

# On Fenchel Mini-Max Learning

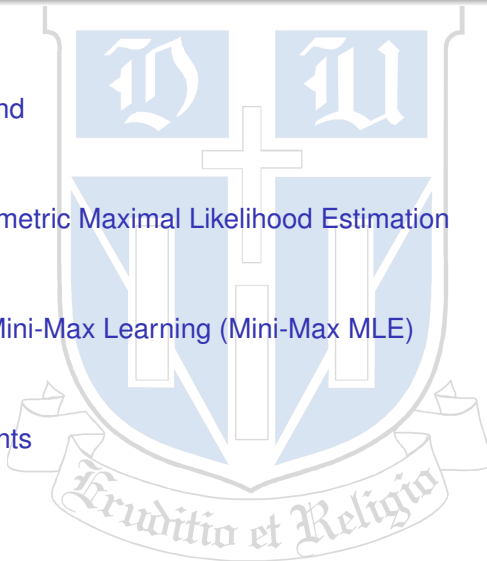
Chenyang Tao<sup>1</sup>, Liqun Chen<sup>1</sup>, Shuyang Dai<sup>1</sup>, Junya Chen<sup>1,2</sup>, Ke Bai<sup>1</sup>,  
Dong Wang<sup>1</sup>, Jianfeng Feng<sup>3</sup>, Wenlian Lu<sup>1</sup>, Georgiy Bobashev<sup>4</sup>,  
Lawrence Carin<sup>1</sup>

<sup>1</sup>Electrical & Computer Engineering, Duke University <sup>2</sup>School of Mathematical  
Sciences, Fudan University <sup>3</sup>ISTBI, Fudan University <sup>4</sup>RTI International

Dec 10, 2019 @ NeurIPS, Vancouver, BC, Canada

# Outline

1. Background
2. Non-Parametric Maximal Likelihood Estimation
3. Fenchel Mini-Max Learning (Mini-Max MLE)
4. Experiments



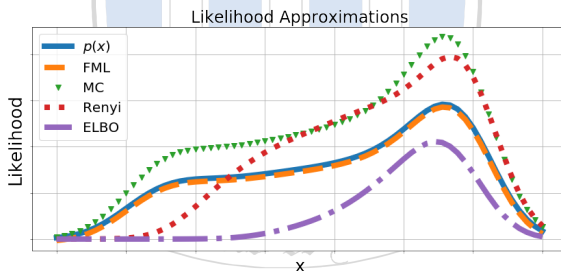
# Our Contributions

We present Fenchel Mini-Max Learning (FML) which highlights

- 1 A **Mini-Max formulation** of *Maximal Likelihood Estimation*
  - 2 **Unbiased** likelihood estimator directly amendable to SGD
  - 3 Amortized estimation with deep neural networks
  - 4 A latent-variable model competitive to variational inference
- We show FML compare favorably to existing alternatives in likelihood-based distribution learning across a wide range of applications, such as
    - density estimation, generative modeling, natural language processing, reinforcement learning, etc.

# A Comparison to Competing Likelihood Estimators

- Naive MC over-estimates:  $\log \hat{p}_{\theta, \text{MC}}(x) \geq \log p_{\theta}(x)$
- Evidence lower bound under-estimates:  $\text{ELBO}_{\theta}(x) \leq \log p_{\theta}(x)$
- Rényi bound approximates:  $\text{RVI}_{\theta}(x) \approx \log p_{\theta}(x)$
- FML estimate is tight:  $\text{FML}_{\theta}(x) = \log p_{\theta}(x)$



# Goals of Probabilistic Modeling

## Four desiderata of probabilistic modeling

- *Estimation* ✓ ⇐ Traditional focuses
    - Find a parameterized model  $\mathcal{M}_\theta$  that best describes the data
  - *Inference* ✓
    - Identify the latent properties (e.g. group identity) of data
  - *Sampling* ✓ ⇐ **Recent interests**
    - Emulating a stochastic rule that generates data
  - *Likelihood evaluation (evidence)* ✗
    - Evaluate the likelihood for a given sample (e.g., outlier)
  - *Other concerns*
    - model scalability, training stability, expressiveness, etc.
- ✓ : frequently used      ✗ : less frequently used

# Goals of Probabilistic Modeling

## Four desiderata of probabilistic modeling

- *Estimation* ✓
  - *Inference* ✓
  - *Sampling* ✓
  - *Likelihood evaluation (evidence)* ✓✗
  - *Other concerns*
    - model scalability, training stability, expressiveness, etc.
- ✓ : frequently used      ✓✗ : less frequently used
- 
- **We want to present a scheme without any major compromise on the above points.**

