

Fixing a Broken ELBO

Ольга Гребенькова

Байесовское муьтимоделирование

24 ноября 2021 г.

Обучение без учителя

Для разных моделей типа VAE основной задачей является нахождения какого-то "полезного" представления данных. То есть мы хотим найти модель следующего вида $p(x, z|\Theta) = p(z|\Theta)p(x|z, \Theta)$.

Классический подход

Обычно находим параметры модели минимизируя KL дивергенцию, что аналогично максимизации правдоподобия модели. Если выражение невычислимо, можем максимизировать нижнюю оценку этого выражения (ELBO) или рассмотреть другие виды дивергенций (RKL).

Проблема

Функции потерь зависят от $p(x|\Theta)$, а не от $p(x, z|\Theta)$. Пример: PixelVAE, где получаем хорошее маргинальное правдоподобие $p(x|\Theta)$. Получение хорошей нижней оценки (а значит и хорошего значения маргинального правдоподобия $p(x|\Theta)$) не достаточно.

x — вектор данных; z — скрытое представление; $e(z|x)$ — выбранный нами энкодер.

Совместное распределение

$$p_e(x, z) = p(x)e(z|x).$$

Маргинальное апостериорное распределение

$$p_e(z) = \int p(x)e(z|x)dx.$$

Условное распределение

$$p_e(x|z) = \frac{p_e(x, z)}{p_e(z)}.$$

$$I_e(X; Z) = \int \int dx dz p_e(x, z) \log \frac{p_e(x, z)}{p(x)p_e(z)} \quad (1)$$

- ❶ $I(X, Z) = I(Z, X)$
- ❷ $I(X, Z) \geq 0$
- ❸ не зависит от репараметризаций
- ❹ показывает как много информации одна случайная величина содержит о другой
- ❺ e — указывает на зависимость от энкодера
- ❻ если X, Y независимы, $I(X, Y) = 0$
- ❼ если $X = Y$, $I(X, Y) = H(X)$
- ❽ скорее не вычислима из-за $p(x)p_e(z)$ (эмпирическое распределение и оценки)

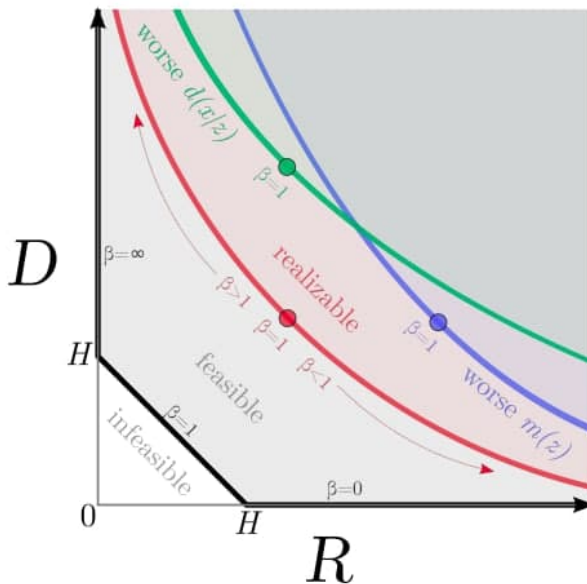
$$H - D \leq I_e(X; Z) \leq R$$

$$H = - \int dx p(x) \log p(x)$$

$$D = - \int dx p(x) \int dz e(z|x) \log d(x|z)$$

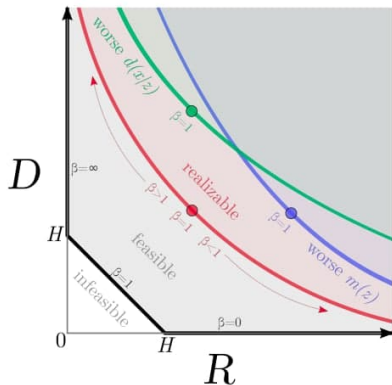
$$R = \int dx p(x) \int dz e(z|x) \log \frac{e(z|x)}{m(z)}$$

- $d(x|z)$ (Decoder) — аппроксимация $p_e(x|z)$
- $m(z)$ (Marginal) — аппроксимация $p_e(z)$
- H — энтропия данных, показывающая сложность датасета
- D — искажение, через кодировщик и декодер, равный обратному логарифму правдоподобия. отрицательная логарифмическая вероятность.
- R — значение, зависящее только от кодировщика и аппроксимации $p_e(z)$. По сути это средняя KL дивергенция между распределением кодировщика и выученным маргинальным распределением.
- Для дискретных данных $H \geq 0, D \geq 0, R \geq 0$.



Вместо поиска оптимального $D(R)$, воспользуемся преобразованием Лежандра и будем искать оптимальные значения D и R для фиксированного значения $\beta = \frac{\partial D}{\partial R}$ минимизируя $\min_{e(z|x), m(z), d(x|z)} D + \beta R$:

$$\min_{e(z|x), m(z), d(x|z)} \int dx p(x) \int dz e(z|x) [-\log d(x|z) + \beta \log \frac{e(z|x)}{m(z)}].$$



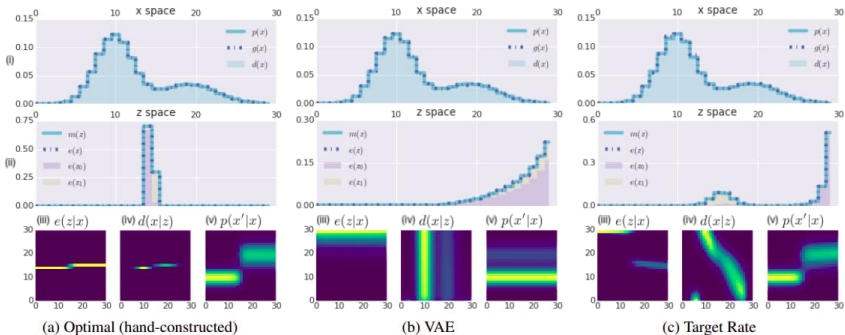


Figure 2. Toy Model illustrating the difference between fitting a model by maximizing ELBO (b) vs minimizing distortion for a fixed rate (c). **Top (i):** Three distributions in data space: the true data distribution, $p^*(x)$, the model's generative distribution, $g(x) = \sum_z m(z)d(x|z)$, and the empirical data reconstruction distribution, $d(x) = \sum_{x'} \sum_z \hat{p}(x')e(z|x')d(x|z)$. **Middle (ii):** Four distributions in latent space: the learned (or computed) marginal $m(z)$, the empirical induced marginal $e(z) = \sum_x \hat{p}(x)e(z|x)$, the empirical distribution over z values for data vectors in the set $\mathcal{X}_0 = \{x_n : z_n = 0\}$, which we denote by $e(z_0)$ in purple, and the empirical distribution over z values for data vectors in the set $\mathcal{X}_1 = \{x_n : z_n = 1\}$, which we denote by $e(z_1)$ in yellow. **Bottom:** Three $K \times K$ distributions: (iii) $e(z|x)$, (iv) $d(x|z)$ and (v) $p(x'|x) = \sum_z e(z|x)d(x|z)$.