# An indirect nonparametric regression method for one-dimensional continuous distributions using warping functions

Zhicheng Chen[*]

**ABSTRACT**

Distributions play a very important role in many applications. Inspired by the newly developed warping transformation of distributions, an indirect nonparametric distribution to distribution regression method is proposed in this article for distribution prediction. Additionally, a hybrid approach by fusing the predictions respectively obtained by the proposed method and the conventional method is further developed for reducing risk when the predictor is contaminated.

**Keywords:** distribution regression, warping transformation, nonparametric regression, hybrid approach

## 1. Introduction

In this article, the correlation between two distribution classes means shapes of their density functions will change simultaneously to some extent (see Fig. 1). The correlation between two distribution classes is distinct from the correlation between two random variables. Two correlated distribution classes do not guarantee random variables respectively follow these two distribution classes are also correlated. For instance, suppose probability density functions $g_t$ and $f_t$ are correlated, consider two time-varying random variables $X_t \sim g_t$ and $Y_t \sim f_t$, if the joint distribution of $(X_t, Y_t)$ can be factorized as $p_{X_t Y_t}(x_t, y_t) = g_t(x_t) f_t(y_t)$, then $X_t$ and $Y_t$ are independent with each other. Similarly, two correlated random variables also do not guarantee the distributions they follow are also correlated.
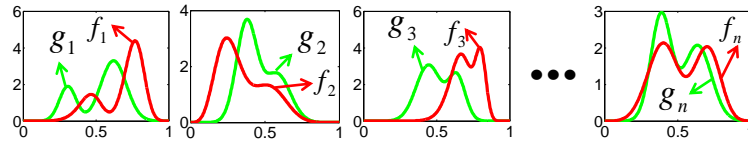


Fig.1    Graphical representation of two correlated distribution classes

Inspired by the seminal work of Dasgupta *et al.* [1], an indirect nonparametric distribution to distribution regression method is proposed in this article for two correlated one-dimensional continuous distribution classes. Other related work includes the conventional distribution to distribution regression (DDR) [2] and distribution to real-value or vector-value regression [3-7], etc.

## 2. Introduction to the warping transformation of distributions

The warping transformation of a distribution is a map that used to transform a distribution to another by deforming the original probability density function with a warping function [1, 8]. The newly reported article by Dasgupta *et al.* [1] has given a very detailed discussion of this transformation.

All distributions in this study are assumed to be continuous with strictly positive support on $[0, 1]$, distributions with general finite supports can be easily tackled by the scale transformation introduced in [1]. Given a probability density function $g(x)$ with strictly positive support on $[0, 1]$, the warping transformation of $g(x)$ by a warping function $\gamma(x)$ defined on $[0, 1]$ is

---

[*] PhD student; Email: 13B933002@hit.edu.cn

$$g_{warp}(x) = g(\gamma(x))\dot{\gamma}(x), \quad x \in [0,1] \tag{1}$$

where $\dot{\gamma}(x) = \dfrac{d}{dx}\gamma(x)$.

Given observed samples from the distribution $f(x)$, $x \in [0,1]$, the optimal warping function used to transform $g(x)$ to get close to $f(x)$, i.e. $f(x) \approx g_{warp}(x)$, can be estimated by the maximum likelihood method proposed in [1].

### 3. The warping transformation-based distribution regression method

For convenience, let $\boldsymbol{\pi}_g$ and $\boldsymbol{\pi}_f$ respectively be two correlated one-dimensional continuous distribution classes with strictly positive support on $[0,1]$. Suppose we have obtained $n$ pairs of probability density functions respectively from $\boldsymbol{\pi}_g$ and $\boldsymbol{\pi}_f$, i.e. $\{g_k, f_k\}_{k=1}^{n}$, given a new density function $g_0$ from $\boldsymbol{\pi}_g$, the task in this section is to develop a nonparametric regression model to predict the corresponding density function $f_0$ from $\boldsymbol{\pi}_f$, i.e. use $\{g_k, f_k\}_{k=1}^{n} \bigcup g_0$ to predict $f_0$. For this purpose, a nonparametric distribution to warping function regression (DWR) is first used to predict the warping function $\gamma_0$, i.e. the mapping relationship from $g_0$ to $f_0$, then use the predicted warping function to transform $g_0$ to obtain a prediction for $f_0$, i.e.

$$\hat{\gamma}_0 = \sum_{k=1}^{n} \hat{\gamma}_k \frac{K(\delta(g_0, g_k)/h)}{\sum_{j=1}^{n} K(\delta(g_0, g_j)/h)} \tag{5a}$$

$$\hat{f}_0(x) = g_0(\hat{\gamma}_0(x))\dot{\hat{\gamma}}_0(x) \tag{5b}$$

where, $\hat{\gamma}_k$ is the estimated warping function from $g_k$ to $f_k$, i.e. $f_k(x) \approx g_k(\hat{\gamma}_k(x))\dot{\hat{\gamma}}_k(x)$. $K(\cdot)$ is the kernel function, $h$ is the bandwidth for the kernel regression, $\delta(g_0, g_k)$ is a metric (such as the $L_1$ distance: $\delta(g_0, g_k) = \int |g_0(\tau) - g_k(\tau)| d\tau$) used to measure the similarity between $g_0$ and $g_k$, $\hat{f}_0(x)$ is the prediction of $f_0(x)$. Note, a convex combination of warping functions is also a warping function [8], thus the regression result in Eq. (5a) being a warping function is guaranteed.

