

A decorative graphic in the top-left corner of the slide, consisting of a grid of colored squares. The grid is 4 squares wide and 4 squares high. The colors of the squares are: Row 1: Teal, Orange, Brown, Teal. Row 2: Orange, Brown, Light Brown, Teal. Row 3: Orange, Teal, Light Brown, Brown. Row 4: Light Brown, Orange, Orange, Brown.

Введение в Байесовские методы

Виктор Кантор

План

Байесовская вероятность и вывод

Байесовский вывод в ML моделях

Принцип наибольшей
обоснованности

Методы оценки обоснованности

Байесовская вероятность и вывод

Частотный подход

Случайность – объективная
неопределенность

- Единственное возможное средство анализа – проведение серии испытаний
- Вероятность – предел частоты наступления события

Байесовский подход

Случайность – мера нашего
незнания. Есть априорное и
апостериорное «незнание»

- Параметры – случайные величины
- При размере выборки в 0 объектов – используем априорное распределение параметра, при размере в 1 и больше – можем получить апостериорное

Критика Байесовского подхода

1. Не предлагаются конструктивные методы вывода априорной вероятности
2. Методы долго работают, часто приводят к вычислительно сложным задачам интегрирования в многомерном пространстве

Что нам дают байесовские методы

1. Способ формализации «здорового смысла»
2. Автоматизация подбора параметров и гиперпараметров (при этом возникают новые гиперпараметры, но более абстрактные)
3. Язык для вывода методов ML из простых базовых предположений

Байесовский вывод

$$\boldsymbol{x} = (x_1, \dots, x_n)$$

Байесовский вывод

$$\mathbf{x} = (x_1, \dots, x_n)$$

$$p(\theta|\mathbf{x}) = \frac{p(\mathbf{x}|\theta)p(\theta)}{p(\mathbf{x})} = \frac{p(\mathbf{x}|\theta)p(\theta)}{\int p(\mathbf{x}|\theta)p(\theta)d\theta}$$

Байесовский вывод

$$\mathbf{x} = (x_1, \dots, x_n)$$

$$p(\theta|\mathbf{x}) = \frac{p(\mathbf{x}|\theta)p(\theta)}{p(\mathbf{x})} = \frac{p(\mathbf{x}|\theta)p(\theta)}{\int p(\mathbf{x}|\theta)p(\theta)d\theta}$$

Процесс получения апостериорной вероятности из априорной и правдоподобия выборки – и есть байесовский вывод

Точечная оценка параметра

$$\hat{\theta}_B = \int \theta p(\theta | \mathbf{x}) d\theta$$

$$\hat{\theta}_{MP} = \arg \max P(\theta | \mathbf{x})$$

Пример: оценка вероятности выпадения орла

Задача – по серии из n подбрасываний монетки и m выпадений орла оценить вероятность q выпадения орла

Пример: оценка вероятности выпадения орла

Задача – по серии из n подбрасываний монетки и m выпадений орла оценить вероятность q выпадения орла

$$p(m|n, q) = C_n^m q^m (1 - q)^{n-m} \sim \mathcal{B}(m|n, q)$$

Пример: оценка вероятности выпадения орла

Задача – по серии из n подбрасываний монетки и m выпадений орла оценить вероятность q выпадения орла

$$p(m|n, q) = C_n^m q^m (1 - q)^{n-m} \sim \mathcal{B}(m|n, q)$$

Вопрос: какую оценку q даст метод максимального правдоподобия?

Пример: оценка вероятности выпадения орла

Байесовская оценка:

$$p(m|n, q) = C_n^m q^m (1 - q)^{n-m} \sim \mathcal{B}(m|n, q)$$

$$p(q|a, b) = \frac{\Gamma(a + b)}{\Gamma(a)\Gamma(b)} q^{a-1} (1 - q)^{b-1} \sim \text{Beta}(q|a, b)$$

Пример: оценка вероятности выпадения орла

Байесовская оценка:

$$p(m|n, q) = C_n^m q^m (1 - q)^{n-m} \sim \mathcal{B}(m|n, q)$$

$$p(q|a, b) = \frac{\Gamma(a + b)}{\Gamma(a)\Gamma(b)} q^{a-1} (1 - q)^{b-1} \sim \text{Beta}(q|a, b)$$

$$p(q|\text{«}m \text{ орлов}\text{»}) \sim \text{Beta}(q|a + m, b + n - m)$$

Пример: оценка вероятности выпадения орла

Байесовская оценка:

$$p(m|n, q) = C_n^m q^m (1 - q)^{n-m} \sim \mathcal{B}(m|n, q)$$

$$p(q|a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} q^{a-1} (1-q)^{b-1} \sim \text{Beta}(q|a, b)$$

$$p(q|\text{«}m \text{ орлов}\text{»}) \sim \text{Beta}(q|a+m, b+n-m)$$

$$\hat{q}_B = \int_0^1 p(q|\text{«}m \text{ орлов}\text{»}) q dq = \frac{m+1}{n+2}$$