# Maximal Likelihood Criteria

## Maximum likelihood estimation (MLE)

- $\{x_i\}_{i=1}^n$  are independent observations of true distribution  $p_d(x)$
- $\{p_\theta(x)\}_{\theta \in \Theta}$  is a family of distributions parameterized by  $\alpha$
- log-likelihood loss:  $\hat{\ell}(\theta) = \sum_i \log p_\theta(x_i)$
- Maximizing expected likelihood:

$$\hat{\theta}_{\text{MLE}} = \arg \max_{\theta \in \mathcal{A}} \{\ell(\theta)\} \quad (1)$$

- This is equivalent to minimizing  $\text{KL}(p_d \parallel p_\theta)$



# Likelihood-based Probabilistic Modeling

## One-step approach

- Direct construction of a parameterized stochastic procedure with tractable likelihoods. (e.g., generative flows)

## Two-step approach

- First estimate an (unnormalized) density with empirical samples (e.g., non-parametric density estimation)
- To draw new samples, use established sampling procedure (e.g., MCMC) to sample from the density estimate
- **This study addresses the challenge of estimation an unnormalized density**



# Background on Estimating unnormalized Statistical Models

## Challenges of MLE with unnormalized statistical models

- In many cases, statistical models are given in the form of unnormalized exponential family  $\tilde{p}(x; \theta) = \exp(-\psi(x; \theta))$ 
  - $\psi(x; \theta)$  is known as the *potential function*.
  - PDF is known up to a multiplicative constant  $Z(\theta)$ ,

$$p(x; \theta) = \frac{1}{Z(\theta)} \tilde{p}(x; \theta), \quad (2)$$

- $Z(\theta) = \int \tilde{p}(x'; \theta) dx'$  is called the *partition function*
- In general,  $Z(\theta)$  is **analytically intractable**

# Density Estimation of Unnormalized Statistical Models

## The challenge of density estimation for unnormalized models

- In MLE training, one optimizes

$$\log p_{\theta}(x) = -\psi_{\theta}(x) - \log Z_{\theta} \quad (3)$$

- Finite sample MC estimate bias the likelihood estimate  $\log \hat{p}_{\psi}(x)$  because the integral is within the log

$$\log \hat{p}_{\theta}(x) = -\psi_{\theta}(x) - \log \hat{Z}_{\theta}, \quad (4)$$

$$\log \hat{Z}_{\theta}^{\text{MC}} = \log \left( \frac{1}{m} \sum_{j=1}^m \exp(-\psi_{\theta}(X'_j)) \right), \quad (5)$$

$$\mathbb{E}_{X_j}[\log \hat{Z}_{\theta}] \leq \log(\mathbb{E}_{X_j}[\hat{Z}_{\theta}]) = \log Z_{\theta}, \text{ (Jensen's ineq)} \quad (6)$$

$$\mathbb{E}_{X_j}[\log \hat{p}_{\theta}(x)] \geq \log p_{\theta}(x). \quad (7)$$

# Density Estimation of Unnormalized Statistical Models

## Existing solutions

### ■ Potential-based

- MCMC-MLE [Geyer, 1991]
- Contrastive divergence (CD), [Hinton, 2002]
- Noise Contrastive Estimation (NCE), [Gutmann, 2010]
- Dynamic dual embedding (DDE), [Dai, 2018]

### ■ Score-based

- Score matching (SM), [Hyvarinen, 2005]
- De-noising auto-encoder (DAE), [Alain, 2014]
- Stein implicit learning (SIL), [Li, 2017]

### ■ Practitioners make trade-offs in terms of the goals they want to achieve, because ...

Spoilers: NO single method hits all bullets : - (

**Table:** Comparison of popular probabilistic modeling procedures.

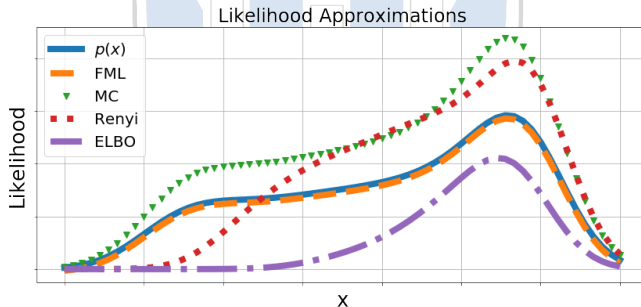
Model	Inference	Sampling	Likelihood	Scalability
<i>MCMC</i>	No	Yes	No	Poor
RBM	Yes	Yes	No	Good
SM	No	No	No	Poor
NCE	No	No	Estimate	Depends
DDE	No	Yes	Exact	Medium
VI	Yes	Yes	Bound	Good
FLOW	No	Yes	Exact	Tricky
SVGD	No	Yes	No	Medium
GAN	No	Yes	No	Good
OT	No	Yes	No	Depends
<b>FML</b>	<b>Yes</b>	<b>Yes</b>	<b>Estimate</b>	<b>Good</b>