This distribution prediction approach can be regarded as an indirect nonparametric distribution to distribution regression. Unlike the conventional distribution to distribution regression proposed by Oliva *et al.* [2], the proposed regression model in Eq. (5) can reflect the information of shape mapping between input and output distributions, thus it has more potential in extrapolating prediction and reducing shape errors.

It is worth noted that the proposed distribution regression method in Eq. (5) is different from the warping transformation-based approach in [1] for estimating the conditional density function of a random variable $Y$ given the observation of another correlated random variable $X$. The process of conditional distribution estimation in [1] can be regarded as real-value to distribution regression (if $X$ is an univariate random variable) or vector-value to distribution regression (if $X$ is a multivariate random variable). The real-value (or vector-value) to distribution regression and the distribution to distribution regression have different application scopes, the former is suited to correlated random variables, while the latter is suited to correlated distributions.

### 4. Limitations of the proposed method and a hybrid approach for reducing risk

The proposed method has more potential in extrapolating prediction as well as reducing shape

errors by borrowing the shape information of $g_0$ with the predicted warping function. However, such approach also has risk when the predictor $g_0$ has been contaminated, which is the main drawback of the proposed method. To reduce the risk, a hybrid approach by fusing predictions obtained by the proposed method and the conventional distribution to distribution regression method [2] should be more preferable. One fusion approach is using the weighted combination, i.e.

$$\hat{f}_0^{\text{fusion}} = \alpha \hat{f}_0^{\text{DWR}} + (1-\alpha) \hat{f}_0^{\text{DDR}}, \quad 0 \leq \alpha \leq 1 \tag{6}$$

where, $\hat{f}_0^{\text{DWR}}$ and $\hat{f}_0^{\text{DDR}}$ are predictions of the target density function $f_0$ respectively obtained by the proposed method (distribution to warping function regression, DWR) and the conventional method (distribution to distribution regression, DDR [2]), $\alpha$ is the combination coefficient.

If a certain amount of test distributions are available, the value of $\alpha$ can be estimated from test distributions. Let $\{g_v, f_v\}_{v=1}^{V}$ be $V$ different pairs of test distributions with $\{g_v\}_{v=1}^{V}$ served as predictors, then the value of $\alpha$ can be estimated by maximizing the combined log likelihood of samples from all test distributions, i.e.

$$\hat{\alpha} = \arg\max_{\alpha \in [0,1]} \left\{ \sum_{v=1}^{V} \left[ \frac{1}{m_v} \sum_{r=1}^{m_v} \log\left( \hat{f}_v^{\text{fusion}}\left( X_{v,r} | \alpha \right) \right) \right] \right\} \tag{7}$$

where, $X_{v,r} (r=1,2,\cdots,m_v)$ are observed samples from distribution $f_v$, $m_v$ is the sample size, The factor $\dfrac{1}{m_v}$ is used to balance the effect caused by the sample size because the number of observed samples from different distributions may be different.

## 5. Conclusions

An indirect nonparametric distribution to distribution regression method is proposed in this article, which can reflect the information of shape mapping between input and output distributions, thus it has more potential in extrapolation and reducing shape errors. However, the proposed method may result in a poor prediction when the density function served as the predictor has been contaminated. To compensate this shortcoming, a hybrid approach by fusing predictions obtained by the proposed method and the conventional distribution to distribution regression method is also proposed.

## 6. Acknowledgments

**References**

[1] Dasgupta S, Pati D and Srivastava A. A Geometric Framework For Density Modeling. arXiv preprint arXiv:1701.05656; 2017.

[2] Oliva JB, Póczos B and Schneider JG. Distribution to Distribution Regression. In ICML (3) 2013; pp: 1049-1057.

[3] Póczos B, Singh A, Rinaldo A and Wasserman LA. Distribution-Free Distribution Regression. In AISTATS 2013; pp: 507-515.

[4] Oliva JB, Neiswanger W, Póczos B, Schneider JG and Xing EP. Fast Distribution To Real Regression. In AISTATS (Vol. 3, p. 2) 2014.

[5] Szabó Z, Gretton A, Póczos B and Sriperumbudur B. Consistent, two-stage sampled distribution regression via mean embedding. arXiv preprint arXiv:1402.1754; 2014.

[6] Szabó Z, Gretton A, Póczos B and Sriperumbudur BK. Two-stage sampled learning theory on distributions. In AISTATS 2015.

[7] Szabó Z, Sriperumbudur B, Póczos B and Gretton, A. Learning theory for distribution regression. Journal of Machine Learning Research 2016; 17(152), 1-40.

[8] Srivastava A and Klassen EP. Functional and Shape Data Analysis. Springer: New York; 2016.