -\frac{1}{N} \left. \frac{\partial^2 \log p(\theta|D)}{\partial \theta \partial \theta^\top} \right|_{\theta_{\text{ML}}}.$$

- Аналогічно можна розкласти априорне розподілення окрестности θ_{ML} , но там не пропадєт член первого порядка, поэтому давайте им и ограничимся:

$$p(\theta) \approx p(\theta_{ML}) + (\theta - \theta_{ML})^\top \nabla_{\theta} p(\theta)|_{\theta_{ML}}.$$

- Итого получается, что

$$p(D) \approx \int \left(p(\theta_{ML}) + (\theta - \theta_{ML})^\top \nabla_{\theta} p(\theta)|_{\theta_{ML}} \right) \times \\ \times e^{\ell(\theta_{ML}) - \frac{N}{2}(\theta - \theta_{ML})^\top J(\theta_{ML})(\theta - \theta_{ML})} d\theta.$$

- Но теперь можно заметить, что

$$\int (\theta - \theta_{\text{ML}}) e^{-\frac{N}{2}(\theta - \theta_{\text{ML}})^\top J(\theta_{\text{ML}})(\theta - \theta_{\text{ML}})} d\theta = 0,$$

потому что это величина, пропорциональная матожиданию $(\theta - \theta_{\text{ML}})$ по гауссиану со средним $(\theta - \theta_{\text{ML}})$ и матрицей ковариаций $J(\theta_{\text{ML}})^{-1}$.

- А значит, наша аппроксимация превращается в

$$p(D) \approx e^{\ell(\theta_{\text{ML}})} p(\theta_{\text{ML}}) \int e^{-\frac{N}{2}(\theta - \theta_{\text{ML}})^\top J(\theta_{\text{ML}})(\theta - \theta_{\text{ML}})} d\theta.$$

- Интеграл теперь можно взять — из него получится нормировочная константа для того же самого гауссиана:

$$p(D) \approx e^{\ell(\theta_{\text{ML}})} p(\theta_{\text{ML}}) (2\pi)^{\frac{d}{2}} N^{-\frac{d}{2}} (\det J(\theta_{\text{ML}}))^{-\frac{1}{2}}, \quad \text{или}$$

$$\log p(D) \approx \ell(\theta_{\text{ML}}) - \frac{d}{2} \log N + \log p(\theta_{\text{ML}}) - \frac{1}{2} \log \det J(\theta_{\text{ML}}) + \frac{d}{2} \log(2\pi).$$

- Выбросим всё, что не растёт с N , и умножим на -2 ; получится *байесовский информационный критерий* (Bayesian information criterion, BIC), он же *критерий Шварца* (Schwartz criterion):

$$\text{BIC}(\mathcal{M}) = -2 \log p(D|\theta_{\text{ML}}, \mathcal{M}) + d \log N,$$

где d — это размерность вектора θ , или число свободных параметров в модели \mathcal{M} .

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 $p_{\text{data}}(\mathbf{x})$:

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

Идея и вывод АИС

- 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 θ_{ML} на том же датасете \mathbf{X} , который используется в этой оценке.

Идея и вывод АИС

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

- θ_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 θ_0 — максимум ожидания логарифма правдоподобия по p_{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_{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

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

Идея и вывод АИС

- 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 $\theta = \theta_0$

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

Идея и вывод AIC

- 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:

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

- Давайте попробуем оценить $b(p_{\text{data}})$:

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

- Будем оценивать слагаемые по отдельности.

- Проще всего оценить B_2 , потому что в нём нет θ_{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 θ_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:

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

- Для оценки B_1 нужно провернуть аналогичный трюк с $\ell(\theta) = \log p(X|\theta)$, разложив его вокруг своего максимума θ_{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_{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:

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

- А если можно всё-таки предполагать, что истинное распределение данных p_{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 в этом разделе 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

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

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

ЭМПИРИЧЕСКИЙ БАЙЕС

- Откуда берутся гиперпараметры?
- Оказывается, их тоже можно оптимизировать!
- У линейной регрессии, например, два гиперпараметра: $\beta = \frac{1}{\sigma^2}$ и α (точность регуляризатора, пусть гребневого).
- Давайте просто попробуем оптимизировать $p(D | \alpha, \beta)$ (marginal likelihood).

