

ECE/CS 559 Lecture 19/20 Th 10/31 ; Th 11/7

Last time: Hebbian updates

(1) Generative Models

- All unsupervised learning is about understanding $P(x)$
- Generative models attempt to create one x :
 - They can do this by modeling $P(x)$ & sampling from it.
 - Or, they can directly model the sampling mechanism
e.g., sample z (standard distribution), then let $x = f(z) + \text{noise}$
 - e.g., $y = \text{caption}$ $x = \text{image}$ (diffusion models)
 - Conditional generative models attend to generation via prompts: $P(x|y)$
 - Examples: $y = \text{caption}$ $x = \text{image}$ (diffusion models)
 $y = \text{question}$ $x = \text{answer}$ (transformers)

(2) Density Estimation

Given Data x , assume i.i.d. from P .

Parametric methods: Let $P(x) = f(x; w)$

Maximize likelihood: $\underset{w}{\operatorname{argmax}} P(\text{Data}) = \prod_{x \in \text{Data}} P(x)$

$$\Leftrightarrow \underset{w}{\operatorname{argmin}} \sum_{x \in \text{Data}} -\log P(x)$$

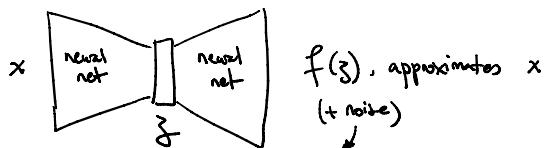
Non-parametric methods: Allow # of parameters w to grow with data size. Example, kernel density estimation



Works well when x discrete or small dimensional.
Otherwise, requires regularization, e.g. smooth densities.

(3) Autoencoders (Variational; VAEs)

- They try to capture $P(x)$ by reducing the "essence" of x to a simple representation z (encoding) then generating x from a sample of z (decoding)



Globally: $P_{\text{enc}}(z|x)$ $P_{\text{prior}}(z)$ $P_{\text{dec}}(x|z)$

- Why sampling $z \sim P_{\text{prior}}(z)$ then letting $x = f(z)$ work?
How do we choose P_{enc} , P_{dec} ?

(4) Derived distributions:

- Say $z, x \in \mathbb{R}^d$. Let $z \sim P(z)$ and let $x = f(z)$.
- What is $P(x)$? Let's assume $P(z)$ are densities, f monotonic
- $P(x) = P(X \leq x) = P(f(z) \leq x) = P(z \leq f^{-1}(x))$
- thus $P(x) = \frac{dP}{dx} = \frac{d}{dx} F_z(f^{-1}(x)) = \frac{d}{dz} F_z(z) \cdot \frac{dz}{dx} = P(f(z)) \cdot \frac{df^{-1}}{dx}$

$$\begin{aligned} \text{Example: } P_z \text{ uniform } [0, 1] \quad x = -\log z \quad f^{-1}(x) = e^{-x} \\ \Rightarrow P_x(x) = \begin{cases} 1 \cdot (-e^{-x}) & e^{-x} \in (0, 1) \\ 0 & \text{otherwise} \end{cases} = \begin{cases} e^{-x} & x < 0 \\ 0 & x \geq 0 \end{cases} \end{aligned}$$

- Can create any new random variable from standard random variables (e.g. uniform, Gaussian, etc.)

(b) From maximum likelihood to ELBO

How do we train an autoencoder?

Idea 1: Fix $P_{\text{prior}}(z)$ and only train the decoder $P_{\text{dec}}(x|z)$

$$\text{by minimizing } \sum_{x \in \text{Data}, z \sim P_{\text{prior}}} -\log P_{\text{dec}}(x|z)$$

Issue: Creates disconnect between x & z (close x 's \nrightarrow close z 's)

