

Lecture 15

Performer, Variational Autoencoders

Rethinking Attention with Performers

Krzysztof Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane, Tamas Sarlos, Peter Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, David Belanger, Lucy Colwell, Adrian Weller

- “We introduce the first Transformer architectures, *Performers*, capable of **provably** accurate and practical estimations of regular (softmax) full rank attention, but of only linear space and time complexity and **not relying on any priors** such as sparsity or low-rankness. Performers use the *Fast Attention Via positive Orthogonal Random features* (FAVOR+) mechanism”

Generalized definition

- Define three different vectors corresponding to each word.
 - Input $x \in \mathbb{R}^d$ $x = [\underline{x_1} \dots x_n]_{d \times n}$
 - Key
 - Query
 - Value

Generalized definition

- Define three different vectors corresponding to each word.

- Input $x \in \mathbb{R}^d$

$$X = [\underline{x_1} \dots x_n]_{d \times n}$$

- Key $k \in \mathbb{R}^p$

$$K = \underbrace{W_k^T}_{p \times d} X_{d \times n}$$

$$K = [\underline{x_1} \dots k_n]_{p \times n}$$

- Query

- Value

Generalized definition

- Define three different vectors corresponding to each word.

- Input $x \in \mathbb{R}^d$ $X = [\underline{x_1} \dots x_n]_{d \times n}$

- Key $k \in \mathbb{R}^p$ $K = \underbrace{W_k^T}_{p \times d} X_{d \times n}$ $K = [\underline{x_1} \dots k_n]_{p \times n}$

- Query $q \in \mathbb{R}^p$ $Q = \underbrace{W_q^T}_{p \times d} X_{d \times n}$ $Q = [q_1 \dots q_n]_{p \times n}$

- Value

Generalized definition

- Define three different vectors corresponding to each word.

- Input $x \in \mathbb{R}^d$ $X = [\underline{x_1} \dots x_n]_{d \times n}$

- Key $k \in \mathbb{R}^p$ $K = \underbrace{W_k^T}_{p \times d} X_{d \times n}$ $K = [\underline{x_1} \dots k_n]_{p \times n}$

- Query $q \in \mathbb{R}^p$ $Q = \underbrace{W_q^T}_{p \times d} X_{d \times n}$ $Q = [q_1 \dots q_n]_{p \times n}$

- Value $v \in \mathbb{R}^m$ $V = \underbrace{W_v^T}_{m \times d} X_{d \times n}$ $V = [\underline{v_1} \dots v_n]_{m \times n}$

Generalized definition

- Define three different vectors corresponding to each word.

- Input $x \in \mathbb{R}^d$ $X = [\underline{x_1} \dots x_n]_{d \times n}$

- Key $k \in \mathbb{R}^p$ $K = \underbrace{W_k^T}_{p \times d} X_{d \times n}$ $K = [\underline{x_1} \dots k_n]_{p \times n}$

- Query $q \in \mathbb{R}^p$ $Q = \underbrace{W_q^T}_{p \times d} X_{d \times n}$ $Q = [q_1 \dots q_n]_{p \times n}$

- Value $v \in \mathbb{R}^m$ $V = \underbrace{W_v^T}_{m \times d} X_{d \times n}$ $V = [\underline{v_1} \dots v_n]_{m \times n}$

$$\mathcal{L} = V \underset{m \times n}{\text{softmax}} \left(\underbrace{\frac{Q^T K}{\sqrt{P}}}_{n \times n} \right)$$

Generalized definition

- Define three different vectors corresponding to each word.

- Input $x \in \mathbb{R}^d$

$$X = [\underline{x_1} \dots x_n]_{d \times n}$$

- Key $k \in \mathbb{R}^p$

$$K = \underbrace{W_k^T}_{p \times d} X_{d \times n}$$

$$K = [\underline{x_1} \dots k_n]_{p \times n}$$

- Query $q \in \mathbb{R}^p$

$$Q = \underbrace{W_q^T}_{p \times d} X_{d \times n}$$

$$Q = [q_1 \dots q_n]_{p \times n}$$

- Value $v \in \mathbb{R}^m$

$$V = \underbrace{W_v^T}_{m \times d} X_{d \times n}$$

$$V = [\underline{v_1} \dots v_n]_{m \times n}$$

$$\mathcal{L} = V \text{softmax} \left(\underbrace{\frac{Q^T K}{\sqrt{P}}}_{n \times n} \right)$$

$m \times n$ $m \times n$ $n \times p$ $p \times n$

$$O(n^2 m)$$

Generalized definition

- Define three different vectors corresponding to each word.