# Mini-Max Likelihood Estimation (Ours)

## General idea

- Inspired by adversarial training, we exploited the Fenchel conjugacy and reformulate MLE into a mini-max game
  - *Min-game* recovers the model likelihood
  - *Max-game* matches model to the data distribution

(See our paper for details on how all four goals are accommodated)



# Fenchel Mini-Max Learning

## Fenchel Conjugacy

- Let  $f(t)$  be a proper convex, lower-semicontinuous function
- The convex conjugate function  $f^*(v)$  is defined as

$$f^*(v) = \sup_{t \in \mathcal{D}(f)} \{tv - f(t)\} \quad (8)$$

- $\mathcal{D}(f)$  denotes the domain of function  $f$
  - $f^*$  is again convex and lower-semicontinuous.
  - Fenchel conj. pair  $(f, f^*)$  are dual to each other  $((f^*)^* = f)$
- Example:  $(-\log(t), -1 - \log(-v))$  is a Fenchel pair
- Change of variable  $v = \exp(-u)$  gives

$$-\log(t) = \max_u \{-u - \exp(-u)t + 1\} \quad (9)$$

# Mini-Max Likelihood Estimation

## Fenchel Mini-Max Formulation

- For unnormalized model distribution  $\tilde{p}_\psi(x) = \exp(-\psi(x))$

$$\hat{\psi}_{\text{MLE}} = \arg \max_{\psi} \left\{ - \min_u \left\{ \sum_{i=1}^n (u_i + e^{-u_i} I(x_i; \psi)) \right\} \right\} \quad (10)$$

- Importance-weighted estimator  $I(x; \psi)$

$$I(x; \psi_\theta) = \int \left( \frac{1}{q(x')} e^{\psi(x) - \psi(x')} \right) q(x') \mathrm{d}x' \quad (11)$$

- $q(x)$  is the proposal distribution
- Auxiliary variable  $u_i$  for each data point  $x_i$  and Min-game returns the normalized log-likelihood, e.g.  $u^* = \log p_\psi(x)$

# Understanding Fenchel Mini-Max Likelihood Estimation

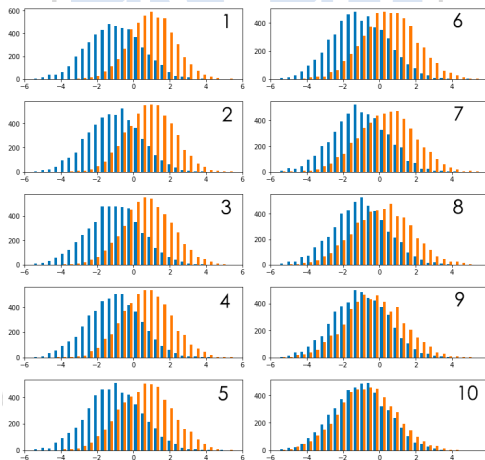
## A gradient analysis

- In standard MLE learning, we have  $\nabla \log p_\theta(x) = \frac{\nabla p_\theta(x)}{p_\theta(x)}$ 
  - The gradient of likelihood is normalized by model evidence
- Insight: while  $\nabla p_\theta(x)$  is difficult to compute (because of  $Z_\theta$ ), unbiased gradient estimate of the inverse likelihood  $\frac{1}{p_\theta(x)}$  is easy to get  $\nabla \left\{ \frac{1}{p_\theta(x)} \right\} = \int \nabla \{ \exp(\psi_\theta(x) - \psi_\theta(x')) \} dx'$
- This connects FML gradient via
 
$$\begin{aligned} \nabla J_\theta(x; \hat{u}_x, \psi) &= -\nabla \left\{ \exp(-\hat{u}_x) \int e^{\psi_\theta(x) - \psi_\theta(x')} dx' \right\} \\ &= -\hat{p}_\theta(x) \nabla \left\{ \frac{1}{p_\theta(x)} \right\} = \frac{\hat{p}_\theta(x)}{p_\theta(x)} \nabla \log p_\theta(x) \approx \nabla \log p_\theta(x), \end{aligned} \quad (12)$$
- It's easy to show when the likelihood ratio is bounded, we can recover the same exact MLE solution



# Learning with FML: Simple Gaussian

(blue: data, orange: model)



Question: Does FML converge to ground-truth?

