

Improved Autoregressive Modeling with Distribution Smoothing

Mikhail Kuznetsov

CMC MSU

February 8, 2021

Let $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ — D -dimensional i.i.d samples from a continuous data distribution $p_{\text{data}}(\mathbf{x})$

Property

An autoregressive model decomposes a joint distribution into univariate conditionals [2]:

$$p(\mathbf{x}) = \prod_{i=1}^D p(x_i \mid \mathbf{x}_{<i})$$

Goal

Find θ — parameters of the model such that $p_{\theta}(\mathbf{x}) \approx p_{\text{data}}(\mathbf{x})$

A commonly used approach for density estimation is maximum likelihood estimation (MLE), i.e., by maximizing

$$L(\theta) \triangleq \frac{1}{N} \sum_{i=1}^N \log p_{\theta}(\mathbf{x}_i)$$

General idea

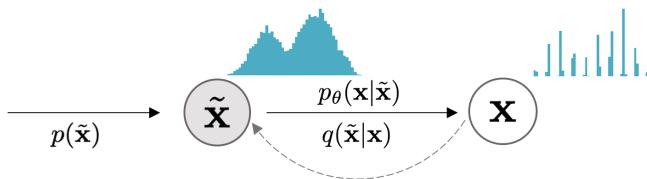


Figure: Overview of the method
($\tilde{\mathbf{x}}$ — the smoothed data, $q(\tilde{\mathbf{x}} | \mathbf{x})$ — the smoothing distribution) [1]

Problem 1: Manifold hypothesis

Many real world data distributions (e.g. natural images) may lie in the vicinity of a low-dimensional manifold and can often have complicated densities with sharp transitions (i.e. high Lipschitz constants) that are difficult to model.

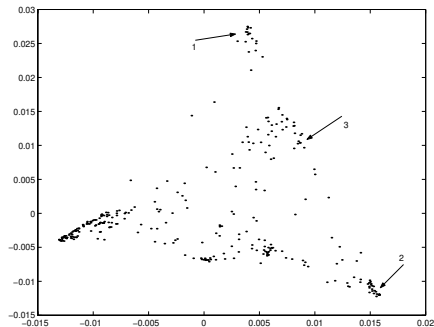


Figure (a): The 300 most frequent words of the Brown corpus represented in the spectral domain.

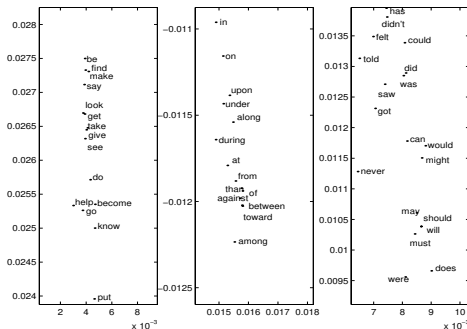
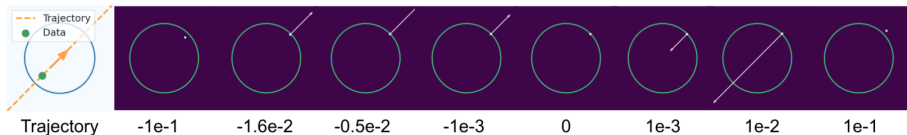


Figure (b): Fragments labeled by arrows: (left) infinitives of verbs, (middle) prepositions, and (right) mostly modal and auxiliary verbs.

Source: (Belkin M., Niyogi P., 2003) [3]

Problem 1: Manifold hypothesis



Problem 2: Compounding errors

$$p(\mathbf{x}) = \prod_{i=1}^D p(x_i \mid \mathbf{x}_{<i})$$

Compounding errors comes from the inaccurate approximation of the conditional distributions.

The reasons for this may be:

- Curse of dimensionality + very limited amount of training data (in some cases)
- The current state is based on the values of the previous states
- Adversarial attacks

Theorem 1 Let:

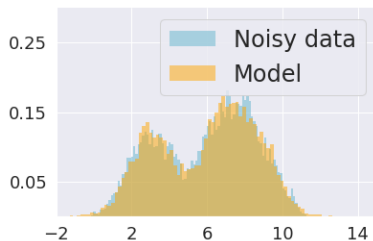
- 1 $p(x)$ — a continuous 1-d distribution that is supported on R
- 2 $q(\tilde{x} | x)$ — 1-d distribution that is:
 - symmetric (i.e. $q(\tilde{x} | x) = q(x | \tilde{x})$)
 - stationary (i.e. translation invariant)
 - $\lim_{x \rightarrow \infty} p(x)q(x | \tilde{x}) = 0$ for any given \tilde{x}

Then: $\text{Lip}(q(\tilde{x})) \leq \text{Lip}(p_{\text{data}}(x))$,

where $q(\tilde{x}) \triangleq \int q(\tilde{x} | x)p(x)dx$ and $\text{Lip}(\cdot)$ denotes the Lipschitz constant of the given 1 - d function.



(a) Data



(b) Smoothed data

Figure: Illustration of the Theorem 1 [1]

Proposition 1 (Informal) Let:

1 $q(\tilde{\mathbf{x}} \mid \mathbf{x})$ is such that:

- symmetric
- stationary
- has small variance
- has negligible higher order moments (i.e. very small)

Then:

$$\mathbf{E}_{p_{\text{data}}(\mathbf{x})} \mathbf{E}_{q(\tilde{\mathbf{x}}|\mathbf{x})} [\log p_{\theta}(\tilde{\mathbf{x}})] \approx \mathbf{E}_{p_{\text{data}}(\mathbf{x})} \left[\log p_{\theta}(\mathbf{x}) + \frac{\eta}{2} \sum_i \frac{\partial^2 \log p_{\theta}}{\partial x_i^2} \right]$$

for some constant η .

Since the samples from p_{data} should be close to a local maximum of the model, this encourages the second order gradients computed at a data point \mathbf{x} to become closer to zero (if it were positive then \mathbf{x} will not be a local maximum), creating a smoothing effect.

- 1 Sing-step recovering
- 2 Obtaining ELBO

- [1] Improved autoregressive modeling with distribution smoothing.
<https://openreview.net/pdf?id=rJA5Pz7lHKb>, 2021.

- [2] Tim Salimans, Andrej Karpathy, Xi Chen, and Diederik P Kingma.
Pixelcnn++: Improving the pixelcnn with discretized logistic mixture likelihood and other modifications.
[arXiv:1701.05517](https://arxiv.org/abs/1701.05517), 2017.

- [3] Mikhail Belkin, Partha Niyogi.
Laplacian eigenmaps for dimensionality reduction and data representation.
<https://www2.imm.dtu.dk/projects/manifold/Papers/Laplacian.pdf>, 2003.