- Input $x \in \mathbb{R}^d$

$$X = [\underline{x_1} \dots x_n]_{d \times n}$$

- Key $k \in \mathbb{R}^p$

$$K = \underbrace{W_k^T}_{p \times d} X_{d \times n}$$

$$K = [\underline{x_1} \dots k_n]_{p \times n}$$

- Query $q \in \mathbb{R}^p$

$$Q = \underbrace{W_q^T}_{p \times d} X_{d \times n}$$

$$Q = [q_1 \dots q_n]_{p \times n}$$

- Value $v \in \mathbb{R}^m$

$$V = \underbrace{W_v^T}_{m \times d} X_{d \times n}$$

$$V = [\underline{v_1} \dots v_n]_{m \times n}$$

$$\mathcal{L} = V \text{ softmax} \left(\underbrace{\frac{Q^T K}{\sqrt{P}}}_{n \times n} \right)$$

$m \times n$ $m \times n$ $n \times p$ $p \times n$

$$X^T \underbrace{W_q W_k}_{p \times p} X$$

$O(n^2 m)$

$$K(x, y) = \phi(x)^T \phi(y)$$

$$\phi : x \rightarrow \phi(x)$$

$$K(x, y) = \phi(x)^T \phi(y)$$

$$\phi: x \rightarrow \phi(x)$$

- most kernels can be approximated by random features

$$K(x, y) = \phi(x)^T \phi(y)$$

$$\phi \quad x \rightarrow \phi(x)$$

- most kernels can be approximated by random features
- Random features has this form:

$$K(x, y) = \phi(x)^T \phi(y)$$

$$\phi: x \rightarrow \phi(x)$$

- most kernels can be approximated by random features
- Random features has this form:

$$\phi(x) = \frac{h(x)}{\sqrt{r}} (f_1(\underline{\omega}_1^T x), f_1(\underline{\omega}_2^T x) \dots f_1(\underline{\omega}_r^T x) \dots f_l(\underline{\omega}_1^T x) \dots f_l(\underline{\omega}_r^T x))$$

$$K(x, y) = \phi(x)^T \phi(y)$$

$$\phi \quad x \rightarrow \phi(x)$$

- most kernels can be approximated by random features
- Random features has this form:

$$\phi(x) = \frac{h(x)}{\sqrt{r}} \underbrace{(f_1(\underline{\omega}_1^T x), f_1(\underline{\omega}_2^T x) \dots f_1(\underline{\omega}_r^T x) \dots f_l(\underline{\omega}_1^T x) \dots f_l(\underline{\omega}_r^T x))}_{l \times r \text{ elements in vector } \phi}$$

$$\phi(y)$$

$$K(x, y) = \phi^T(x) \phi(y)$$

For example:

$$h(x) = 1 \quad f_1 = \sin \quad f_2 = \cos \quad l = 2$$

$$\underline{\omega} \sim N(0, I_l)$$

For example:

$$h(\underline{x}) = 1 \quad f_1 = \sin \quad f_2 = \cos \quad l = 2 \quad \underline{\omega} \sim N(0, I_l)$$

$$\phi(\underline{x}) = \frac{1}{\sqrt{r}} (\sin(\underline{\omega}_1^T \underline{x}) \sin(\underline{\omega}_2^T \underline{x}) \dots \sin(\underline{\omega}_1^T \underline{x}) \cos(\underline{\omega}_1^T \underline{x}) \dots \cos(\underline{\omega}_r^T \underline{x}))$$

