

# Noise Contrastive Estimation

Sargur N. Srihari  
srihari@cedar.buffalo.edu

# Topics in Partition Function

- Definition of Partition Function
  1. The log-likelihood gradient
  2. Stochastic maximum likelihood and contrastive divergence
  3. Pseudolikelihood
  4. Score matching and Ratio matching
  5. Denoising score matching
  6. Noise-contrastive estimation
  7. Estimating the partition function

# Intuition

- Noise contrastive estimation is based on the idea that a good generative model should be able to distinguish data from noise.
- A closely related idea is that a good generative model should be able to generate samples that no classifier can distinguish from data.
- This idea yields generative adversarial networks

# Noise Contrastive Estimation

- The probability distribution estimated by the model is represented explicitly as

$$\log p_{\text{model}}(\mathbf{x}) = \log \tilde{p}_{\text{model}}(\mathbf{x}; \boldsymbol{\theta}) + c,$$

- where  $c$  is explicitly introduced as an approximation of  $-\log Z(\boldsymbol{\theta})$
- Rather than estimating only  $\boldsymbol{\theta}$ , the noise contrastive estimation procedure treats  $c$  as just another parameter and estimates  $\boldsymbol{\theta}$  and  $c$  simultaneously, using the same algorithm for both

# Result of NCE

- The resulting  $\log p_{\text{model}}(\mathbf{x})$  thus may not correspond exactly to a valid probability distribution, but will become closer and closer to being valid as the estimate of  $c$  improves
- Such an approach would not be possible using maximum likelihood as the criterion for the estimator
  - The maximum likelihood criterion would choose to set  $c$  arbitrarily high, rather than setting  $c$  to create a valid probability distribution.

# NCE converts unsupervised to supervised learning

- NCE works by reducing the unsupervised learning problem of estimating  $p(\mathbf{x})$  to that of learning a probabilistic binary classifier in which one of the categories corresponds to the data generated by the model
- This supervised learning problem is constructed in such a way that maximum likelihood estimation in this supervised learning problem defines an asymptotically consistent estimator of the original problem.

# Noise Distribution

- Introduce a second distribution, the noise distribution  $p_{\text{noise}}(\mathbf{x})$ 
  - The noise distribution should be tractable to evaluate and to sample from
- We can now construct a model over both  $\mathbf{x}$  and a new, binary class variable  $y$ 
  - In the new joint model, we specify that

$$\begin{aligned} p_{\text{joint}}(y = 1) &= \frac{1}{2}, \\ p_{\text{joint}}(\mathbf{x} \mid y = 1) &= p_{\text{model}}(\mathbf{x}), \\ p_{\text{joint}}(\mathbf{x} \mid y = 0) &= p_{\text{noise}}(\mathbf{x}). \end{aligned}$$

- In other words,  $y$  is a switch variable that determines whether we will generate  $\mathbf{x}$  from the model or from the noise distribution

# Joint model for training data

- We can construct a similar joint model of training data
- In this case, the switch variable determines whether we draw  $\mathbf{x}$  from the data or from the noise distribution

$$p_{\text{train}}(y = 1) = \frac{1}{2}, p_{\text{train}}(\mathbf{x} \mid y = 1) = p_{\text{data}}(\mathbf{x}),$$

and

$$p_{\text{train}}(\mathbf{x} \mid y = 0) = p_{\text{noise}}(\mathbf{x}).$$

# Maximum Likelihood Training

- We can now just use standard maximum likelihood learning on the supervised learning problem of fitting  $p_{\text{joint}}$  to  $p_{\text{train}}$

$$\theta, c = \arg \max_{\theta, c} \mathbb{E}_{\mathbf{x}, y \sim p_{\text{train}}} \log p_{\text{joint}}(y | \mathbf{x})$$

# Role of Logistic Regression

- The distribution  $p_{\text{joint}}$  is essentially a logistic regression model applied to the difference in log probabilities of the model and the noise distribution:

$$\begin{aligned} p_{\text{joint}}(y = 1 \mid \mathbf{x}) &= \frac{p_{\text{model}}(\mathbf{x})}{p_{\text{model}}(\mathbf{x}) + p_{\text{noise}}(\mathbf{x})} \\ &= \frac{1}{1 + \frac{p_{\text{noise}}(\mathbf{x})}{p_{\text{model}}(\mathbf{x})}} \\ &= \frac{1}{1 + \exp\left(\log \frac{p_{\text{noise}}(\mathbf{x})}{p_{\text{model}}(\mathbf{x})}\right)} \\ &= \sigma\left(-\log \frac{p_{\text{noise}}(\mathbf{x})}{p_{\text{model}}(\mathbf{x})}\right) \\ &= \sigma(\log p_{\text{model}}(\mathbf{x}) - \log p_{\text{noise}}(\mathbf{x})). \end{aligned}$$

# Ease of using NCE

- NCE is thus simple to apply so long as
  - $-\log \tilde{p}_{\text{model}}$  is easy to back-propagate through, and,
  - as specified above,  $p_{\text{noise}}$  is easy to evaluate (in order to evaluate  $p_{\text{joint}}$ ) and sample from (in order to generate the training data).

# Application in Neural Language Models

- Applied to modeling the conditional distribution over a word given the context of the word
- Though the word may be drawn from a large vocabulary, there is only one word