

Вариационные приближения

Сергей Николенко

НИУ ВШЭ — Санкт-Петербург

30 мая 2020 г.

Random facts:

- 30 мая — Всемирный день черепахи, учреждённый в 2000 году по инициативе Американского общества спасения черепах
- 30 мая 1416 г. инквизиция сожгла чешского реформатора Иеронима Пражского, а 30 мая 1431 г. в Руане сожгли Жанну д'Арк
- 30 мая 1816 г. доктор Хосе Гаспар Родригес де Франсия, «Великий правитель» и «Отец парагвайской нации», был провозглашён вечным диктатором
- 30 мая 1896 г. в Нью-Йорке произошло первое в мире ДТП с участием автомобиля, в котором сломал ногу велосипедист Ивлин Томас; в Москве в тот же день произошла давка на Ходынском поле, в которой погибли 1389 человек
- 30 мая 1953 г. в *Nature* вышла статья Джеймса Уотсона и Френсиса Крика, в которой они предложили двойную спираль как структурную модель ДНК

ЕМ в общем виде

- Часто нужно оценивать $p(Z | X)$ для латентных переменных Z и данных X .
- Но это слишком сложно! Один вариант — сэмплировать из $p(Z | X)$.
- Другой вариант — лапласовские приближения, но тоже нечасто работают.
- Давайте решать в общем виде.

Обоснование алгоритма EM

- Вспомним сначала формальное обоснование алгоритма EM.
- Мы решаем задачу максимизации правдоподобия по данным $\mathbf{X} = \{x_1, \dots, x_N\}$.

$$L(\theta | \mathbf{X}) = p(\mathbf{X} | \theta) = \prod p(x_i | \theta)$$

или, что то же самое, максимизации $\ell(\theta | \mathbf{X}) = \log L(\theta | \mathbf{X})$.

- EM может помочь, если этот максимум трудно найти аналитически.

Обоснование алгоритма EM

- Давайте предположим, что в данных есть *скрытые компоненты*, такие, что если бы мы их знали, задача была бы проще.
- Замечание: совершенно не обязательно эти компоненты должны иметь какой-то физический смысл. :) Может быть, так просто удобнее.
- В любом случае, получается набор данных $Z = (X, Y)$ с совместной плотностью

$$p(z \mid \theta) = p(x, y \mid \theta) = p(y \mid x, \theta)p(x \mid \theta).$$

- Получается полное правдоподобие $L(\theta \mid Z) = p(X, Y \mid \theta)$. Это случайная величина (т.к. Y неизвестно).

Обоснование алгоритма EM

- Заметим, что настоящее правдоподобие $L(\theta) = E_Y [p(X, \mathcal{Y} \mid \theta) \mid X, \theta]$.
- E-шаг алгоритма EM вычисляет условное ожидание (логарифма) полного правдоподобия при условии X и текущих оценок параметров θ_n :

$$Q(\theta, \theta_n) = E [\log p(X, \mathcal{Y} \mid \theta) \mid X, \theta_n] .$$

- Здесь θ_n – текущие оценки, а θ – неизвестные значения (которые мы хотим получить в конечном счёте); т.е. $Q(\theta, \theta_n)$ – это функция от θ .

Обоснование алгоритма EM

- E-шаг алгоритма EM вычисляет условное ожидание (логарифма) полного правдоподобия при условии X и текущих оценок параметров θ :

$$Q(\theta, \theta_n) = E [\log p(X, \mathcal{Y} \mid \theta) \mid X, \theta_n] .$$

- Условное ожидание – это

$$E [\log p(X, \mathcal{Y} \mid \theta) \mid X, \theta_n] = \int_y \log p(X, y \mid \theta) p(y \mid X, \theta_n) dy ,$$

где $p(y \mid X, \theta_n)$ – маргинальное распределение скрытых компонентов данных.

- EM лучше всего применять, когда это выражение легко подсчитать, может быть, даже аналитически.
- Вместо $p(y \mid X, \theta_n)$ можно подставить $p(y, X \mid \theta_n) = p(y \mid X, \theta_n)p(X \mid \theta_n)$, от этого ничего не изменится.

Обоснование алгоритма EM

- В итоге после E-шага алгоритма EM мы получаем функцию $Q(\theta, \theta_n)$.
- На M-шаге мы максимизируем

$$\theta_{n+1} = \arg \max_{\theta} Q(\theta, \theta_n).$$

- Затем повторяем процедуру до сходимости.
- В принципе, достаточно просто находить θ_{n+1} , для которого $Q(\theta_{n+1}, \theta_n) > Q(\theta_n, \theta_n)$ – Generalized EM.
- Осталось понять, что значит $Q(\theta, \theta_n)$ и почему всё это работает.