For example:

$$h(x) = 1 \quad f_1 = \sin \quad f_2 = \cos \quad l = 2 \quad \underline{\omega} \sim N(0, I_l)$$

$$\phi(\underline{x}) = \frac{1}{\sqrt{r}} (\sin(\underline{\omega}_1^T x) \sin(\underline{\omega}_2^T x) \dots \sin(\underline{\omega}_1^T x) \cos(\underline{\omega}_1^T x) \dots \cos(\underline{\omega}_r^T x))$$

$$K(x, y) = \phi^T \phi(y) = e^{\frac{-|x-y|^2}{\gamma}}$$

Gaussian

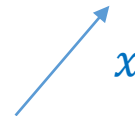
For example:

$$h(x) = 1 \quad f_1 = \sin \quad f_2 = \cos \quad l = 2 \quad \underline{\omega} \sim N(0, I_l)$$

$$\phi(\underline{x}) = \frac{1}{\sqrt{r}} (\sin(\underline{\omega}_1^T x) \sin(\underline{\omega}_2^T x) \dots \sin(\underline{\omega}_1^T x) \cos(\underline{\omega}_1^T x) \dots \cos(\underline{\omega}_r^T x))$$

$$K(x, y) = \phi^T \phi(y) = e^{\frac{-|x-y|^2}{\gamma}}$$

Gaussian



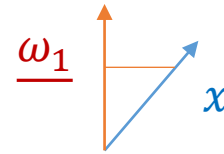
For example:

$$h(x) = 1 \quad f_1 = \sin \quad f_2 = \cos \quad l = 2 \quad \underline{\omega} \sim N(0, I_l)$$

$$\phi(\underline{x}) = \frac{1}{\sqrt{r}} (\sin(\underline{\omega}_1^T x) \sin(\underline{\omega}_2^T x) \dots \sin(\underline{\omega}_1^T x) \cos(\underline{\omega}_1^T x) \dots \cos(\underline{\omega}_r^T x))$$

$$K(x, y) = \phi^T \phi(y) = e^{\frac{-|x-y|^2}{r}}$$

Gaussian



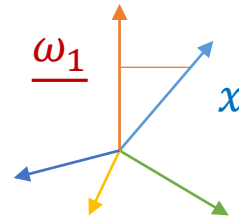
For example:

$$h(x) = 1 \quad f_1 = \sin \quad f_2 = \cos \quad l = 2 \quad \underline{\omega} \sim N(0, I_l)$$

$$\phi(\underline{x}) = \frac{1}{\sqrt{r}} (\sin(\underline{\omega}_1^T x) \sin(\underline{\omega}_2^T x) \dots \sin(\underline{\omega}_1^T x) \cos(\underline{\omega}_1^T x) \dots \cos(\underline{\omega}_r^T x))$$

$$K(x, y) = \phi^T \phi(y) = e^{\frac{-|x-y|^2}{r}}$$

Gaussian



For example:

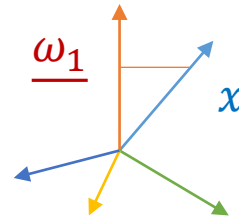
$$h(x) = 1 \quad f_1 = \sin \quad f_2 = \cos \quad l = 2 \quad \underline{\omega} \sim N(0, I_l)$$

$$\phi(\underline{x}) = \frac{1}{\sqrt{r}} (\sin(\underline{\omega}_1^T x) \sin(\underline{\omega}_2^T x) \dots \sin(\underline{\omega}_1^T x) \cos(\underline{\omega}_1^T x) \dots \cos(\underline{\omega}_r^T x))$$

$$K(x, y) = \phi^T \phi(y) = e^{\frac{-|x-y|^2}{\gamma}}$$

Gaussian

$$\text{softmax} \left(\frac{Q^T K}{\sqrt{P}} \right)$$



$$\sigma(\underline{s})_i = a_i = \frac{e^{s_i}}{\sum_j e^{s_j}}$$

$$\sigma(\underline{s})_i = a_i = \frac{e^{s_i}}{\sum_j e^{s_j}}$$

$$A = \exp\left(\frac{Q^T K}{\sqrt{P}}\right)$$

$$\underbrace{\begin{bmatrix} \cdots & \vdots & \vdots \\ \vdots & \ddots & \vdots \\ \cdots & \cdots & \vdots \end{bmatrix}}_{A} \quad n \times n$$

