

RVM и EM-алгоритм

Сергей Николенко

Академия MADE — Mail.Ru

6 мая 2020 г.

Random facts:

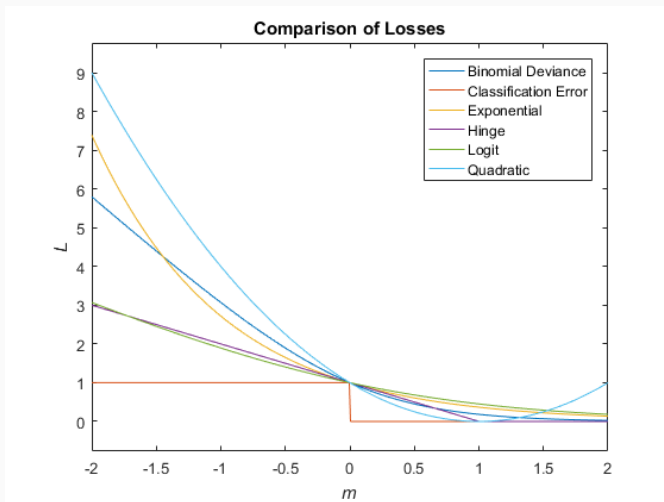
- 6 мая 1626 г. голландский колонист Петер Минейт (Peter Minuit) за горсть бижутерии общей стоимостью 25 долларов приобрёл у индейцев остров Манхэттен
- 6 мая 1840 г. был выпущен в обращение «Чёрный пенни»
- 6 мая 1890 г. Церковь Иисуса Христа Святых последних дней официально отказалась от многожёнства, и в том же году Юта получила статус штата
- 6 мая 1949 г. в Кембридже свою первую программу запустил EDSAC (Electronic delay storage automatic calculator), первый настоящий практический компьютер на архитектуре фон Неймана; он вычислил списки квадратов и простых чисел
- 6 мая 2012 г. в Москве прошёл «Марш миллионов» для выражения протеста против инаугурации Владимира Путина

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 ν , или ещё вдобавок ϵ), 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:

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

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

- 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:

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

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

- Апостериорное распределение нам знакомо:

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

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

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

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 α и β ? 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 $\gamma_i = 1 - \alpha_i \Sigma_{ii}$.

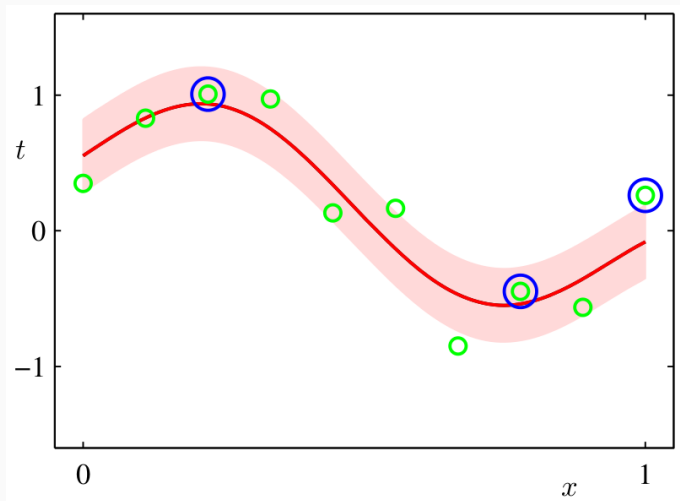
- Теперь можно просто итеративно пересчитывать α, β из \mathbf{m}, Σ , потом наоборот, потом опять наоборот, и до сходимости.

- В результате получается обычно, что большинство α_i неограниченно растут, и соответствующие веса можно считать нулевыми.
- Оставшиеся называются *relevance vectors*, их обычно мало.
- Если теперь мы найдём α^*, β^* , 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 $\sigma^2(\mathbf{x}) = (\beta^*)^{-1} + \phi(\mathbf{x})^\top \mathbf{\Sigma} \phi(\mathbf{x})$.

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 α_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 $\nabla \ln p(\mathbf{w} \mid \mathbf{t}, \boldsymbol{\alpha})$, и получится

