

# ВАРИАЦИОННЫЕ ПРИБЛИЖЕНИЯ II

---

Сергей Николенко

СПбГУ — Санкт-Петербург

19 февраля 2021 г.

---

## *Random facts:*

- 19 февраля — Всемирный день китов; именно 19 февраля 1986 года вступил в силу запрет International Whaling Commission на любую коммерческую добычу китов; Япония присоединилась к запрету, но оставила за собой возможность «отлова китов в научных целях»; после завершения научных работ китовое мясо не вполне известным науке образом оказывается в японских ресторанах
- 19 февраля в Армении — День дарения книг, отмечающийся с 2008 года в день рождения Ованеса Туманяна (19 февраля 1869 г.)
- 19 февраля 356 г. Константин закрыл все языческие храмы в Римской империи
- 19 февраля 1855 г. Урбен Леверье представил Парижской академии наук первую в мире карту погоды — прогноз для Европы на утро того же дня
- 19 февраля 1954 г. указом Президиума Верховного Совета СССР Крымская область была передана из состава РСФСР в состав УССР

# ВАРИАЦИОННОЕ ПРИБЛИЖЕНИЕ ДЛЯ ГАУССИАНА

---

# Одномерный гауссиан

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

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

- Вводим сопряжённые априорные распределения:

$$p(\mu | \tau) = N(\mu | \mu_0, (\lambda_0 \tau)^{-1}),$$
$$p(\tau) = \text{Gamma}(\tau | a_0, b_0).$$

- 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

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

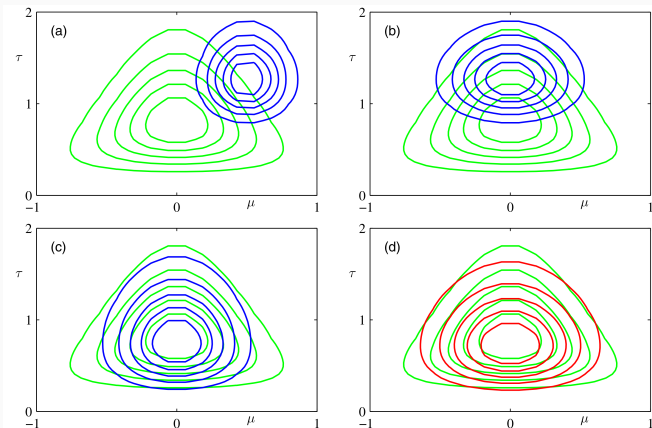
- А  $q_\tau(\tau)$  – гамма-распределение с параметрами

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

- 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  $\mathbb{E}[\tau]$ :

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

- Автоматически получили несмещённую оценку дисперсии!

# ВАРИАЦИОННОЕ ПРИБЛИЖЕНИЕ ДЛЯ СМЕСИ ГАУССИАНОВ

---

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

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

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

- Выберем сопряжённые априорные распределения:

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

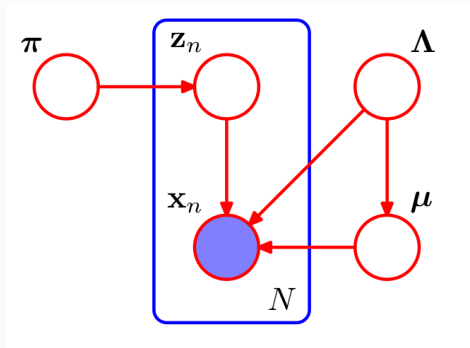
$$p(\mu, \Lambda) = p(\mu | \Lambda) p(\Lambda)$$

$$= \prod_{k=1}^K N(\mu_k | \mathbf{m}_0, (\beta_0 \Lambda_k)^{-1}) W(\Lambda_k | W_0, \nu_0).$$



# СМЕСЬ ГАУССИАНОВ

- Вот такая графическая модель:



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

- Теперь вариационное приближение. Сначала сама факторизация:

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

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

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

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

- Сначала  $q^*(\mathbf{Z})$ :

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

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

- Нормируем:

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

- 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:

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

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

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

- То же самое происходило и в EM-алгоритме.

