
Estimating Densities with Non-Parametric Exponential Families

Lin Yuan
yuanl@purdue.edu
School of ECE
Purdue University

Sergey Kirshner
skirshne@purdue.edu
Dept. of Statistics
Purdue University

Robert Givan
givan@purdue.edu
School of ECE
Purdue University

Abstract

We propose a novel approach for density estimation with exponential families for the case when the true density may not fall within the chosen family. Our approach augments the sufficient statistics with features designed to accumulate probability mass in the neighborhood of the observed points, resulting in a non-parametric model similar to kernel density estimators. We show that under mild conditions, the resulting model uses only the sufficient statistics if the density is within the chosen exponential family, and asymptotically, it approximates densities outside of the chosen exponential family.

1 Non-parametric Exponential Family

Suppose \mathbf{X} is a vector of random variables with support $\mathcal{X} \subseteq \mathbb{R}^m$. A distribution for \mathbf{X} belongs to the exponential family of distributions with sufficient statistics $\mathbf{t} : \mathcal{X} \rightarrow \mathcal{H} \subseteq \mathbb{R}^d$, if its probability density has a functional form:¹ $f^E(\mathbf{x}|\boldsymbol{\lambda}) = \frac{1}{Z(\boldsymbol{\lambda})} q(\mathbf{x}) \exp\langle \boldsymbol{\lambda}, \mathbf{t}(\mathbf{x}) \rangle$, where $Z(\boldsymbol{\lambda}) = \int_{\mathcal{X}} q(\mathbf{x}) \exp\langle \boldsymbol{\lambda}, \mathbf{t}(\mathbf{x}) \rangle d\mathbf{x} < \infty$ is a partition function, $\boldsymbol{\lambda}$ is a vector of canonical parameters, $q : \mathcal{X} \rightarrow \mathbb{R}$ is a base measure, and $\langle \cdot, \cdot \rangle$ denotes the Euclidean inner product. Assuming q is fixed, let $\mathcal{EF}_{\mathbf{t}}$ denote the set of all possible distributions of the form f^E with the set of sufficient statistics \mathbf{t} .

Given samples $\mathbf{x}^{1:n} \triangleq (\mathbf{x}^1, \dots, \mathbf{x}^n) \stackrel{i.i.d.}{\sim} f$ where $f : \mathcal{X} \rightarrow \mathbb{R}$ is an unknown density with the same support as q . Let $\hat{f}_n : \mathcal{X} \rightarrow \mathbb{R}$ be the empirical distribution for $\mathbf{x}^{1:n}$, $\hat{f}_n(\mathbf{x}|\mathbf{x}^{1:n}) = \frac{1}{n} \sum_{i=1}^n \delta(\mathbf{x}^i)$ where $\delta(\mathbf{x})$ is a Dirac delta function. Exponential families can be obtained as a solution to the optimization problem of minimizing the relative entropy subject to matching the moment constraints of the empirical and the estimated distributions.

However, if the true distribution does not fall within the chosen exponential family, $f \notin \mathcal{EF}_{\mathbf{t}}$, the estimated model may provide a poor approximation to the true density (e.g. putting most probability mass on the mean when f is multi-modal). Kernel density estimators (KDEs), on the other hand, ensure a small portion of probability mass around the observed data points \mathbf{x}^i by a weighted combination of kernel functions: $f_n^{\text{KDE}}(\mathbf{x}|\mathbf{x}^{1:n}) = \frac{1}{n} \sum_{i=1}^n K_{\mathbf{H}}(\mathbf{x}; \mathbf{x}^i)$, where $K_{\mathbf{H}}(\mathbf{x}; \mathbf{x}^i) = |\mathbf{H}|^{-\frac{1}{2}} K\left(\mathbf{H}^{-\frac{1}{2}}(\mathbf{x} - \mathbf{x}^i)\right)$. K is a multivariate kernel function, a bounded probability density function on \mathbb{R}^m . $K_{\mathbf{H}}$ is a multivariate kernel function with a symmetric positive definite bandwidth matrix \mathbf{H} ; in this paper, we assume $\mathbf{H} = h^2 \mathbf{I}_d$ (assuming $\mathbf{x} \in \mathbb{R}^d$). Typically, kernel functions are probability density functions peaked around each data point \mathbf{x}^i preserving certain probability mass. For example, the uniform kernel is an indicator function $\mathcal{I}_{\mathbf{x}^i}(\mathbf{x}) = 1$ if and only if $\|\mathbf{x} - \mathbf{x}^i\|_2 < \frac{h}{2}$, and $E_f[\mathcal{I}_{\mathbf{x}^i}]$ is the probability mass inside the $\frac{h}{2}$ -ball around \mathbf{x}^i . Thus,

¹For notational convenience, we denote $\mathbf{X} = \mathbf{x}$ by \mathbf{x} .