Сопряженные распределения

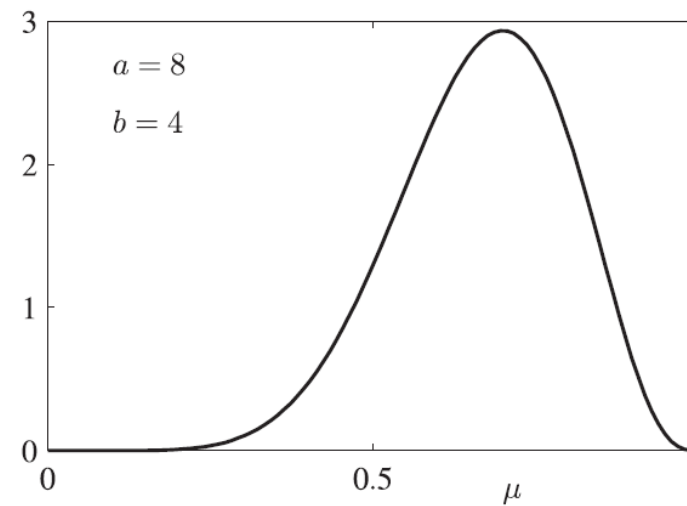
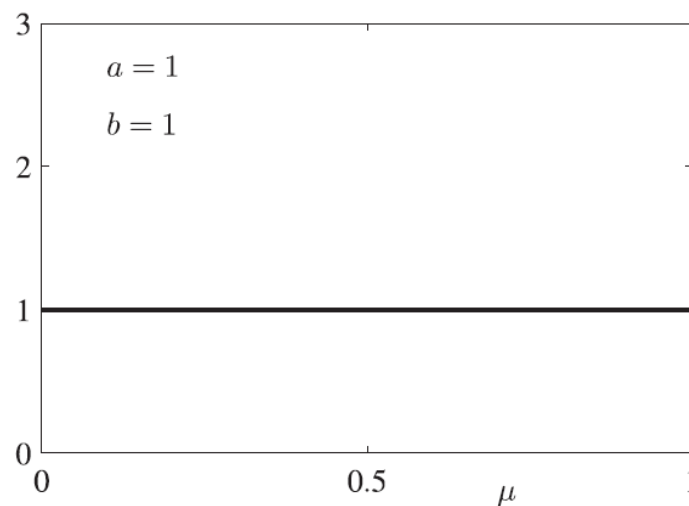
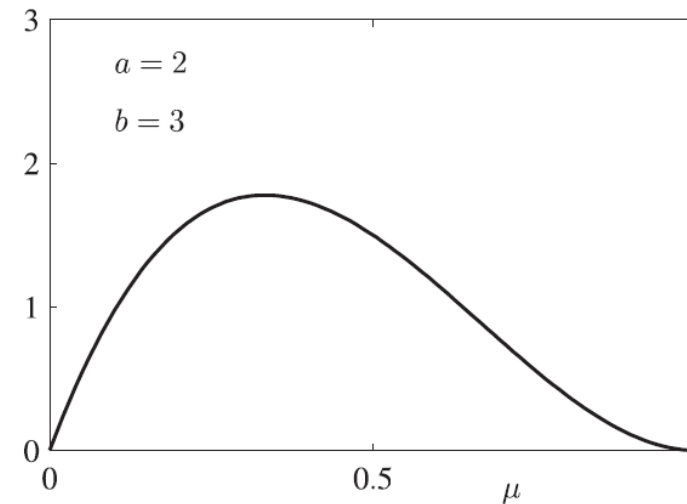
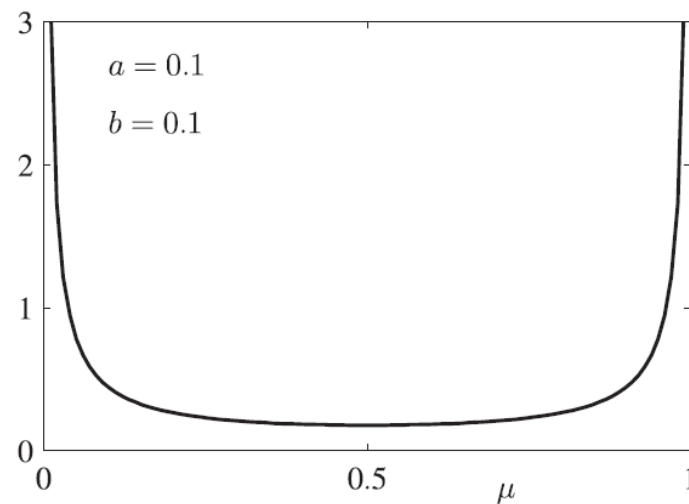
Если из $p(\theta) \sim \mathcal{A}(\alpha_0)$ и $p(\boldsymbol{x}|\theta) \sim \mathcal{B}(\beta)$

следует $p(\theta|\boldsymbol{x}) \sim \mathcal{A}(\alpha_1)$, то говорят, что распределение

\mathcal{A} сопряженное к \mathcal{B}

Бета-распределение

Пример:
распределение,
сопряженное к
биномиальному



Еще про сопряженные распределения

- Сопряженное распределение можно выписать явно для любого распределения из экспоненциального семейства:

$$p(\boldsymbol{x}|\boldsymbol{\alpha}) = h(\boldsymbol{x})g(\boldsymbol{\alpha}) \exp(\boldsymbol{\alpha}^T u(\boldsymbol{x}))$$

- К этому семейству относятся нормальное, гамма-, бета-, равномерное, Бернулли, Дирихле, хи-квадрат, Пуассоновское и другие распределения
- Именно сопряженное распределение лучше брать как априорное

Байесовский вывод в ML моделях

Оценка параметров модели по MLE

$$\theta_{ML} = \operatorname{argmax}_{\theta} p_{\theta}(X) = \operatorname{argmax}_{\theta} \sum_{i=1}^n \ln p_{\theta}(x_i)$$

Байесовская оценка параметров модели

$$\begin{aligned}\theta_{MP} &= \operatorname{argmax}_{\theta} p(\theta|X) = \operatorname{argmax}_{\theta} p(X|\theta)p(\theta) = \\ &\operatorname{argmax}_{\theta} \sum_{i=1}^n \ln p(x_i|\theta) + \ln p(\theta)\end{aligned}$$

Байесовская оценка параметров модели

$$\begin{aligned}\theta_{MP} &= \operatorname{argmax}_{\theta} p(\theta|X) = \operatorname{argmax}_{\theta} p(X|\theta)p(\theta) = \\ &\operatorname{argmax}_{\theta} \sum_{i=1}^n \ln p(x_i|\theta) + \boxed{\ln p(\theta)}\end{aligned}$$