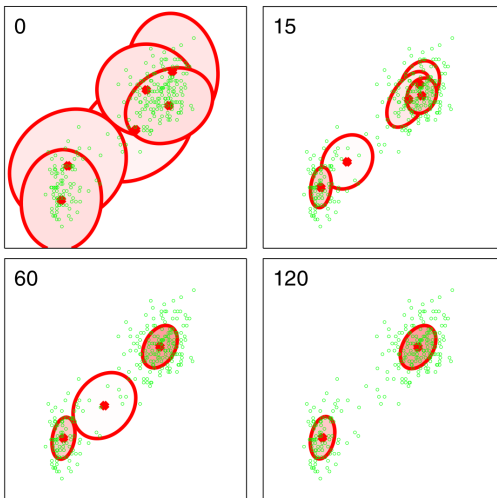
- Теперь  $q^*(\pi, \mu, \Lambda)$ :

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

- Вот уже получилось, что  $q^*(\pi, \mu, \Lambda)$  раскладывается в  $q^*(\pi)q^*(\mu, \Lambda)$ , опять же без предположений.
- Более того,  $q^*(\mu, \Lambda) = \prod_{k=1}^K q(\mu_k, \Lambda_k)$ .
- И теперь можно по отдельности посчитать (упражнение), получится типичный M-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.

- Есть другие примеры вариационных приближений.
- Обращение матриц; например, для линейной регрессии надо посчитать  $\beta^* = C^{-1}\mathbf{b}$ :

$$J(\beta) = \frac{1}{2}(\beta^* - \beta)^\top C(\beta^* - \beta) = \dots = \text{Const} - \beta^\top \mathbf{b} + \frac{1}{2}\beta^\top C\beta,$$

и теперь можно решать такую задачу выпуклой оптимизации.

- Метод конечных элементов – для уравнения Пуассона  $-u''(x) = f(x)$ ,  $x \in (a, b)$ :

$$J(u) = \frac{1}{2} \int_a^b (u'(x) - u^{*'}(x))^2 dx = \dots = \text{Const} - \int_a^b u(x)f(x)dx + \frac{1}{2} \int_a^b u'(x)^2 dx,$$

и если ищем в подпространстве  $\tilde{u}(x) = \sum_{i=1}^k \alpha_i \phi_i(x)$ , то опять

$$\tilde{J}(\alpha) = \alpha^\top \mathbf{b} + \frac{1}{2}\alpha^\top C\alpha.$$

- В графических моделях – теория среднего поля (mean field theory). Пусть дано  $p(\mathbf{x})$ ,  $\mathbf{x} = (\mathbf{x}_v, \mathbf{x}_h)$ , и надо найти

$$\log p(\mathbf{x}_v) = \log \sum_{\mathbf{x}_h} p(\mathbf{x}_v, \mathbf{x}_h), \quad p(\mathbf{x}_h | \mathbf{x}_v) = p(\mathbf{x}_h, \mathbf{x}_v) / p(\mathbf{x}_v).$$

- Опять делаем тот же трюк:

$$J(q) = \log p(\mathbf{x}_v) - \text{KL}(q_{\mathbf{x}_h} \| p_{\mathbf{x}_h | \mathbf{x}_v}) = \log p(\mathbf{x}_v) - \sum_{\mathbf{x}_h} q(\mathbf{x}_h) \log \frac{q(\mathbf{x}_h)}{p(\mathbf{x}_h | \mathbf{x}_v)}$$

$$\dots = H(q) + \mathbb{E}_q [\log p(\mathbf{x}_h, \mathbf{x}_v)] = H(q) + \sum_{C \in \mathcal{C}} \sum_{\mathbf{x}_{C \cap h}} q(\mathbf{x}_{C \cap h}) \log \Psi_C(\mathbf{x}_C),$$

где  $q(\mathbf{x}_{C \cap h})$  – маргинальная вероятность по скрытым переменным из клики  $C$ .

