

# Tutorial on variational Autoencoders

## variational Autoencoders (vae's)

↳ circa 2016 emerged as one of the most popular approaches to unsupervised learning built on top of standard function approximators (neural networks)

type of "generative" modelling

goal of generative models is to learn some distribution  $p(\cdot)$  that is as close as possible to  $p_{\text{gt}}(\cdot)$  (generator)

VAE

- ↳ Assumptions are weak (this is good)
- ↳ Backprop train is "fast"

---

### 1.1 Preliminaries: Latent variable models

W/ gen models, the more complicated the dependencies between dimensions = harder to train

A gen model performs well when it decides which character [0-9] to draw before determining any pixels

↳ this is latent variable



## 2 Variational Autoencoders

Mathematically, VAEs have little to do w/ classical Autoencoders

trad Autoencoder has encoder-decoder

In training VAEs resemble this structure

to solve train equation VAEs have 2 problems to solve

① how to define latent variables,  $z$   
(decide what info they hold)

② how to deal w/ integral over  $z$

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

### ① latent vars

for digits need to not just encode the digit in its latent dimensions but its abstract stylistic properties

we want to avoid deciding by hand the info each dimension of  $z$  encodes & the dimensions dependencies

VAEs assume there is no simple interp of the dimensions of  $z$



(3)

instead they assert that ~~these~~ sample  $q$ ,  $z$  can be drawn from a simple dist  
 $\rightarrow N(0, I)$

$I = \text{identity matrix}$

how?

A distribution in  $d$  dimensions can be generated by taking a set of  $d$  variables w/ a norm distributed & mapping them through a sufficiently complicated function

example:

$\rightarrow$  norm dist points of 2D graph

$\rightarrow$  eg:  $g(z) = z/10 + z/\|z\|$

$\rightarrow$  ring shape of same points on 2D graph

provided powerful function approx (NN)  
we can simply learn a function which maps our indep, norm dist  $z$  values to whatever latent variables are needed for the model

$\hookrightarrow$  then map latent vars to  $x$

$$N(x | f(z, 0))$$

$\hookrightarrow$   $f$  can be an NN

within the NN layers = ① learn latent

② map to pixels to render



don't need to hand craft the latent structure - or make sure it exists

if there exists a latent structure that maximizes the likelihood (ability to reproduce) then the NN layers will uncover it