

# Метод релевантных векторов и снова про EM-алгоритм

---

Сергей Николенко

Samsung AI Center — Москва

12 мая 2020 г.

---

## *Random facts:*

- 12 мая 1311 г. во Франции были казнены 54 видных тамплиера
- 12 мая 1364 г. в Кракове был создан Ягеллонский университет, а 12 мая 1551 г. появился первый университет в обеих Америках — университет Сан-Маркос в Лиме
- 12 мая 1935 в городке Акрон, штат Огайо, был основан первый клуб анонимных алкоголиков
- 12 мая 1976 г. была создана Московская Хельсинкская группа; организатором и первым руководителем был физик, член-корр АН Армянской ССР Юрий Орлов, а о создании было объявлено на пресс-конференции, проведённой в московской квартире академика Андрея Сахарова
- 12 мая 2004 г. в Англиканской церкви была введена должность веб-священника, а 12 мая 2010 г. был запущен первый кириллический домен .рф

# Эмпирический байес

---

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

- Получается:

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

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

- Выделяем полный квадрат так же, как раньше:

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

$$\mathbf{m}_N = \beta A^{-1}\mathbf{X}^T\mathbf{y}.$$

- Теперь

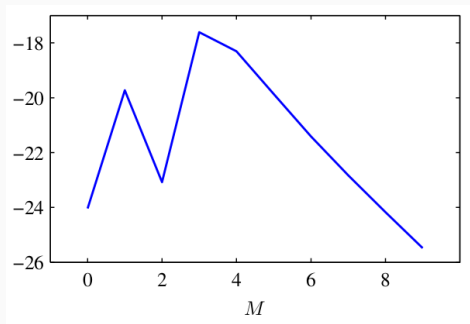
$$\int e^{-\frac{1}{2}(\mathbf{w}-\mathbf{m}_N)^T A(\mathbf{w}-\mathbf{m}_N)} d\mathbf{w} = (2\pi)^{\frac{d}{2}} \sqrt{\det A^{-1}}.$$

- Получается:

$$\begin{aligned} \ln p(D | \alpha, \beta) = \\ \frac{d}{2} \ln \alpha + \frac{N}{2} \ln \beta - \frac{\beta}{2} \|\mathbf{y} - X\mathbf{m}_N\|^2 - \frac{\alpha}{2} \mathbf{m}_N^T \mathbf{m}_N - \frac{1}{2} \ln \det A - \frac{N}{2} \ln(2\pi). \end{aligned}$$

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

- Пример графика по числу параметров:



RVM

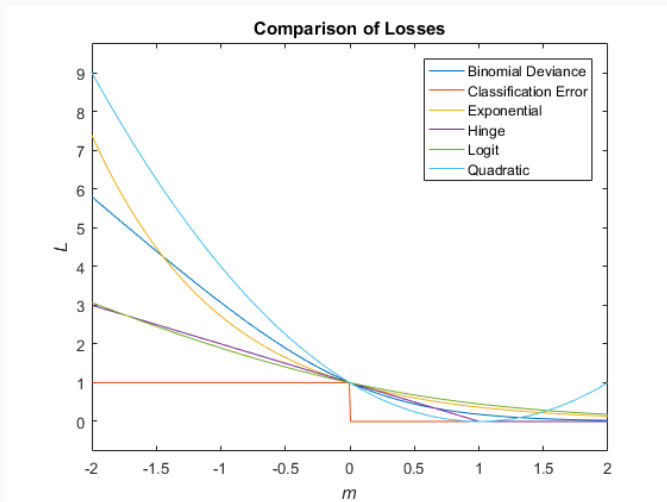
---

# Функции потерь в классификации

- Сначала – ещё один важный другой взгляд: разные методы классификации отличаются друг от друга тем, какую функцию ошибки они оптимизируют.
- У классификации проблема с «правильной» функцией ошибки, то есть ошибкой собственно классификации:
  - она и не везде дифференцируема,
  - и производная её никому не нужна.
- Давайте посмотрим на разные функции потерь (loss functions); мы уже несколько видели, ещё кое-что осталось.