$$\sigma(\underline{s})_i = a_i = \frac{e^{s_i}}{\sum_j e^{s_j}}$$

$$A = \exp\left(\frac{Q^T K}{\sqrt{P}}\right)$$

$$\underbrace{q_{l \times p}^T K_{p \times n}}_{l \times n}$$

$$\underbrace{\begin{bmatrix} \cdots & \vdots & \vdots \\ \vdots & \ddots & \vdots \\ \cdots & \vdots & \vdots \end{bmatrix}}_{A} \Big]_{n \times n}$$

$$\sigma(\underline{s})_i = a_i = \frac{e^{s_i}}{\sum_j e^{s_j}}$$

$$A = \exp\left(\frac{Q^T K}{\sqrt{P}}\right)$$

$$\underbrace{q_{l \times p}^T K_{p \times n}}_{l \times n}$$

$$\underbrace{\begin{bmatrix} \dots & & \\ \vdots & \ddots & \vdots \\ \dots & & \end{bmatrix}}_{A} \quad n \times n$$

$$\begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \quad n \times 1$$

$$\begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix} \quad n \times 1$$

$$\sigma(\underline{s})_i = a_i = \frac{e^{s_i}}{\sum_j e^{s_j}}$$

$$A = \exp\left(\frac{Q^T K}{\sqrt{P}}\right)$$

$$\underbrace{q_{l \times p}^T K_{p \times n}}_{l \times n}$$

$$\underbrace{\begin{bmatrix} \dots & & \\ \vdots & \ddots & \vdots \\ \dots & & \end{bmatrix}}_{A} \quad n \times n$$

$$\begin{matrix} 1 \\ \vdots \\ 1 \end{matrix} \quad n \times 1$$

$$\begin{matrix} x \\ \vdots \\ x \end{matrix} \quad n \times 1$$

$$\text{diag}(\underbrace{A \underline{1}}_D) = \begin{bmatrix} x & 0 & 0 & 0 & - \\ 0 & x & 0 & 0 & - \\ 0 & 0 & x & 0 & - \\ \dots & \dots & \dots & \dots & \dots \end{bmatrix}$$

$$\sigma(\underline{s})_i = a_i = \frac{e^{s_i}}{\sum_j e^{s_j}}$$

$$A = \exp\left(\frac{Q^T K}{\sqrt{P}}\right)$$

$$\underbrace{q_{l \times p}^T K_{p \times n}}_{l \times n}$$

$$\underbrace{\begin{bmatrix} \dots & & \\ \vdots & \ddots & \vdots \\ \dots & & \end{bmatrix}}_{A} \quad n \times n$$

$$\begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \quad n \times 1$$

$$\begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix} \quad n \times 1$$

$$\text{diag}(\underbrace{A \underline{1}}_D) = \begin{bmatrix} x & 0 & 0 & 0 & - \\ 0 & x & 0 & 0 & - \\ 0 & 0 & x & 0 & - \\ \dots & \dots & \dots & \dots & \dots \end{bmatrix}$$

$$\text{softmax}\left(\frac{Q^T K}{\sqrt{P}}\right) = AD^{-1}$$

$$\sigma(\underline{s})_i = a_i = \frac{e^{s_i}}{\sum_j e^{s_j}}$$

$$A = \exp\left(\frac{Q^T K}{\sqrt{P}}\right)$$

$$\underbrace{q_{l \times p}^T K_{p \times n}}_{l \times n}$$

$$\underbrace{\begin{bmatrix} \dots & & \\ \vdots & \ddots & \vdots \\ \dots & & \end{bmatrix}}_{A} \quad n \times n$$

$$\begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \quad n \times 1$$

$$\begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix} \quad n \times 1$$

$$\text{diag}(\underbrace{A \underline{1}}_D) = \begin{bmatrix} x & 0 & 0 & 0 & - \\ 0 & x & 0 & 0 & - \\ 0 & 0 & x & 0 & - \\ \dots & \dots & \dots & \dots & \dots \end{bmatrix}$$

