

# Statistical Archetypal Analysis

Chenyue Wu<sup>a,\*</sup>, Esteban G. Tabak<sup>a</sup>

<sup>a</sup>*Department of Mathematics, Courant Institute of Mathematical Sciences, New York University, 251 Mercer Street, New York, NY 10012, United States*

---

## Abstract

Statistical Archetypal Analysis (SAA) is introduced for the dimensional reduction of a collection of probability distributions known via samples. Applications include medical diagnosis from clinical data in the form of distributions (such as distributions of blood pressure or heart rates from different patients), the analysis of climate data such as temperature or wind speed at different locations, and the study of bifurcations in stochastic dynamical systems. Distributions can be embedded into a Hilbert space with a suitable metric, and then analyzed similarly to feature vectors in Euclidean space. However, most dimensional reduction techniques –such as Principal Component Analysis– are not interpretable for distributions, as neither the components nor the reconstruction of input data by components are themselves distributions. To obtain an interpretable result, Archetypal Analysis (AA) is extended to distributions, requiring the components to be mixtures of the input distributions and approximating the input distributions by mixtures of components.

*Keywords:* archetypal analysis, dimension reduction, energy distance, kernel embedding, principal component analysis

---

## 1. Introduction

Finite collections of probability distributions appear naturally in a variety of settings, often as conditional distributions  $\rho(x|z)$  where  $z$  adopts a discrete set of values. For instance,  $x$  may represent a collection of clinical variables such as body temperature, blood pressure and cholesterol level, and  $z$  may stand for covariates such as sex, age group or medical treatment. In an example that this paper analyzes in some detail,  $x$  is the atmospheric temperature measured at ground level and  $z$  stands for the station where the measurements are performed.

It is therefore a natural extension of data analysis to use as either labels or features, probability distributions instead of the more conventional discrete-valued variables, continuum scalars or vectors. Thus one might want to predict

---

\*Corresponding author

*Email addresses:* [chenyue@cims.nyu.edu](mailto:chenyue@cims.nyu.edu) (Chenyue Wu), [tabak@cims.nyu.edu](mailto:tabak@cims.nyu.edu) (Esteban G. Tabak)

12 not the temperature at a particular location and time but its probability distri-  
13 bution, or cluster populations for medical purposes according to the probability  
14 distributions of a group of clinical variables.

15 A basic quantity that permeates data analysis is the distance between data  
16 points. There are several statistical distances in the literature that measure  
17 the dissimilarity between two probability distributions. Some are based on ana-  
18 logues of the Euclidean distance, some on information theory, some on optimal  
19 transport. Typically, each sheds a different light on what makes two distri-  
20 butions different. In this article, we use the energy distance as a measure of  
21 dissimilarity among distributions, as it is easy to evaluate efficiently from sam-  
22 ple points and can be derived from an inner product, thus rendering accessible  
23 many data analysis tools.

24 We study the problem of dimensional reduction of sets of distributions. Af-  
25 ter being equipped with a metric and embedded into a Hilbert space, distri-  
26 butions can be analyzed similarly to conventional feature vectors. However,  
27 there is a gap between the dimensional reduction of distributions and vectors:  
28 interpretability. Traditional dimension reduction techniques, such as principal  
29 components analysis, lack interpretability when applied to probability distribu-  
30 tions, as the projection of each distribution onto the low dimensional subspace  
31 found is almost surely not a probability distribution: even though probability  
32 distributions can be embedded into a Hilbert space, almost all elements in this  
33 space are not probability distributions, since these are constrained by positivity  
34 and normalization.

35 To overcome this difficulty in interpretation, we use the tools of archetypal  
36 analysis. Archetypal analysis finds a small number of “archetypes” that are  
37 convex combinations of the original data points, and approximates the origi-  
38 nal data points again via convex combinations of these archetypes. A convex  
39 combination can be interpreted as a mixture of probability distributions, so  
40 the archetypes found by archetypal analysis are mixtures of the original dis-  
41 tributions and the original distributions are approximated within the family of  
42 mixtures of the archetypes.

43 This paper is arranged as follows: Section 2 gives a review of archetypal anal-  
44 ysis, of the algorithms for archetypal analysis in the general case and specifically  
45 for energy distance. Section 3 reviews reproducing kernel Hilbert space, energy  
46 distance, describes how distributions equipped with the energy distance can be  
47 embedded into a Hilbert space, and describes algorithms to evaluate the energy  
48 distance from samples. Section 4 introduces statistical archetypal analysis for  
49 the dimensional reduction of probability distributions and includes applications  
50 with numerical experiment.

## 51 **2. Archetypal Analysis**

52 Archetypal analysis approximates data points by convex combination of pro-  
53 totypes, where these prototypes, denoted “archetypes”, are themselves convex  
54 combinations of the data points.

55 Archetypal analysis was introduced in [1] –see also [2]– as a dimensional  
 56 reduction method alternative to principal components analysis (PCA), yielding  
 57 more interpretable results. It originated in the study of a dataset consisting of  
 58 6 head dimensions for 200 soldiers, with the goal of designing face masks for the  
 59 Swiss Army. For this dataset, PCA found principal components that did not  
 60 resemble a head shape. To have patterns resembling “pure types” in the data,  
 61 each entry in the dataset was approximated by a mixture of the patterns. To  
 62 make patterns resemble the data, each pattern itself was a mixture of the data  
 63 points.

For a data matrix  $X = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$  representing  $n$  observations, each of  
 dimension  $m$ , Archetypal Analysis seeks  $k \ll n$   $m$ -dimensional archetypes  $Z =$   
 $(\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_k)$ , such that each  $\mathbf{x}_i$  can be approximated by a convex combination  
 of the  $\mathbf{z}_k$ :

$$x_i \approx a_{1i}\mathbf{z}_1 + a_{2i}\mathbf{z}_2 + \dots + a_{ki}\mathbf{z}_k, \quad a_{ji} \geq 0, \quad \sum_j a_{ji} = 1,$$

where the  $\mathbf{z}_j$  themselves are convex combinations of the data:

$$\mathbf{z}_j = b_{1j}\mathbf{x}_1 + b_{2j}\mathbf{x}_2 + \dots + b_{nj}\mathbf{x}_n, \quad b_{ij} \geq 0, \quad \sum_i b_{ij} = 1.$$

After setting a number of archetypes  $k$ , the coefficients  $a$  and  $b$  arise from the  
 optimization problem

$$\min_{a_{ji}, b_{lj}} \sum_{i=1}^n \left\| \mathbf{x}_i - \sum_{j=1}^k a_{ji} \sum_{l=1}^n b_{lj} \mathbf{x}_l \right\|^2, \quad (1)$$

with constraints

$$a_{ji} \geq 0, \quad \sum_j a_{ji} = 1, \quad b_{ij} \geq 0, \quad \sum_i b_{ij} = 1,$$

or, in terms of the matrices  $A = (a_{ji})_{k \times n}$  and  $B = (b_{ij})_{n \times k}$ ,

$$\min_{A, B} \|X - XBA\|_F^2 \quad (2)$$

64 under the same constraints, with  $\|\cdot\|_F$  denoting the Frobenius norm

$$\|M\|_F = \left( \sum_{i=1}^p \sum_{j=1}^q |m_{ij}|^2 \right)^{\frac{1}{2}}. \text{ Alternatively, this can be written as:}$$

$$A, B = \operatorname{argmin} \operatorname{tr} [(I_n - BA)^\top G (I_n - BA)], \quad (3)$$

65 where  $G = X^\top X$  is the Gram matrix of data. This restatement is particularly  
 66 convenient, as it will allow us to formulate the problem in terms of inner products  
 67 among the data points instead of the points themselves, which in our problem  
 68 are distributions. Thus we need a norm for distributions that derive from an  
 69 inner product, for which we will adopt the energy distance.