Регуляризатор для
функции правдоподобия

Пример: априор в линейной регрессии

Линейная регрессия:

$$(X, \mathbf{t}) = \{\mathbf{x}_i, t_i\}_{i=1}^n \quad t = f(\mathbf{x}) + \varepsilon \quad \varepsilon \sim \mathcal{N}(\varepsilon|0, \sigma^2)$$

$$f(\mathbf{x}) = \sum_{j=1}^m w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$$

Пример: априор в линейной регрессии

MLE:
$$p(\mathbf{t}|X, \mathbf{w}) = \prod_{i=1}^n p(t_i|\mathbf{x}_i, \mathbf{w}) = \prod_{i=1}^n \mathcal{N}(t_i|f(\mathbf{x}_i, \mathbf{w}), \sigma^2) =$$

$$\prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(t_i - \mathbf{w}^T \phi(\mathbf{x}_i))^2}{2\sigma^2}\right) = \frac{1}{(\sqrt{2\pi}\sigma)^n} \exp\left(-\frac{\sum_{i=1}^n (t_i - \mathbf{w}^T \phi(\mathbf{x}_i))^2}{2\sigma^2}\right)$$

$$-\frac{1}{2\sigma^2} \sum_{i=1}^n (t_i - \mathbf{w}^T \phi(\mathbf{x}_i))^2 = -\frac{1}{2\sigma^2} (\mathbf{t} - \Phi \mathbf{w})^T (\mathbf{t} - \Phi \mathbf{w}) \rightarrow \max_{\mathbf{w}}$$

$$\mathbf{w}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

Пример: априор в линейной регрессии

Байес:

$$w_{MP} = \arg \max_w p(\mathbf{w}|X, \mathbf{t}) = \arg \max_w p(\mathbf{t}|X, \mathbf{w})p(\mathbf{w})$$

$$p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|0, \alpha^{-1}I)$$

$$-\frac{1}{2\sigma^2} \sum_{i=1}^n (t_i - \phi(\mathbf{x}_i))^2 - \frac{\alpha}{2} \|\mathbf{w}\|^2 \rightarrow \max_{\mathbf{w}}$$

$$\mathbf{w}_{MP} = (\sigma^{-2}\Phi^T\Phi + \alpha I)^{-1}\sigma^{-2}\Phi^T\mathbf{t}$$

Принцип наибольшей обоснованности

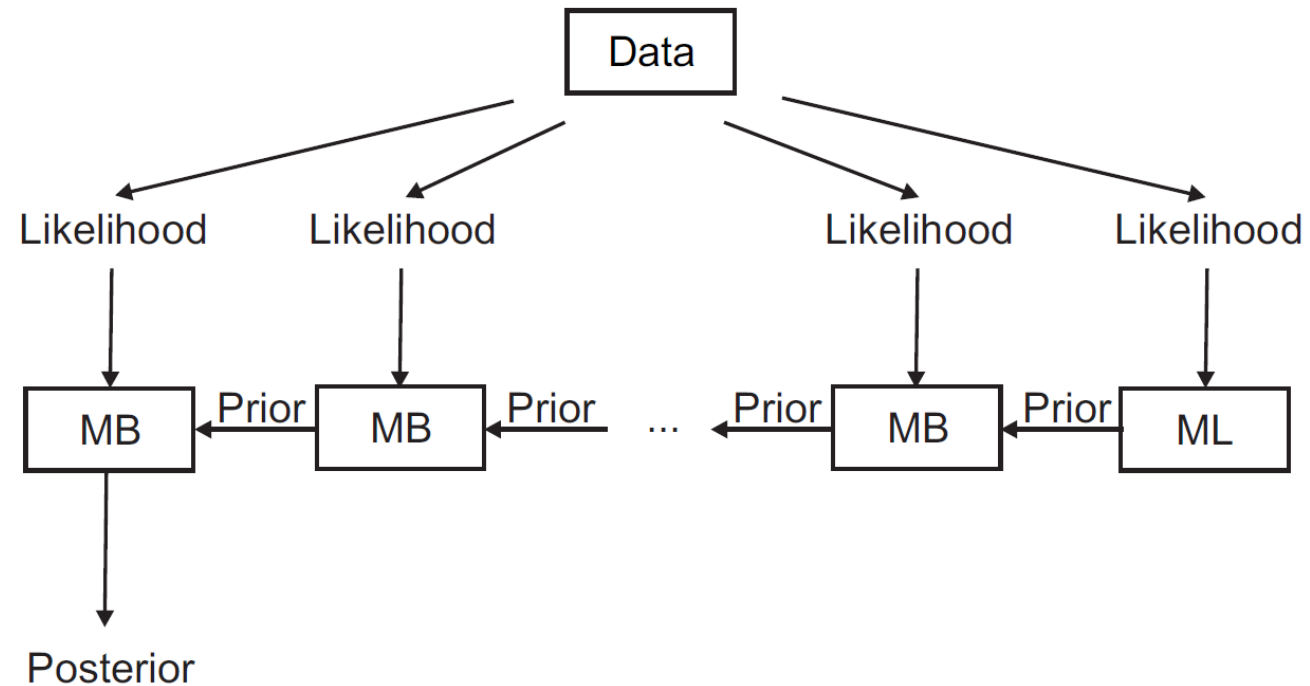
Иерархическая схема Байеса

Распределение параметра тоже можно задать в параметрическом виде: $p(\theta) = p(\theta|\alpha)$

$$p(\alpha|\mathbf{x}) = \frac{p(\mathbf{x}|\alpha)p(\alpha)}{p(\mathbf{x})} = \frac{p(\mathbf{x}|\alpha)p(\alpha)}{\int p(\mathbf{x}|\alpha)p(\alpha)d\alpha}$$

$$p(\mathbf{x}|\alpha) = \int p(\mathbf{x}|\theta)p(\theta|\alpha)d\theta \text{ - обоснованность (evidence) модели}$$

Иерархическая схема Байеса



Обычно ограничиваются двухуровневым выводом

Пример: генератор случайных чисел

- Есть генератор случайных чисел от 1 до N и два числа, сгенерированные им: 6 и 8.
- Определить матожидание генерируемых чисел
- Определить, $N = 10$ или $N = 100$

Пример: генератор случайных чисел

- Есть генератор случайных чисел от 1 до N и два числа, сгенерированные им: 6 и 8.
- Определить матожидание генерируемых чисел
- Определить, $N = 10$ или $N = 100$

Какой ответ на первый вопрос даст классическая максимизация правдоподобия?