# Функции потерь в классификации



- SVM – отличный метод. Но и у него есть недостатки.
  1. Выходы SVM – решения, а апостериорные вероятности непонятно как получить.
  2. SVM – для двух классов, обобщить на несколько проблематично.
  3. Есть параметр  $C$  (или  $\nu$ , или ещё вдобавок  $\epsilon$ ), который надо подбирать.
  4. Предсказания – линейные комбинации ядер, которым необходимо быть положительно определёнными и которые центрированы на точках из датасета.
- Сейчас мы рассмотрим байесовский аналог SVM – *relevance vector machines* (RVM).

- RVM удобнее сразу формулировать для регрессии.
- Вспомним обычную нашу линейную модель:

$$p(t | \mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t | y(\mathbf{x}), \beta^{-1}), \text{ где}$$

$$y(\mathbf{x}) = \sum_{i=1}^M w_i \phi_i(\mathbf{x}) = \mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}).$$

- RVM – это вариант такой модели, который старается работать как SVM.
- Рассмотрим

$$y(\mathbf{x}) = \sum_{n=1}^N w_n k(\mathbf{x}, \mathbf{x}_n) + b.$$

- Т.е. мы сразу ищем решение в форме линейной комбинации значений ядра (вспомним «эквивалентное ядро» для линейной регрессии), но, в отличие от SVM, теперь ядро никак не ограничивается.

- Для  $N$  наблюдений вектора  $\mathbf{x}$  (обозначим через  $\mathbf{X}$ ) со значениями  $\mathbf{t}$  получим правдоподобие

$$p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N p(t_n | \mathbf{x}_n, \mathbf{w}, \beta^{-1}).$$

- Априорное распределение тоже будет нормальное, но вместо единого гиперпараметра для всех весов мы введём отдельный гиперпараметр для каждого:

$$p(\mathbf{w} | \boldsymbol{\alpha}) = \prod_{i=1}^M \mathcal{N}(w_i | 0, \alpha_i^{-1}).$$

- Отдельные гиперпараметры:

$$p(\mathbf{w} | \boldsymbol{\alpha}) = \prod_{i=1}^M \mathcal{N}(w_i | 0, \alpha_i^{-1}).$$

- Идея здесь в том, что при максимизации апостериорной вероятности большая часть  $\alpha_i$  просто уйдёт на бесконечность, и соответствующие веса будут нулевыми.
- Сейчас увидим, как это получается.

- Апостериорное распределение нам знакомо:

$$p(\mathbf{w} \mid \mathbf{t}, \mathbf{X}, \boldsymbol{\alpha}, \beta) = \mathcal{N}(\mathbf{w} \mid \mathbf{m}, \boldsymbol{\Sigma}), \text{ где}$$

$$\mathbf{m} = \beta \boldsymbol{\Sigma} \boldsymbol{\Phi}^T \mathbf{t},$$

$$\boldsymbol{\Sigma} = \left( \mathbf{A} + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi} \right)^{-1},$$

где  $\mathbf{A} = \text{diag}(\alpha_1, \dots, \alpha_M)$ , а  $\boldsymbol{\Phi}$  в нашем случае – это  $\mathbf{K}$ , симметрическая матрица с элементами  $k(\mathbf{x}_n, \mathbf{x}_m)$ .

- Как найти  $\alpha$  и  $\beta$ ? Нужно максимизировать маргинальное правдоподобие датасета

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

- Это свёртка двух гауссианов, тоже гауссиан:

$$\begin{aligned} \ln p(\mathbf{t} | \mathbf{X}, \alpha, \beta) &= \ln \mathcal{N}(\mathbf{t} | \mathbf{0}, \mathbf{C}) = \\ &= -\frac{1}{2} [N \ln(2\pi) + \ln |\mathbf{C}| + \mathbf{t}^\top \mathbf{C}^{-1} \mathbf{t}], \text{ где } \mathbf{C} = \beta^{-1} \mathbf{I} + \Phi \mathbf{A}^{-1} \Phi^\top. \end{aligned}$$

- Как это оптимизировать?