70 **3. Energy Distance**

71 The energy distance is a metric defined on probability measures ([3, 4]),  
 72 which we will use to measure dissimilarity among probability distributions.

**Definition 1** (Energy Distance). For probability measures  $\mu, \nu$  on  $\mathbb{R}^d$ , random vectors  $X, X' \sim \mu(x), Y, Y' \sim \nu(y)$ ,  $\mathbb{E}\|X\| < \infty, \mathbb{E}\|Y\| < \infty$ , the energy distance between  $\mu$  and  $\nu$ ,  $D(\mu, \nu)$ , is defined by

$$D^2(\mu, \nu) = 2\mathbb{E}\|X - Y\| - \mathbb{E}\|X - X'\| - \mathbb{E}\|Y - Y'\|, \quad (4)$$

73 where  $\|\cdot\|$  is the Euclidean norm on  $\mathbb{R}^d$ , and  $X, X', Y$  and  $Y'$  are pairwise  
 74 independent.

The energy distance as defined above is a metric on distributions ([5], [6]). It can be viewed as the metric induced by kernel embedding ([7]) with kernel

$$k(x, y) = \|x - x_0\| + \|y - y_0\| - \|x - y\|, \quad (5)$$

where  $x_0$  is a fixed value in  $\mathbb{R}^d$ , whose choice does not affect the induced metric. The kernel induces an inner product between distributions  $P$  and  $Q$ :

$$\langle P, Q \rangle = \mathbb{E}_{X, Y} k(X, Y) \quad (6)$$

where  $X \sim P, Y \sim Q$ , with  $X$  and  $Y$  independent. The corresponding square-distance is given by

$$\begin{aligned} \gamma_k^2(P, Q) &= \langle P, P \rangle + \langle Q, Q \rangle - 2\langle P, Q \rangle \\ &= \mathbb{E}_{X, X'} k(X, X') + \mathbb{E}_{Y, Y'} k(Y, Y') - 2\mathbb{E}_{X, Y} k(X, Y), \end{aligned} \quad (7)$$

where the random vectors  $X, X' \sim P(x), Y, Y' \sim Q(y)$  are pairwise independent (conditions for kernels to yield a metric can be found in [5, 8]). In terms of the kernel in (5),

$$\begin{aligned} \gamma_k^2(\mu, \nu) &= 2\mathbb{E}\|X - x_0\| - \mathbb{E}\|X - X'\| + 2\mathbb{E}\|Y - x_0\| - \mathbb{E}\|Y - Y'\| \\ &\quad - 2\mathbb{E}\|X - x_0\| - 2\mathbb{E}\|Y - x_0\| + 2\mathbb{E}\|X - Y\| = D^2(\mu, \nu). \end{aligned}$$

75 A number of distances for distributions is available in the literature of statis-  
 76 tics, probability and information theory, such as the Kullback-Leibler divergence  
 77 ([9, 10]) and the  $p$ -Wasserstein metric between two probability measures  $\mu(x)$   
 78 and  $\nu(x)$  on a metric space  $(M, d)$  ([11]). We chose the energy distance be-  
 79 cause it can be estimated efficiently from samples and it embeds the probability  
 80 measures into a Hilbert space, which facilitates further analysis.

81 *3.1. Estimating the Energy Distance from Data*

In calculating the energy distance between two distributions  $\mu$  and  $\nu$  given independent random vectors  $X \sim \mu, Y \sim \nu$  and their i.i.d. copies  $X', Y'$ ,

$$D(\mu, \nu) = \sqrt{2\mathbb{E}\|X - Y\| - \mathbb{E}\|X - X'\| - \mathbb{E}\|Y - Y'\|},$$

82 one needs to evaluate three expectations:  $\mathbb{E}\|X - Y\|$ ,  $\mathbb{E}\|X - X'\|$  and  $\mathbb{E}\|Y - Y'\|$ .  
 83 If we only have samples of  $\mu$  and  $\nu$ , these expectation can be approximated by  
 84 their empirical means.

Specifically, when we have samples  $\{x_i\}_{i=1}^{n_X}$  of  $\mu$  and  $\{y_j\}_{j=1}^{n_Y}$  of  $\nu$ , we can estimate the energy distance between  $\mu$  and  $\nu$  by the energy distance between their corresponding empirical distributions  $\hat{\mu}$  and  $\hat{\nu}$ :

$$D(\hat{\mu}, \hat{\nu}) = \sqrt{2\mathbb{E}\|\hat{X} - \hat{Y}\| - \mathbb{E}\|\hat{X} - \hat{X}'\| - \mathbb{E}\|\hat{Y} - \hat{Y}'\|}. \quad (8)$$

In the equations above,

$$\mathbb{E}\|\hat{X} - \hat{Y}\| = \frac{1}{n_X n_Y} \sum_{i,j=1}^{i=n_X, j=n_Y} \|x_i - y_j\| \quad (9)$$

is the empirical mean of  $\mathbb{E}\|X - Y\|$ . For  $\hat{X}'$ , we use the same samples available for  $X$ ,

$$\mathbb{E}\|\hat{X} - \hat{X}'\| = \frac{1}{n_X n_X} \sum_{i,i'=1}^{i=n_X, i'=n_X} \|x_i - x_{i'}\|. \quad (10)$$

Similarly,

$$\mathbb{E}\|\hat{Y} - \hat{Y}'\| = \frac{1}{n_Y n_Y} \sum_{j,j'=1}^{j=n_Y, j'=n_Y} \|y_j - y_{j'}\|. \quad (11)$$

85 According to the formulations above for estimating energy distance from  
 86 samples, if we have  $n_X$  sample points for  $\mu$  and  $n_Y$  sample points for  $\nu$ , the  
 87 time complexity of estimating their energy distance is  $O(n_X n_Y + n_X^2 + n_Y^2)$ .