# Convergence Guarantees

## Proposition (Convergence of FML with stochastic gradient descent)

- Let  $\{\eta_t\}$  be the scheduled learning rate and  $\{\xi_t\}$  be the approximation error. If the generalized learning rate  $\tilde{\eta}_t = \eta_t \xi_t$  satisfies  $\sum_t \mathbb{E}[\tilde{\eta}_t] = \infty$  and  $\sum_t \mathbb{E}[\tilde{\eta}_t^2] < \infty$ , we have
  - (*Convex setting*) if the likelihood loss has a unique equilibrium point and is asymptotically stable (e.g. strictly convex), then under standard Robbins-Monro regularity conditions FML-SGD converges to MLE solution  $\hat{\theta}_{\text{MLE}}$  with probability 1 from any  $\theta_0$ .
  - (*Non-convex setting*) if the likelihood loss has Lipschitz-continuous, then FML-SGD will converge to a stationary point of  $f(\theta)$  with probability 1, i.e.,  $\|\nabla_{\theta} f(\theta)\| \rightarrow 0$  as  $t \rightarrow \infty$ .



# Fenchel Mini-Max Learning for Latent Variable Models

## Applying the same trick to latent variable model

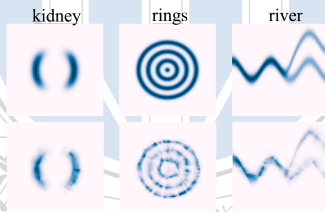
$$\arg \max_{\alpha, \beta} \left\{ \min_u \left\{ \sum_{i=1}^n \left( u_i + e^{-u_i} \int q_{\beta}(z|x) \left\{ \frac{p_{\alpha}^{\tau_t}(x, z)}{q_{\beta}(z|x)} \right\} dz \right) \right\} \right\} \quad (13)$$

- $q_{\beta}(z|x)$  resembles the approximate posterior in VI
  - In IW-VAE it also functions as proposal distribution
- $\tau_t$  is the annealing factor
- *In contrast to VI, we optimize an estimate of the log-likelihood rather than a lower bound*



**Table:** Quantitative evaluation on toy models.

Model	Parameter estimation error $\uparrow \downarrow$					Likelihood consistency score $\uparrow$				
	banana	kidney	rings	river	wave	banana	kidney	rings	river	wave
MC	3.46	3.9	4.71	1.71	1.78	0.961	0.881	0.508	0.702	0.619
SM	7.79	2.75	3.62	1.64	2.61	×	×	×	×	×
NCE	3.88	2.5	4.81	2.85	<b>1.20</b>	0.968	0.882	0.557	0.721	0.759
KEF	×	×	×	×	×	0.973	0.755	0.183	0.436	0.265
DDE	6.59	7.31	24.9	29.1	25.7	0.944	0.830	0.426	0.520	0.186
FML (ours)	<b>3.05</b>	<b>1.9</b>	<b>2.59</b>	<b>1.13</b>	1.27	<b>0.974</b>	<b>0.901</b>	<b>0.562</b>	<b>0.731</b>	<b>0.782</b>



**Figure:** FML predicted likelihood using nonparametric potentials.

**Table:** Results on language models, with the example synthesized text representative of typical results.

	PPL ↓	BLEU-2 ↑	BLEU-3 ↑	BLEU-4 ↑	BLEU-5 ↑
<i>EMNLP WMT news</i>					
VAE	12.5	76.1	46.8	23.1	11.6
FML	<b>11.6</b>	<b>77.2</b>	<b>47.4</b>	<b>24.3</b>	<b>12.2</b>
<i>MS COCO</i>					
VAE	9.5	82.1	60.7	38.9	24.8
FML	<b>8.6</b>	<b>84.2</b>	<b>64.4</b>	<b>40.3</b>	<b>25.2</b>
Sampled sentences from respective models on <i>WMT news</i>					
VAE	<i>"China's economic crisis, the number of US exports, which is still in recent years of the UK's population."</i>				
FML	<i>"In addition, police officials have also found a new investigation into the area where they could take a further notice of a similar investigation into."</i>				

# The End

Thank you.

