

Title: Reparameterization trick

URL: https://en.wikipedia.org/wiki/Reparameterization_trick

PageID: 77939986

Categories: Category:Machine learning, Category:Stochastic optimization

Source: Wikipedia (CC BY-SA 4.0).

The reparameterization trick (aka "reparameterization gradient estimator") is a technique used in statistical machine learning , particularly in variational inference , variational autoencoders , and stochastic optimization . It allows for the efficient computation of gradients through random variables, enabling the optimization of parametric probability models using stochastic gradient descent , and the variance reduction of estimators .

It was developed in the 1980s in operations research , under the name of "pathwise gradients", or "stochastic gradients". [1] [2] Its use in variational inference was proposed in 2013. [3]

Mathematics

Let z be a random variable with distribution $q_{\phi}(z)$, where ϕ is a vector containing the parameters of the distribution.

REINFORCE estimator

Consider an objective function of the form: $L(\phi) = \mathbb{E}_{z \sim q_{\phi}(z)}[f(z)]$. Without the reparameterization trick, estimating the gradient $\nabla_{\phi} L(\phi)$ can be challenging, because the parameter appears in the random variable itself. In more detail, we have to statistically estimate: $\nabla_{\phi} L(\phi) = \nabla_{\phi} \int dz q_{\phi}(z) f(z)$. The REINFORCE estimator, widely used in reinforcement learning and especially policy gradient , [4] uses the following equality: $\nabla_{\phi} L(\phi) = \int dz q_{\phi}(z) \nabla_{\phi} (\ln q_{\phi}(z)) f(z) = \mathbb{E}_{z \sim q_{\phi}(z)}[\nabla_{\phi} (\ln q_{\phi}(z)) f(z)]$. This allows the gradient to be estimated: $\nabla_{\phi} L(\phi) \approx \frac{1}{N} \sum_{i=1}^N \nabla_{\phi} (\ln q_{\phi}(z_i)) f(z_i)$. The REINFORCE estimator has high variance, and many methods were developed to reduce its variance . [5]

Reparameterization estimator

The reparameterization trick expresses z as: $z = g_{\phi}(\epsilon)$, $\epsilon \sim p(\epsilon)$. Here, g_{ϕ} is a deterministic function parameterized by ϕ , and ϵ is a noise variable drawn from a fixed distribution $p(\epsilon)$. This gives: $L(\phi) = \mathbb{E}_{\epsilon \sim p(\epsilon)}[f(g_{\phi}(\epsilon))]$. Now, the gradient can be estimated as: $\nabla_{\phi} L(\phi) = \mathbb{E}_{\epsilon \sim p(\epsilon)}[\nabla_{\phi} f(g_{\phi}(\epsilon))]$. This allows the gradient to be estimated: $\nabla_{\phi} L(\phi) \approx \frac{1}{N} \sum_{i=1}^N \nabla_{\phi} f(g_{\phi}(\epsilon_i))$.

Examples

For some common distributions, the reparameterization trick takes specific forms:

Normal distribution : For $z \sim \mathcal{N}(\mu, \sigma^2)$, we can use: $z = \mu + \sigma \epsilon$, $\epsilon \sim \mathcal{N}(0, 1)$.

Exponential distribution : For $z \sim \text{Exp}(\lambda)$, we can use: $z = -\frac{1}{\lambda} \log(\epsilon)$, $\epsilon \sim \text{Uniform}(0, 1)$.

$\epsilon \sim \text{Uniform}(0,1)$ Discrete distribution can be reparameterized by the Gumbel distribution (Gumbel-softmax trick or "concrete distribution"). [6]

In general, any distribution that is differentiable with respect to its parameters can be reparameterized by inverting the multivariable CDF function, then apply the implicit method. See [1] for an exposition and application to the Gamma Beta , Dirichlet , and von Mises distributions .

Applications

Variational autoencoder

In Variational Autoencoders (VAEs), the VAE objective function, known as the Evidence Lower Bound (ELBO), is given by:

$$\text{ELBO}(\phi, \theta) = \mathbb{E}_{z \sim q_{\phi}(z|x)} [\log p_{\theta}(x|z)] - D_{\text{KL}}(q_{\phi}(z|x) || p(z)) \quad \{\text{ELBO}\}(\phi, \theta) = \mathbb{E}_{z \sim q_{\phi}(z|x)} [\log p_{\theta}(x|z)] - D_{\text{KL}}(q_{\phi}(z|x) || p(z))$$