The corresponding inner product between distributions  $\mu$  and  $\nu$ , given independent random vectors  $X \sim \mu$ ,  $Y \sim \nu$  and  $X', Y'$ , is

$$\langle \mu, \nu \rangle = \mathbb{E}\|X - x_0\| + \mathbb{E}\|Y - x_0\| - \mathbb{E}\|X - Y\|, \quad (12)$$

where  $x_0$  is a fixed point. Similarly, calculation of this inner product involves three expectations:  $\mathbb{E}\|X - x_0\|$ ,  $\mathbb{E}\|Y - x_0\|$ ,  $\mathbb{E}\|X - Y\|$ . When  $\mu$  and  $\nu$  are known via their samples  $\{x_i\}_{i=1}^{n_X}$  and  $\{y_j\}_{j=1}^{n_Y}$ , their inner product  $\langle \mu, \nu \rangle$  is estimated by

$$\langle \hat{\mu}, \hat{\nu} \rangle = \mathbb{E}\|\hat{X} - x_0\| + \mathbb{E}\|\hat{Y} - x_0\| - \mathbb{E}\|\hat{X} - \hat{Y}\|, \quad (13)$$

where

$$\mathbb{E}\|\hat{X} - x_0\| = \frac{1}{n_X} \sum_{i=1}^{n_X} \|x_i - x_0\|, \quad (14)$$

$$\mathbb{E}\|\hat{Y} - x_0\| = \frac{1}{n_Y} \sum_{j=1}^{n_Y} \|y_j - x_0\|, \quad (15)$$

$$\mathbb{E}\|\hat{X} - \hat{Y}\| = \frac{1}{n_X n_Y} \sum_{i,j=1}^{i=n_X, j=n_Y} \|x_i - y_j\|, \quad (16)$$

---

**Algorithm 1** Generic algorithm for estimating the energy distance

---

**Input:** Samples  $x_i, y_j$  of  $X, Y$  respectively.

**Output:** Empirical estimation of  $\mathbb{E}\|X - Y\|$ .

```
1: procedure ENERGY( $\{x_i\}, \{y_j\}$ )
2:   sum = 0
3:   for all  $x_i$  do
4:     for all  $y_j$  do
5:       sum = sum +  $\|x_i - y_j\|$ 
6:     end for
7:   end for
8:   return  $\frac{\text{sum}}{n_X n_Y}$ 
9: end procedure
```

---

88 and the time complexity of estimating this inner product is  $O(n_X n_Y)$ . If we  
89 have  $n$  sample points for both  $\mu$  and  $\nu$ , the time complexity is  $O(n^2)$ .

90 Notice that estimating the energy distance and the corresponding inner prod-  
91 uct from  $n$  sample points of both distributions have time complexities  $O(n^2)$ ,  
92 which becomes computationally expensive when using a large number of sample  
93 points. In the following section, a fast algorithm for energy distance between  
94 one-dimensional distributions is introduced, making the application of energy  
95 distance much more efficient.

### 96 3.2. Fast Algorithm in One Dimension

According to (9), (10), (11), (14), (15) and (16), both the data-based com-  
putations of energy distance (8) and corresponding inner product (13) have the  
same complexity of evaluating

$$\mathbb{E}\|\hat{X} - \hat{Y}\| = \frac{1}{n_X n_Y} \sum_{i=1}^{n_X} \sum_{j=1}^{n_Y} \|x_i - y_j\|, \quad (17)$$

97 where the  $(x_i, y_j)$  are  $(n_X, n_Y)$  samples of  $(X, Y)$ . Generally, the time com-  
98 plexity of evaluating (17) is  $O(n^2)$  via Algorithm 1, which simply takes the  
99 arithmetic mean of  $\|x_i - y_j\|$ .

100

In one-dimensional space, however, the fact that  $\|\cdot\| = |\cdot|$  enables us to use

the identity  $|x - y| = \mathbf{1}_{x-y>0}(x - y) - \mathbf{1}_{x-y\leq 0}(x - y)$  to obtain

$$\begin{aligned}
& \mathbb{E}\|\hat{X} - \hat{Y}\| \\
&= \frac{1}{n_X n_Y} \sum_{i=1}^{n_X} \sum_{j=1}^{n_Y} |x_i - y_j| \\
&= \frac{1}{n_X n_Y} \sum_{i=1}^{n_X} \sum_{j=1}^{n_Y} \mathbf{1}_{\{x_i - y_j > 0\}}(x_i - y_j) - \mathbf{1}_{\{x_i - y_j \leq 0\}}(x_i - y_j) \\
&= \frac{1}{n_X} \sum_{i=1}^{n_X} \frac{\#\{j|y_j < x_i\} - \#\{j|y_j \geq x_i\}}{n_Y} x_i + \frac{1}{n_Y} \sum_{j=1}^{n_Y} \frac{\#\{i|x_i \leq y_j\} - \#\{i|x_i > y_j\}}{n_X} y_j,
\end{aligned}$$

101 where  $\#\{\dots\}$  denote the number of elements in a set.

102 If  $\{x_i\}_{i=1}^{n_X}$  and  $\{y_j\}_{j=1}^{n_Y}$  are sorted arrays, the latter expression can be calcu-  
103 lated in the linear time  $O(n_X + n_Y)$ , since each of  $\#\{j|y_j < x_i\}$ ,  $\#\{j|y_j \geq x_i\}$ ,  
104  $\#\{i|x_i \leq y_j\}$  and  $\#\{i|x_i > y_j\}$  can be calculated in linear time by merging  
105  $\{x_i\}_{i=1}^{n_X}$  and  $\{y_j\}_{j=1}^{n_Y}$  into one sorted array (Algorithm 2.)

106 If given unsorted samples, we need to sort them before applying Algorithm  
107 2. Feasible sorting algorithms are quick sort, which has an  $O(n \log n)$  average  
108 complexity and an  $O(n^2)$  worst case complexity, heap sort and merge sort, which  
109 have an  $O(n \log n)$  worst case complexity. Therefore even for unsorted samples,  
110 the complexity of estimating the energy distance can be bounded by  $O(n \log n)$ .

## 111 4. Statistical Archetypal Analysis

### 112 4.1. Dimensional Reduction

113 In this section, we study the dimensional reduction of probability distri-  
114 butions, mapping a collection of distributions to a low-dimensional space with  
115 minimal loss of information. Probability distributions have infinite dimension;  
116 when they are known via samples, they can be said to have a dimensionality of  
117 the order of the number of samples points. Our dimensional reduction on this  
118 high-dimensional dataset consists of two steps: we embed the distributions into  
119 an Euclidean space, and then use dimensional reduction methods developed for  
120 Euclidean spaces.

121 Probability distributions equipped with the energy distance form a convex  
122 subset of a Hilbert space. Therefore a collection of  $N$  distributions  $\mu_i$  can be  
123 naturally embedded into an  $N$ -dimensional Euclidean space, since every finite  
124 dimension subspace of a Hilbert space is isometric to an Euclidean space.

