

QUESTIONS

LOGITS

$$f(\mathbf{x}) = \text{softmax}(\mathbf{W} \mathbf{x})$$

Are the values of the model before the softmax.

CROSS ENTROPY FORMULA

The **cross-entropy** loss is defined for two vectors $\mathbf{y}, \hat{\mathbf{y}} \in \Delta_c$ as:

$$\text{CE}(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_i y_i \log(\hat{y}_i) .$$

$$\text{CE}(\mathbf{y}, \hat{\mathbf{y}}) = - \log(\hat{y}_t)$$

One-hot encoding -> just one y_i is 1

RNN FORMULA

$$\mathbf{H}_t = \phi(\mathbf{X}_t \mathbf{W}_{xh} + \mathbf{H}_{t-1} \mathbf{W}_{hh} + \mathbf{b}_h)$$

$$\mathbf{O} = \mathbf{H} \mathbf{W}_{hq} + \mathbf{b}_q$$

GRU

$$\mathbf{R}_t = \sigma(\mathbf{X}_t \mathbf{W}_{xr} + \mathbf{H}_{t-1} \mathbf{W}_{hr} + \mathbf{b}_r)$$

$$\mathbf{Z}_t = \sigma(\mathbf{X}_t \mathbf{W}_{xz} + \mathbf{H}_{t-1} \mathbf{W}_{hz} + \mathbf{b}_z)$$

$$\tilde{\mathbf{H}}_t = \tanh(\mathbf{X}_t \mathbf{W}_{xh} + (\mathbf{R}_t \odot \mathbf{H}_{t-1}) \mathbf{W}_{hh} + \mathbf{b}_h)$$

$$\mathbf{H}_t = \mathbf{Z}_t \odot \mathbf{H}_{t-1} + (1 - \mathbf{Z}_t) \odot \tilde{\mathbf{H}}_t$$

LSTM

$$\mathbf{I}_t = \sigma(\mathbf{X}_t \mathbf{W}_{xi} + \mathbf{H}_{t-1} \mathbf{W}_{hi} + \mathbf{b}_i),$$

$$\mathbf{F}_t = \sigma(\mathbf{X}_t \mathbf{W}_{xf} + \mathbf{H}_{t-1} \mathbf{W}_{hf} + \mathbf{b}_f)$$

$$\mathbf{O}_t = \sigma(\mathbf{X}_t \mathbf{W}_{xo} + \mathbf{H}_{t-1} \mathbf{W}_{ho} + \mathbf{b}_o),$$

$$\tilde{\mathbf{C}}_t = \tanh(\mathbf{X}_t \mathbf{W}_{xc} + \mathbf{H}_{t-1} \mathbf{W}_{hc} + \mathbf{b}_c)$$

$$\mathbf{C}_t = \mathbf{F}_t \odot \mathbf{C}_{t-1} + \mathbf{I}_t \odot \tilde{\mathbf{C}}_t$$

$$\mathbf{H}_t = \mathbf{O}_t \odot \tanh(\mathbf{C}_t)$$

LOSS GAN:

$$J_D = -\frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \log D(\mathbf{x}) - \frac{1}{2} \mathbb{E}_{\mathbf{z} \sim p_G(\mathbf{z})} \log(1 - D(G(\mathbf{z})))$$

$$J_G = -J_D$$

MINMAX The generator minimizes the log-probability of the discriminator being correct.

$$J_D = -\frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \log D(\mathbf{x}) - \frac{1}{2} \mathbb{E}_{\mathbf{z} \sim p_G(\mathbf{z})} \log(1 - D(G(\mathbf{z})))$$

$$J_G = -\frac{1}{2} \mathbb{E}_{\mathbf{z} \sim p_G(\mathbf{z})} \log D(G(\mathbf{z}))$$

NON-SATURATING GAME the generator maximizes the log-probability of the discriminator being mistaken.

