

Latent Variable Models

Distributions $P(X)$

Dashboard X in high dimensional space \mathcal{X}

Generative model task



X distributed according to unknown $P_{gt}(X)$



A model P is as similar to P_{gt}

Situations:

Every dashboard X in the data set, there is one or more settings of the latent variables which causes the model to generate something similar to X



A vector of latent variable in high dimensional space Z , which can be easily sampled according to probability density function $P(z)$ defined over Z

→ shows $P = \int P(z) dz$ ↗ variable

↓ probability distribution ↓ density function

Then, create a family of deterministic functions $f(z; \theta)$, parametered by θ in some space Θ where $f: Z \times \Theta \rightarrow \mathcal{X}$

Then $f(z; \theta)$ is a random variable in the space \mathcal{X}

From all of above, we want to optimize θ such that we can sample z from $P(z)$ and, with high probability, $f(z; \theta)$ will be like the X 's in our dataset

maximize the probability of each x according to

$$P(x) = \int f(z; \theta) P(z) dz$$

↓

$f(z; \theta)$ can be replaced by a distribution $P(x|z; \theta)$, which allows us to make the dependence of x on z

maximum likelihood \Rightarrow
 1. training sets samples
 2. Similar
 2. dissimilar

The output distribution is often Gaussian in VAE

$$P(x|z; \theta) = \mathcal{N}(x | f(z; \theta), \sigma^2 * I)$$

means it has mean $f(z; \theta)$
 covariance = $I * \sigma^2$

why use Gaussian
 ↓
 use gradient descent to increase $P(x)$
 ↓
 by making $f(z; \theta)$ approach x for the variable z

$$P(z) = \mathcal{N}(z | 0, I)$$

setting up the objective.

we need to introduce a new function $Q(Z|X)$ which can take a value of X and give us a distribution over Z values that likely to produce X .
which means the space of Z values under Q will be much smaller than $P(Z)$
why? \rightarrow for most of Z , $P(X|Z)$ will be zero, hence contribute nothing to our estimate of $P(X)$

Kullback-Leibler divergence (KL divergence or D)

$$D[Q(Z)||P(Z|X)] = E_{Z \sim Q} [\log Q(Z) - \log P(Z|X)]$$

\downarrow by applying Bayes rule
to $P(Z|X)$

$$D[Q(Z)||P(Z|X)] = E_{Z \sim Q} [\log Q(Z) - \log P(X|Z) - \log P(Z)] + \log P(X)$$

\downarrow $\log P(X)$ is not depend on Z
rearrange the equation

$$\log P(X) - D[Q(Z)||P(Z|X)] = E_{Z \sim Q} [\log P(X|Z)] - D[Q(Z)||P(Z)]$$

$$Q(Z) = Q(Z|X)$$

$$\min D[Q(Z|X)||P(Z|X)]$$

\downarrow $\log P(X)$ plus an error term (which make's Q produce Z 's that can reproduce a given X ; this term should be small, if Q is high capacity)

\downarrow something we can optimise
using stochastic descent

The framework:

Q is 'encoding' X into Z

P is 'decoding' Z into X

$P(Z|X)$ is not something we can compute analytically: it describe the values of Z that are likely to give rise to a sample like X because $Q(Z|X)$ is pulling to match $P(Z|X)$, make intractable $P(Z|X)$ tractable

Optimizing the objective.

how to use stochastic gradient descent on the right hand side? $\rightarrow Q(Z|X)$ need to be specific

$\therefore Q(Z|X) = \mathcal{N}(Z|U(X); \mathcal{Q}, \Sigma(X; \mathcal{Q}))$ U, Σ are arbitrary deterministic functions with parameter \mathcal{Q} \rightarrow can be learned from data.

Σ is constrained to diagonal matrix to make it computational and clear.