- Можно подсчитать производные и получить

$$\alpha_i = \frac{\gamma_i}{m_i^2},$$
$$\beta^{-1} = \frac{\|\mathbf{t} - \Phi \mathbf{m}\|^2}{N - \sum_i \gamma_i},$$

где  $\gamma_i = 1 - \alpha_i \Sigma_{ii}$ .

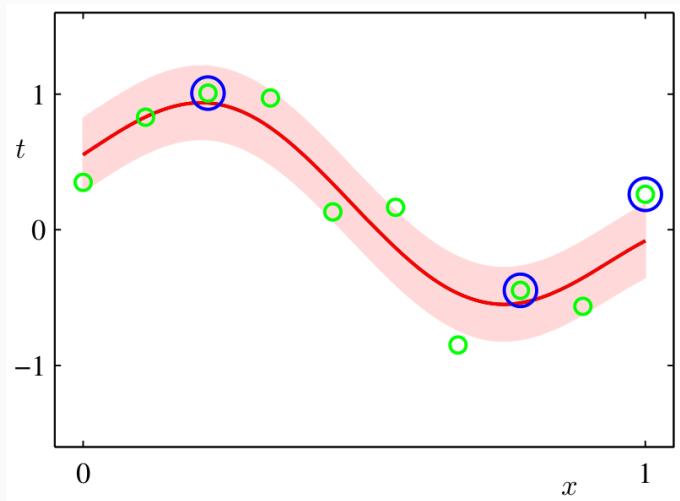
- Теперь можно просто итеративно пересчитывать  $\alpha, \beta$  из  $\mathbf{m}, \Sigma$ , потом наоборот, потом опять наоборот, и до сходимости.

- В результате получается обычно, что большинство  $\alpha_i$  неограниченно растут, и соответствующие веса можно считать нулевыми.
- Оставшиеся называются *relevance vectors*, их обычно мало.
- Если теперь мы найдём  $\alpha^*, \beta^*$ , то предсказывать в новых точках можно как

$$\begin{aligned} p(t | \mathbf{x}, \mathbf{X}, \mathbf{t}, \alpha^*, \beta^*) &= \int p(t | \mathbf{x}, \mathbf{w}, \beta^*) p(\mathbf{w} | \mathbf{X}, \mathbf{t}, \alpha^*, \beta^*) d\mathbf{w} = \\ &= \mathcal{N}(t | \mathbf{m}^\top \phi(\mathbf{x}), \sigma^2(\mathbf{x})), \end{aligned}$$

где  $\sigma^2(\mathbf{x}) = (\beta^*)^{-1} + \phi(\mathbf{x})^\top \Sigma \phi(\mathbf{x})$ .

# RVM для регрессии



# RVM для классификации

- Можно сделать то же самое и для классификации.  
Рассмотрим классификацию с двумя классами,  $t \in \{0, 1\}$ :

$$y(\mathbf{x}, \mathbf{w}) = \sigma(\mathbf{w}^\top \phi(\mathbf{x})).$$

- И добавим сюда, опять же, априорное распределение с разными  $\alpha_i$  для каждого веса:

$$p(\mathbf{w} \mid \boldsymbol{\alpha}) = \prod_{i=1}^M \mathcal{N}(w_i \mid 0, \alpha_i^{-1}).$$

- Идея: инициализируем  $\boldsymbol{\alpha}$ , считаем лапласовское приближение к апостериорному распределению, максимизируем, получаем новое  $\boldsymbol{\alpha}$ , и т.д.

