

# Вариационные приближения

---

Сергей Николенко

Академия MADE — Mail.Ru

15 мая 2021 г.

---

## *Random facts:*

- 15 мая 1157 г. «великий любитель женщин, сладкой пищи и питья» князь Юрий Долгорукий был отравлен на пиру у киевского боярина Петрилы; а 15 мая 1682 г. его потомок князь Юрий Долгоруков был убит взбунтовавшимися стрельцами
- 15 мая 1703 г. милостью короля Георга II первым премьер-министром Великобритании стал лидер партии вигов, лорд-канцлер Роберт Уолпол
- 15 мая 1942 г. на аэродроме Кольцово (Свердловск) Григорий Бахчиванджи испытал первый советский ракетный самолёт БИ-1; а 15 мая 1954 г. из ворот завода в Рентоне (Вашингтон) вышел самолёт-прототип, позже названный Boeing 707
- 15 мая 2001 г. произошёл «CSX 8888 incident»: никем не управляемый поезд #8888 с 42 вагонами (в том числе с опасными химикатами) в течение двух часов двигался со скоростью около 80 км/ч; в 2010 г. об этом вышел фильм *Unstoppable*
- 15 мая 2009 г. президент Таджикистана Эмомали Рахмон запретил чиновникам появляться на одних с ним портретах

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

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

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

# Обоснование алгоритма EM

- E-шаг алгоритма EM вычисляет условное ожидание (логарифма) полного правдоподобия при условии  $\mathbf{X}$  и текущих оценок параметров  $\theta$ :

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

- Условное ожидание – это

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

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

- EM лучше всего применять, когда это выражение легко подсчитать, может быть, даже аналитически.
- Вместо  $p(y \mid \mathbf{X}, \theta_n)$  можно подставить  $p(y, \mathbf{X} \mid \theta_n) = p(y \mid \mathbf{X}, \theta_n)p(\mathbf{X} \mid \theta_n)$ , 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  $\theta_{n+1}$ , 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.



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

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

- Получили

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

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

# 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:

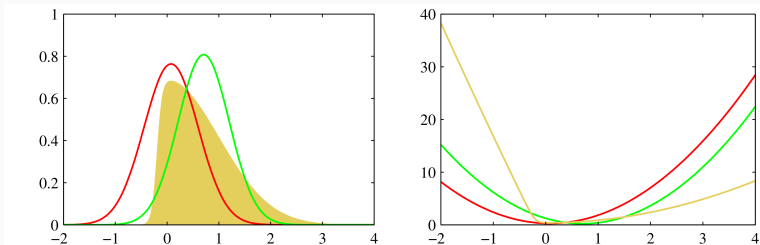
$$\ln p(X) = \mathcal{L}(q) + \text{KL}(q\|p), \text{ где}$$

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

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

- $\mathcal{L}(q)$  — 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.

- 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

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

- 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:

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

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

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

- ...получится, что

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

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

- ...получится

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

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

- Аналогично,

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

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

- 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

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

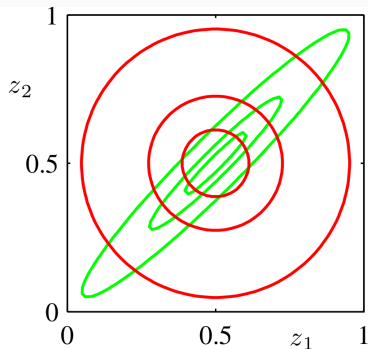
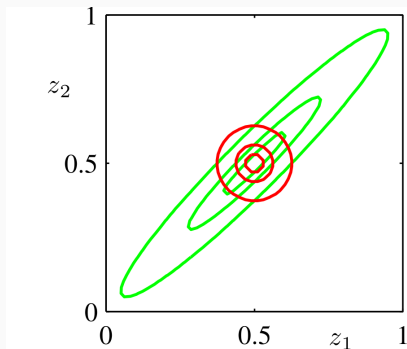
и можно оптимизировать по отдельности:

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

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

# Разные KL-дивергенции

- Разные дисперсии ответа:

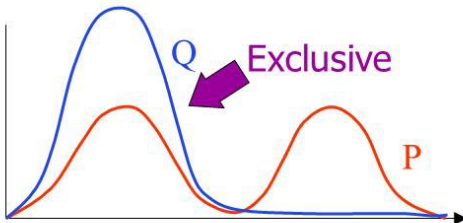


# Разные KL-дивергенции

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

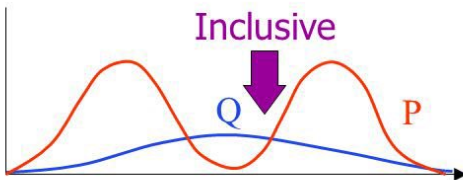
Minimising  
 $KL(Q||P)$

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



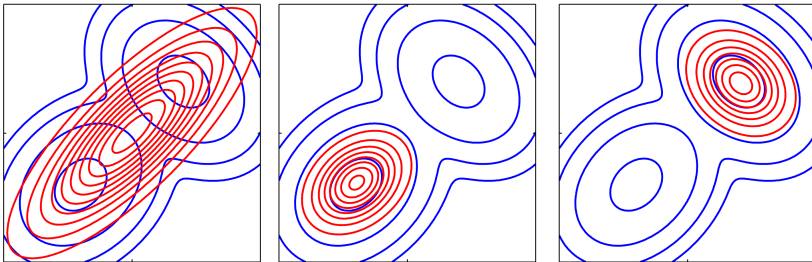
Minimising  
 $KL(P||Q)$

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



# 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...

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

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

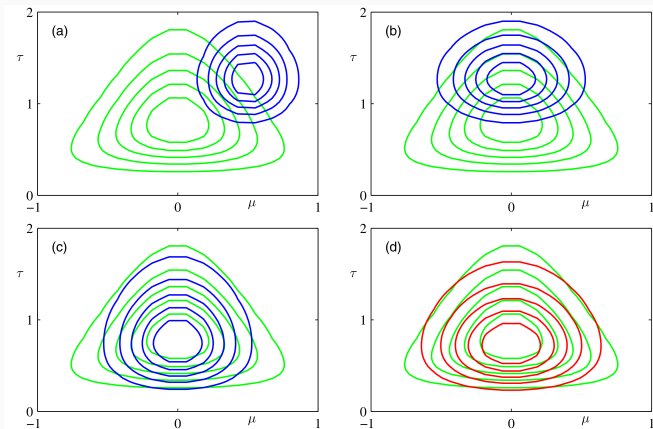
- А  $q_\tau(\tau)$  – гамма-распределение с параметрами

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

- Всё получилось как надо, но без предположений о форме  $q_\tau$  и  $q_\mu$ .

# Одномерный гауссиан

- Вот такой вывод в пространстве  $(\mu, \tau)$ :



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

- Получатся моменты для  $\mu$

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

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

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

# Вариационное приближение для смеси гауссианов

---

# Смесь гауссианов

- Смесь гауссианов:  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ ,  $Z = \{z_1, \dots, z_n\}$ ,

$$p(Z | \boldsymbol{\pi}) = \prod_{n=1}^N \prod_{k=1}^K \pi_k^{z_{nk}},$$

$$p(X | Z, \boldsymbol{\mu}, \Lambda) = \prod_{n=1}^N \prod_{k=1}^K \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \Lambda_k^{-1}).$$

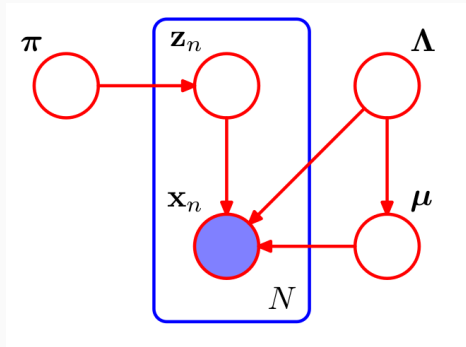
- Выберем сопряжённые априорные распределения:

$$p(\boldsymbol{\pi}) = \text{Dir}(\boldsymbol{\pi} | \boldsymbol{\alpha}_0) = C(\boldsymbol{\alpha}_0) \prod_{k=1}^K \pi_k^{\alpha_0 - 1},$$

$$\begin{aligned} p(\boldsymbol{\mu}, \Lambda) &= p(\boldsymbol{\mu} | \Lambda) p(\Lambda) \\ &= \prod_{k=1}^K \mathcal{N}(\boldsymbol{\mu}_k | \mathbf{m}_0, (\beta_0 \Lambda_k)^{-1}) \mathcal{W}(\Lambda_k | \mathbf{W}_0, \nu_0). \end{aligned}$$

# Смесь гауссианов

- Вот такая графическая модель:



- Распределение Дирихле пусть будет симметричное для простоты; часто ещё  $\mathbf{m}_0 = \mathbf{0}$ .
- Заметьте разницу между латентными переменными и параметрами модели.



# Вариационное приближение

- Теперь вариационное приближение. Сначала сама факторизация:

$$p(X, Z, \pi, \mu, \Lambda) = p(X | Z, \mu, \Lambda) p(Z | \pi) p(\pi) p(\mu | \Lambda) p(\Lambda).$$

- Мы наблюдаем только  $X$ , остальное всё надо как-то оценить.
- Интересно, что единственное предположение про наше вариационное приближение выглядит так:

$$q(Z, \pi, \mu, \Lambda) = q(Z) q(\pi, \mu, \Lambda).$$

- И всё! Дальше всё само собой получится. Но не сразу...

