

Machine Learning for Applied Physics and High Energy Physics 1

Lecture 22 Generative Adversarial Neural Networks (GANs).

In this series of lectures we will discuss two generative model frameworks that have gained wide appeal in the last years: Generative Adversarial Neural Networks (GANs) and Variational Autoencoders (VAEs). These models are based on differential neural networks, therefore they will be trained by backpropagation as implemented, e.g., in high-level python packages such as Tensor.

22.1 Limitations of maximising the likelihood.

The Kullback-Leibler (KL) divergence plays a central role in many generative models. It measures the similarity between two probability distributions $p(\vec{x})$ and $q(\vec{x})$. Given two distributions, there are two distinct KL divergences

$$D_{KL}(p \parallel q) = \int d\vec{x} p(\vec{x}) \ln \frac{p(\vec{x})}{q(\vec{x})}$$

$$D_{KL}(q \parallel p) = \int d\vec{x} q(\vec{x}) \ln \frac{q(\vec{x})}{p(\vec{x})}$$

These can be combined to construct the symmetric square metric

$$D_{JS}(p, q) = \frac{1}{2} \left[D_{KL}\left(p \parallel \frac{p+q}{2}\right) + D_{KL}\left(q \parallel \frac{p+q}{2}\right) \right]$$

which is called Jensen-Shannon divergence. An important property of the KL divergence is its positivity

$$D_{KL}(p \parallel q) \geq 0 \quad (\text{with equality if and only if } p(\vec{x}) = q(\vec{x}))$$

In ML, the two distributions we are interested in are 2

- the model distribution $p_g(\vec{x})$
- the data distribution $p_{\text{data}}(\vec{x})$

We want $p_g(\vec{x})$ to be as similar as possible to $p_{\text{data}}(\vec{x})$.

Example: let us consider two Gaussian distributions (1D)

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma_p} e^{-(x-x_p)^2/2\sigma_p^2} \quad q(x) = \frac{1}{\sqrt{2\pi}\sigma_q} e^{-(x-x_q)^2/2\sigma_q^2}$$

$$D_{KL}(p(x) \| q(x)) = \ln \frac{\sigma_q}{\sigma_p} + \frac{1}{2\sigma_q^2} \left\{ (x_q - x_p)^2 + (\sigma_p^2 - \sigma_q^2) \right\}$$

The two distributions are alike if the difference between their central values and standard deviations is small.

Maximising the log-likelihood of the data under the model is the same as minimising the KL divergence between the data distribution and the model distribution $D_{KL}(p_{\text{data}} \| p_g)$.

Let us rewrite

- ENTROPY

$$\begin{aligned} D_{KL}(p_{\text{data}} \| p_g) &= \int d\vec{x} p_{\text{data}}(\vec{x}) \ln p_{\text{data}}(\vec{x}) - \int d\vec{x} p_{\text{data}}(\vec{x}) \ln p_g(\vec{x}) \\ &= -S[p_{\text{data}}] - \langle \ln p_g(\vec{x}) \rangle_{\text{data}} \end{aligned}$$

We have

$$\langle \ln p_g(\vec{x}) \rangle_{\text{data}} = -S[p_{\text{data}}] - D_{KL}(p_{\text{data}} \| p_g)$$

The equivalence follows from the positivity of the KL divergence and the fact that the entropy of the data distribution is constant. The original formulation of GANs (Goodfellow) minimises an upper bound on the Jensen-Shannon divergence between the model distribution $p_g(\vec{x})$ and the

data distribution $p_{\text{data}}(\bar{x})$.

Let us compare more closely the behaviour of the two KL divergences $D_{\text{KL}}(p_{\text{data}} \| p_{\theta})$ and $D_{\text{KL}}(p_{\theta} \| p_{\text{data}})$. Both these KL-divergences measure similarities between p_{data} and p_{θ} , yet they are sensitive to very different things.

- $D_{\text{KL}}(p_{\theta} \| p_{\text{data}})$ is insensitive to setting $p_{\theta} \approx 0$ even when $p_{\text{data}} \neq 0$ whereas $D_{\text{KL}}(p_{\text{data}} \| p_{\theta})$ punishes this harshly.
- $D_{\text{KL}}(p_{\text{data}} \| p_{\theta})$ is insensitive to placing weight in the model distribution in regions where $p_{\text{data}} \approx 0$ whereas $D_{\text{KL}}(p_{\theta} \| p_{\text{data}})$ punishes this harshly.

In other words, $D_{\text{KL}}(p_{\text{data}} \| p_{\theta})$ prefers model that have a high probability in regions with lots of training data points whereas $D_{\text{KL}}(p_{\theta} \| p_{\text{data}})$ punishes models for putting high probability where there is no data.

This suggests that the way likelihood-based methods can fail is by improperly "filling in" any low-probability density regions between peaks in the data distribution. The Jensen-Shannon divergence is sensitive to both placing weight where there is data and not placing weight where there is no data.

Now, $D_{\text{KL}}(p_{\theta} \| p_{\text{data}})$ is impossible to compute because we do not know $p_{\text{data}}(\bar{x})$. Conversely, $D_{\text{KL}}(p_{\text{data}} \| p_{\theta})$ can be computed easily from the data using sampling. The idea of adversarial learning is to circumvent the aforementioned difficulty by using an adversarial learning procedure. It consists in training a discriminator network to distinguish between real data points and samples generated from the model.

By punishing the model for generating points that can be easily ⁴ discriminated from the data, adversarial learning decreases the weight of regions in the model space that are far away from data points (regions that inevitably arise when maximising the likelihood).

22.2 Generative models and adversarial learning

The central idea of GANs is to construct two differentiable neural networks. The first neural network, usually a de-convolutional network approximates a generator function

$G(\bar{z}, \mathcal{D}_G)$ that takes as input a \bar{z} sampled from some prior in the latent space, and outputs a \bar{x} from the model. The second network approximates a discriminator function $D(\bar{x}, \mathcal{D}_D)$ that is designed to distinguish between \bar{x} from the data and samples generated by the model: $\bar{x} = G(\bar{z}, \mathcal{D}_G)$. The scalar $D(\bar{x})$ represents the probability that \bar{x} came from the data rather than the model $p_{\mathcal{D}_G}$. We train D to distinguish actual data points from synthetic examples and the generative network to fool the discriminative network.

To define the cost function for training, it is useful to define the functional

$$V(D, G) = \mathbb{E}_{x \sim p_{\text{data}}} (-\log D(\bar{x})) \\ + \mathbb{E}_{z \sim p_{\text{prior}}} (\log [1 - D(G(\bar{z}))])$$

We take the cost function for the discriminator and generator to be $C^{(G)} = -C^{(D)} = \frac{1}{2} V(D, G)$. This choice of cost functions corresponds to what is called a zero-sum game. Since the discriminator is maximised, we can write a cost function for the generator as

$$C(G) = \max_D V(G, D)$$

which can be rewritten as

$$C(G) = -\ln 4 + 2 D_{JS}(p_{data}, p_G)$$