Обоснование алгоритма EM

- Мы хотим перейти от θ_n к θ , для которого $\ell(\theta) > \ell(\theta_n)$.

$$\begin{aligned}\ell(\theta) - \ell(\theta_n) &= \\&= \log \left(\int_y p(\mathbf{X} | y, \theta) p(y | \theta) dy \right) - \log p(\mathbf{X} | \theta_n) = \\&= \log \left(\int_y p(y | \mathbf{X}, \theta_n) \frac{p(\mathbf{X} | y, \theta) p(y | \theta)}{p(y | \mathbf{X}, \theta_n)} dy \right) - \log p(\mathbf{X} | \theta_n) \geq \\&\geq \int_y p(y | \mathbf{X}, \theta_n) \log \left(\frac{p(\mathbf{X} | y, \theta) p(y | \theta)}{p(y | \mathbf{X}, \theta_n)} \right) dy - \log p(\mathbf{X} | \theta_n) = \\&= \int_y p(y | \mathbf{X}, \theta_n) \log \left(\frac{p(\mathbf{X} | y, \theta) p(y | \theta)}{p(\mathbf{X} | \theta_n) p(y | \mathbf{X}, \theta_n)} \right) dy.\end{aligned}$$

- Получили

$$\begin{aligned}\ell(\theta) &\geq l(\theta, \theta_n) = \\ &= \ell(\theta_n) + \int_y p(y | \mathbf{X}, \theta_n) \log \left(\frac{p(\mathbf{X} | y, \theta)p(y | \theta)}{p(\mathbf{X} | \theta_n)p(y | \mathcal{X}, \theta_n)} \right) dy.\end{aligned}$$

- Мы нашли нижнюю оценку на $\ell(\theta)$ везде, касание происходит в точке θ_n .

Вариационные приближения

- Вариационный вывод: функционалы, производные по функциям... в общем, можно оптимизировать функционалы.
- Для нас это значит, что можно оптимизировать приближение q из какого-то класса к заданному p .
- Пусть есть $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ и $\mathbf{Z} = \{\mathbf{z}_1, \dots, \mathbf{z}_N\}$.
- Мы знаем $p(\mathbf{X}, \mathbf{Z})$ из модели, хотим найти приближение для $p(\mathbf{Z} | \mathbf{X})$ и $p(\mathbf{X})$.

- Как и в EM:

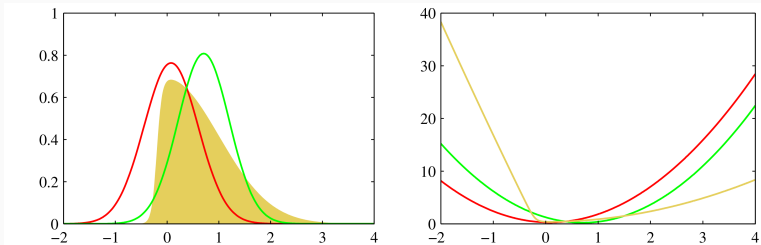
$$\ln p(X) = \mathcal{L}(q) + \text{KL}(q\|p), \text{ где}$$

$$\mathcal{L}(q) = \int q(Z) \ln \frac{p(X, Z)}{q(Z)} dZ,$$

$$\text{KL}(q\|p) = - \int q(Z) \ln \frac{p(Z | X)}{q(Z)} dZ.$$

- $\mathcal{L}(q)$ — это вариационная нижняя оценка, её можно теперь максимизировать, и KL будет автоматически минимизироваться.

- Пример – сравним с лапласовским:



- Если $q(\mathcal{Z})$ произвольное, то мы просто получим $q(\mathcal{Z}) = p(\mathcal{Z} | \mathbf{X})$; но это вряд ли получится.
- Будем ограничивать.

Факторизуемые распределения

- Главный частный случай — пусть $Z = Z_1 \sqcup \dots \sqcup Z_M$, и

$$q(Z) = \prod_{i=1}^M q_i(Z_i).$$

- Но больше никаких предположений! В этом прелесть — оптимизируем сразу функции!
- Это соответствует теории среднего поля в физике (mean field theory).

Факторизуемые распределения

- Тогда:

$$\begin{aligned}\mathcal{L}(q) &= \int \prod_i q_i \left(\ln p(\mathbf{X}, \mathbf{Z}) - \sum_i \ln q_i \right) d\mathbf{Z} \\ &= \int q_j \left[\int \ln p(\mathbf{X}, \mathbf{Z}) \prod_{i \neq j} q_i d\mathbf{Z}_i \right] d\mathbf{Z}_j - \int q_j \ln q_j d\mathbf{Z}_j + \text{const} \\ &= \int q_j \ln \tilde{p}(\mathbf{X}, \mathbf{Z}_j) d\mathbf{Z}_j - \int q_j \ln q_j d\mathbf{Z}_j + \text{const},\end{aligned}$$

где $\ln \tilde{p}(\mathbf{X}, \mathbf{Z}_j) = \mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{Z})] + \text{const}$.