# Вариационное приближение

- Сначала  $q^*(Z)$ :

$$\begin{aligned}\ln q^*(Z) &= \mathbb{E}_{\pi, \mu, \Lambda} [\ln p(X, Z, \pi, \mu, \Lambda)] + \text{const} \\ &= \mathbb{E}_{\pi, \mu, \Lambda} [\ln p(Z | \pi)] + \mathbb{E}_{\mu, \Lambda} [\ln p(X | Z, \mu, \Lambda)] + \text{const} \\ &= \dots = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \ln \rho_{nk} + \text{const},\end{aligned}$$

$$\begin{aligned}\text{где } \ln \rho_{nk} &= \mathbb{E}[\ln \pi_k] + \frac{1}{2} \mathbb{E}[\ln |\Lambda_k|] - \frac{D}{2} \ln(2\pi) - \\ &\quad - \frac{1}{2} \mathbb{E}_{\mu_k, \Lambda_k} [(\mathbf{x}_n - \mu_k)^\top \Lambda_k (\mathbf{x}_n - \mu_k)].\end{aligned}$$

- Нормируем:

$$q^*(Z) = \prod_{n=1}^N \prod_{k=1}^K r_{nk}^{z_{nk}}, \text{ где } r_{nk} = \frac{\rho_{nk}}{\sum_{j=1}^K \rho_{nj}}.$$

# Вариационное приближение

- Теперь  $E[z_{nk}] = r_{nk}$ , т.е.  $r_{nk}$  – то, насколько точка  $\mathbf{x}_n$  принадлежит кластеру  $k$ .
- Можно определить статистики с их учётом, как обычно:

$$N_k = \sum_{n=1}^N r_{nk},$$

$$\bar{\mathbf{x}}_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} \mathbf{x}_n,$$

$$\mathbf{S}_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} (\mathbf{x}_n - \bar{\mathbf{x}}_k)(\mathbf{x}_n - \bar{\mathbf{x}}_k)^\top.$$

- То же самое происходило и в EM-алгоритме.

# Вариационное приближение

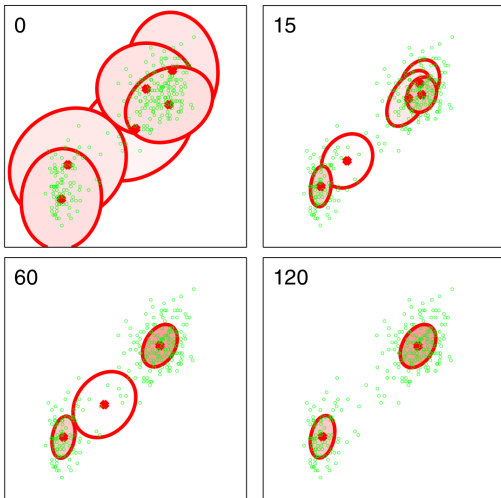
- Теперь  $q^*(\boldsymbol{\pi}, \boldsymbol{\mu}, \Lambda)$ :

$$\begin{aligned}\ln q^*(\boldsymbol{\pi}, \boldsymbol{\mu}, \Lambda) &= \ln p(\boldsymbol{\pi}) + \sum_{k=1}^K \ln p(\boldsymbol{\mu}_k, \Lambda_k) + \mathbb{E}_Z[\ln p(Z \mid \boldsymbol{\pi})] \\ &\quad + \sum_{k=1}^K \sum_{n=1}^N \mathbb{E}[Z_{nk}] \ln \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_k \Lambda_k^{-1}) + \text{const.}\end{aligned}$$

- Вот уже получилось, что  $q^*(\boldsymbol{\pi}, \boldsymbol{\mu}, \Lambda)$  раскладывается в  $q^*(\boldsymbol{\pi})q^*(\boldsymbol{\mu}, \Lambda)$ , опять же без предположений.
- Более того,  $q^*(\boldsymbol{\mu}, \Lambda) = \prod_{k=1}^K q(\boldsymbol{\mu}_k, \Lambda_k)$ .
- И теперь можно по отдельности посчитать (упражнение), получится типичный М-шаг.
- Причём распределения останутся той же формы (т.к. были сопряжённые).

# Вариационное приближение

- Теперь даже model selection автоматически получается, просто у некоторых компонент  $N_k \approx 0$ :



- Никакого оверфиттинга или коллапса компонент.

Спасибо!

Спасибо за внимание!