The non-saturating game is heuristically motivated. Generator can still learn even when discriminator successfully rejects all generator samples.

BATCH NORMALIZATION

$$\mathbf{X}' = \frac{\mathbf{X} - \mu}{\sqrt{\sigma^2}}$$

$$\tilde{\mu}_j = \frac{1}{b} \sum_i [\mathbf{H}]_{i,j}, \quad \tilde{\sigma}_j^2 = \frac{1}{b} \sum_i ([\mathbf{H}]_{i,j} - \tilde{\mu}_j)^2$$

Empirical mean and variance

$$\mathbf{H}' = \frac{\mathbf{H} - \tilde{\mu}}{\sqrt{\tilde{\sigma}^2 + \epsilon}}$$

BN in convolution layers

$$\tilde{\mu}_z = \frac{1}{bhw} \sum_{i,j,k} [\mathbf{H}]_{i,j,k,z}, \quad \tilde{\sigma}_z^2 = \frac{1}{bhw} \sum_{i,j,k} ([\mathbf{H}]_{i,j,k,z} - \tilde{\mu}_z)^2$$

Where h and w are the height and width respectively of the feature map

LAYER NORMALIZATION

$$\tilde{\mu}_i = \frac{1}{f} \sum_j [\mathbf{H}]_{i,j}, \quad \tilde{\sigma}_i^2 = \frac{1}{f} \sum_j ([\mathbf{H}]_{i,j} - \tilde{\mu}_i)^2$$

CONVOLUTION

$$h[i,j] = u[i,j] + \sum_{k_1, k_2} W[i,j,k_1,k_2] \cdot x[i+k_1, j+k_2]$$

Translation Invariance-> a shift in the inputs \mathbf{x} should lead to a shift in h -> W and u do not depend on i and j

$$h[i,j] = u + \sum_{k_1, k_2} W[k_1, k_2] \cdot x[i+k_1, j+k_2]$$

Locality principle This means that outside some range $|k_1|, |k_2| > \Delta$, we set $W[k_1, k_2] = 0$:

$$h[i,j] = u + \sum_{k_1=-\Delta}^{\Delta} \sum_{k_2=-\Delta}^{\Delta} W[k_1, k_2] \cdot x[i+k_1, j+k_2]$$

QUESTIONS

DILATED CONVOLUTION

$$y[i, j] = \sum_{k_1=0}^{M_1-1} \sum_{k_2=0}^{M_2-1} h[k_1, k_2] x[i - d \cdot k_1, j - d \cdot k_2]$$

which expands the receptive field without loss of resolution. The expansion of the filter allows to increase its dimensions by filling the empty positions with zeros.

CAUSAL CONVOLUTION

$$[\mathbf{H}]_{i,d} = \phi \left(\sum_{i'=i-k}^{+k} \sum_{z=1}^c [W]_{i'+k+1,z,d} [\mathbf{X}]_{i+i',z} \right)$$

(n, c')

n temporal seq, c' output features, n temporal seq, c input features.

i temporal indices through n, d indices through c', i' temporal kernel indices, k dim of window convolution, z features, s filter dimensions

$$[\mathbf{H}]_{i,d} = \phi \left(\sum_{i'=0}^k \sum_{z=1}^c [W]_{i'+1,z,d} [\mathbf{X}]_{i-i',z} \right)$$

(n, c')

i-i' stands for only the left values of the x

SELF ATTENTION LAYER

$$\mathbf{h}_i = \sum_{j=-k}^k \mathcal{W}_j \mathbf{x}_{i+j}$$

1D Convolution->
Remove Locality

$$\mathbf{h}_i = \sum_{j=-i+1}^{n-i} g(i+j) \mathbf{x}_j$$

Continuous
Convolution

$$\mathbf{h}_i = \sum_{j=1}^n \alpha(\mathbf{x}_i, \mathbf{x}_j) \mathbf{x}_j$$

Remove regular
Dependencies
Non local nn

$$\mathbf{h}_i = \sum_{j=1}^n \text{softmax}_j \left(\frac{1}{\sqrt{d}} \mathbf{x}_i^\top \mathbf{x}_j \right) \mathbf{x}_j$$