$$\text{softmax}\left(\frac{Q^T K}{\sqrt{P}}\right) = AD^{-1}$$

$$\begin{bmatrix} \frac{1}{x} & 0 & 0 & 0 & - \\ 0 & \frac{1}{x} & 0 & 0 & - \\ 0 & 0 & \frac{1}{x} & 0 & - \\ \dots & \dots & \dots & \dots & \dots \end{bmatrix}$$

$$e^{-|x-y|^2} = e^{-(x-y)^T(x-y)} = e^{-[x^T x + y^T y - 2x^T y]} = e^{-x^T x} \cdot e^{-y^T y} \cdot e^{2x^T y}$$

$$\underbrace{e^{-|x-y|^2}} = e^{\frac{-(x-y)^T(x-y)}{2}} = e^{\frac{-[x^T x + y^T y - 2x^T y]}{2}} = e^{-x^T x} \cdot e^{-y^T y} \cdot \underbrace{e^{2x^T y}}$$

$$\underbrace{e^{-|x-y|^2}} = e^{\frac{-(x-y)^T(x-y)}{2}} = e^{\frac{-[x^T x + y^T y - 2x^T y]}{2}} = \underbrace{e^{\frac{-x^T x}{2}} \cdot e^{\frac{-y^T y}{2}} \cdot e^{\frac{2x^T y}{2}}}$$

$$\underbrace{e^{-|x-y|^2}} = e^{\frac{-(x-y)^T(x-y)}{2}} = e^{\frac{-[x^T x + y^T y - 2x^T y]}{2}} = \underbrace{e^{\frac{-x^T x}{2}} \cdot e^{\frac{-y^T y}{2}} \cdot e^{\frac{2x^T y}{2}}}$$

$$K_{gauss} \cdot e^{\frac{x^T x}{2}} \cdot e^{\frac{y^T y}{2}}$$

$$\underbrace{e^{-|x-y|^2}} = e^{\frac{-(x-y)^T(x-y)}{2}} = e^{\frac{-[x^T x + y^T y - 2x^T y]}{2}} = \underbrace{e^{\frac{-x^T x}{2}} \cdot e^{\frac{-y^T y}{2}} \cdot e^{\frac{2x^T y}{2}}}$$

$$\underbrace{K_{gauss} \cdot e^{\frac{x^T x}{2}} \cdot e^{\frac{y^T y}{2}}}_{K_{SM} \leftarrow e^{x^T y}}$$

$$\underbrace{e^{-|x-y|^2}} = e^{\frac{-(x-y)^T(x-y)}{2}} = e^{\frac{-[x^T x + y^T y - 2x^T y]}{2}} = \underbrace{e^{\frac{-x^T x}{2}} \cdot e^{\frac{-y^T y}{2}} \cdot e^{\frac{2x^T y}{2}}}$$

$$\underbrace{K_{\text{gauss}} \cdot e^{\frac{x^T x}{2}} \cdot e^{\frac{y^T y}{2}}}_{K_{SM} \leftarrow e^{x^T y}}$$

$$h(x) = \frac{x^T x}{2}$$

$$\underbrace{e^{-|x-y|^2}} = e^{\frac{-(x-y)^T(x-y)}{2}} = e^{\frac{-[x^T x + y^T y - 2x^T y]}{2}} = \underbrace{e^{\frac{-x^T x}{2}} \cdot e^{\frac{-y^T y}{2}} \cdot e^{\frac{2x^T y}{2}}}$$

$$\underbrace{K_{gauss} \cdot e^{\frac{x^T x}{2}} \cdot e^{\frac{y^T y}{2}}}_{K_{SM} \leftarrow e^{x^T y}}$$

$$h(x) = \frac{x^T x}{2}$$

$$\phi(x) = \frac{h(x)}{\sqrt{r}} (f_1(\underline{\omega}_1^T x), f_1(\underline{\omega}_2^T x) \dots f_1(\underline{\omega}_r^T x) \dots f_l(\underline{\omega}_1^T x) \dots f_l(\underline{\omega}_r^T x))$$

$$\underbrace{e^{-|x-y|^2}} = e^{\frac{-(x-y)^T(x-y)}{2}} = e^{\frac{-[x^T x + y^T y - 2x^T y]}{2}} = \underbrace{e^{\frac{-x^T x}{2}} \cdot e^{\frac{-y^T y}{2}} \cdot e^{\frac{2x^T y}{2}}}$$