- Как максимизировать теперь $\mathcal{L}(q)$ по q_j ?

Факторизуемые распределения

- Надо заметить, что мы получили там KL-дивергенцию между $q_j(Z_j)$ и $\tilde{p}(X, Z_j)$.
- Т.е. оптимальное решение получится при

$$\ln q_j^*(Z_j) = \mathbb{E} [\ln p(X, Z)] + \text{const.}$$

- Константа здесь просто для нормализации.
- Оказывается, достаточно взять ожидание от логарифма совместного распределения.
- Но явных формул не получается, потому что ожидание считается по остальным $q_i^*, i \neq j$.
- И всё-таки обычно что-то можно сделать; давайте рассмотрим примеры.

- Первый пример — приблизим двумерный гауссиан одномерными:

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z} \mid \boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}),$$

$$\boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \quad \boldsymbol{\Lambda} = \begin{pmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{21} & \Lambda_{22} \end{pmatrix}.$$

- И мы хотим приблизить $q(\mathbf{z}) = q_1(z_1)q_2(z_2)$.
- Повычисляем...

- ...получится, что

$$\ln q_1^*(z_1) = -\frac{1}{2}z_1^2\Lambda_{11} + z_1\mu_{11}\Lambda_{11} - z_1\Lambda_{12}(E[z_2] - \mu_2) + \text{const.}$$

- Чудесным образом получился гауссиан! Сам собой, без предположений.
- Найдём его среднее и дисперсию...

- ...получится

$$q_1^*(z_1) = \mathcal{N}(z_1 \mid m_1, \Lambda_{11}^{-1}), \text{ где}$$

$$m_1 = \mu_1 - \Lambda_{11}^{-1} \Lambda_{12} (E[z_2] - \mu_2).$$

- Аналогично,

$$q_2^*(z_2) = \mathcal{N}(z_2 \mid m_2, \Lambda_{22}^{-1}), \text{ где}$$

$$m_2 = \mu_2 - \Lambda_{22}^{-1} \Lambda_{21} (E[z_1] - \mu_1).$$

- Какое решение у этой системы?

- Да просто $E[z_1] = m_1 = \mu_1$, $E[z_2] = m_2 = \mu_2$.
- А если бы мы минимизировали $KL(p||q)$, получилось бы

$$KL(p||q) = - \int p(\mathbf{Z}) \left[\sum_i \ln q_i(\mathbf{Z}_i) \right] d\mathbf{Z} + \text{const},$$

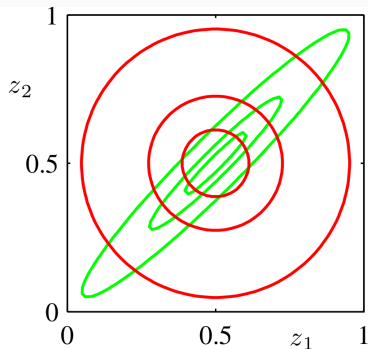
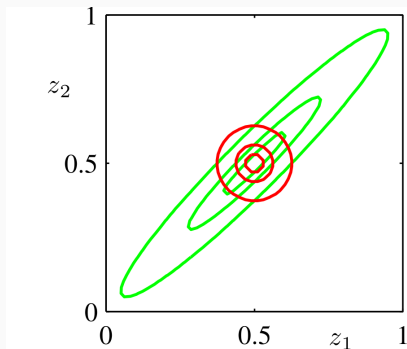
и можно оптимизировать по отдельности:

$$q_j^*(\mathbf{Z}_j) = \int p(\mathbf{Z}) \prod_{i \neq j} d\mathbf{Z}_i = p(\mathbf{Z}_j).$$

- Т.е. просто маргинализация.
- Почему бы так и не сделать? В чём разница?

Разные KL-дивергенции

- Разные дисперсии ответа:

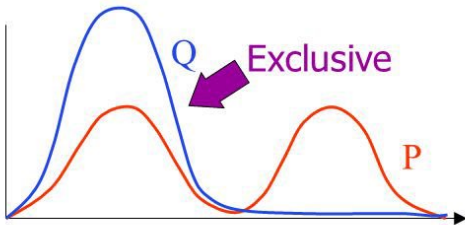


Разные KL-дивергенции

- Минимизация $KL(p||q)$ накрывает всю p , а $KL(q||p)$ строит всю q в регионе больших p :