Пример: генератор случайных чисел

Байесовский взгляд:

$$D(\mathbf{q}|\boldsymbol{\alpha}) = \frac{1}{B(\boldsymbol{\alpha}^0)} q_1^{\alpha_1^0-1} \dots q_N^{\alpha_N^0-1}, \quad \sum_{i=1}^N q_i = 1, \quad q_i \geq 0 \quad p(x_1, x_2|\mathbf{q}) = q_8 q_6$$

По формуле Байеса:

$$p(\mathbf{q}|x_1, x_2) = \frac{1}{Z} q_1^0 \dots q_5^0 q_6^1 q_7^0 q_8^1 q_9^0 \dots q_N^0 = D(\mathbf{q}|\boldsymbol{\alpha}^1)$$

Пример: генератор случайных чисел

$$\mathbb{E}q_i = \frac{\alpha_i}{\sum_{j=1}^N \alpha_j}$$

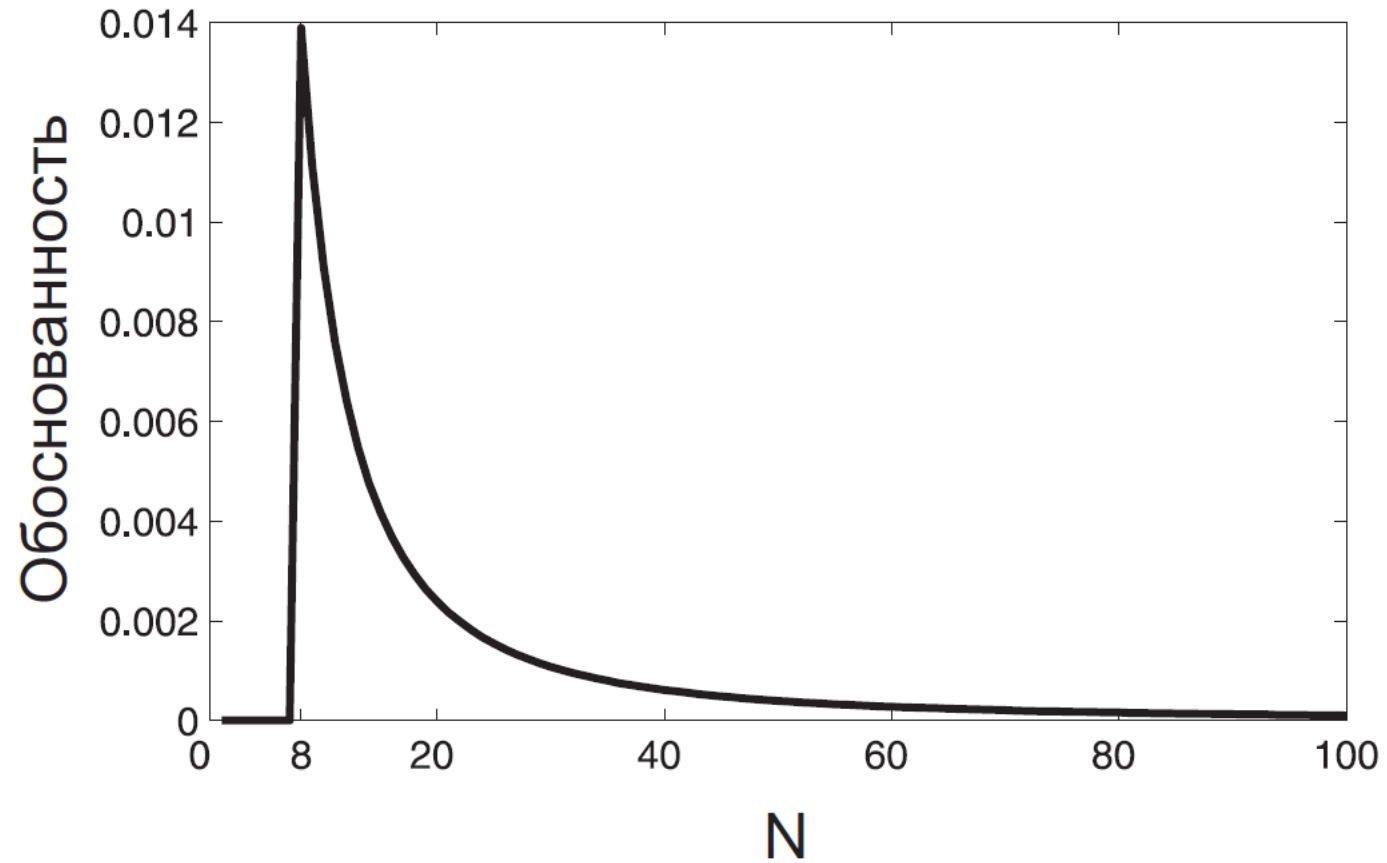
$$q_6 = q_8 = \frac{2}{12} = \frac{1}{6} \approx 0.16, \quad q_i = \frac{1}{12} \approx 0.08 \quad \forall i \neq 6, 8$$