Assume that  $x_i \in \mathbb{R}^N$ ,  $i \in [1, 2, \dots, N]$ , are points in  $\mathbb{R}^N$  such that  $\|x_i - x_j\| = D(\mu_i, \mu_j)$  where  $D(\cdot)$  is the energy distance (In other words,  $x_i$  is the image of  $\mu_i$  under the embedding into an Euclidean space.) Principal Components Analysis (PCA) solves the following optimization problem for centered  $x_i$ :

$$\min_{z_j, a_{ji}} \sum_{i=1}^N \left\| x_i - \sum_{j=1}^K a_{ji} z_j \right\|^2 \tag{18}$$

---

**Algorithm 2** Fast algorithm for estimating the energy distance in 1D

---

**Input:** Sorted samples  $x_i, y_j$  of 1D random variable  $X, Y$  respectively

**Output:** Empirical estimation of  $\mathbb{E}\|X - Y\|$

```

1: procedure FASTENERGY( $\{x_i\}, \{y_j\}$ )
2:    $\text{sum}_X = 0, \text{sum}_Y = 0, i = 1, j = 1$ 
3:   while  $i \leq n_X$  and  $j \leq n_Y$  do
4:     if  $x_i \leq y_j$  then
5:        $\text{sum}_X = \text{sum}_X + \frac{(j-1) - [n_Y - (j-1)]}{n_Y} x_i$ 
6:        $i = i + 1$ 
7:     else
8:        $\text{sum}_Y = \text{sum}_Y + \frac{(i-1) - [n_X - (i-1)]}{n_X} y_j$ 
9:        $j = j + 1$ 
10:    end if
11:  end while
12:  if  $i > n_X$  then
13:     $\text{sum}_Y = \text{sum}_Y + \sum_{k=j}^{n_Y} y_k$ 
14:  else
15:     $\text{sum}_X = \text{sum}_X + \sum_{k=i}^{n_X} x_k$ 
16:  end if
17:  return  $\text{sum}_X/n_X + \text{sum}_Y/n_Y$ 
18: end procedure

```

---

125 under the constraints that the  $z_j$  are orthonormal vectors. Thus PCA maps  
126 each data point to the closest point in the vector space spanned by the  $z_j$ . Here  
127  $K$  is the dimension of the low dimensional space sought, and  $\sum_{j=1}^K a_{ji} z_j$  is the  
128 image of  $x_i$  under this dimensional reduction.

129 PCA and other mainstream dimensional reduction techniques are not ap-  
130 propriate for probability distributions from two perspectives: 1) the  $z_j$  in (18)  
131 is generally not a probability distribution, neither are almost all points in the  
132 space spanned by the  $z_j$ . 2) the coefficients  $a_{ji}$  for each  $x_i$  in (18) may be nega-  
133 tive, so they cannot clearly express how each  $z_j$  contributes to the representation  
134 of  $x_i$ . We will use instead archetypal analysis for the dimensional reduction of  
135 distributions, which does not suffer from this lack of interpretability.

#### 136 4.2. Statistical Archetypal Analysis

As seen in Section 2, archetypal analysis has a formulation similar to (18),  
except that it requires the optimization of

$$\min_{z_j, a_{ji}} \sum_{i=1}^N \left\| x_i - \sum_{j=1}^K a_{ji} z_j \right\|^2 \quad (19)$$



137 under the constraints that  $a_{ji} \geq 0$ ,  $\sum_{j=1}^K a_{ji} = 1$ , and each  $z_j$  is a convex combi-  
 138 nation of the  $x_i$ , i.e.  $z_j = \sum_{l=1}^N b_{lj} x_l$ , with  $b_{lj} \geq 0$ ,  $\sum_{l=1}^N b_{lj} = 1$ .

Switching from vectors  $x_i$ ,  $z_j$  to distributions  $\mu_i$ ,  $\nu_j$  and using energy distance instead of Euclidean distance, statistical archetypal analysis adopts the form

$$\min_{\nu_j, a_{ji}} \sum_{i=1}^N \left\| \mu_i - \sum_{j=1}^K a_{ji} \nu_j \right\|^2, \quad \nu_j = \sum_{l=1}^N b_{lj} \mu_l, \quad (20)$$

139 with the same constraints over the  $a$  and  $b$ , which now adopt the natural inter-  
 140 pretation that the  $\nu_j$  are mixtures of the  $\mu_i$  and the latter are well-approximated  
 141 by mixtures of the  $\nu_j$ .  $\|\cdot\|$  in (20) is the energy distance, but can naturally be  
 142 extended to any metric induced by kernel embedding as discussed in Section 3.

Since the energy distance, that we shall use for the norm in (20), derives from an inner product, statistical archetypal analysis can be rewritten as in (3):

$$\operatorname{argmin}_{A, B} \operatorname{tr}[(I_n - BA)^\top G(I_n - BA)], \quad (21)$$

where each column of  $A$  represents one archetype as a convex combination of the original distributions, and each column of  $B$  contains the coefficients for the approximate reconstruction of each original distribution from the archetypes.  $G$  is the Gram matrix of pairwise inner products among the distributions,

$$G_{ij} = \mathbb{E}k(\mathbf{X}_i, \mathbf{X}_j)$$

143 for independent  $\mathbf{X}_i \sim \mu_i$  and  $\mathbf{X}_j \sim \mu_j$  and kernel  $k$ .

144 When each  $\mu_i$  is known via samples  $\{y_m^{(i)}\}_{m=1}^{M_i}$  of size  $M_i$ , we can replace  $\mu_i$   
 145 by its empirical distribution at data points  $y_m^{(i)}$  with weights  $\frac{1}{M_i}$ . In this setting,  
 146  $\nu_j$  becomes an empirical distribution concentrated at the union of the  $y_m^{(l)}$  over  
 147  $l = 1, 2, \dots, N$ , with weights  $\frac{b_{lj}}{M_l}$  for all  $m$ . The resulting number of samples of  
 148  $\nu_j$  appears large, since it contains the support of every empirical distribution  
 149  $\mu_i$ . However, since the solution of (19) is sparse, most entries in  $b_{lj}$  are zero, so  
 150 we only need to keep those data points  $y_m^{(l)}$  for  $\nu_j$  where  $b_{lj}$  is non-zero.

151 Statistical archetypal analysis overcomes the two difficulties in interpreta-  
 152 tion when applying dimension reduction on probability distribution. Archetypes  
 153  $\{\nu_i\}_{i=1}^k$  found by archetypal analysis, which are mixtures of the  $\{\mu_i\}_{i=1}^n$ , are all  
 154 probability distributions. The low-dimensional space used to capture informa-  
 155 tion of the dataset of distributions in this case is the convex hull of all archetypes,  
 156 i.e. the family of mixtures of all archetypes. Each coefficient  $a_{ji}$  in (19) stands  
 157 for the contribution of the  $j^{\text{th}}$  archetype  $\nu_j$  to  $\mu_i$ .

### 158 4.3. Numerical Examples

#### 159 4.3.1. Synthetic Data

160 In our first example, we simulate 100 probability distributions  $\{\mu_i\}_{i=1}^{100}$ , each  
 161 a Gaussian mixture  $\mu_i = \lambda_i \mathcal{N}(-6, 2) + (1 - \lambda_i) \mathcal{N}(6, 1)$ , where each  $\lambda_i$  is drawn

162 independently from the uniform distribution in  $[0, 1]$ .

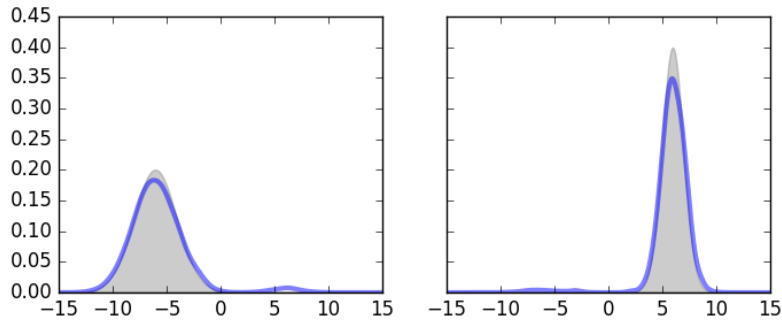


Figure 1: Archetypes of synthetic data for  $k=2$ . The curves are the found archetypes and the shadows are the two components  $\mathcal{N}(-6, 2)$  and  $\mathcal{N}(6, 1)$  in the mixture family respectively.