$$\underbrace{K_{gauss} \cdot e^{\frac{x^T x}{2}} \cdot e^{\frac{y^T y}{2}}}_{K_{SM} \leftarrow e^{x^T y}}$$

$$h(x) = \frac{x^T x}{2}$$

$$\phi(x) = \frac{h(x)}{\sqrt{r}} (f_1(\underline{\omega}_1^T x), f_1(\underline{\omega}_2^T x) \dots f_1(\underline{\omega}_r^T x) \dots f_l(\underline{\omega}_1^T x) \dots f_l(\underline{\omega}_r^T x))$$

$$\phi(\underline{x}) = \frac{x^T x}{2\sqrt{r}} (\sin(\underline{\omega}_1^T x) \sin(\underline{\omega}_2^T x) \dots \sin(\underline{\omega}_1^T x) \cos(\underline{\omega}_1^T x) \dots \cos(\underline{\omega}_r^T x))$$

$$V \text{ softmax} \left(\frac{Q^T K}{\sqrt{P}} \right)$$

$$\underbrace{V \operatorname{softmax} \left(\frac{Q^T K}{\sqrt{P}} \right)}_{V Q'^T K'}$$

$$\underbrace{V \operatorname{softmax} \left(\frac{Q^T K}{\sqrt{P}} \right)}_{V Q^T K'}$$

$n \times p$
 $p \times n$

$$\underbrace{V \operatorname{softmax} \left(\frac{Q^T K}{\sqrt{P}} \right)}_{V Q^T K'}$$

$\xrightarrow{\hspace{10em}} n \times p$
 $\xrightarrow{\hspace{10em}} p \times n$

$$O(mn^2)$$

$$V \operatorname{softmax} \left(\frac{Q^T K}{\sqrt{P}} \right)$$

$V (Q'^T K')$

$$O(mn^2)$$

$$\underbrace{m \times n \quad n \times r'}_{m \times r'}$$

$$r' \times n$$

$$r' \times n$$

$$O(mr'n^2)$$

$$V \operatorname{softmax} \left(\frac{Q^T K}{\sqrt{P}} \right)$$

$V (Q'^T K')$

$\xrightarrow{\quad} n \times p$
 $\xrightarrow{\quad} p \times n$

$$O(mn^2)$$

$$\underbrace{m \times n \quad n \times r'}_{m \times r'} \quad \begin{matrix} r' \times n \\ r' \times n \end{matrix}$$

$$O(mr'n^2)$$

$$\begin{bmatrix} \vdots & \dots & \vdots \\ \vdots & \ddots & \vdots \\ \vdots & \dots & \vdots \end{bmatrix}$$

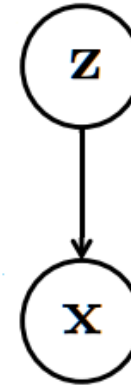
Variational Auto encoder (VEA)

Variational Inference

- Problem Definition

- Observable Data: $x = \{x_1, x_2, \dots, x_n\}$

- Hidden Variable: $z = \{z_1, z_2, \dots, z_n\}$

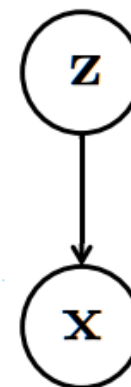


Variational Inference

- Problem Definition

- Observable Data: $x = \{x_1, x_2, \dots, x_n\}$

- Hidden Variable: $z = \{z_1, z_2, \dots, z_n\}$



$$p(z|x) = \frac{p(z, x)}{p(x)} = \frac{p(x|z)p(z)}{\int p(x|z)p(z)dz}$$

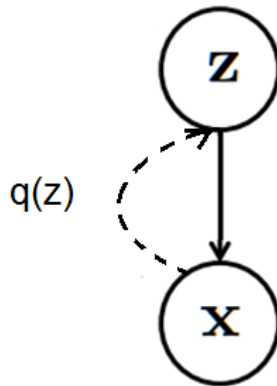
Variational Inference

- Solutions
 - Monte Carlo Sampling
 - Metropolis Hasting
 - Gibbs Sampling
 - Variational Inference