$$\mu_{MP}(N = 10) = 5.75$$

$$q_6 = q_8 = \frac{2}{102} = \frac{1}{51} \approx 0.02, \quad q_i = \frac{1}{102} \approx 0.01 \quad \forall i \neq 6, 8$$

$$\mu_{MP}(N = 100) \approx 49.65$$

Пример: генератор случайных чисел



Метод релевантных векторов в регрессии

Более сложное априорное распределение, чем было в Ridge Regression + Байесовский вывод:

$$p(\boldsymbol{w}|\boldsymbol{\alpha}) = \mathcal{N}(0, A^{-1})$$

$$p(\boldsymbol{t}|\boldsymbol{\alpha}, \sigma^2) = \int p(\boldsymbol{t}|X, \boldsymbol{w}, \sigma^2)p(\boldsymbol{w}|\boldsymbol{\alpha})d\boldsymbol{w} \rightarrow \max_{\boldsymbol{\alpha}, \sigma^2}$$

Метод релевантных векторов в регрессии

Обоснованность:

$$p(\boldsymbol{t}|\boldsymbol{\alpha}, \sigma^2) = \int p(\boldsymbol{t}|X, \boldsymbol{w}, \sigma^2)p(\boldsymbol{w}|\boldsymbol{\alpha})d\boldsymbol{w} = \int Q(\boldsymbol{w})d\boldsymbol{w}$$

Метод релевантных векторов в регрессии

Обоснованность:

$$p(\mathbf{t}|\boldsymbol{\alpha}, \sigma^2) = \int p(\mathbf{t}|X, \mathbf{w}, \sigma^2)p(\mathbf{w}|\boldsymbol{\alpha})d\mathbf{w} = \int Q(\mathbf{w})d\mathbf{w}$$

$$L(\mathbf{w}) = \log Q(\mathbf{w})$$

$$L(\mathbf{w}) = L(\mathbf{w}_{MP}) + (\nabla_{\mathbf{w}}L(\mathbf{w}_{MP}))^T(\mathbf{w} - \mathbf{w}_{MP}) + \frac{1}{2}(\mathbf{w} - \mathbf{w}_{MP})^T H(\mathbf{w} - \mathbf{w}_{MP})$$

$$\mathbf{w}_{MP} = \arg \max_{\mathbf{w}} L(\mathbf{w}) \Rightarrow \nabla_{\mathbf{w}}L(\mathbf{w}_{MP}) = 0$$

$$H = \nabla \nabla L(\mathbf{w}_{MP})$$

Метод релевантных векторов в регрессии

$$L(\boldsymbol{w}) = L(\boldsymbol{w}_{MP}) + (\nabla_{\boldsymbol{w}} L(\boldsymbol{w}_{MP}))^T (\boldsymbol{w} - \boldsymbol{w}_{MP}) + \frac{1}{2} (\boldsymbol{w} - \boldsymbol{w}_{MP})^T H (\boldsymbol{w} - \boldsymbol{w}_{MP})$$

$$\boldsymbol{w}_{MP} = \arg \max_{\boldsymbol{w}} L(\boldsymbol{w}) \Rightarrow \nabla_{\boldsymbol{w}} L(\boldsymbol{w}_{MP}) = 0$$

$$H = \nabla \nabla L(\boldsymbol{w}_{MP})$$

тогда:

$$\begin{aligned} \int Q(\boldsymbol{w}) d\boldsymbol{w} &= \int \exp \left(L(\boldsymbol{w}_{MP}) + \frac{1}{2} (\boldsymbol{w} - \boldsymbol{w}_{MP})^T H (\boldsymbol{w} - \boldsymbol{w}_{MP}) \right) d\boldsymbol{w} = \\ &= Q(\boldsymbol{w}_{MP}) \sqrt{(2\pi)^m} \sqrt{\det((-H)^{-1})} = \sqrt{(2\pi)^m} \frac{Q(\boldsymbol{w}_{MP})}{\sqrt{\det(-H)}} \end{aligned}$$

Метод релевантных векторов в регрессии

$$L(\mathbf{w}) = L(\mathbf{w}_{MP}) + (\nabla_{\mathbf{w}} L(\mathbf{w}_{MP}))^T (\mathbf{w} - \mathbf{w}_{MP}) + \frac{1}{2} (\mathbf{w} - \mathbf{w}_{MP})^T H (\mathbf{w} - \mathbf{w}_{MP})$$

$$\mathbf{w}_{MP} = \arg \max_{\mathbf{w}} L(\mathbf{w}) \Rightarrow \nabla_{\mathbf{w}} L(\mathbf{w}_{MP}) = 0$$

$$H = \nabla \nabla L(\mathbf{w}_{MP})$$

тогда:

$$\begin{aligned} \int Q(\mathbf{w}) d\mathbf{w} &= \int \exp \left(L(\mathbf{w}_{MP}) + \frac{1}{2} (\mathbf{w} - \mathbf{w}_{MP})^T H (\mathbf{w} - \mathbf{w}_{MP}) \right) d\mathbf{w} = \\ &= Q(\mathbf{w}_{MP}) \sqrt{(2\pi)^m} \sqrt{\det((-H)^{-1})} = \sqrt{(2\pi)^m} \frac{Q(\mathbf{w}_{MP})}{\sqrt{\det(-H)}} \end{aligned}$$

«Взяли» интеграл малыми усилиями, осталось
посчитать Q в точке оптимума и оценить
нормировочную константу