KDEs match the mass around each \mathbf{x}^i to that of \hat{f}_n , therefore are able to approximate any f as $n \rightarrow \infty$. Inspired by KDEs, our proposed non-parametric exponential family adds the kernel functions $K_{\mathbf{H}}(\mathbf{x}; \mathbf{x}^i) \triangleq t_a^i(\mathbf{x})$ ² to the set of sufficient statistics $\mathbf{t}(\mathbf{x})$. We then approximately match $E_f[t_a^i]$ with $E_{\hat{f}_n}[t_a^i]$: $|E_{\hat{f}_n}[t_a^i] - E_f[t_a^i]| < \beta_i$. In addition to the canonical parameters $\boldsymbol{\lambda}$ for sufficient statistics, non-parametric exponential families have *augmented* parameters $\boldsymbol{\lambda}_a$ for the augmented statistics $\mathbf{t}_a(\mathbf{x}) \triangleq [t_a^1 \dots t_a^n]$.

The proposed density approximation ($f_n^{NE}(\mathbf{x})$) is a solution to

$$\begin{aligned} f_n^{NE}(\mathbf{x}|\mathbf{x}^{1:n}) &= \arg \min_{f^{NE} \in \mathcal{F}} KL(f^{NE} \parallel q) \text{ subj to} \\ E_{f_n^{NE}(\mathbf{x}|\mathbf{x}^{1:n})}[\mathbf{t}(\mathbf{x})] &= E_{\hat{f}_n(\mathbf{x}|\mathbf{x}^{1:n})}[\mathbf{t}(\mathbf{x})], \\ |E_{f_n^{NE}}[t_a^i(\mathbf{x})] - E_{\hat{f}_n}[t_a^i(\mathbf{x})]| &\leq \beta_i, i = 1, \dots, n. \end{aligned} \quad (1)$$

f_n^{NE} falls within the generalized MaxEnt framework [1]:

$$\begin{aligned} f(\mathbf{x}) &= \frac{1}{Z(\boldsymbol{\lambda}, \boldsymbol{\lambda}_a)} q(\mathbf{x}) \exp[\langle \boldsymbol{\lambda}, \mathbf{t}(\mathbf{x}) \rangle + \langle \boldsymbol{\lambda}_a, \mathbf{t}_a(\mathbf{x}) \rangle] \\ Z(\boldsymbol{\lambda}, \boldsymbol{\lambda}_a) &= \int_{\mathcal{X}} q(\mathbf{x}) \exp[\langle \boldsymbol{\lambda}, \mathbf{t}(\mathbf{x}) \rangle + \langle \boldsymbol{\lambda}_a, \mathbf{t}_a(\mathbf{x}) \rangle] d\mathbf{x}. \end{aligned} \quad (2)$$

Let $\mathbf{s}(\mathbf{x}) \triangleq (\mathbf{t}(\mathbf{x}), \mathbf{t}_a(\mathbf{x}))$ and $\boldsymbol{\theta} \triangleq (\boldsymbol{\lambda}, \boldsymbol{\lambda}_a)$ be a combined set of statistics and parameters, respectively, for the augmented model. A specific set of parameter values for the distribution in (2) satisfying the constraints in (1) can be found by maximizing the penalized log-likelihood

$$l(\boldsymbol{\theta}) = \left\langle \boldsymbol{\theta}, \frac{1}{n} \sum_{i=1}^n \mathbf{s}(\mathbf{x}^i) \right\rangle - \ln Z(\boldsymbol{\theta}) - \sum_{i=1}^n \beta_i |\lambda_a^i|. \quad (3)$$

1.1 Theoretical Properties

The proofs appear in [2].