- Апостериорное распределение:

$$\begin{aligned}\ln p(\mathbf{w} \mid \mathbf{t}, \boldsymbol{\alpha}) &= \ln(p(\mathbf{t} \mid \mathbf{w})p(\mathbf{w} \mid \boldsymbol{\alpha})) - \ln p(\mathbf{t} \mid \boldsymbol{\alpha}) = \\ &= \sum_{n=1}^N [t_n \ln y_n + (1 - t_n) \ln(1 - y_n)] - \frac{1}{2} \mathbf{w}^\top \mathbf{A} \mathbf{w} + \text{const.}\end{aligned}$$

- Мы уже обсуждали, как его максимизировать – IRLS; для этого подсчитаем

$$\begin{aligned}\nabla \ln p(\mathbf{w} \mid \mathbf{t}, \boldsymbol{\alpha}) &= \boldsymbol{\Phi}^\top (\mathbf{t} - \mathbf{y}) - \mathbf{A} \mathbf{w}, \\ \nabla \nabla \ln p(\mathbf{w} \mid \mathbf{t}, \boldsymbol{\alpha}) &= - \left( \boldsymbol{\Phi}^\top \mathbf{B} \boldsymbol{\Phi} + \mathbf{A} \right),\end{aligned}$$

где  $\mathbf{B}$  – диагональная матрица с элементами  $b_n = y_n(1 - y_n)$ .

- Лапласовское приближение получится из  $\nabla \ln p(\mathbf{w} | \mathbf{t}, \boldsymbol{\alpha})$ , и получится

$$\mathbf{w}^* = \mathbf{A}^{-1} \boldsymbol{\Phi}^T (\mathbf{t} - \mathbf{y}),$$
$$\boldsymbol{\Sigma} = \left( \boldsymbol{\Phi}^T \mathbf{B} \boldsymbol{\Phi} + \mathbf{A} \right)^{-1},$$

и распределение для предсказания получится

$$p(\mathbf{t} | \boldsymbol{\alpha}) = \int p(\mathbf{t} | \mathbf{w}) p(\mathbf{w} | \boldsymbol{\alpha}) d\mathbf{w} \approx$$
$$\approx p(\mathbf{t} | \mathbf{w}^*) p(\mathbf{w}^* | \boldsymbol{\alpha}) (2\pi)^{M/2} |\boldsymbol{\Sigma}|^{1/2}.$$

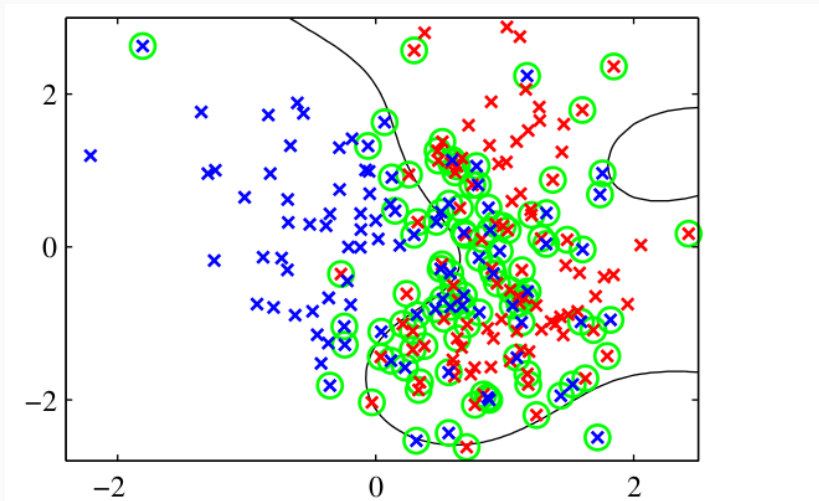
- $p(\mathbf{t} | \boldsymbol{\alpha}) = \int p(\mathbf{t} | \mathbf{w})p(\mathbf{w} | \boldsymbol{\alpha})d\mathbf{w} \approx p(\mathbf{t} | \mathbf{w}^*)p(\mathbf{w}^* | \boldsymbol{\alpha})(2\pi)^{M/2}|\boldsymbol{\Sigma}|^{1/2}.$

- Теперь мы оптимизируем это по  $\boldsymbol{\alpha}$ : берём производную, получаем