Метод релевантных векторов в регрессии

Получение МР оценки на \mathbf{w} :

$$\beta = \sigma^{-2}$$

$$\begin{aligned} L(\mathbf{w}) = & -\frac{1}{2}\beta(\mathbf{t} - \Phi\mathbf{w})^T(\mathbf{t} - \Phi\mathbf{w}) - \frac{1}{2}\mathbf{w}^T A \mathbf{w} - \frac{n}{2} \log(2\pi) - \\ & - \frac{m}{2} \log(2\pi) + \frac{1}{2} \log \det(A) = -\frac{1}{2}\beta[\mathbf{t}^T \mathbf{t} - 2\mathbf{w}^T \Phi^T \mathbf{t} + \mathbf{w}^T \Phi^T \Phi \mathbf{w}] + C \end{aligned}$$

$$\nabla L(\mathbf{w}) = -\frac{1}{2}\beta(-2\Phi^T \mathbf{t} + 2\Phi^T \Phi \mathbf{w}) - A \mathbf{w} = 0 \Rightarrow \boxed{\mathbf{w}_{MP} = (\beta\Phi^T \Phi + A)^{-1} \beta\Phi^T \mathbf{t}}$$

Метод релевантных векторов в регрессии

Итоговое выражение

$$p(\mathbf{t}|X, \boldsymbol{\alpha}, \sigma^2) = \int p(\mathbf{t}|X, \sigma^2)p(\mathbf{w}|\boldsymbol{\alpha})d\mathbf{w} = \sqrt{(2\pi)^m} \frac{Q(\mathbf{w}_{MP})}{\sqrt{\det(-H)}} =$$
$$\frac{1}{\sqrt{(2\pi)^n \det(\beta^{-1}I + \Phi A^{-1}\Phi^T)^{1/2}}} \exp\left(-\frac{1}{2}\mathbf{t}^T (\beta^{-1}I + \Phi A^{-1}\Phi^T)^{-1}\mathbf{t}\right)$$

Метод релевантных векторов в регрессии

Итерационные формулы для параметров (получим приравняв производные обоснованности по ним к нулю):

$$\alpha_i^{new} = \frac{\gamma_i}{w_{MP,i}^2} \quad \gamma_i = 1 - \alpha_i^{old} \Sigma_{ii}$$

$$(\sigma^2)^{new} = \frac{\|\mathbf{t} - \Phi \mathbf{w}\|^2}{n - \sum_{i=1}^m \gamma_i}$$

$$\Sigma = (\beta \Phi^T \Phi + A)^{-1} \quad \mathbf{w}_{MP} = \beta \Sigma \Phi^T \mathbf{t}$$

Метод релевантных векторов в регрессии

Получение прогноза:

$$p(t_*|\mathbf{x}_*, \mathbf{t}, X) = \int p(t_*|\mathbf{x}_*, \mathbf{w}, \sigma_{MP}^2) p(\mathbf{w}|\mathbf{t}, X, \boldsymbol{\alpha}_{MP}, \sigma_{MP}^2) d\mathbf{w} = \mathcal{N}(t_*|y_*, \sigma_*^2)$$

$$y_* = \mathbf{w}_{MP}^T \boldsymbol{\phi}(\mathbf{x}_*)$$

$$\sigma_*^2 = \sigma_{MP}^2 + \boldsymbol{\phi}(\mathbf{x}_*)^T \Sigma \boldsymbol{\phi}(\mathbf{x}_*)$$

Метод релевантных векторов в регрессии

Алгоритм 1: Метод релевантных векторов для задачи регрессии

Вход: Обучающая выборка $\{\mathbf{x}_i, t_i\}_{i=1}^n$ $\mathbf{x}_i \in \mathbb{R}^d$, $t_i \in \mathbb{R}$; Матрица обобщенных признаков $\Phi = \{\phi_j(\mathbf{x}_i)\}_{i,j=1}^{n,m}$;

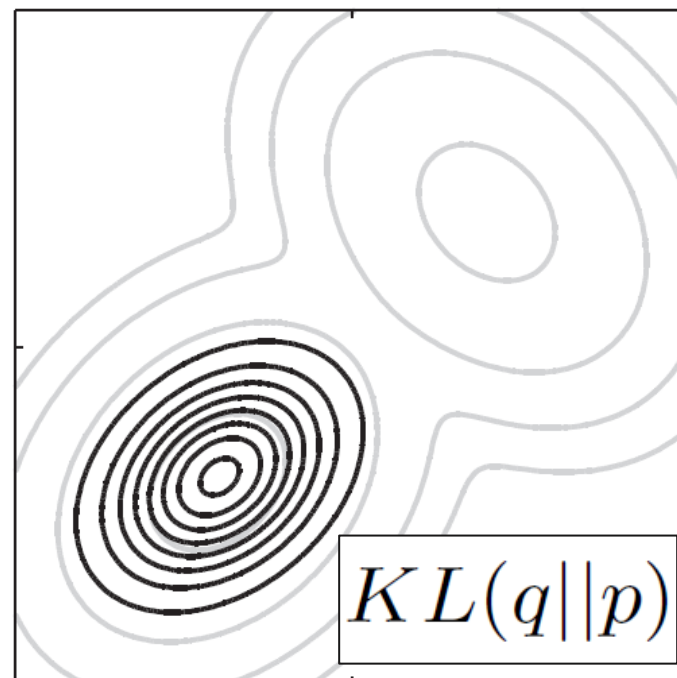
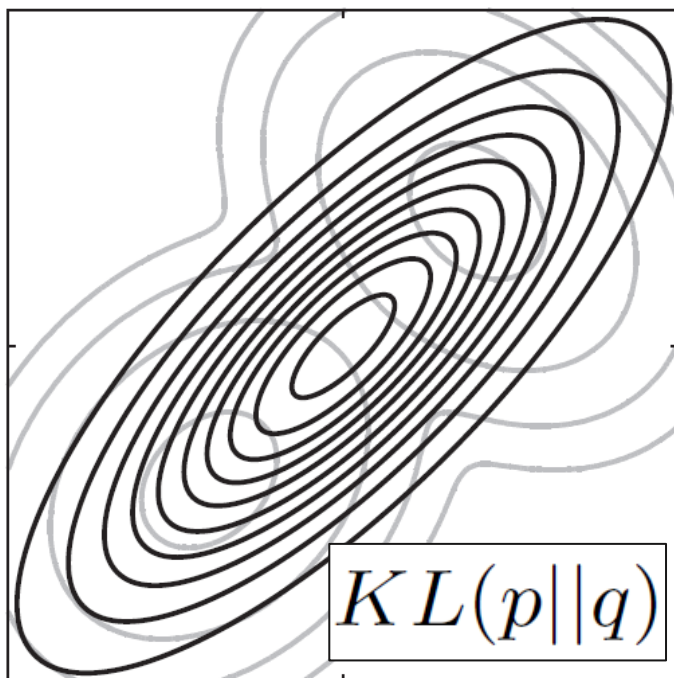
Выход: Набор весов \mathbf{w} , матрица Σ и оценка дисперсии шума β^{-1} для решающего правила $t_*(\mathbf{x}) = \sum_{j=1}^m w_j \phi_j(\mathbf{x})$, $\sigma_*^2(\mathbf{x}) = \beta^{-1} + \phi^T(\mathbf{x}_*) \Sigma \phi(\mathbf{x}_*)$;