Theorem 1.1. *Suppose a vector of random variables \mathbf{X} with support on \mathcal{X} has a density $f \in \mathcal{EF}_{\mathbf{t}}$ with features $\mathbf{t} : \mathcal{X} \rightarrow \mathcal{H} \subseteq \mathbb{R}^d$ and a vector of canonical parameters $\boldsymbol{\lambda} \in \mathcal{C} \subseteq \mathbb{R}^d$. Suppose $\mathbf{x}^1, \dots, \mathbf{x}^n \triangleq \mathbf{x}^{1:n}$ is a sequence of i.i.d. random vectors drawn from f . Let $f_n^{NE}(\mathbf{x}|\hat{\boldsymbol{\theta}}_n, \mathbf{x}^{1:n}) \in \mathcal{NEF}_{\mathbf{s}}$ be the MLE solution of (1), $\hat{\boldsymbol{\theta}}_n = (\tilde{\boldsymbol{\lambda}}_n, \tilde{\boldsymbol{\lambda}}_{a,n})$, with all $\beta_i = \beta > 0, i = 1, \dots, n$. Assuming 1. \mathcal{X} is compact, 2. \mathbf{t} is continuous, 3. $\mathcal{EF}_{\mathbf{t}}$ is a family of uniformly equicontinuous functions w.r.t \mathbf{x} , 4. Kernel K has bounded variation and has a bandwidth parameter \mathbf{H} such that the series $\sum_{n=1}^{\infty} e^{-\gamma n|\mathbf{H}|}$ converges for every positive value of γ , then as $n \rightarrow \infty, \tilde{\lambda}_{a,n}^i \xrightarrow{P} 0, \forall i = 1, \dots, n$ and $\tilde{\boldsymbol{\lambda}}_n \xrightarrow{P} \boldsymbol{\lambda}$.*

Theorem 1.1 shows that if the true distribution falls within the exponential family, then as sample size increases, the estimated density from the non-parametric exponential family will have vanishing reliance on the augmented parameters.

Theorem 1.2. *Given a probability density function $f(\mathbf{x}) : \mathcal{X} \rightarrow \mathbb{R}$, let $f_n^{NE}(\mathbf{x}|\hat{\boldsymbol{\theta}}_n, \mathbf{x}^{1:n}) \in \mathcal{NEF}_{\mathbf{s}}$ be a solution satisfying (1). If 1. f is uniformly continuous on \mathcal{X} , 2. $K_{\mathbf{H}}(\mathbf{x})$ is uniformly continuous on \mathcal{X} , 3. $\sup_{\mathbf{x} \in \mathcal{X}} K_{\mathbf{H}}(\mathbf{x}) < \infty$, 4. $\lim_{\|\mathbf{x}\| \rightarrow \infty} \prod_{i=1}^m \mathbf{x}_i = 0$, 5. $\lim_{n \rightarrow \infty} |\mathbf{H}|^{\frac{1}{2}} = 0$, 6. $\lim_{n \rightarrow \infty} n|\mathbf{H}|^{\frac{1}{2}} = \infty$, then $f_n^{NE}(\mathbf{x}|\hat{\boldsymbol{\theta}}_n, \mathbf{x}^{1:n}) \xrightarrow{P} f(\mathbf{x})$ pointwise on \mathcal{X} .*

Theorem 1.2 indicates the weak consistency of the non-parametric exponential family density estimator. Thus our proposed non-parametric approach can be used to approximate densities which are not from exponential families.

²We omit \mathbf{H} (h in univariate kernels) from t_a^i for the simplicity of notation. It is a tuning parameter that may be set globally for all $i = 1, \dots, n$.

Algorithm 1 Non-Parametric Exponential Family Coordinate Descent

INPUT: Samples $\mathbf{x}^1, \dots, \mathbf{x}^n \in \mathbb{R}^d$, sufficient statistics $\mathbf{t} : \mathcal{X} \rightarrow \mathcal{H}$, augmented features $t_a^i : \mathcal{H} \rightarrow \mathbb{R}$, $i = 1, \dots, n$, ℓ_1 regularization parameters β

OUTPUT: MLE $\boldsymbol{\theta} = (\lambda^1, \dots, \lambda^d, \lambda_a^1, \dots, \lambda_a^n)$

Initialize $\boldsymbol{\theta}^{(0)}$

Compute the sufficient statistics $E_{\hat{f}_n(\mathbf{x}|\mathbf{x}^{1:n})}[\mathbf{t}(\mathbf{x})]$

repeat

 iteration $k = k + 1$, $\boldsymbol{\theta}^{(k)} = \boldsymbol{\theta}^{(k-1)}$

for $i = 1, \dots, d$ **do**

$g_i^{(k)} = E_{\hat{f}_n} [t^i(\mathbf{x})] - E_{f_n^{NE}(\mathbf{x}|\boldsymbol{\theta}^{(k)})} [t^i(\mathbf{x})]$