- Теория среднего поля – это когда  $q(\mathbf{x}_h) = \prod_{i \in h} q_i(x_i)$ .



# ЗАДАЧИ ОБРАБОТКИ ЕСТЕСТВЕННОГО ЯЗЫКА

---

- Работы, связанные с естественным языком, — это одна из ключевых задач для создания искусственного интеллекта.
- А.А. Марков-старший, 1913: численные эксперименты с текстом «Евгения Онегина», языковые модели через марковские цепи
- Ранний оптимизм: 1950-е, Ноам Хомский, Дартмутский семинар.
- Georgetown-IBM experiment, 1954: «через 3-5 лет машинный перевод будет решён».
- Но оказалось, что всё сложно; первая зима нейронных сетей — это отчасти fail проекта по машинному переводу.
- И до сих пор мы, хотя умеем решать связанные с языком задачи всё лучше и лучше, очень далеки от истинного понимания.

- Почему сложно? Во многом из-за модели окружающего мира, commonsense reasoning.
- Пример — разрешение *анафоры*:
  - мама вымыла раму, и теперь она блестит;
  - мама вымыла раму, и теперь она устала.
- Это пример хорошо определённой задачи, фактически задачи классификации, для которой легко набрать датасет, но она в нетривиальных случаях всё равно очень, очень сложная.
- Какие ещё бывают задачи в обработке естественного языка?

## (1) Синтаксические задачи:

- *частеречная разметка* (part-of-speech tagging): разметить в заданном тексте слова по частям речи (существительное, глагол, прилагательное...) и, возможно, по морфологическим признакам (род, падеж...);

## (1) Синтаксические задачи:

- *частеречная разметка* (part-of-speech tagging);
- *морфологическая сегментация* (morphological segmentation):  
разделить слова в заданном тексте на *морфемы*, т.е. синтаксические единицы вроде приставок, суффиксов и окончаний; для некоторых языков (например, английского) это не очень актуально, но в русском языке морфологии очень много;

## (1) Синтаксические задачи:

- *частеречная разметка* (part-of-speech tagging);
- *морфологическая сегментация* (morphological segmentation);
- другой вариант задачи о морфологии отдельных слов — *стемминг* (stemming), в котором требуется выделить основы слов, или *лемматизация* (lemmatization), в которой слово нужно привести к базовой форме (например, форме единственного числа мужского рода);

## (1) Синтаксические задачи:

- *частеречная разметка* (part-of-speech tagging);
- *морфологическая сегментация* (morphological segmentation);
- *стемминг* (stemming)
- *выделение границ предложения* (sentence boundary disambiguation): разбить заданный текст на предложения; задача непростая даже в русском и английском, а в языках вроде китайского весьма нетривиальной становится даже задача *пословной сегментации* (word segmentation), потому что поток иероглифов без пробелов может делиться на слова по-разному;

## (1) Синтаксические задачи:

- *частеречная разметка* (part-of-speech tagging);
- *морфологическая сегментация* (morphological segmentation);
- *стемминг* (stemming)
- *выделение границ слов и предложений*;
- *распознавание именованных сущностей* (named entity recognition): найти в тексте собственные имена людей, географических и прочих объектов, разметив их по типам сущностей (люди, места, названия компаний и т.п.);



## (1) Синтаксические задачи:

- *частеречная разметка* (part-of-speech tagging);
- *морфологическая сегментация* (morphological segmentation);
- *стемминг* (stemming)
- *выделение границ слов и предложений*;
- *распознавание именованных сущностей* (named entity recognition);
- *разрешение смысла слов* (word sense disambiguation): выбрать, какой из омонимов, какой из разных смыслов одного и того же слова используется в данном отрывке текста;