163 We set number of archetypes  $k$  to 2 and perform archetypal analysis on the  
164 synthetic data. The two archetypes found are shown and compared to  $\mathcal{N}(-6, 2)$   
165 and  $\mathcal{N}(6, 1)$ , the two components in the mixture family, in Figure 1. Both of  
166 them are close to the components except at the center and tail part. This is due  
167 to the definition of archetypes, which is a mixture of input distributions. Unless  
168 we have exactly these two components as input, the archetypes will always have  
169 a heavier tail.

#### 170 4.3.2. Temperature Data

171 We work with ground temperature data from United States Climate Refer-  
172 ence Network (USCRN) Quality Controlled Datasets ([dataset][12], [13], [14]).  
173 The temperature data are measured hourly in 43 cities across the United States.  
174 We operate on data from which the diurnal and seasonal signal has been removed  
175 using the optimal-transport based methodology in [15]. In addition, this dataset  
176 has missing values, which are filled using a low rank approximation to the data  
177 matrix. Figure 2 shows the 43 cities on the map with Table 1 a complete list of  
178 the cities.

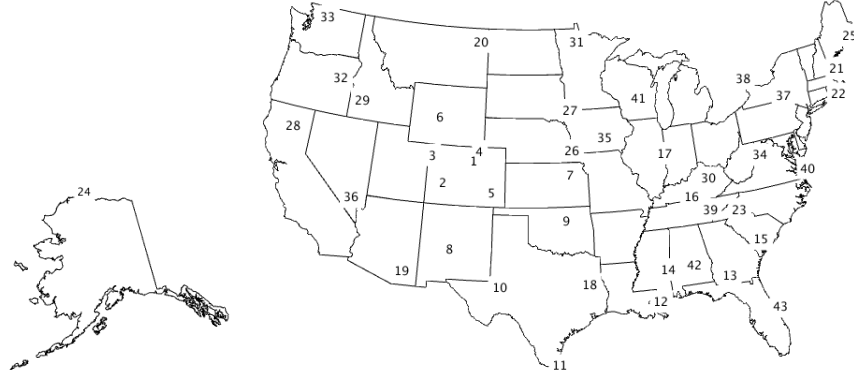


Figure 2: Locations on the map where the data were collected.

Table 1: Locations on the map where the data were collected

Index	City	Index	City	Index	City
1	Boulder	16	KY-Bowling Green	31	MN-Goodridge
2	Montrose	17	IL-Champaign	32	OR-John Day
3	Dinosaur	18	TX-Palestine	33	WA-Darrington
4	Nunn	19	AZ-Tucson	34	WV-Elkins
5	LaJunta	20	MT-Wolf Point	35	IA-Des Moines
6	Lander	21	NH-Durham	36	NV-Mercury
7	ManhattanKs	22	RI-Kingston	37	NY-Ithaca
8	Socorro	23	NC-Asheville	38	ON-Egbert
9	Stillwater	24	AK-Barrow	39	TN-Crossville
10	Monahans	25	ME-Old Town	40	VA-Cape Charles
11	Edinburg	26	NE-Lincoln	41	WI-Necedah
12	Lafayette	27	SD-Sioux Falls	42	AL-Selma
13	Newton	28	CA-Redding	43	FL-Titusville
14	MsNewton	29	ID-Murphy		
15	SC-Blackville	30	KY-Versailles		

We choose alternatively  $K = 3, 5$  as the number of archetypes. For  $K = 3$ , the resulting archetypes are shown in Figure 3; the corresponding mixtures are as follows:

$$\begin{aligned}
 \text{Archetype 1: } & 0.66875 \times \text{MN-Goodridge} \\
 & + 0.01916 \times \text{NY-Ithaca} + 0.31209 \times \text{WV-Elkins}, \\
 \text{Archetype 2: } & 0.22233 \times \text{Edinburg} + 0.77767 \times \text{Lafayette}, \\
 \text{Archetype 3: } & 0.01292 \times \text{VA-Cape Charles} + 0.98707 \times \text{WA-Darrington}.
 \end{aligned}$$

<sup>179</sup> The main difference between these three archetypes is how much they are spread.

180 The first archetype has the heaviest tail among the three while the last archetype  
 181 has the largest peak at center. The second archetype also has a marked asym-  
 182 metry.

183 Figure 4 shows the plane spanned by these three archetypes. The bottom  
 184 left cross is the first archetype, the bottom right cross is the second archetype  
 185 and the top cross is the third archetype, which consists almost exclusively of the  
 186 distribution at WA-Darrington. Each point represents the best approximation  
 187 within the convex hull to its corresponding distribution for one city.

188 The approximation of distributions at each station by mixtures of archetypes  
 189 are shown in Figures 5–9. We can see that, except for the distribution at  
 190 Titusville, FL, the distributions at all 43 stations can be well approximated by  
 191 mixtures of just three archetypes. These results indicate strongly that there is  
 192 a low dimension structure underlying this dataset.

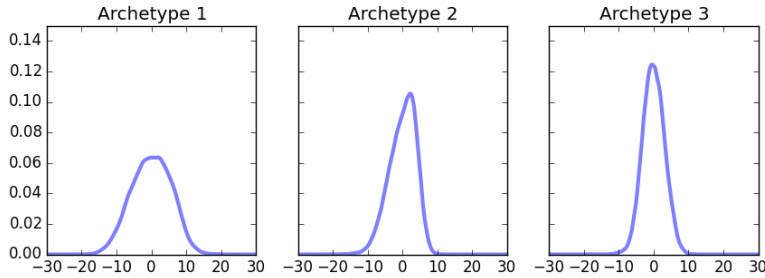


Figure 3: Archetypes of temperature data for  $k=3$ .

For  $K = 5$ , the archetypes are shown in Figure 10; the corresponding mix-  
 tures are:

- Archetype 1:  $0.07329 \times \text{AK-Barrow} + 0.71803 \times \text{NH-Durham}$   
 $+ 0.20867 \times \text{RI-Kingston},$
- Archetype 2:  $0.68181 \times \text{Dinosaur} + 0.31819 \times \text{Lafayette},$
- Archetype 3:  $0.09892 \times \text{Edinburg} + 0.90108 \times \text{FL-Titusville},$
- Archetype 4:  $0.44146 \times \text{MN-Goodridge} + 0.41827 \times \text{MT-Wolf Point}$   
 $+ 0.14027 \times \text{ManhattanKs},$
- Archetype 5:  $0.01306 \times \text{VA-Cape Charles} + 0.98694 \times \text{WA-Darrington}.$

193 When the number of archetypes  $K$  is increased from 3 to 5, the archetypes  
 194 found for  $K = 3$  are not the same as for  $K = 5$ : only the last archetypes  
 195 for  $K = 3$  and  $K = 5$  are close. This is due to the fact that in archetypal  
 196 analysis, when the number of archetypes is increased, the shape of convex hull  
 197 of archetypes changes so as to be as close to the data points as possible.