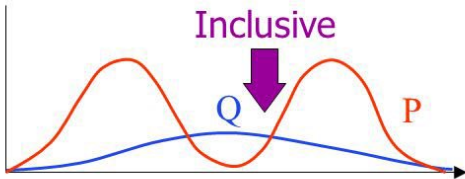
Minimising
 $KL(Q||P)$

$$= \sum_H Q(H) \ln \frac{Q(H)}{P(H|V)}$$



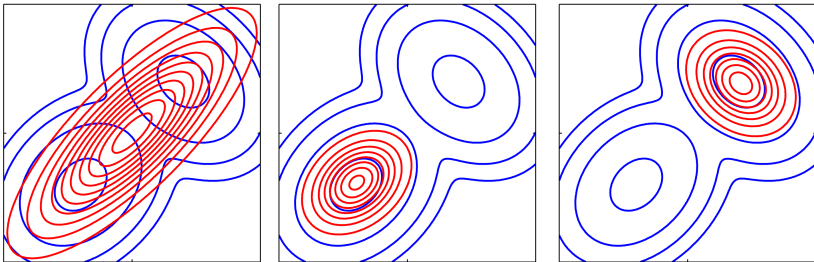
Minimising
 $KL(P||Q)$

$$= \sum_H P(H|V) \ln \frac{P(H|V)}{Q(H)}$$



Разные KL-дивергенции

- Например, для двумерного гауссиана:



- В машинном обучении гораздо интереснее, конечно, пик найти.

Вариационное приближение для гауссиана

Одномерный гауссиан

- И ещё пример: давайте найдём параметры одномерного гауссиана по точкам $\mathbf{X} = \{x_1, \dots, x_N\}$. Правдоподобие:

$$p(\mathbf{X} \mid \mu, \tau) = \left(\frac{\tau}{2\pi}\right)^{N/2} e^{-\frac{\tau}{2} \sum_{n=1}^N (x_n - \mu)^2}.$$

- Вводим сопряжённые априорные распределения:

$$\begin{aligned} p(\mu \mid \tau) &= \mathcal{N}(\mu \mid \mu_0, (\lambda_0 \tau)^{-1}), \\ p(\tau) &= \text{Gamma}(\tau \mid a_0, b_0). \end{aligned}$$

- Мы это только что подсчитали точно, но давайте приблизим теперь апостериорное распределение как

$$q(\mu, \tau) = q_\mu(\mu) q_\tau(\tau).$$

- На самом деле так не раскладывается!
- Это то, что мы делали для $q(\mathbf{Z}) = \prod_{i=1}^M q_i(\mathbf{Z}_i)$. Посчитаем...

Одномерный гауссиан

- ... $q_\mu(\mu)$ – гауссиан с параметрами

$$\mu_N = \frac{\lambda_0 \mu_0 + N \bar{x}}{\lambda_0 + N}, \quad \lambda_N = (\lambda_0 + N) \mathbb{E}[\tau].$$

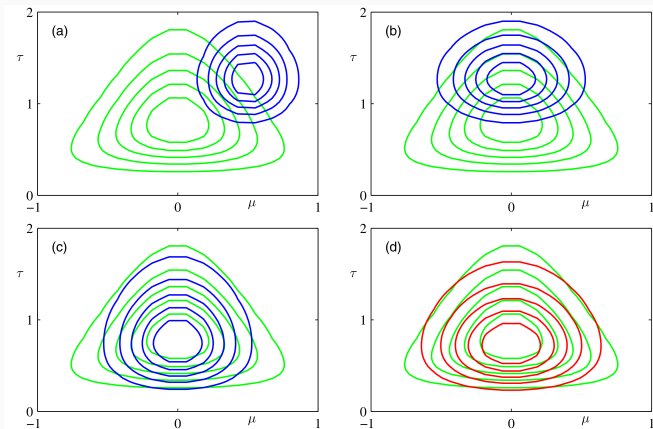
- А $q_\tau(\tau)$ – гамма-распределение с параметрами

$$a_N = a_0 + \frac{N+1}{2}, \quad b_N = b_0 + \frac{1}{2} \mathbb{E}_\mu \left[\sum_n (x_n - \mu)^2 + \lambda_0 (\mu - \mu_0)^2 \right].$$

- Всё получилось как надо, но без предположений о форме q_τ и q_μ .

Одномерный гауссиан

- Вот такой вывод в пространстве (μ, τ) :



- А для $\mu_0 = a_0 = b_0 = \lambda_0 = 0$ (non-informative priors) можно и точно посчитать...

- Получатся моменты для μ

$$E[\mu] = \bar{x}, \quad E[\mu^2] = \bar{x}^2 + \frac{1}{NE[\tau]}.$$

- Это можно подставить и найти $E[\tau]$:

$$\frac{1}{E[\tau]} = \frac{1}{N} \sum_{n=1}^N (x_n - \bar{x})^2.$$

Спасибо!

Спасибо за внимание!