 Perform line search along $g_i^{(k)}$ to update $\lambda^{i,(k)}$

end for

for $j = 1 \dots n$ **do**

 Solve two equations for λ_a^j (denote the solutions $\lambda_a^{j,-}$ and $\lambda_a^{j,+}$, respectively):

$E_{f_n^{NE}(\mathbf{x}|\boldsymbol{\theta}^{(k)})} [t_a^j(\mathbf{x})] = E_{\hat{f}_n} [t_a^j(\mathbf{x})] - \beta_j$ and $E_{f_n^{NE}(\mathbf{x}|\boldsymbol{\theta}^{(k)})} [t_a^j(\mathbf{x})] = E_{\hat{f}_n} [t_a^j(\mathbf{x})] + \beta_j$

 set $\lambda_a^{j,(k)} = \lambda_a^{j,-}$ if $\lambda_a^{j,-} > 0$, set $\lambda_a^{j,(k)} = \lambda_a^{j,+}$ if $\lambda_a^{j,+} < 0$, and set $\lambda_a^{j,(k)} = 0$ otherwise

end for

until convergence

return $\boldsymbol{\theta}^{(k)}$

1.2 Estimating Parameters for Non-Parametric Exponential Families

Recently there have been a number of methods developed for optimization of convex non-smooth functions, some of them specifically aimed at log-linear problems such as (3) [e.g., 3–5]. We employed a coordinate descent algorithm similar to the SUMMET algorithm of [1] (see Algorithm 1), primarily, due to its simplicity. Other possible approaches can be employed as well and may end up more efficient for this formulation.

The proposed algorithm iterates between optimizing canonical parameters $\boldsymbol{\lambda}$ (by setting $E_{\hat{f}_n}[\mathbf{t}(\mathbf{x})] = E_{f_n^{NE}(\mathbf{x}|\boldsymbol{\theta}^{(k)})}[\mathbf{t}(\mathbf{x})]$) and sequentially optimizing the augmented parameters $\boldsymbol{\lambda}_a$ so that the Karush-Kuhn-Tucker conditions [e.g., 6] are satisfied:

$$E_{\hat{f}_n} [t_a^i(\mathbf{x})] - E_{f_n^{NE}(\mathbf{x}|\boldsymbol{\theta})} [t_a^i(\mathbf{x})] \in \begin{cases} \{\beta_i\} & \lambda_a^i > 0, \\ \{-\beta_i\} & \lambda_a^i < 0, \\ (-\beta_i, \beta_i) & \beta_i = 0. \end{cases}$$

Algorithm 1 belies the inherent difficulty of: (1) calculation of the partial derivative $g_i^{(k)}$, and (2) an implicit search procedure to update $\lambda_a^{j,(k)}$, both involve calculating intractable integrals. If the support is low-dimensional and the mass is contained in a small volume, then the partition function (and thus the gradient) can be computed by numerical integration (quadrature). Alternatively, a common approach to MLE with an intractable partition function $Z(\boldsymbol{\theta})$ is Markov Chain Monte Carlo MLE [MCMC-MLE, 7]. For example, the time complexity at each iteration k is $O(Sn^2)$, where S is the number of Monte-Carlo samples we choose to use. However, we believe developments in optimization [e.g. 3, 5] will help us find an efficient solution.

2 Experimental Evaluation

We illustrate the behavior of the proposed non-parametric density estimator matching first and second order moment constraints (NPGaussian, i.e. $\mathbf{t}(x) = (x, x^2)$) in the univariate setting. Normal density (in $\mathcal{E}\mathcal{F}$, $\mathcal{N}(0, 1)$), t-distribution (not in $\mathcal{E}\mathcal{F}$, $\text{df}=6$) and a mixture of two normals (not in $\mathcal{E}\mathcal{F}$, $\frac{1}{2}\mathcal{N}(-3, 1) + \frac{1}{2}\mathcal{N}(3, 1)$) are used for simulating i.i.d samples. We vary the sample size from 10 to 1000 for training and compute the out-of-sample likelihood with an evaluation set of 100000 samples for testing. Adaptive Gauss-Lobatto quadrature is used for numerical integration. We compared the performance of our non-parametric approach, the model from the true functional family,

