

Обзор статьи
Opening the Black Box of Deep Neural Networks
via Information

Глеб Пособин

9 октября 2017 г.

Взаимная информация

Взаимная информация для двух случайных величин X, Y :

$$I(X; Y) = D_{kl}(p(x, y) || p(x)p(y)) = \sum_{x, y} p(x, y) \log \frac{p(x, y)}{p(x)p(y)}$$

Инвариантность к обратимым преобразованиям:

$$I(X; Y) = I(\phi(X); \psi(Y))$$

Data processing inequality: для любых трёх случайных величин X, Y, Z , таких что $X \rightarrow Y \rightarrow Z$, верно:

$$I(X; Y) \geq I(X; Z)$$

Достаточные статистики

Две случайные величины X, Y .

Достаточная статистика $S(X)$:

$$I(S(X); Y) = I(X; Y)$$

Минимальная достаточная статистика $T(X)$:

$$T(X) = \arg \min_{S(X): I(S(X); Y) = I(X; Y)} I(S(X); X)$$

Точные минимальные достаточные статистики очень редко существуют.

Вероятностные достаточные статистики

T — случайная величина, зависящая от X . Разрешим $I(X; Y) \neq I(T; Y)$.

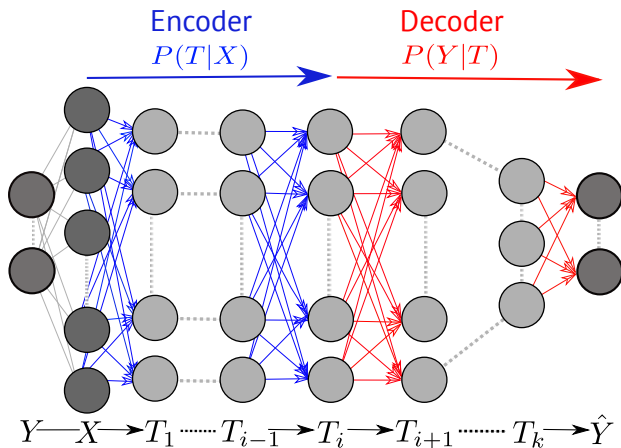
Information Bottleneck trade-off:

$$\min_{p(t|x), p(y|t), p(t)} I(X; T) - \beta I(T; Y)$$

Решения удовлетворяют системе:

$$\left\{ \begin{array}{l} p(t|x) = \frac{p(t)}{Z(x; \beta)} \exp(-\beta D_{KL}[p(y|x) \parallel p(y|t)]) \\ p(t) = \sum_x p(t|x) p(x) \\ p(y|t) = \sum_x p(y|x) p(x|t) \end{array} \right.$$

Нейронные сети как марковская цепь



$$I(X; Y) \geq I(T_1; Y) \geq \dots \geq I(T_k; Y) \geq I(\hat{Y}, Y)$$

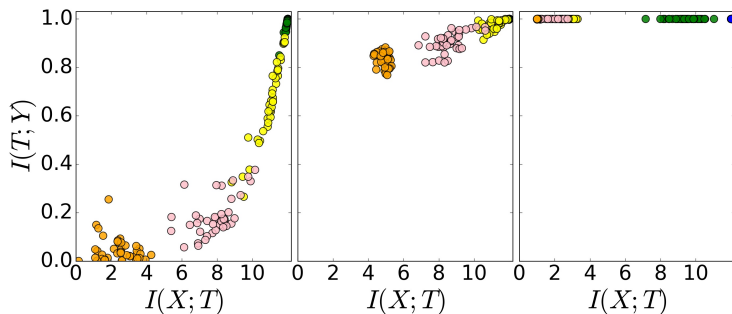
$$H(X) \geq I(X; T_1) \geq \dots \geq I(X; T_k) \geq I(X; \hat{Y})$$

Эксперимент

Сеть: ≤ 7 полносвязных скрытых слоёв с 12-10-7-5-4-3-2 нейронами, активации — \tanh , обучаем обычным SGD.

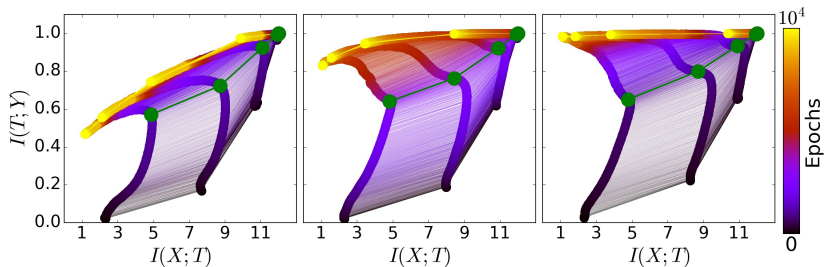
Датасет: искусственно сгенерированный, $X \in \{0, 1\}^{12}$, соответствует 12 точкам на двумерной сфере, $Y \in \{0, 1\}$.