\mathcal{Q} is omitted for easy handwriting

So, the KL divergence $D[Q(Z|X) || P(Z)]$ can be expressed by two multivariate Gaussian distributions

$$D(\mathcal{N}(U_0, \Sigma_0) || \mathcal{N}(U_1, \Sigma_1)) = \frac{1}{2} \left(\text{tr}(\Sigma_1^{-1} \Sigma_0) + (U_1 - U_0)^T \Sigma_1^{-1} (U_1 - U_0) - k + \log \left(\frac{\det \Sigma_1}{\det \Sigma_0} \right) \right)$$

where k is the dimensionality of the distribution

\downarrow it can be simplified to.

$$D[\mathcal{N}(U(X), \Sigma(X)) || \mathcal{N}(0, I)] = \frac{1}{2} (\text{tr} \Sigma(X) + (U(X))^T (U(X)) - k - \log \det(\Sigma(X)))$$

For equation

$$\log P(X) - D[Q(Z|X) || P(Z|X)] = E_{Z \sim Q} [\log P(X|Z)] - D[Q(Z|X) || P(Z)]$$

Tricks:

for $E_{Z \sim Q} [\log P(X|Z)]$, if we use sampling to estimate it, getting a good estimate requires passing many samples of Z through f , it's expensive.

So, in std, we take one sample of Z and treat $\log P(X|Z)$ as an approximation of $E_{Z \sim Q} [\log P(X|Z)]$. After all, we already doing std over different values of X sampled from dataset.

So the equation we need to optimize is

$$E_{X \sim D} [\log P(X) - D[Q(Z|X) || P(Z|X)]] = E_{X \sim D} [E_{Z \sim Q} [\log P(X|Z)] - D[Q(Z|X) || P(Z)]]$$

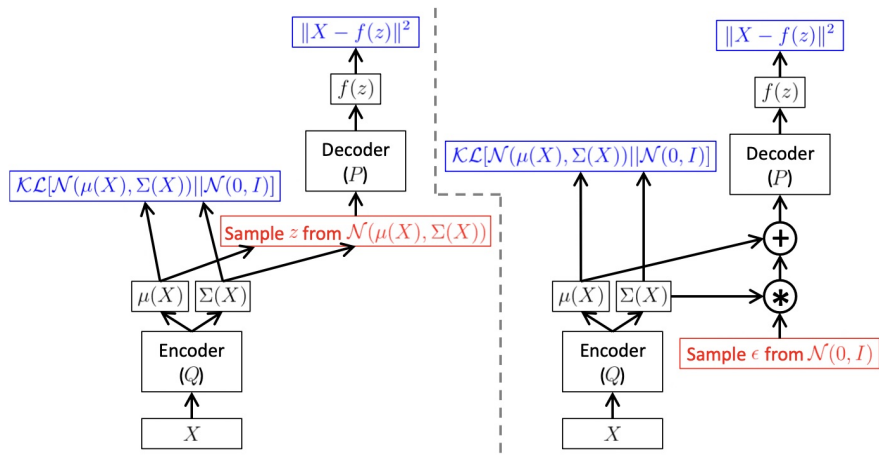


Figure 4: A training-time variational autoencoder implemented as a feed-forward neural network, where $P(X|z)$ is Gaussian. Left is without the “reparameterization trick”, and right is with it. Red shows sampling operations that are non-differentiable. Blue shows loss layers. The feedforward behavior of these networks is identical, but backpropagation can be applied only to the right network.

Probability density function

Definition

value at any given sample (or point) in the sample space (the set of possible values taken by the random variable) can be interpreted as providing a relative likelihood that the value of random variable would be equal to that sample.

probability density is the probability per unit length



So, absolute likelihood for a particular value is 0



Since an infinite set of possible random values.

Absolutely continuous univariate distribution

$$Pr[a \leq X \leq b] = \int_a^b f(x) dx \quad \text{A random variable has density } f_x$$

And F_x is the cumulative distribution function of X

$$F_x = \int_{-\infty}^x f(u) du.$$

if f_x is continuous at (x)

$$f_x(x) = \frac{d}{dx}(F_x(x))$$

countable intersection

$$X \cap A, X \cup A, C_n(C \cap X) = U \setminus (C \cap X)$$

Formal definition

A random variable X with values in a measurable space (X, A) (usually \mathbb{R}^n with the Borel sets as subsets) has as probability distribution the measure X_*P on (X, A) : the density of X with respect to a reference measure u on (X, A) is:

$$f = \frac{dX_*P}{du}$$

That is, f is any measurable function

$$Pr[X \in A] = \int_{X^{-1}A} dP = \int_A f du$$

Covariance

$$\text{cov}(x, y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n} \rightarrow \text{for population}$$

$$\text{Cov}(X, Y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n-1} \rightarrow \text{for a sample.}$$

↓ just show the direction, while correlation can show the strength of relationship

$$\rho(x, y) = \frac{\text{Cov}(X, Y)}{\sigma_x \sigma_y} \rightarrow \text{standard deviation.}$$