Variational Inference

- Approximate $p(z|x)$ by $q(z)$
- Minimize the KL Divergence:

$$D_{KL} [q(z) || p(z|x)] = - \int q(z) \log \frac{p(z|x)}{q(z)} dz$$



Variational Lower Bound

$$D_{KL} [q(z) || p(z|x)] = - \int q(z) \log \frac{p(z|x)}{q(z)} dz$$

Variational Lower Bound

$$\begin{aligned}D_{KL} [q(z) || p(z|x)] &= - \int q(z) \log \frac{p(z|x)}{q(z)} dz \\ &= - \int q(z) \log \frac{p(z, x)}{q(z)p(x)} dz\end{aligned}$$

Variational Lower Bound

$$\begin{aligned}D_{KL} [q(z) || p(z|x)] &= - \int q(z) \log \frac{p(z|x)}{q(z)} dz \\ &= - \int q(z) \log \frac{p(z, x)}{q(z)p(x)} dz \\ &= - \int q(z) \log \frac{p(z, x)}{q(z)} dz + \int q(z) \log (p(x)) dz\end{aligned}$$

Variational Lower Bound

$$\begin{aligned}D_{KL} [q(z) || p(z|x)] &= - \int q(z) \log \frac{p(z|x)}{q(z)} dz \\&= - \int q(z) \log \frac{p(z, x)}{q(z)p(x)} dz \\&= - \int q(z) \log \frac{p(z, x)}{q(z)} dz + \int q(z) \log(p(x)) dz \\&= - \int q(z) \left(\log(p(z, x)) - \log(q(z)) \right) dz + \log(p(x))\end{aligned}$$

Variational Lower Bound

$$\begin{aligned}D_{KL} [q(z) || p(z|x)] &= - \int q(z) \log \frac{p(z|x)}{q(z)} dz \\&= - \int q(z) \log \frac{p(z, x)}{q(z)p(x)} dz \\&= - \int q(z) \log \frac{p(z, x)}{q(z)} dz + \int q(z) \log(p(x)) dz \\&= - \int q(z) \left(\log(p(z, x)) - \log(q(z)) \right) dz + \log(p(x)) \\&= - \underbrace{\left(E_{q(z)} \left[\log(p(z, x)) \right] - E_{q(z)} \left[\log(q(z)) \right] \right)}_{\text{Evidence Lower Bound (ELBO)}} + \log(p(x))\end{aligned}$$

Variational Lower Bound

$$D_{KL}[q(z)||p(z|x)] = - \underbrace{(E_{q(z)}[\log(p(z, x))] - E_{q(z)}[\log(q(z))])}_{\text{Evidence Lower Bound (ELBO)}} + \log(p(x))$$

Variational Lower Bound

$$D_{KL} [q(z) || p(z|x)] = - \underbrace{(E_{q(z)} [\log(p(z, x))] - E_{q(z)} [\log(q(z))])}_{\text{Evidence Lower Bound (ELBO)}} + \log(p(x))$$

$$D_{KL} [q(z) || p(z|x)] = -L [q(z)] + \log(p(x))$$

Variational Lower Bound

$$D_{KL} [q(z) || p(z|x)] = - \underbrace{(E_{q(z)} [\log(p(z, x))] - E_{q(z)} [\log(q(z))])}_{\text{Evidence Lower Bound (ELBO)}} + \log(p(x))$$

$$D_{KL} [q(z) || p(z|x)] = -L [q(z)] + \log(p(x))$$

$$\log(p(x)) = D_{KL} [q(z) || p(z|x)] + L [q(z)]$$

Variational Lower Bound

$$D_{KL} [q(z) || p(z|x)] = - \underbrace{(E_{q(z)} [\log(p(z, x))] - E_{q(z)} [\log(q(z))])}_{\text{Evidence Lower Bound (ELBO)}} + \log(p(x))$$

$$D_{KL} [q(z) || p(z|x)] = -L [q(z)] + \log(p(x))$$

$$\log(p(x)) = D_{KL} [q(z) || p(z|x)] + L [q(z)]$$

Minimizing $D_{KL} [q(z) || p(z|x)]$
is equal to Maximizing $L [q(z)]$