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

и распределение для предсказания получится

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

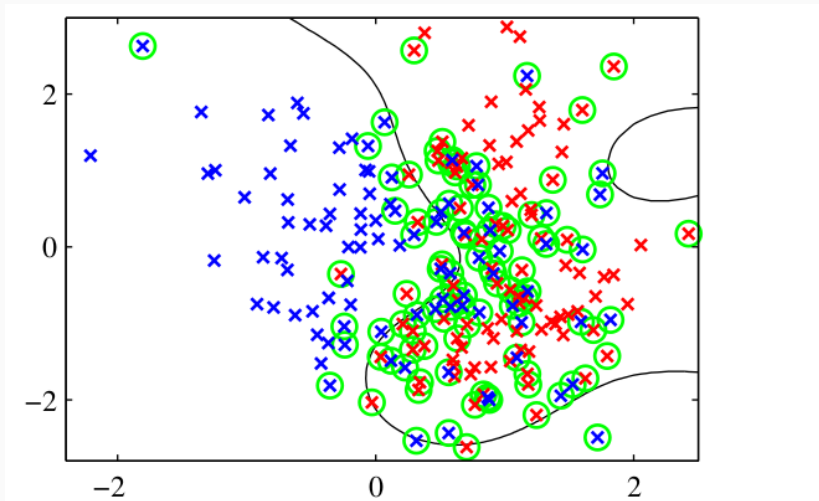
- $p(\mathbf{t} | \boldsymbol{\alpha}) = \int p(\mathbf{t} | \mathbf{w})p(\mathbf{w} | \boldsymbol{\alpha})d\mathbf{w} \approx p(\mathbf{t} | \mathbf{w}^*)p(\mathbf{w}^* | \boldsymbol{\alpha})(2\pi)^{M/2}|\boldsymbol{\Sigma}|^{1/2}.$
- Теперь мы оптимизируем это по $\boldsymbol{\alpha}$: берём производную, получаем

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

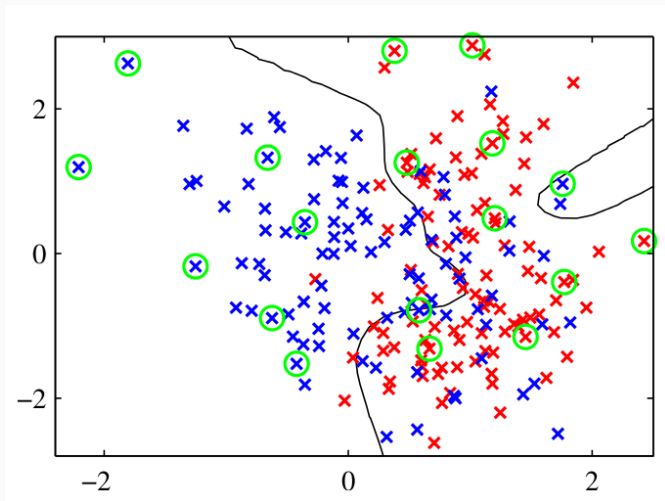
$$\alpha_i = \frac{\gamma_i}{(w_i^*)^2}, \quad \gamma_i = 1 - \alpha_i \Sigma_{ii}.$$

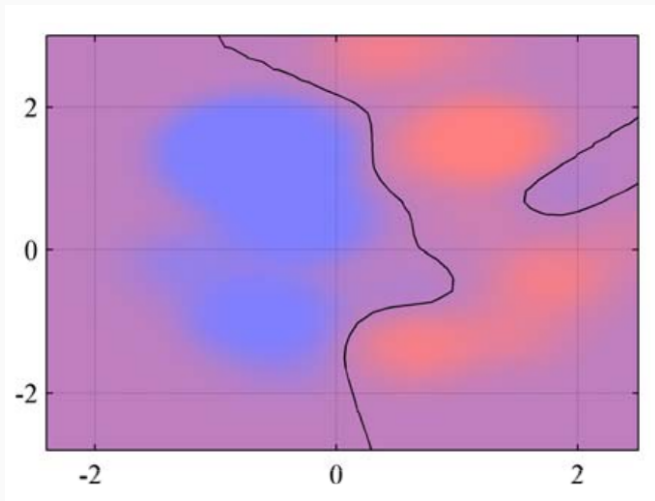
- Т.е. формула получилась точно такая же, как в случае регрессии.

Было: SVM



Стало: RVM





- 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.

Сравнение SVM и RVM

- 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

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

Быстрый алгоритм обучения RVM

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

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

где $\boldsymbol{\varphi}_i$ – i -я строка Φ (ϕ_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 φ_i перекрывается с остальными векторами модели.
- Quality – то, насколько φ_i сонаправлен с ошибкой между \mathbf{t} и \mathbf{y}_{-i} (ошибкой модели без φ_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 β , φ_1 , $\alpha_1 = s_1^2/(q_1^2 - s_1)$, остальные $\alpha_j = \infty$;
 2. вычислить Σ , m , q_i и s_i для всех i ;
 3. выбрать i , проапдейтить α_i , проапдейтить β ;
 4. goto 2 и так пока не сойдётся.