- Получается:

$$p(D | \alpha, \beta) = \int p(\mathbf{w})p(D | \mathbf{w})d\mathbf{w},$$

$$\ln p(D | \alpha, \beta) = \left(\frac{\beta}{2\pi}\right)^{\frac{N}{2}} \left(\frac{\alpha}{2\pi}\right)^{\frac{d}{2}} \int e^{-\frac{\beta}{2}\|\mathbf{y}-X\mathbf{w}\|^2 - \frac{\alpha}{2}\mathbf{w}^\top\mathbf{w}}d\mathbf{w}.$$

- Выделяем полный квадрат так же, как раньше:

$$A = \beta X^\top X + \alpha \mathbf{I},$$

$$\mu_N = \beta A^{-1} X^\top \mathbf{y}.$$

- Теперь

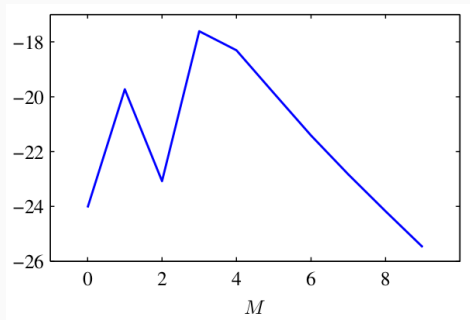
$$\int e^{-\frac{1}{2}(\mathbf{w}-\mu_N)^\top A(\mathbf{w}-\mu_N)} d\mathbf{w} = (2\pi)^{\frac{d}{2}} \sqrt{\det A^{-1}}.$$

- Получается:

$$\ln p(D | \alpha, \beta) = \frac{d}{2} \ln \alpha + \frac{N}{2} \ln \beta - \frac{\beta}{2} \|\mathbf{y} - X\mu_N\|^2 - \frac{\alpha}{2} \mu_N^\top \mu_N - \frac{1}{2} \ln \det A - \frac{N}{2} \ln(2\pi).$$

- Это теперь надо максимизировать по α и β , а можно и разные d перебирать, если речь идёт о том, как выбрать оптимальное число признаков.

- Пример графика по числу параметров:



- А как оптимизировать?

- Обозначим через λ_i собственные числа матрицы $\beta\mathbf{X}^\top\mathbf{X}$:

$$(\beta\mathbf{X}^\top\mathbf{X}) \mathbf{u}_i = \lambda_i \mathbf{u}_i$$

для некоторых собственных векторов \mathbf{u}_i .

- Тогда \mathbf{u}_i будут являться и собственными векторами матрицы \mathbf{A} , с собственными числами $\alpha + \lambda_i$:

$$\mathbf{A}\mathbf{u}_i = (\beta\mathbf{X}^\top\mathbf{X} + \alpha\mathbf{I}) \mathbf{u}_i = \lambda_i \mathbf{u}_i + \alpha \mathbf{u}_i.$$

- Теперь будем оптимизировать логарифм маргинального правдоподобия $\log p(\mathbf{y}|\mathbf{X}, \alpha, \beta)$, взяв производную по α :

$$\frac{\partial \log \det \mathbf{A}}{\partial \alpha} = \frac{\partial \log \prod_i (\alpha + \lambda_i)}{\partial \alpha} = \sum_i \frac{\partial \log(\alpha + \lambda_i)}{\partial \alpha} = \sum_i \frac{1}{\alpha + \lambda_i},$$

и вся производная будет равна

$$\frac{\partial \log p(\mathbf{y}|\mathbf{X}, \alpha, \beta)}{\partial \alpha} = \frac{M}{2\alpha} - \frac{1}{2} \mu_N^\top \mu_n - \frac{1}{2} \sum_i \frac{1}{\alpha + \lambda_i}.$$

- Приравняем производную нулю:

$$\alpha \mu_N^\top \mu_N = M - \alpha \sum_i \frac{1}{\alpha + \lambda_i} = \sum_i \frac{\lambda_i}{\alpha + \lambda_i}.$$

- Получается итеративный процесс

$$\alpha^{(k+1)} = \frac{1}{\mu_N(\alpha^{(k)})^\top \mu_N(\alpha^{(k)})} \sum_i \frac{\lambda_i}{\alpha^{(k)} + \lambda_i}.$$

- Аналогично по β : можно заметить, что

$$\frac{\partial \lambda_i}{\partial \beta} = \frac{\lambda_i}{\beta}.$$

- Значит,

$$\frac{\partial \log \det \mathbf{A}}{\partial \beta} = \frac{\partial \log \prod_i (\alpha + \lambda_i)}{\partial \beta} = \sum_i \frac{\partial \log(\alpha + \lambda_i)}{\partial \beta} = \frac{1}{\beta} \sum_i \frac{\lambda_i}{\alpha + \lambda_i}, \text{ и}$$

$$\frac{\partial \log p(\mathbf{y}|\mathbf{X}, \alpha, \beta)}{\partial \beta} = \frac{N}{2\beta} - \frac{1}{2} \sum_{n=1}^N (t_n - \mu_n^\top \mathbf{x}_n)^2 - \frac{1}{2\beta} \sum_i \frac{\lambda_i}{\alpha + \lambda_i}.$$

- Итого по β получаем

$$\beta^{(k+1)} = \frac{N - \sum_i \frac{\lambda_i}{\alpha^{(k)} + \lambda_i}}{\sum_{n=1}^N \left(t_n - \mu_n \left(\alpha^{(k)}, \beta^{(k)} \right)^\top \mathbf{x}_n \right)^2}.$$

Спасибо за внимание!