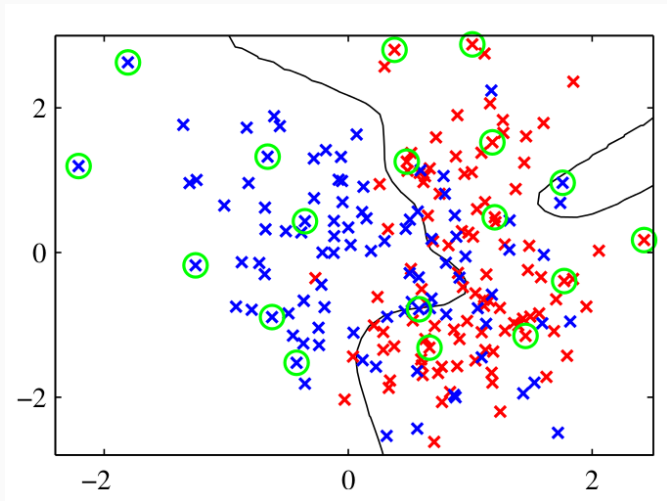
$$-\frac{1}{2}(w_i^*)^2 + \frac{1}{2\alpha_i} - \frac{1}{2}\Sigma_{ii} = 0, \text{ т.е.}$$

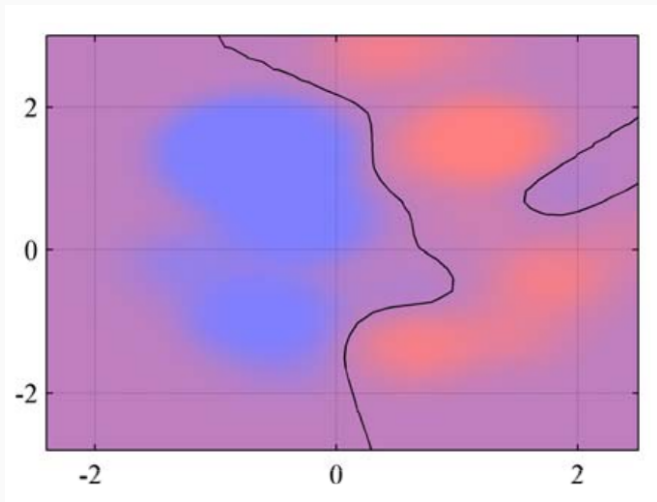
$$\alpha_i = \frac{\gamma_i}{(w_i^*)^2}, \quad \gamma_i = 1 - \alpha_i \Sigma_{ii}.$$

- Т.е. формула получилась точно такая же, как в случае регрессии.









- На несколько классов теперь обобщается естественным образом:

$$a_k = \mathbf{w}_k^\top \mathbf{x}, \quad y_k(\mathbf{x}) = \frac{e^{a_k}}{\sum_j e^{a_j}}.$$

- И дальше всё то же самое.

- RVM как-то получилось лучше со всех сторон.
- Главный минус – в RVM обучение гораздо дольше (хотя есть алгоритмы и побыстрее, чем мы рассматривали, но всё равно дольше).
- Но даже это не то чтобы минус, потому что в SVM нужна кросс-валидация для подбора параметров, т.е. на самом деле обучение SVM дольше, чем кажется.
- Есть и (даже более важный) плюс с точки зрения скорости – в RVM гораздо быстрее применение модели к новым точкам, потому что опорных векторов гораздо меньше.

- В RVM для регрессии получается правдоподобие

$$\begin{aligned}\ln p(\mathbf{t} | \mathbf{X}, \boldsymbol{\alpha}, \beta) &= \ln \mathcal{N}(\mathbf{t} | \mathbf{0}, \mathbf{C}) = \\ &= -\frac{1}{2} [N \ln(2\pi) + \ln |\mathbf{C}| + \mathbf{t}^T \mathbf{C}^{-1} \mathbf{t}], \text{ где } \mathbf{C} = \beta^{-1} \mathbf{I} + \boldsymbol{\Phi} \mathbf{A}^{-1} \boldsymbol{\Phi}^T.\end{aligned}$$

- Выделим вклад в  $\mathbf{C}$  одного компонента  $\alpha_i$ :

$$\begin{aligned}\mathbf{C} &= \beta^{-1}\mathbf{I} + \sum_{j \neq i} \alpha_j^{-1} \boldsymbol{\varphi}_j \boldsymbol{\varphi}_j^T + \alpha_i^{-1} \boldsymbol{\varphi}_i \boldsymbol{\varphi}_i^T = \\ &= \mathbf{C}_{-i} + \alpha_i^{-1} \boldsymbol{\varphi}_i \boldsymbol{\varphi}_i^T,\end{aligned}$$

где  $\boldsymbol{\varphi}_i$  –  $i$ -я строка  $\Phi$  ( $\phi_n$  был  $n$ -м столбцом).