Алгоритм EM

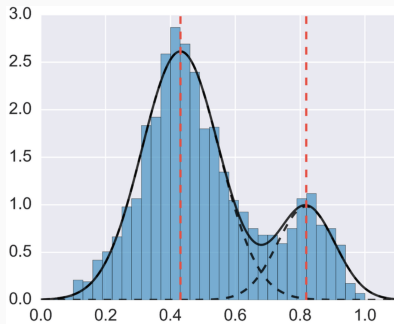
- Часто возникает ситуация, когда в имеющихся данных некоторые переменные присутствуют, а некоторые — отсутствуют.
- Даны результаты сэмплирования распределения вероятностей с несколькими параметрами, из которых известны не все.

- Эти неизвестные параметры тоже расцениваются как случайные величины.
- Задача — найти наиболее вероятную гипотезу, то есть ту гипотезу h , которая максимизирует

$$E[\ln p(D|h)].$$

Частный случай

- Построим один из простейших примеров применения алгоритма ЕМ. Пусть случайная переменная y сэмплируется из суммы двух нормальных распределений. Дисперсии даны (одинаковые), нужно найти только средние μ_1, μ_2 .



- Какое тут правдоподобие? Как его оптимизировать?

- Нельзя понять, какие y_i были порождены каким распределением — классический пример *скрытых переменных*.
- Один тестовый пример полностью описывается как тройка $\langle y_i, z_{i1}, z_{i2} \rangle$, где $z_{ij} = 1$ iff y_i был сгенерирован j -м распределением.

- Сгенерировать какую-нибудь гипотезу $h = (\mu_1, \mu_2)$.
- Пока не дойдем до локального максимума:
 - Вычислить ожидание $E(z_{ij})$ в предположении текущей гипотезы (E-шаг).
 - Вычислить новую гипотезу $h' = (\mu'_1, \mu'_2)$, предполагая, что z_{ij} принимают значения $E(z_{ij})$ (M-шаг).

В примере с гауссианами

- В примере с гауссианами:

$$\begin{aligned} E(z_{ij}) &= \frac{p(y = y_i | \mu = \mu_j)}{p(y = y_i | \mu = \mu_1) + p(y = y_i | \mu = \mu_2)} = \\ &= \frac{e^{-\frac{1}{2\sigma^2}(y_i - \mu_j)^2}}{e^{-\frac{1}{2\sigma^2}(y_i - \mu_1)^2} + e^{-\frac{1}{2\sigma^2}(y_i - \mu_2)^2}}. \end{aligned}$$

- Мы подсчитываем эти ожидания, а потом подправляем гипотезу:

$$\mu_j \leftarrow \frac{1}{m} \sum_{i=1}^m E(z_{ij}) y_i.$$

- Звучит логично, но с какой стати это всё работает?

Обоснование алгоритма EM

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

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

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

- 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

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

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

Обоснование алгоритма EM

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

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

- 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 θ :

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

- Условное ожидание — это

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

где $p(\mathcal{Z} \mid \mathcal{Y}, \theta_n)$ — маргинальное распределение скрытых компонентов данных.

- EM лучше всего применять, когда это выражение легко подсчитать, может быть, даже аналитически.
- Вместо $p(\mathcal{Z} \mid \mathcal{Y}, \theta_n)$ можно подставить $p(\mathcal{Z}, \mathcal{Y} \mid \theta_n) = p(\mathcal{Z} \mid \mathcal{Y}, \theta_n) p(\mathcal{Y} \mid \theta_n)$, 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.

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

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

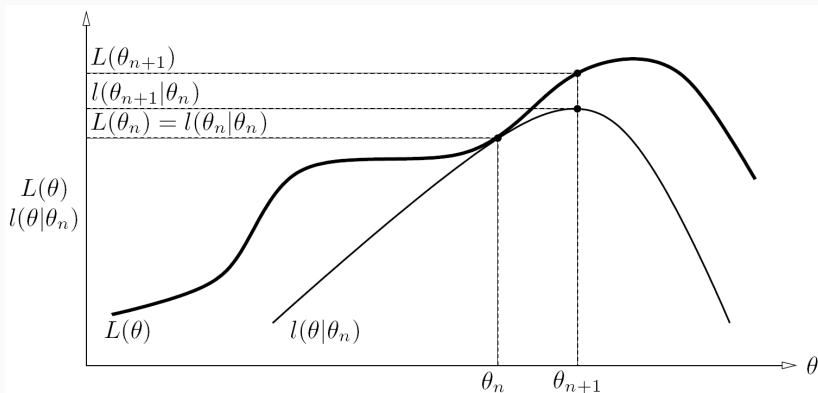
- Получили

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

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 θ_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 *ММ-алгоритм* (minorization-maximization). Мы к ним ещё вернёмся.

Обоснование алгоритма EM



Обоснование алгоритма EM

- Осталось только понять, что максимизировать можно Q .

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

а остальное от θ не зависит. Вот и получился EM.

Спасибо!

Спасибо за внимание!