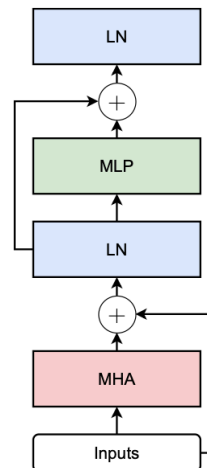
Normdot prod
Fast efficient
parallelize

$$\mathbf{Q} = \mathbf{XW}_q, \quad \mathbf{K} = \mathbf{XW}_k, \quad \mathbf{V} = \mathbf{XW}_v$$

(n, q) (n, q) (n, v)

$$\mathbf{H} = \text{softmax} \left(\frac{\mathbf{QK}^\top}{\sqrt{q}} \right) \mathbf{V}$$

(n, v)



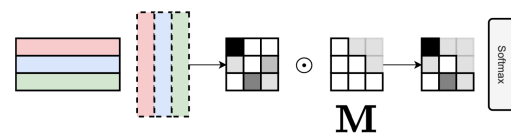
(a) Post-normalized block

$$\mathbf{X}' = [\mathbf{X} \parallel \mathbf{E}] \quad \text{or} \quad \mathbf{X}' = \mathbf{X} + \mathbf{E}$$

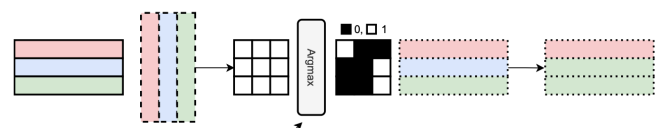
CAUSAL TRANSFORMER

$$\mathbf{H} = \text{softmax} \left(\frac{\mathbf{QK}^\top \odot \mathbf{M}}{\sqrt{q}} \right) \mathbf{V}$$

$$M_{ij} = \begin{cases} 1 & \text{if } j \leq i \\ -\infty & \text{otherwise} \end{cases}$$



HARD ATTENTION



EMBEDDINGS

Permutations

$$\text{MHA}(\mathbf{PX} \parallel \mathbf{E}) \neq \mathbf{P} \cdot \text{MHA}(\mathbf{X} \parallel \mathbf{E})$$

QUESTIONS

LINEAR TRANSFORMERS

$$\mathbf{h}_i = \frac{\sum_j \phi(\mathbf{q}_i)^\top \phi(\mathbf{k}_j) \mathbf{v}_j}{\sum_j \phi(\mathbf{q}_i)^\top \phi(\mathbf{k}_j)} = \frac{\phi(\mathbf{q}_i)^\top \overbrace{\sum_j \phi(\mathbf{k}_j) \mathbf{v}_j}^{\mathbf{S}}}{\underbrace{\phi(\mathbf{q}_i)^\top \sum_j \phi(\mathbf{k}_j)}_{\mathbf{Z}}}.$$

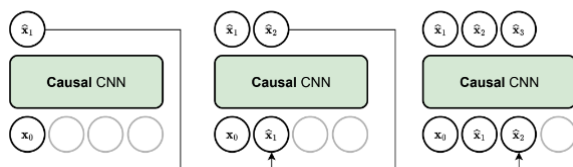
$O(n^2) \rightarrow O(n)$

RELATIVE POSITION EMBEDDING ATTENTION

$$\alpha(\mathbf{x}_i, \mathbf{x}_j, i - j) = \mathbf{x}_i^\top \mathbf{x}_j + b_{ij}$$