198 The approximation to the original distributions by mixtures of archetypes  
 199 are shown in Figures 11–15. In this example, we find that when the number of

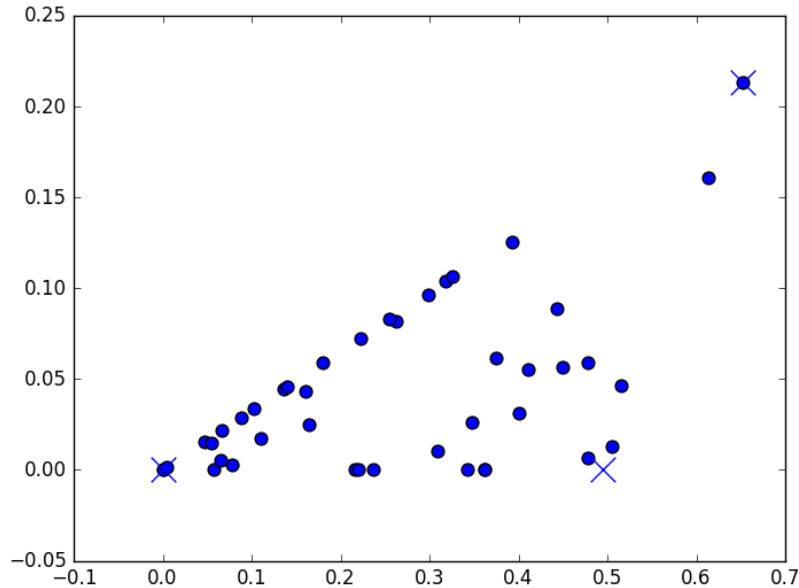


Figure 4: Convex hull spanned by archetypes of temperature data for  $k=3$ . A cross stands for one archetype and a point for the distribution at each city.

200 archetypes is increased to 5, the mixtures of archetypes offer an almost perfect  
 201 approximation to the distributions for all the 43 cities.

## 202 5. Conclusions

203 This article develops statistical archetypal analysis for dimension reduction  
 204 of probability distributions. Archetypal analysis constrains the archetypes –  
 205 analogues of principal components– to convex combinations of the data, and  
 206 approximates the data as convex combinations of these archetypes, hence providing  
 207 an interpretable fit for distributions, with patterns that can be interpreted  
 208 as mixtures of distributions.

209 In order to perform archetypal analysis on distributions, one needs a metric  
 210 and a linear structure. A natural way to introduce these is through an embedding  
 211 of the distributions into a Hilbert space, for which we have used the energy  
 212 distance (one of the many choices provided by the theory of reproducing kernel  
 213 Hilbert spaces for distributions.)

214 As a proof of concept, statistical archetypal analysis was applied to both  
 215 synthetic and temperature data. Statistical archetypal analysis recovers the  
 216 components of a mixture family used to generate synthetic data, and reveals a

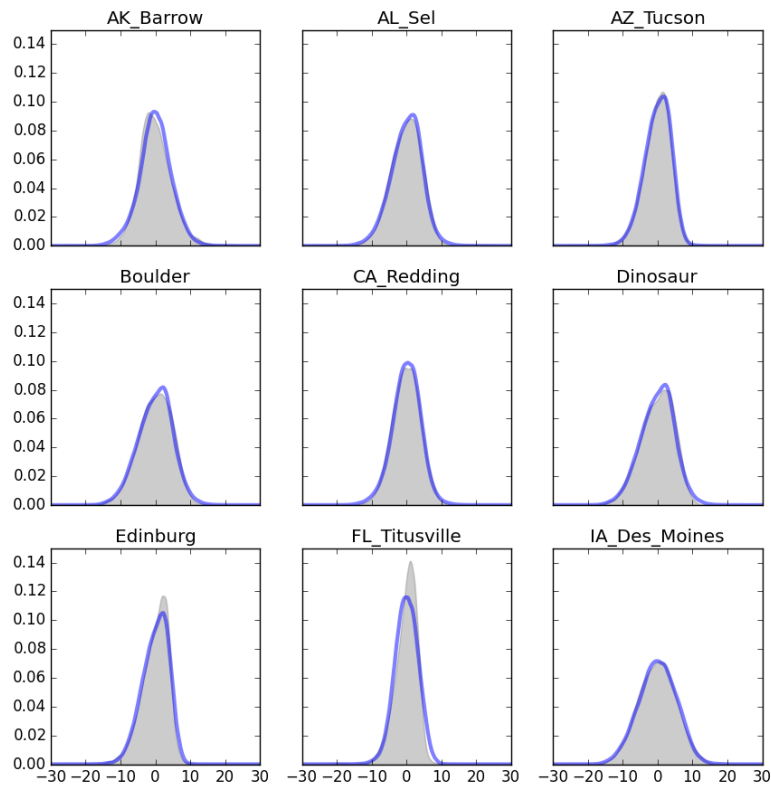


Figure 5: Reconstruction of distribution at each city by 3 archetypes. The original distributions are depicted as shadows; their approximation by mixtures of archetypes as solid curves.

217 low dimensional structure in the distributions of temperature data across the  
 218 United States.

## 219 References

- 220 [1] A. Cutler, L. Breiman, Archetypal analysis, *Technometrics* 36 (4) (1994)  
 221 338–347. doi:10.1080/00401706.1994.10485840.
- 222 [2] J. Friedman, T. Hastie, R. Tibshirani, *The elements of statistical learning*,  
 223 Vol. 1, Springer series in statistics Springer, Berlin, 2001.
- 224 [3] M. L. Rizzo, G. J. Székely, *Energy distance*, Wiley Interdisciplinary Re-

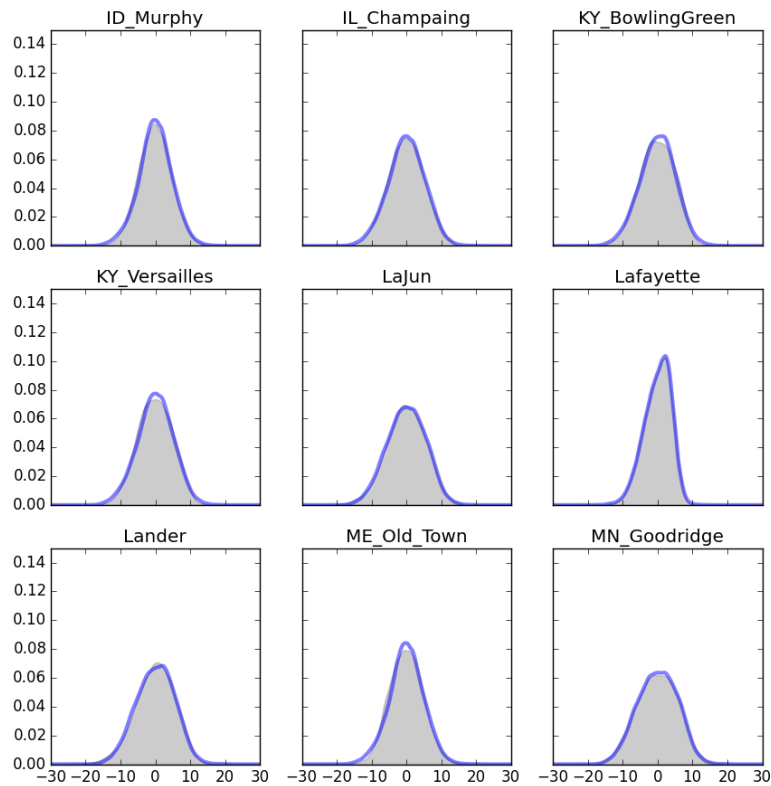


Figure 6: Reconstruction of distribution at each city by 3 archetypes. The original distributions are depicted as shadows; their approximation by mixtures of archetypes as solid curves.