# Быстрый алгоритм обучения RVM

- $\mathbf{C} = \mathbf{C}_{-i} + \alpha_i^{-1} \boldsymbol{\varphi}_i \boldsymbol{\varphi}_i^\top$ .
- Верны следующие тождества для определителя и обратной матрицы:

$$|\mathbf{C}| = |\mathbf{C}_{-i}| |1 + \alpha_i^{-1} \boldsymbol{\varphi}_i^\top \mathbf{C}_{-i}^{-1} \boldsymbol{\varphi}_i|,$$
$$\mathbf{C}^{-1} = \mathbf{C}_{-i}^{-1} - \frac{\mathbf{C}_{-i}^{-1} \boldsymbol{\varphi}_i \boldsymbol{\varphi}_i^\top \mathbf{C}_{-i}^{-1}}{\alpha_i + \boldsymbol{\varphi}_i^\top \mathbf{C}_{-i}^{-1} \boldsymbol{\varphi}_i}.$$

**Упражнение.** Докажите это.

- Значит,  $L(\boldsymbol{\alpha})$  можно переписать в виде

$$L(\boldsymbol{\alpha}) = L(\boldsymbol{\alpha}_{-i}) + \lambda(\alpha_i), \text{ где}$$

$$\lambda(\alpha_i) = \frac{1}{2} \left[ \ln \alpha_i - \ln(\alpha_i + s_i) + \frac{q_i^2}{\alpha_i + s_i} \right].$$

- Здесь

$$\begin{aligned} s_i &= \boldsymbol{\varphi}_i^\top \mathbf{C}_{-i}^{-1} \boldsymbol{\varphi}_i && \text{sparsity } \boldsymbol{\varphi}_i \\ q_i &= \boldsymbol{\varphi}_i^\top \mathbf{C}_{-i}^{-1} \mathbf{t} && \text{quality } \boldsymbol{\varphi}_i. \end{aligned}$$



# Быстрый алгоритм обучения RVM

- $s_i = \varphi_i^\top \mathbf{C}_{-i}^{-1} \varphi_i$ ,  $q_i = \varphi_i^\top \mathbf{C}_{-i}^{-1} \mathbf{t}$ .
- Sparsity – то, насколько  $\varphi_i$  перекрывается с остальными векторами модели.
- Quality – то, насколько  $\varphi_i$  сонаправлен с ошибкой между  $\mathbf{t}$  и  $\mathbf{y}_{-i}$  (ошибкой модели без  $\varphi_i$ ).
- Чем больше sparsity и чем меньше quality, тем более вероятно, что этот базисный вектор из модели исключат (т.е.  $\alpha_i \rightarrow \infty$ ).

# Быстрый алгоритм обучения RVM

- $\lambda(\alpha_i) = \frac{1}{2} \left[ \ln \alpha_i - \ln(\alpha_i + s_i) + \frac{q_i^2}{\alpha_i + s_i} \right]$ .
- Возьмём производную, приравняем нулю, получим (т.к.  $\alpha_i \geq 0$ )

$$\alpha_i = \begin{cases} \infty, & q_i^2 \leq s_i, \\ \frac{s_i^2}{q_i^2 - s_i}, & q_i^2 > s_i. \end{cases}$$

- Как мы и ожидали.