- 1: инициализация: $\alpha_i := 1$, $i = 1, \dots, m$, $\beta := 1$, AlphaBound $:= 10^{12}$, WeightBound $:= 10^{-6}$, NumberOfIterations $:= 100$;
 - 2: **для** $k = 1, \dots, \text{NumberOfIterations}$
 - 3: $A := \text{diag}(\alpha_1, \dots, \alpha_m)$;
 - 4: $\Sigma := (\beta \Phi^T \Phi + A)^{-1}$;
 - 5: $\mathbf{w}_{MP} := \Sigma \beta \Phi^T \mathbf{t}$;
 - 6: **для** $j = 1, \dots, m$
 - 7: **если** $w_{MP,j} < \text{WeightBound}$ или $\alpha_j > \text{AlphaBound}$ **то**
 - 8: $w_{MP,j} := 0$, $\alpha_j := +\infty$, $\gamma_j := 0$;
 - 9: **иначе**
 - 10: $\gamma_j := 1 - \alpha_j \Sigma_{jj}$, $\alpha_j := \frac{\gamma_j}{w_{MP,j}^2}$;
 - 11: $\beta := \frac{n - \sum_{j=1}^m \gamma_j}{\|\mathbf{t} - \Phi \mathbf{w}_{MP}\|^2}$
-

Методы оценки обоснованности

Дивергенция Кульбака-Лейблера

$$KL(q||p) = - \int q(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x}$$



Задача оценки обоснованности

$$p(Z|X) = \frac{p(X, Z)}{p(X)}$$

$p(X)$ - обоснованность выбранной модели

Прямое интегрирование обычно невозможно, поэтому ограничиваются приближением $p(Z|X)$ некоторым распределением $q(Z)$

Разложение обоснованности

$$\begin{aligned}\log p(X) &= \log p(X) \int q(Z) dZ = \int \log p(X) q(Z) dZ = \\ &\int \log \frac{p(X, Z)}{p(Z|X)} q(Z) dZ = \int \log \frac{p(X, Z) q(Z)}{q(Z) p(Z|X)} q(Z) dZ = \\ &\int \log \frac{p(X, Z)}{q(Z)} q(Z) dZ - \int \log \frac{p(Z|X)}{q(Z)} q(Z) dZ = \mathcal{L}(q) + KL(q||p)\end{aligned}$$

Факторизованное приближение

Будем использовать для приближения $q(Z) = \prod_{i=1}^k q_i(z_i)$

$$\begin{aligned}\mathcal{L}(q) &= \int \prod_i q_i \left(\log p(X, Z) - \sum_i \log q_i \right) dZ = \\ &= \int q_j \left(\int \log p(X, Z) \prod_{i \neq j} q_i dz_i \right) dz_j - \int q_j \log q_j dz_j + C\end{aligned}$$

Обозначим $\log \tilde{p}(X, z_j) = \mathbb{E}_{i \neq j} \log p(X, Z) = \int \log p(X, Z) \prod_{i \neq j} q_i dz_i$

Факторизованное приближение

Будем использовать для приближения $q(Z) = \prod_{i=1}^k q_i(z_i)$

$$\begin{aligned}\mathcal{L}(q) &= \int \prod_i q_i \left(\log p(X, Z) - \sum_i \log q_i \right) dZ = \\ &= \int q_j \left(\int \log p(X, Z) \prod_{i \neq j} q_i dz_i \right) dz_j - \int q_j \log q_j dz_j + C\end{aligned}$$

Обозначим $\log \tilde{p}(X, z_j) = \mathbb{E}_{i \neq j} \log p(X, Z) = \int \log p(X, Z) \prod_{i \neq j} q_i dz_i$

$$\mathcal{L}(q) = \int q_j \log \frac{\tilde{p}(X, z_j)}{q_j} dz_j + C = -KL(q || \tilde{p}) + C$$

Факторизованное приближение

Обозначим $\log \tilde{p}(X, z_j) = \mathbb{E}_{i \neq j} \log p(X, Z) = \int \log p(X, Z) \prod_{i \neq j} q_i dz_i$

$$\mathcal{L}(q) = \int q_j \log \frac{\tilde{p}(X, z_j)}{q_j} dz_j + C = -KL(q||\tilde{p}) + C$$

$$\log q_j^*(z_j) = \mathbb{E}_{i \neq j} \log p(X, Z) + C$$

Пример: вариационная линейная регрессия

$$p(\mathbf{t}|\mathbf{w}) = \prod_{i=1}^n \mathcal{N}(t_i | \mathbf{w}^T \phi(\mathbf{x}_i), \beta^{-1}) \quad p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w} | 0, \alpha^{-1} I)$$

$$p(\alpha) = \mathcal{G}(\alpha | a_0, b_0)$$

Будем искать приближение $p(\mathbf{w}, \alpha | \mathbf{t})$ в виде $q(\mathbf{w}, \alpha) = q(\mathbf{w})q(\alpha)$