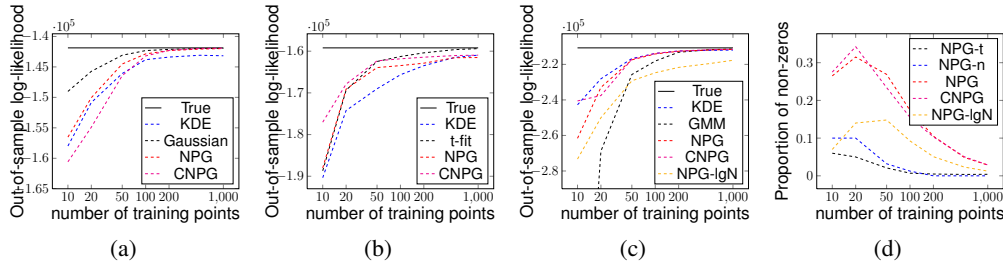


Figure 1: Estimating simple one dimensional densities. Results are averaged over 20 runs. The x axis is in log scale. (a) Normal distribution (b) t distribution (c) Mixed normal distribution (d) Number of non-zero λ_a s. In Figure(a,b,c), NPG: NPGaussian with $O(1/\sqrt{n})$ schedule, NPG-IgN: NPGaussian with $O(1/\log(n))$ schedule, CNPG: constrained NPGaussian with true global moment statistics. In (d), NPG-t: number of non-zeros for estimating normal distribution with NPG, NPG-n: number of non-zeros for estimating mixed normal distribution with NPG, CNPG: number of non-zeros for estimating mixed normal distribution with CNPG, NPG-IgN: number of non-zeros for estimating mixed normal distribution with NPG-IgN.

and another non-parametric approach (KDE). There are two sets of tuning parameters, bandwidth h and the box constraint parameter β , assumed to be the same for all $i = 1, \dots, n$. β was set according to a fixed schedule $\beta(n) = O(1/\sqrt{n})$. h (both for KDE and for our approach) was determined based on cross-validated log-likelihood.³ Gaussian kernel function is used for NPGaussian and for KDE. When estimating normal density, the NPGaussian model (NPG) quickly converges to the normal density with the augmented parameters vanishing as suggested by Theorem 1.1 (Figure 1(a,d)). We also consider the case when the true sufficient statistics are given to us (constrained NPGaussian, CNPG). The CNPG model shows improvement over NPGaussian for small n . However, as the training sample size increases, both CNPG and NPGaussian show similar performance as the moment constraints $t(x)$ are more accurately approximated. NPGaussian gives comparable performance with KDE when estimating densities $\notin \mathcal{EF}_t$ (Figure 1(b,c)). We also experimented with $O(1/\log(n))$ regularization schedule (NPG-IgN) for β s when estimating the mixed normal distribution. As n increase, NPG-IgN becomes sparser than NPGaussian (Fig. 1(d)) and gives a worse performance than KDE (Figure 1(c)), which agrees with Theorem 1.2.

References

- [1] M. Dudik, S. Phillips, and R. Schapire. Maximum entropy density estimation with generalized regularization and an application to species distribution modeling. *Journal of Machine Learning Research*, 8: 1217–1260, Jun 2007.
- [2] L. Yuan, S. Kirshner, and R. Givan. Estimating Densities with Non-Parametric Exponential Families. *ArXiv e-prints*, June 2012. URL <http://arxiv.org/abs/1206.5036>.
- [3] F. Bach, R. Jenatton, J. Mairal, and G. Obozinski. Convex optimization with sparsity-inducing norms. In S. Sra, S. Nowozin, and S. J. Wright., editors, *Optimization for Machine Learning*, chapter 2. MIT press, 2011.
- [4] T. Wu and K. Lange. Coordinate descent algorithms for lasso penalized regression. *The Annals of Applied Statistics*, 2(1):224–244, March 2008.
- [5] S. Shalev-Shwartz and A. Tewari. Stochastic methods for l_1 regularized loss minimization. *Journal of Machine Learning Research*, 12:1865–1892, June 2011.
- [6] J. Nocedal and S. J. Wright. *Numerical Optimization*. Springer Series in Operations Research. Springer, 2nd edition, 2006.
- [7] C. J. Geyer and E. A. Thompson. Constrained monte carlo maximum likelihood for dependent data. *Journal of the Royal Statistical Society*, 54(3):pp. 657–699, 1992.

³Interestingly, unlike KDE, cross-validation implies the choice of h does not impact NPGaussian much.