

Probabilistic models: $p(x)$ → unsupervised learning (learn the distribution) aka "generative".
 $p(y|x)$ → supervised learning

Key user of generative modelling: - rep. learning - actually generating things - same uses as language modeling.

"language modelling" images is possible - Pixel RNNs.
 → how should we read pixels?

In paper: scanline order, ← issue: pixels right above are considered "far away".

$$p(x) = \prod_{i=1}^n p(x_i | x_1, \dots, x_{i-1})$$

256-way softmax

→ generate one color channel at a time.

practical: Pixel CNN (must be faster, only uses masked neighborhood of pixel)

robust: Pixel Transformer

- ↳ all pixels are equally "close" for self-attention.
- ↳ to reduce computation, only compute attention for nearby pixels ("modified" masked attention)

Autoregressive generative models

Training:

1. Divide x into x_1, \dots, x_n dimensions
2. Discretize each $x_i \rightarrow k$ values
3. Model $p(x)$ using chain rule
4. Use sequence model to model $p(x)$.

Using:

- ↳ Sampling (x_1, x_2, \dots , etc.)
- ↳ Completion (feed in known values)
- ↳ Representation: like GPT-2/BERT.

Conditional autoregressive models - generate something conditioned on another thing.

- ↳ pass in context thru encoder prior to "kicking off" sequence model.

General design for generative models:

image → nontrivial model → loss → same image. (RNN, Pixel CNN/RNN, Transformer)

Autoencoders: train network to encode image, then decode same image.

image → encoder → hidden state → decoder → recreated image.

- ↳ force a "bottleneck" structure to compress initially (classic) → non-linear dimensionality reduction: Antiquated: Dencel.
- ↳ force a sparse, conditioned hidden state
- ↳ corrupt the AE input with noise

widely used: denoising autoencoder (Ch1 in the Stanfords)

- ↳ (sort of) analogous to BERT.

↳ simple to implement, but many ad-hoc choices. image → encoder → HS → decoder → image

Early days: autoencoders used for layerwise pretraining

- ↳ train deep networks w/o end-to-end mess
- ↳ advent of RELU, BN, etc. put a stop to this

Today: less widely used

- ↳ representation: VAEs, contrastive learning
- ↳ generations: GANs, VAEs, autoregressive models
- ↳ viable option for "quick and dirty" rep. learning
- ↳ Problem: sampling is hard.

Representing LVMs?

- ↳ easy: FCNN
- ↳ better: transpose convolutions (upscaling, unpooling)

Latent Variable Models: A cum of GMMs $\sum_z p(x|z)p(z)$

→ assume $p(x)$ is super complex

→ assume $p(z)$ is simple (Gaussian)

$$p(x|z) = \mathcal{N}(\mu_{nn}(z), \Sigma_{nn}(z))$$

$$p(x) = \int p(x|z)p(z) dz$$

simple simple

simplest

LVM in DL: basically, turn random numbers into valid samples

- ↳ VAEs, GANs, normalizing flows, etc.

$p(x)$ categorical variable.

$$p(y|x) = \sum_z p(y|x,z)p(z)$$

To train LVMs:

→ model $p_\theta(x)$

→ data $\mathcal{D} = \{x_1, \dots, x_n\}$

→ max likelihood fit:

$$\theta = \arg \max_{\theta} \frac{1}{N} \sum_i \log p_\theta(x_i)$$

Expected log-likelihood:

$$\theta = \arg \max_{\theta} \frac{1}{N} \sum_i \mathbb{E}_{z \sim p(z|x_i)} [\log p_\theta(x_i|z)]$$

calculate through probabilistic inference.

↳ "guess" most likely z given x_i

combine a bunch of trained AEs, reduce dim, speed up backprop

if x is discrete, discrete logistic is common

intractable... no closed-form solution