Idea 2: Fix x $P_{\text{enc}}(z|x)$ and train $P_{\text{prior}}(z) \circ P_{\text{dec}}(x|z)$

$$\text{by minimizing } \sum_{x \in \text{Data}, z \sim P_{\text{enc}}(1|x)} -\log P_{\text{prior}}(z) \cdot P_{\text{dec}}(x|z)$$

$$\text{approximates } \mathbb{E}_x \mathbb{E}_{z \sim P_{\text{enc}}} \underbrace{-\log P(z)}_{\text{implies } P(x) P(z|x)} P_{\text{dec}}(x|z) \Leftrightarrow$$

$$\mathbb{E}_x \left[\log \frac{1}{P_{\text{gen}}(x)} \right] + \mathbb{E}_x \left[\mathbb{E}_{z \sim P_{\text{enc}}} \left[\log \frac{1}{P_{\text{gen}}(z|x)} \right] \right] - \mathbb{E}_x \left[\mathbb{E}_z \left[\log \frac{1}{P_{\text{gen}}(z|x)} \right] \right]$$

min at $P_{\text{gen}}(z|x) = P_{\text{enc}}(z|x)$ max at $P_{\text{gen}}(z|x) = P_{\text{dec}}(x|z)$ - Entropy (P_{enc})

Issue: If the neural net of the decoder cannot easily "invert" the choice of P_{enc} , it won't work well.

Idea 3: Train P_{enc} and P_{dec} jointly.

Issue: P_{enc} in the loss is just a sampling distribution.

Hard to optimize: bad local minima / P_{enc} can degenerate by following w/ P_{dec}

Idea 4: Include P_{enc} in the loss: $\mathbb{E} \left[-\log \frac{P_{\text{prior}}(z) P_{\text{dec}}(x|z)}{P_{\text{enc}}(z|x)} \right]$

"max-entropy" effect

Evidence Lower Bound (ELBO) Loss $< \log P(x)$

③ Example: Gaussian Variational Autoencoder

$$P_{\text{prior}}(\mathbf{z}) \propto \exp\left(-\frac{\|\mathbf{z}\|^2}{2}\right) \quad \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$P_{\text{encoder}}(\mathbf{z}|\mathbf{x}) \propto \exp\left(-\frac{\|\mathbf{z} - g(\mathbf{x}; \mathbf{w})\|^2}{2}\right) \quad \mathcal{N}(g(\mathbf{x}; \mathbf{w}), \mathbf{I})$$

$$P_{\text{decoder}}(\mathbf{x}|\mathbf{z}) \propto \exp\left(-\frac{\|\mathbf{x} - f(\mathbf{z}; \mathbf{w})\|^2}{2}\right) \quad \mathcal{N}(f(\mathbf{z}; \mathbf{w}), \mathbf{I})$$

$$\text{ELBO has explicit form: } \mathbb{E}\left[\frac{\|\mathbf{z}\|^2}{2} - \frac{\|g(\mathbf{x}; \mathbf{w})\|^2}{2}\right] = \frac{1}{2}\mathbb{E}\left[\|g(\mathbf{x}; \mathbf{w})\|^2\right]$$

(minimize representation norm)

$$+ \frac{1}{2}\mathbb{E}\left[\|\mathbf{x} - \underbrace{f(g(\mathbf{x}; \mathbf{w}) + \mathbf{N}; \mathbf{w})}_{\text{encode}}\|_2^2\right] \quad \mathbf{N} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

(minimize reconstruction error)

This looks almost like k-means and it's not a coincidence! The similarity is closer if we quantize \mathbf{z} .

④ Hierarchical VAEs

- Instead of one representation stage, have multiple:

$$\mathbf{x} \xrightarrow{\text{encode}} \mathbf{z}_1 \xrightarrow{\text{encode}} \mathbf{z}_2 \xrightarrow{\text{encode}} \mathbf{z}_3 \xrightarrow{\text{decode}} \mathbf{z}_2 \xrightarrow{\text{decode}} \mathbf{z}_1 \xrightarrow{\text{decode}} \mathbf{x}$$

- Why? Allows for gradual tweaks, making it easier to train.

Example: Diffusion models use encoders that just add noise.

$$\text{ELBO: } \log \frac{P_{\text{prior}}(\mathbf{z}_1) P_{\text{dec}}(\mathbf{z}_2|\mathbf{z}_1) P_{\text{dec}}(\mathbf{z}_3|\mathbf{z}_2) P_{\text{dec}}(\mathbf{x}|\mathbf{z}_1)}{P_{\text{enc}}(\mathbf{z}_3|\mathbf{z}_2) P_{\text{enc}}(\mathbf{z}_2|\mathbf{z}_1) P_{\text{enc}}(\mathbf{z}_1|\mathbf{x})}$$