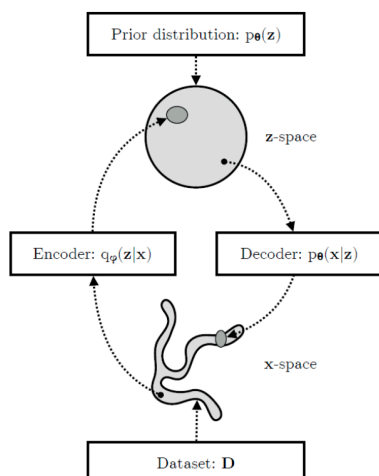
Depend on the relative distance

AUTOREGRESSIVE MODEL FOR TIME SERIES



train the model to predict the next time step by shifting the target sequence and teaching the model to predict future values based only on past data. This enables autoregressive generation, where the model predicts one step at a time, using previous predictions to generate the next step.

VARIATIONAL AUTOENCODER



$$q_\phi(\mathbf{z}|\mathbf{x}) \approx p_\theta(\mathbf{z}|\mathbf{x}) \quad p_\theta(\mathbf{x}) = \int p_\theta(\mathbf{x}, \mathbf{z}) d\mathbf{z}$$

$$p_\theta(\mathbf{x}_1, \dots, \mathbf{x}_M) = \prod_{j=1}^M p_\theta(\mathbf{x}_j | A(\mathbf{x}_j))$$

$$\boldsymbol{\eta} = f(A(\mathbf{x}))$$

$$p_\theta(\mathbf{x} | A(\mathbf{x})) = p_\theta(\mathbf{x} | \boldsymbol{\eta})$$

Parametric inference model- joint distribution

$$\begin{aligned} \log p_\theta(\mathbf{x}) &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \{\log p_\theta(\mathbf{x})\} \\ &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left\{ \log \left(\frac{p_\theta(\mathbf{x}, \mathbf{z})}{p_\theta(\mathbf{z}|\mathbf{x})} \right) \right\} \\ &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left\{ \log \left(\frac{p_\theta(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \frac{q_\phi(\mathbf{z}|\mathbf{x})}{p_\theta(\mathbf{z}|\mathbf{x})} \right) \right\} \\ &= \underbrace{\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left\{ \log \left(\frac{p_\theta(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right) \right\}}_{\mathcal{L}_{\theta, \phi}(\mathbf{x})} + \underbrace{\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left\{ \log \left(\frac{q_\phi(\mathbf{z}|\mathbf{x})}{p_\theta(\mathbf{z}|\mathbf{x})} \right) \right\}}_{\mathcal{D}_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x}) || p_\theta(\mathbf{z}|\mathbf{x}))} \end{aligned}$$

ELBO - Evidence Lower Bound

A lower bound on the log-likelihood of the data, that if it is maximized Maximize the marginal likelihood, $p_\theta(\mathbf{z}|\mathbf{x})$ which improves the generation. Minimize the KL Divergence improve the approximation of

$$q_\phi(\mathbf{z}|\mathbf{x}) \text{ to } p_\theta(\mathbf{z}|\mathbf{x}).$$

Kullback - Leibler distance distance from 2 probabilities distribution

Calculating the gradient

$$\begin{aligned} \nabla_{\theta} \mathcal{L}_{\theta, \phi}(\mathbf{x}) &= \nabla_{\theta} \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \{\log p_\theta(\mathbf{x}, \mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x})\} \\ &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \{\nabla_{\theta} (\log p_\theta(\mathbf{x}, \mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x}))\} \\ &\simeq \nabla_{\theta} (\log p_\theta(\mathbf{x}, \mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x})) \\ &= \nabla_{\theta} (\log p_\theta(\mathbf{x}, \mathbf{z})) \end{aligned}$$

QUESTIONS

$$\begin{aligned}\nabla_{\phi} \mathcal{L}_{\theta, \phi}(\mathbf{x}) &= \nabla_{\phi} \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \{ \log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x}) \} \\ &\neq \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \{ \nabla_{\theta} (\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})) \}\end{aligned}$$

$$\log d_{\phi}(\mathbf{x}, \epsilon) = \log \left| \det \left(\frac{\partial \mathbf{z}}{\partial \epsilon} \right) \right| = \sum_i \log \sigma_i$$