## (1) Синтаксические задачи:

- *частеречная разметка* (part-of-speech tagging);
- *морфологическая сегментация* (morphological segmentation);
- *стемминг* (stemming)
- *выделение границ слов и предложений*;
- *распознавание именованных сущностей* (named entity recognition);
- *разрешение смысла слов* (word sense disambiguation);
- *синтаксический парсинг* (syntactic parsing): по заданному предложению (и, возможно, его контексту) построить его синтаксическое дерево, прямо по Хомскому;

## (1) Синтаксические задачи:

- *частеречная разметка* (part-of-speech tagging);
- *морфологическая сегментация* (morphological segmentation);
- *стемминг* (stemming)
- *выделение границ слов и предложений*;
- *распознавание именованных сущностей* (named entity recognition);
- *разрешение смысла слов* (word sense disambiguation);
- *синтаксический парсинг* (syntactic parsing);
- *разрешение кореференций* (coreference resolution):  
определить, к каким объектам или другим частям текста относятся те или иные слова и обороты; частный случай этой задачи — то самое разрешение анафоры, которое мы обсуждали выше.

- (1) Синтаксические задачи.
- (2) Хорошо определённые семантические задачи:
  - *языковые модели* (language models): по заданному отрывку текста предсказать следующее слово или следующий символ; эта задача очень важна, например, для распознавания речи;

- (1) Синтаксические задачи.
- (2) Хорошо определённые семантические задачи:
  - *языковые модели* (language models);
  - *анализ тональности* (sentiment analysis): определить по тексту его тональность, т.е. позитивное ли отношение несёт этот текст или негативное;

(1) Синтаксические задачи.

(2) Хорошо определённые семантические задачи:

- *языковые модели* (language models);
- *анализ тональности* (sentiment analysis);
- *выделение отношений* или *фактов* (relationship extraction, fact extraction): выделить из текста хорошо определённые отношения или факты об упоминающихся там сущностях; например, кто с кем находится в родственных отношениях, в каком году основана упоминающаяся в тексте компания и т.д.;

- (1) Синтаксические задачи.
- (2) Хорошо определённые семантические задачи:
  - *языковые модели* (language models);
  - *анализ тональности* (sentiment analysis);
  - *выделение отношений* или *фактов* (relationship extraction, fact extraction);
  - *ответы на вопросы* (question answering): дать ответ на заданный вопрос; в зависимости от постановки это может быть или чистая классификация (выбор из вариантов ответа, как в тесте), или классификация с очень большим числом классов (ответы на фактологические вопросы вроде «кто» или «в каком году»), или даже порождение текста (если отвечать на вопросы нужно в рамках естественного диалога).

- (1) Синтаксические задачи.
- (2) Хорошо определённые семантические задачи.
- (3) Хуже определённые семантические задачи:
  - собственно *порождение текста* (text generation);



- (1) Синтаксические задачи.
- (2) Хорошо определённые семантические задачи.
- (3) Хуже определённые семантические задачи:
  - собственно *порождение текста* (text generation);
  - *автоматическое реферирование* (automatic summarization): по тексту породить его краткое содержание, abstract, так сказать; это можно рассмотреть как задачу классификации, если просить модель выбрать из текста готовые предложения, лучше всего отражающие смысл всего текста, а можно как задачу порождения, если краткое содержание нужно написать с нуля;

- (1) Синтаксические задачи.
- (2) Хорошо определённые семантические задачи.
- (3) Хуже определённые семантические задачи:
  - собственно *порождение текста* (text generation);
  - *автоматическое реферирование* (automatic summarization);
  - *машинный перевод* (machine translation): по тексту на одном языке породить соответствующий текст на другом языке;

- (1) Синтаксические задачи.
- (2) Хорошо определённые семантические задачи.
- (3) Хуже определённые семантические задачи:
  - собственно *порождение текста* (text generation);
  - *автоматическое реферирование* (automatic summarization);
  - *машинный перевод* (machine translation);
  - *диалоговые модели* (dialog and conversational models):  
поддержать разговор с человеком.

Спасибо за внимание!