- 225 views: Computational Statistics 8 (1) (2016) 27–38. doi:10.1002/wics.  
 226 1375.
- 227 [4] G. J. Székely, M. L. Rizzo, Energy statistics: A class of statistics based  
 228 on distances, Journal of statistical planning and inference 143 (8) (2013)  
 229 1249–1272. doi:10.1016/j.jspi.2013.03.018.
- 230 [5] L. B. Klebanov, A class of probability metrics and its statistical  
 231 applications, in: Statistical Data Analysis Based on the L1-Norm  
 232 and Related Methods, Springer, 2002, pp. 241–252. doi:10.1007/  
 233 978-3-0348-8201-9\_20.
- 234 [6] G. J. Székely, M. L. Rizzo, A new test for multivariate normality, Journal

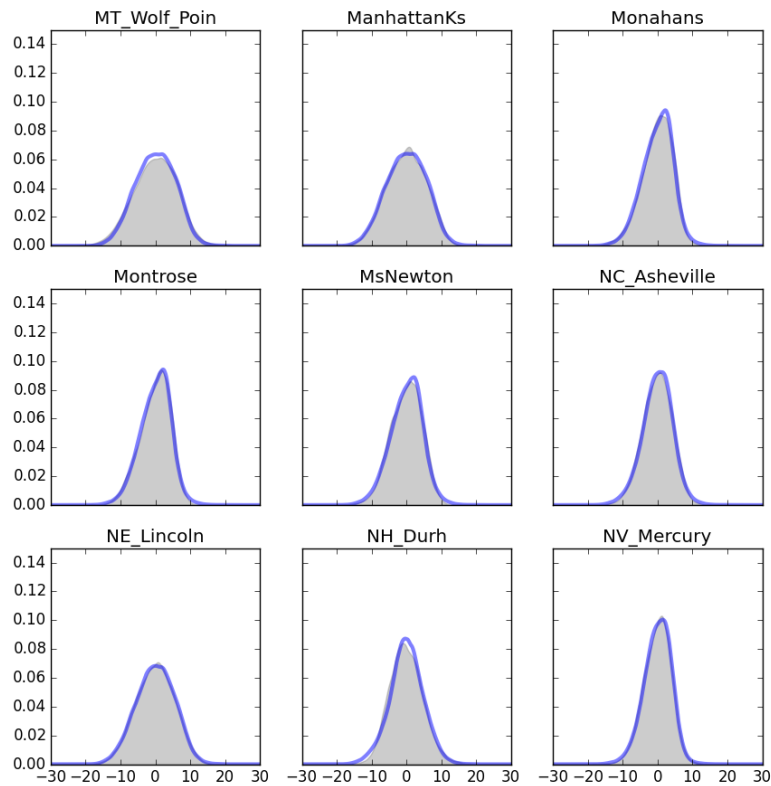


Figure 7: Reconstruction of distribution at each city by 3 archetypes. The original distributions are depicted as shadows; their approximation by mixtures of archetypes as solid curves.

- 235 of Multivariate Analysis 93 (1) (2005) 58–80. doi:10.1016/j.jmva.2003.  
 236 12.002.
- 237 [7] S. T. Rachev, L. Klebanov, S. V. Stoyanov, F. Fabozzi, The methods of  
 238 distances in the theory of probability and statistics, Springer Science &  
 239 Business Media, 2013. doi:10.1007/978-1-4614-4869-3.
- 240 [8] B. K. Sriperumbudur, A. Gretton, K. Fukumizu, B. Schölkopf, G. R. Lanck-  
 241 riet, Hilbert space embeddings and metrics on probability measures, Jour-  
 242 nal of Machine Learning Research 11 (Apr) (2010) 1517–1561.
- 243 [9] S. Kullback, Information theory and statistics, Courier Corporation, 1968.
- 244 [10] C. M. Bishop, Pattern recognition and machine learning, Springer, 2006.



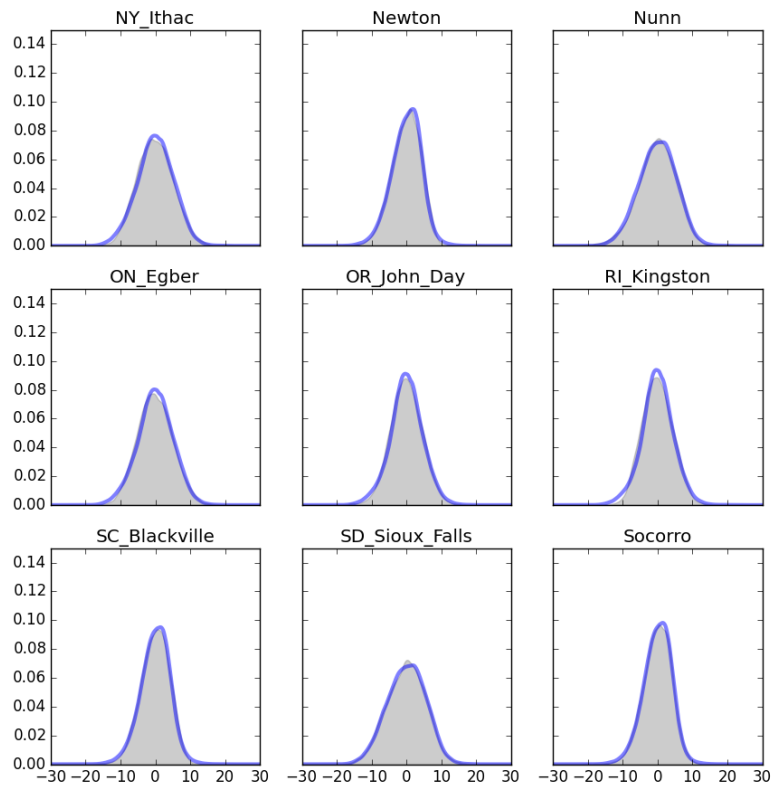


Figure 8: Reconstruction of distribution at each city by 3 archetypes. The original distributions are depicted as shadows; their approximation by mixtures of archetypes as solid curves.