- И алгоритм теперь получается такой:
  1. инициализировать  $\beta$ ,  $\varphi_1$ ,  $\alpha_1 = s_1^2/(q_1^2 - s_1)$ , остальные  $\alpha_j = \infty$ ;
  2. вычислить  $\Sigma$ ,  $\mathbf{m}$ ,  $q_i$  и  $s_i$  для всех  $i$ ;
  3. выбрать  $i$ , проапдейтить  $\alpha_i$ , проапдейтить  $\beta$ ;
  4. goto 2 и так пока не сойдётся.

# EM-алгоритм

---

# Обоснование алгоритма EM

- Дадим формальное обоснование алгоритма EM.
- Мы решаем задачу максимизации правдоподобия по данным  $\mathcal{Y} = \{y_1, \dots, y_N\}$ .

$$L(\boldsymbol{\theta} | \mathcal{Y}) = p(\mathcal{Y} | \boldsymbol{\theta}) = \prod p(y_i | \boldsymbol{\theta})$$

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

- 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  $\mathcal{X} = (\mathcal{Y}, \mathcal{Z})$  с совместной плотностью

$$p(\mathbf{x} | \boldsymbol{\theta}) = p(\mathbf{y}, \mathbf{z} | \boldsymbol{\theta}) = p(\mathbf{z} | \mathbf{y}, \boldsymbol{\theta})p(\mathbf{y} | \boldsymbol{\theta}).$$

- Получается полное правдоподобие  $L(\boldsymbol{\theta} | \mathcal{X}) = p(\mathcal{Y}, \mathcal{Z} | \boldsymbol{\theta})$ . Это случайная величина (т.к.  $\mathcal{Z}$  неизвестно).

# Обоснование алгоритма EM

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

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}_n) = E [\log p(\mathcal{Y}, \mathcal{Z} | \boldsymbol{\theta}) | \mathcal{Y}, \boldsymbol{\theta}_n].$$

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

# Обоснование алгоритма EM

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

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

- Условное ожидание — это

$$E[\log p(\mathcal{Y}, \mathcal{Z} | \theta) | \mathcal{Y}, \theta_n] = \int_{\mathbf{z}} \log p(\mathcal{Y}, \mathbf{z} | \theta) p(\mathbf{z} | \mathcal{Y}, \theta_n) d\mathbf{z},$$

где  $p(\mathbf{z} | \mathcal{Y}, \theta_n)$  — маргинальное распределение скрытых компонентов данных.

- EM лучше всего применять, когда это выражение легко подсчитать, может быть, даже аналитически.
- Вместо  $p(\mathbf{z} | \mathcal{Y}, \theta_n)$  можно подставить  $p(\mathbf{z}, \mathcal{Y} | \theta_n) = p(\mathbf{z} | \mathcal{Y}, \theta_n)p(\mathcal{Y} | \theta_n)$ , от этого ничего не изменится.



# Обоснование алгоритма EM

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

$$\boldsymbol{\theta}_{n+1} = \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}_n).$$

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

# Обоснование алгоритма EM

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

$$\begin{aligned}\ell(\theta) - \ell(\theta_n) &= \\ &= \log \left( \int_{\mathbf{z}} p(\mathcal{Y} | \mathbf{z}, \theta) p(\mathbf{z} | \theta) d\mathbf{z} \right) - \log p(\mathcal{Y} | \theta_n) = \\ &= \log \left( \int_{\mathbf{z}} p(\mathbf{z} | \mathcal{Y}, \theta_n) \frac{p(\mathcal{Y} | \mathbf{z}, \theta) p(\mathbf{z} | \theta)}{p(\mathbf{z} | \mathcal{Y}, \theta_n)} d\mathbf{z} \right) - \log p(\mathcal{Y} | \theta_n) \geq \\ &\geq \int_{\mathbf{z}} p(\mathbf{z} | \mathcal{Y}, \theta_n) \log \left( \frac{p(\mathcal{Y} | \mathbf{z}, \theta) p(\mathbf{z} | \theta)}{p(\mathbf{z} | \mathcal{Y}, \theta_n)} \right) d\mathbf{z} - \log p(\mathcal{Y} | \theta_n) = \\ &= \int_{\mathbf{z}} p(\mathbf{z} | \mathcal{Y}, \theta_n) \log \left( \frac{p(\mathcal{Y} | \mathbf{z}, \theta) p(\mathbf{z} | \theta)}{p(\mathcal{Y} | \theta_n) p(\mathbf{z} | \mathcal{Y}, \theta_n)} \right) d\mathbf{z}.\end{aligned}$$