Пример: вариационная линейная регрессия

$$p(\mathbf{t}|\mathbf{w}) = \prod_{i=1}^n \mathcal{N}(t_i | \mathbf{w}^T \phi(\mathbf{x}_i), \beta^{-1}) \quad p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w} | 0, \alpha^{-1} I)$$

$$p(\alpha) = \mathcal{G}(\alpha | a_0, b_0) \quad q(\mathbf{w}, \alpha) = q(\mathbf{w})q(\alpha)$$

$$\begin{aligned} \log q^*(\alpha) &= \mathbb{E}_{\mathbf{w}} \log p(\mathbf{t}, \mathbf{w}, \alpha) = \\ \mathbb{E}_{\mathbf{w}} (\log p(\mathbf{w}|\alpha)p(\alpha)) + C &= \mathbb{E}_{\mathbf{w}} \log p(\mathbf{w}|\alpha) + \log p(\alpha) + C = \\ \frac{m}{2} \log \alpha - \frac{\alpha}{2} \mathbb{E} \mathbf{w}^T \mathbf{w} &+ (a_0 - 1) \log \alpha - b_0 \alpha + C_1 \end{aligned}$$

Пример: вариационная линейная регрессия

$$\alpha \sim \mathcal{G}(\alpha | a_n, b_n) \qquad a_n = a_0 + \frac{m}{2}, \quad b_n = b_0 + \frac{1}{2} \mathbb{E} \mathbf{w}^T \mathbf{w}$$

Аналогично для \mathbf{w} :

$$\begin{aligned} \log q^*(\mathbf{w}) &= \mathbb{E}_\alpha \log p(\mathbf{t}, \mathbf{w}, \alpha) = \mathbb{E}_\alpha \log (p(\mathbf{t} | \mathbf{w}) p(\mathbf{w} | \alpha) p(\alpha)) = \log p(\mathbf{t} | \mathbf{w}) + \mathbb{E}_\alpha \log p(\mathbf{w} | \alpha) + C = \\ &= -\frac{\beta}{2} \sum_{i=1}^n (\mathbf{w}^T \phi(\mathbf{x}_i) - t_i)^2 - \frac{1}{2} \mathbb{E}_\alpha \cdot \mathbf{w}^T \mathbf{w} + C_1 = -\frac{1}{2} \mathbf{w}^T (\mathbb{E}_\alpha I + \beta \Phi^T \Phi) \mathbf{w} + \beta \mathbf{w}^T \Phi^T \mathbf{t} + C_2. \end{aligned}$$

В итоге:

$$\mathbf{w} \sim \mathcal{N}(\mathbf{w} | \boldsymbol{\mu}_n, S_n) \qquad \boldsymbol{\mu}_n = \beta S_n \Phi^T \mathbf{t}, \quad S_n = (\mathbb{E}_\alpha I + \beta \Phi^T \Phi)^{-1}$$

Другой подход: методы Монте-Карло

$$\int_a^b f(x)dx \approx \frac{b-a}{n} \sum_{i=1}^n f(x_i) = \hat{f}, \quad x_i \sim U[a, b]$$

Другой подход: методы Монте-Карло

$$\int_a^b f(x)dx \approx \frac{b-a}{n} \sum_{i=1}^n f(x_i) = \hat{f}, \quad x_i \sim U[a, b]$$

$$\int_D f(\mathbf{x})p(\mathbf{x})d\mathbf{x} \approx \frac{|D|}{n} \sum f(\mathbf{x}_i)p(\mathbf{x}_i), \quad \mathbf{x} \sim U(D)$$

Другой подход: методы Монте-Карло

$$\int_a^b f(x)dx \approx \frac{b-a}{n} \sum_{i=1}^n f(x_i) = \hat{f}, \quad x_i \sim U[a, b]$$

$$\int_D f(\mathbf{x})p(\mathbf{x})d\mathbf{x} \approx \frac{|D|}{n} \sum f(\mathbf{x}_i)p(\mathbf{x}_i), \quad \mathbf{x} \sim U(D)$$

Такой подход для вероятностных интегралов оказывается не слишком эффективным (как вы думаете, почему?)

Оценка математического ожидания сэмплированием

$$\int f(\boldsymbol{x})p(\boldsymbol{x})d\boldsymbol{x} \approx \frac{1}{n} \sum f(\boldsymbol{x}_i), \quad \boldsymbol{x} \sim p(\boldsymbol{x})$$

Оценка матожидания сэмплированием

$$\int f(\boldsymbol{x})p(\boldsymbol{x})d\boldsymbol{x} \approx \frac{1}{n} \sum f(\boldsymbol{x}_i), \quad \boldsymbol{x} \sim p(\boldsymbol{x})$$

Главная наша проблема теперь – как сэмплировать из распределения p или хотя бы чего-то похожего

Модификации МС методов

- Сэмплирование из другого распределения

$$\mathbb{E}_p f = \int f(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} = \frac{1}{Z_p} \int f(\mathbf{x}) \frac{\tilde{p}(\mathbf{x})}{q(\mathbf{x})} q(\mathbf{x}) d\mathbf{x} \approx$$

$$\frac{1}{n Z_p} \sum_{i=1}^n f(\mathbf{x}_i) \frac{\tilde{p}(\mathbf{x}_i)}{q(\mathbf{x}_i)} = \frac{1}{n \sum_{i=1}^n r_i} \sum_{i=1}^n f(\mathbf{x}_i) r_i, \quad \mathbf{x} \sim q(\mathbf{x})$$

- МСМС (Markov Chain Monte Carlo)
- Гибридные МС методы