Графики



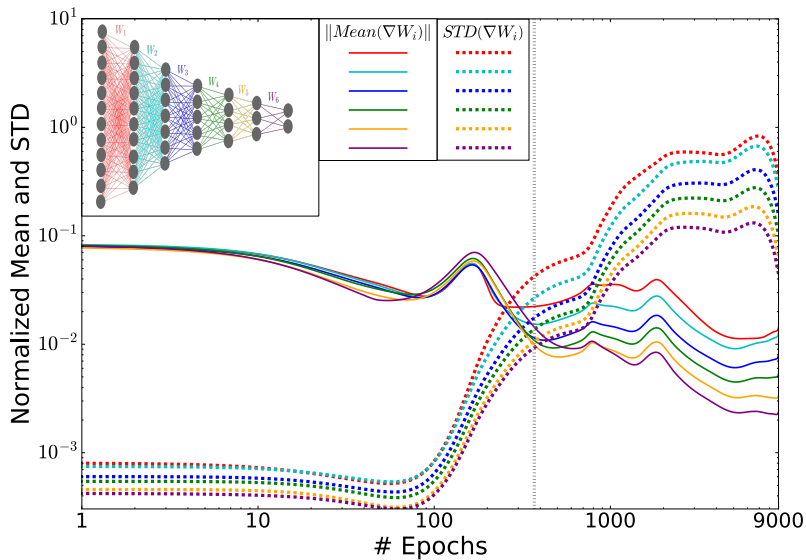
Слева — до обучения, по центру — после 400 эпох, справа — после 9000 эпох. Разные цвета — разные слои.

Графики



Слева — на 5% данных, по центру — на 45%, справа — 85% данных.

Графики



Две фазы SGD

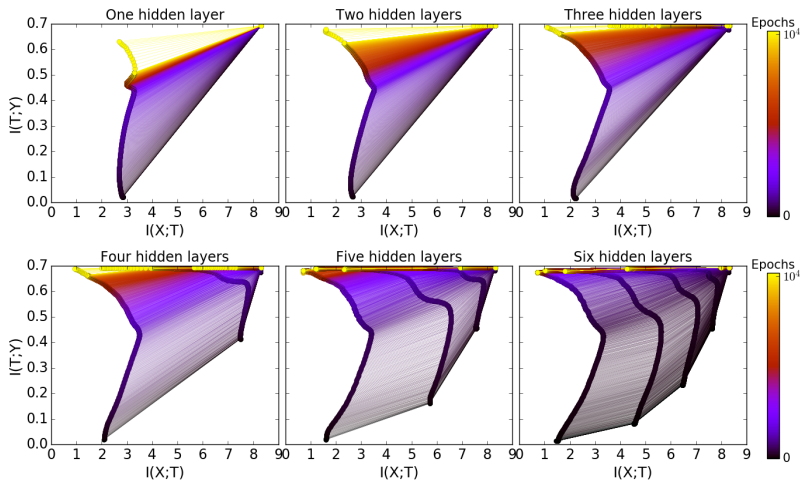
Первая фаза (дрейф): большие средние и маленькие дисперсии градиентов.

Вторая фаза (диффузия): средние значения градиентов гораздо меньше дисперсий, веса фактически совершают случайное блуждание.

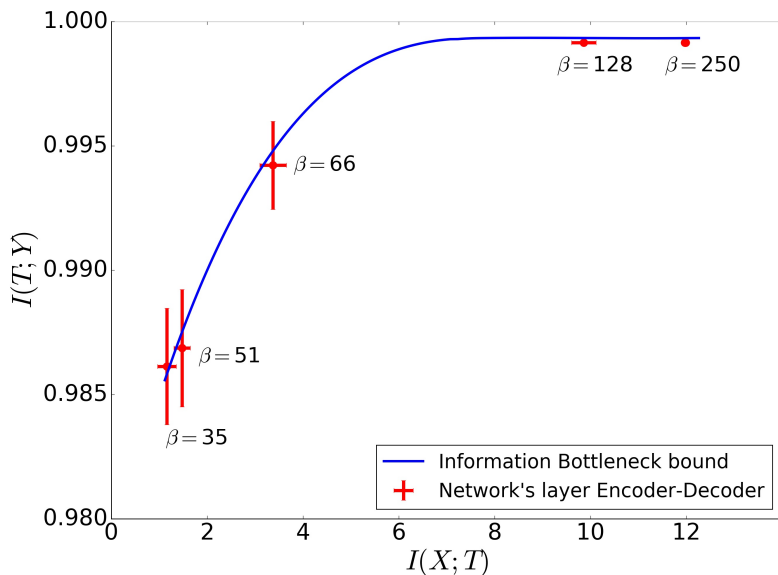
Теоретически подобные случайные блуждания максимизируют $H(X|T_i) = H(X) - I(X; T_i)$, то есть минимизируют $I(X; T_i)$.

Значит, на второй фазе градиентный спуск «сжимает» выученные представления объектов, стараясь не ухудшать $I(T; Y)$.

Влияние числа слоёв на обучение



Близость к Information Bottleneck bound



Результаты

Самое важное: возможно, SGD работает именно из-за диффузии весов сети во второй фазе, когда происходит «сжатие» представлений.

Слои сходятся к Information Bottleneck bound.

Большее число слоёв ускоряет сходимость.

Плюсы/минусы

Плюсы:

Хорошая идея

Интересные результаты, хоть и на искусственном примере и для маленькой сети

Минусы:

Нигде не показаны результаты на отложенной выборке

Маленькая сеть с только полносвязными слоями

Странный искусственный датасет

⇒ непонятно, верны ли выводы на практике, нужно дополнительно проверять

Ссылки

Ravid Shwartz-Ziv, Naftali Tishby:

Opening the Black Box of Deep Neural Networks via Information

<http://arxiv.org/abs/1703.00810>

Naftali Tishby, Fernando C. Pereira, William Bialek:

The information bottleneck method

<http://arxiv.org/abs/physics/0004057>