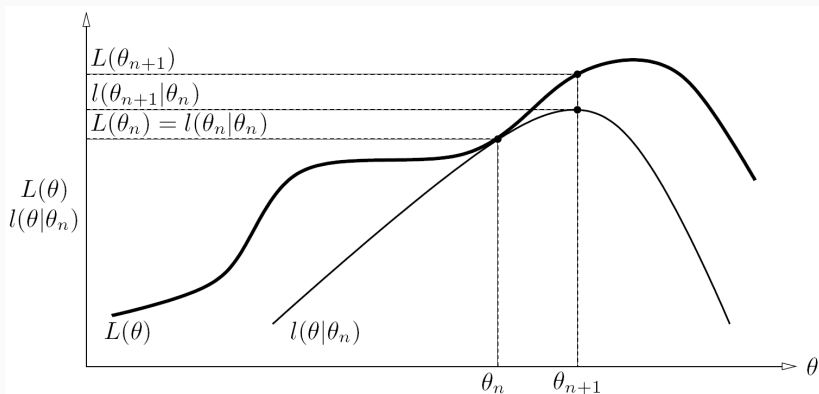
- Получили

$$\begin{aligned}\ell(\boldsymbol{\theta}) &\geq \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\theta}_n) = \\ &= \ell(\boldsymbol{\theta}_n) + \int_{\mathbf{z}} p(\mathbf{z} | \mathcal{Y}, \boldsymbol{\theta}_n) \log \left( \frac{p(\mathcal{Y} | \mathbf{z}, \boldsymbol{\theta}) p(\mathbf{z} | \boldsymbol{\theta})}{p(\mathcal{Y} | \boldsymbol{\theta}_n) p(\mathbf{z} | \mathcal{Y}, \boldsymbol{\theta}_n)} \right) d\mathbf{z}.\end{aligned}$$

**Упражнение.** Докажите, что  $\mathcal{L}(\boldsymbol{\theta}_n, \boldsymbol{\theta}_n) = \ell(\boldsymbol{\theta}_n)$ .

- Иначе говоря, мы нашли нижнюю оценку на  $\ell(\theta)$  везде, касание происходит в точке  $\theta_n$ .
- Т.е. мы нашли нижнюю оценку для правдоподобия и смещаемся в точку, где она максимальна (или хотя бы больше текущей).
- Такая общая схема называется *MM-алгоритм* (minorization-maximization). Мы к ним ещё вернёмся.

# Обоснование алгоритма EM



- Осталось только понять, что максимизировать можно  $Q$ .

$$\begin{aligned}\theta_{n+1} &= \arg \max_{\theta} \ell(\theta, \theta_n) = \arg \max_{\theta} \left\{ \ell(\theta_n) + \right. \\ &\quad \left. + \int_{\mathcal{Z}} f(y | \mathcal{X}, \theta_n) \log \left( \frac{p(\mathcal{X} | y, \theta) f(y | \theta)}{p(\mathcal{X} | \theta_n) f(y | \mathcal{X}, \theta_n)} \right) dz \right\} = \\ &= \arg \max_{\theta} \left\{ \int_{\mathcal{Z}} p(z | \mathcal{X}, \theta_n) \log (p(\mathcal{X} | y, \theta) p(z | \theta)) dz \right\} = \\ &= \arg \max_{\theta} \left\{ \int_{\mathcal{Z}} p(z | \mathcal{X}, \theta_n) \log p(\mathcal{X}, y | \theta) dz \right\} = \\ &= \arg \max_{\theta} \{Q(\theta, \theta_n)\},\end{aligned}$$

а остальное от  $\theta$  не зависит. Вот и получился EM.

Спасибо!

Спасибо за внимание!