REPARAMETRIZATION TRICK

Can't be approximated, difficult to obtain -> Changes of variables of invertible and differentiable

$$\mathbf{z} = g(\epsilon, \phi, \mathbf{x})$$

$$\begin{aligned}\nabla_{\phi} \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \{ f(\mathbf{z}) \} &= \nabla_{\phi} \mathbb{E}_{p(\epsilon)} \{ f(\mathbf{z}) \} \\ &= \mathbb{E}_{p(\epsilon)} \{ \nabla_{\phi} f(\mathbf{z}) \} \\ &\simeq \nabla_{\phi} f(\mathbf{z}).\end{aligned}$$

$$\begin{aligned}\mathcal{L}_{\theta, \phi}(\mathbf{x}) &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \{ \log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x}) \} \\ &= \mathbb{E}_{p(\epsilon)} \{ \log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x}) \}\end{aligned}$$

Monte Carlo estimator

$$\begin{aligned}\epsilon &\sim p(\epsilon) \\ \mathbf{z} &= g(\phi, \mathbf{x}, \epsilon) \\ \tilde{\mathcal{L}}_{\theta, \phi}(\mathbf{x}) &= \log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})\end{aligned}$$

$$\log q_{\phi}(\mathbf{z}|\mathbf{x}) = \log p(\epsilon) - \log d_{\phi}(\mathbf{x}, \epsilon)$$

$$\log d_{\phi}(\mathbf{x}, \epsilon) = \log \left| \det \left(\frac{\partial \mathbf{z}}{\partial \epsilon} \right) \right|$$

$$q_{\phi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}^2))$$

Factorize Gaussian Encoder

$$\begin{aligned}\epsilon &\sim \mathcal{N}(0, \mathbf{I}) \\ (\boldsymbol{\mu}, \log \boldsymbol{\sigma}) &= f_{\phi}(\mathbf{x}) \\ \mathbf{z} &= \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \epsilon\end{aligned}$$

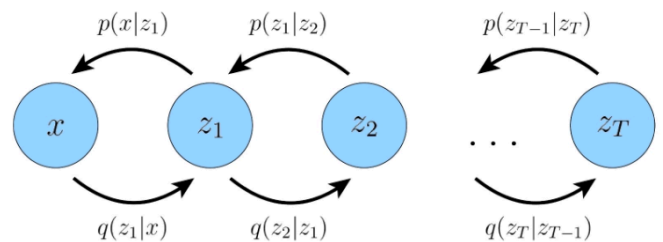
$$\log q_{\phi}(\mathbf{z}|\mathbf{x}) = \log p(\epsilon) - \log d_{\phi}(\mathbf{x}, \epsilon) = \sum_i \log \mathcal{N}(\epsilon_i; 0, 1) - \log \sigma_i.$$

HIERARCHICAL VARIATIONAL AUTOENCODER

A generalization of VAE that extends to multiple hierarchies over latent variables, "more abstract latents"

MARKOV HVAE

Each Transition down hierarchy is Markovian, where decoding each latent z_t only condition on previous latent z_{t-1}



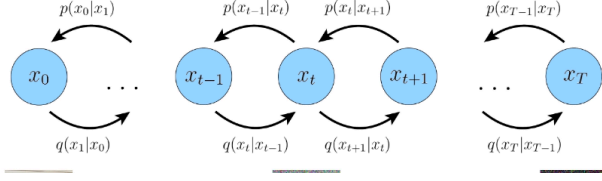
$$p(\mathbf{x}, \mathbf{z}_{1:T}) = p(\mathbf{z}_T) p_{\theta}(\mathbf{x}|\mathbf{z}_1) \prod_{t=2}^T p_{\theta}(\mathbf{z}_{t-1}|\mathbf{z}_t)$$

$$q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x}) = q_{\phi}(\mathbf{z}_1|\mathbf{x}) \prod_{t=2}^T q_{\phi}(\mathbf{z}_t|\mathbf{z}_{t-1})$$

$$\log p(\mathbf{x}) \geq \mathbb{E}_{q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}_{1:T})}{q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x})} \right]$$

QUESTIONS

VARIATIONAL DIFFUSION MODELS



Latent dimension is exactly equal to the data dimensions.