where $q_{\phi}(z|x)$ is the encoder (recognition model), $p_{\theta}(x|z)$ is the decoder (generative model), and $p(z)$ is the prior distribution over latent variables. The gradient of ELBO with respect to θ is simply $\mathbb{E}_{z \sim q_{\phi}(z|x)} [\nabla_{\theta} \log p_{\theta}(x|z)] \approx \frac{1}{L} \sum_{l=1}^L \nabla_{\theta} \log p_{\theta}(x|z_l)$ but the gradient with respect to ϕ requires the trick. Express the sampling operation $z \sim q_{\phi}(z|x)$ as: $z = \mu_{\phi}(x) + \sigma_{\phi}(x) \odot \epsilon$, $\epsilon \sim N(0, I)$ where $\mu_{\phi}(x)$ and $\sigma_{\phi}(x)$ are the outputs of the encoder network, and \odot denotes element-wise multiplication. Then we have $\nabla_{\phi} \text{ELBO}(\phi, \theta) = \mathbb{E}_{z \sim q_{\phi}(z|x)} [\nabla_{\phi} \log p_{\theta}(x|z) + \nabla_{\phi} \log q_{\phi}(z|x) - \nabla_{\phi} \log p(z)]$ where $z = \mu_{\phi}(x) + \sigma_{\phi}(x) \odot \epsilon$. This allows us to estimate the gradient using Monte Carlo sampling: $\nabla_{\phi} \text{ELBO}(\phi, \theta) \approx \frac{1}{L} \sum_{l=1}^L [\nabla_{\phi} \log p_{\theta}(x|z_l) + \nabla_{\phi} \log q_{\phi}(z_l|x) - \nabla_{\phi} \log p(z_l)]$ where $z_l = \mu_{\phi}(x) + \sigma_{\phi}(x) \odot \epsilon_l$ and $\epsilon_l \sim N(0, I)$ for $l = 1, \dots, L$.

This formulation enables backpropagation through the sampling process, allowing for end-to-end training of the VAE model using stochastic gradient descent or its variants.

Variational inference

More generally, the trick allows using stochastic gradient descent for variational inference. Let the variational objective (ELBO) be of the form: $\text{ELBO}(\phi) = \mathbb{E}_{z \sim q_{\phi}(z)} [\log p(x, z) - \log q_{\phi}(z)]$. Using the reparameterization trick, we can estimate the gradient of this objective with respect to ϕ : $\nabla_{\phi} \text{ELBO}(\phi) \approx \frac{1}{L} \sum_{l=1}^L \nabla_{\phi} [\log p(x, g_{\phi}(\epsilon_l)) - \log q_{\phi}(g_{\phi}(\epsilon_l))]$, $\epsilon_l \sim p(\epsilon)$ where $g_{\phi}(\epsilon) = \mu_{\phi}(x) + \sigma_{\phi}(x) \odot \epsilon$ and $\epsilon \sim N(0, I)$ for $l = 1, \dots, L$.

Dropout

The reparameterization trick has been applied to reduce the variance in dropout, a regularization technique in neural networks. The original dropout can be reparameterized with Bernoulli distributions: $y = (W \odot \epsilon) x$, $\epsilon_{ij} \sim \text{Bernoulli}(\alpha_{ij})$ where W is the weight matrix, x is the input, and α_{ij} are the (fixed) dropout rates.

More generally, other distributions can be used than the Bernoulli distribution, such as the gaussian noise: $y_i = \mu_i + \sigma_i \epsilon_i$, $\epsilon_i \sim \mathcal{N}(0, 1)$ where $\mu_i = \mathbf{m}_i^T \mathbf{x}$ and $\sigma_i^2 = \mathbf{v}_i^T \mathbf{x}^2$, with \mathbf{m}_i and \mathbf{v}_i being the mean and variance of the i -th output neuron. The reparameterization trick can be applied to all such cases, resulting in the variational dropout method. [7]

See also

Variational autoencoder

Stochastic gradient descent

Variational inference

References

Further reading

Ruiz, Francisco R.; AUEB, Titsias RC; Blei, David (2016). "The Generalized Reparameterization Gradient". *Advances in Neural Information Processing Systems*. 29. arXiv : 1610.02287 . Retrieved September 23, 2024 .

Zhang, Cheng; Butepage, Judith; Kjellstrom, Hedvig; Mandt, Stephan (2019-08-01). "Advances in Variational Inference". *IEEE Transactions on Pattern Analysis and Machine Intelligence*. 41 (8): 2008– 2026. arXiv : 1711.05597 . Bibcode : 2019ITPAM..41.2008Z . doi : 10.1109/TPAMI.2018.2889774 . ISSN 0162-8828 . PMID 30596568 .

Mohamed, Shakir (October 29, 2015). "Machine Learning Trick of the Day (4): Reparameterisation Tricks". *The Spectator* . Retrieved September 23, 2024 .