The structure of the latent encoder at each timestep is not learned but pre-defined as a Linear Gaussian model

Parameters of latent encoder vary over time such as way that the distribution of the latent at final time step T is a standard gaussian

$$q(\mathbf{x}_{1:T}|\mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})$$

$$\mu_t(\mathbf{x}_t) = \sqrt{\alpha_t}\mathbf{x}_{t-1} \quad \text{and} \quad \Sigma_t(\mathbf{x}_t) = (1 - \alpha_t)\mathbf{I},$$

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t}\mathbf{x}_{t-1}, (1 - \alpha_t)\mathbf{I})$$

$$p(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)$$

where $p(\mathbf{x}_T) = \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$.

$$\begin{aligned} \log p(\mathbf{x}) &= \log \int p(\mathbf{x}_{0:T}) d\mathbf{x}_{1:T} \\ &\geq \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] \\ &= \underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)]}_{\text{reconstruction term}} - \underbrace{\mathbb{E}_{q(\mathbf{x}_{T-1}|\mathbf{x}_0)} [\mathcal{D}_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_{T-1}) || p(\mathbf{x}_T))]}_{\text{prior matching term}} \\ &\quad - \sum_{t=1}^{T-1} \underbrace{\mathbb{E}_{q(\mathbf{x}_{t-1}, \mathbf{x}_{t+1}|\mathbf{x}_0)} [\mathcal{D}_{\text{KL}}(q(\mathbf{x}_t|\mathbf{x}_{t-1}) || p_{\theta}(\mathbf{x}_t|\mathbf{x}_{t+1}))]}_{\text{consistency term}} \end{aligned}$$

DROPOUT

M is a binary matrix drawn from a Bernoulli matrix with probability (1-p). Dropout in the training .

$$\tilde{\mathbf{H}} = \mathbf{H} \odot \mathbf{M},$$

MONTECARLO DROPOUT

Present Inside the inference.

The output is computed averaging the result with different dropout masks

$$\hat{y} = \mathbb{E}_{p(\mathbf{M})} [f(\mathbf{x}; \mathbf{M})] \approx \frac{1}{N} \sum_{i=1}^N f(\mathbf{x}; \mathbf{M}_i)$$

WASSERSTEIN GAN

We want a mass transportation of the distribution to the preferred one.

We want to calculate “Earth Mover Distance” minim amount of “Work” to this transportation

$$\max_{w \in \mathcal{W}} \mathbb{E}_{x \sim p(x)} \{D_w(x)\} - \mathbb{E}_{\tilde{x} \sim p(\tilde{x})} \{D_w(\tilde{x})\}$$

$$\leq \sup_{\|f\|_L \leq 1} \mathbb{E}_{x \sim p(x)} \{D(x)\} - \mathbb{E}_{\tilde{x} \sim p(\tilde{x})} \{D(\tilde{x})\} = K \cdot \mathcal{W}(\cdot)$$

We replace the discriminator with the critic. Its goal is to output scores for real and fake data that allow the generator to minimize the real data distribution and the generated distribution

$$\mathcal{W}(p(x), p(\tilde{x})) = \sup_{\|f\|_L \leq 1} \mathbb{E}_{x \sim p(x)} \{f(x)\} - \mathbb{E}_{\tilde{x} \sim p(\tilde{x})} \{f(\tilde{x})\}$$