- 245 [11] C. R. Givens, R. M. Shortt, et al., A class of wasserstein metrics for prob-  
 246 ability distributions, *The Michigan Mathematical Journal* 31 (2) (1984)  
 247 231–240. doi:10.1307/mmj/1029003026.
- 248 [12] M. Palecki, J. Lawrimore, R. Leeper, J. Bell, S. Emblar, N. Casey, U.s. cli-  
 249 mate reference network products [hourly], [https://www1.ncdc.noaa.gov/  
 250 pub/data/uscrn/products/hourly02](https://www1.ncdc.noaa.gov/pub/data/uscrn/products/hourly02) (2013). doi:10.7289/V5H13007.
- 251 [13] H. J. Diamond, T. R. Karl, M. A. Palecki, C. B. Baker, J. E. Bell, R. D.  
 252 Leeper, D. R. Easterling, J. H. Lawrimore, T. P. Meyers, M. R. Helfert,  
 253 et al., Us climate reference network after one decade of operations: Status  
 254 and assessment, *Bulletin of the American Meteorological Society* 94 (4)  
 255 (2013) 485–498. doi:10.1175/BAMS-D-12-00170.1.

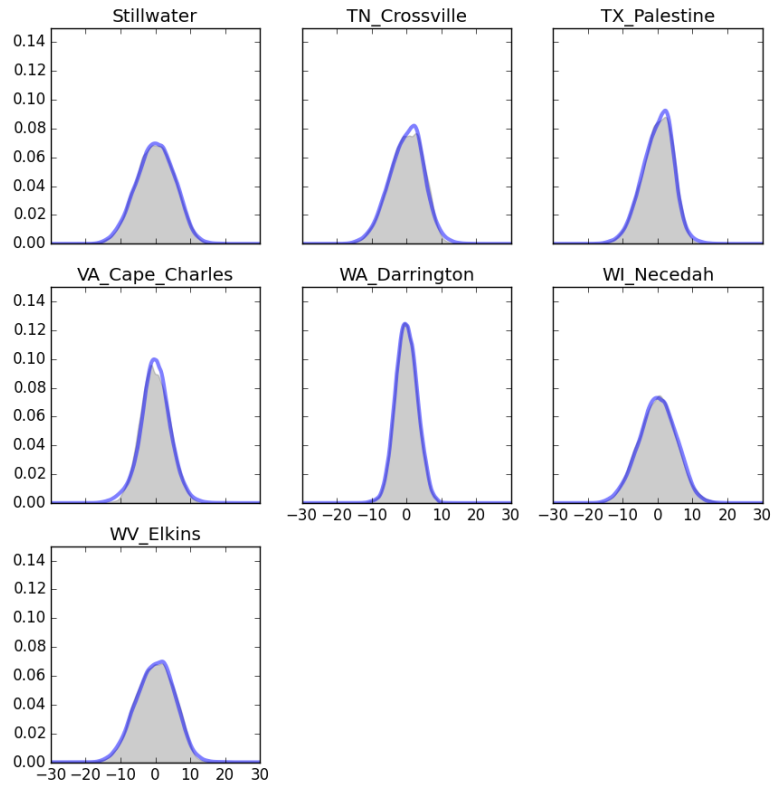


Figure 9: Reconstruction of distribution at each city by 3 archetypes. The original distributions are depicted as shadows; their approximation by mixtures of archetypes as solid curves.

- 256 [14] J. E. Bell, M. A. Palecki, C. B. Baker, W. G. Collins, J. H. Lawrimore,  
 257 R. D. Leeper, M. E. Hall, J. Kochendorfer, T. P. Meyers, T. Wilson,  
 258 et al., Us climate reference network soil moisture and temperature ob-  
 259 servations, *Journal of Hydrometeorology* 14 (3) (2013) 977–988. doi:  
 260 10.1175/JHM-D-12-0146.1.
- 261 [15] E. G. Tabak, G. Trigila, Explanation of variability and removal of confound-  
 262 ing factors from data through optimal transport, accepted for publication  
 263 in CPAM (2017).

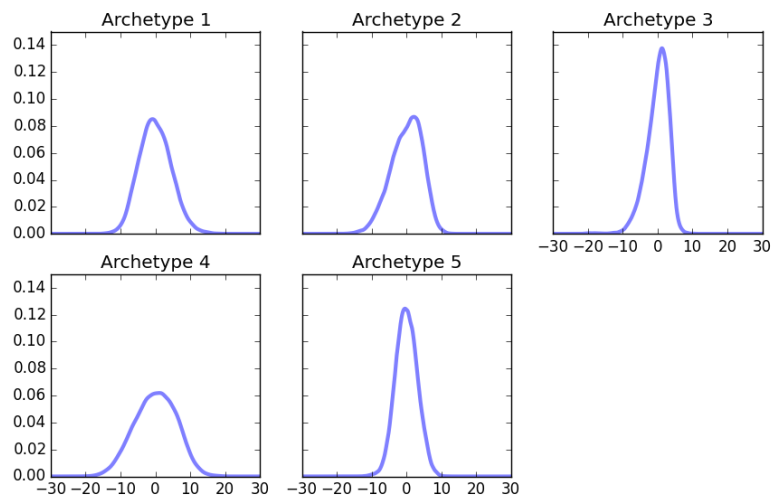


Figure 10: Archetypes of temperature data for  $k=5$ .

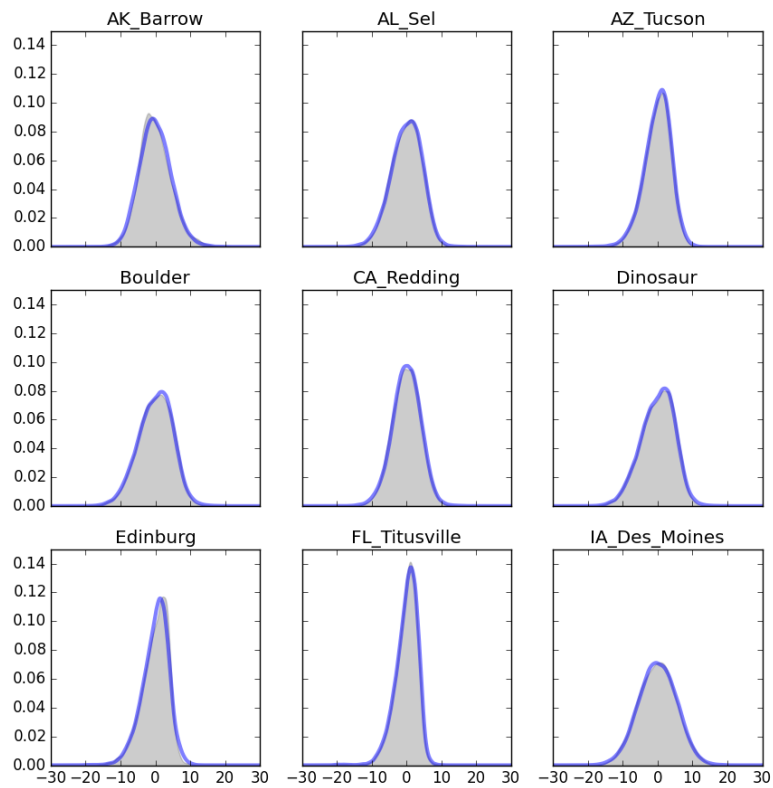


Figure 11: Reconstruction of distribution at each city by 5 archetypes. The original distributions are depicted as shadows; their approximation by mixtures of archetypes as solid curves.

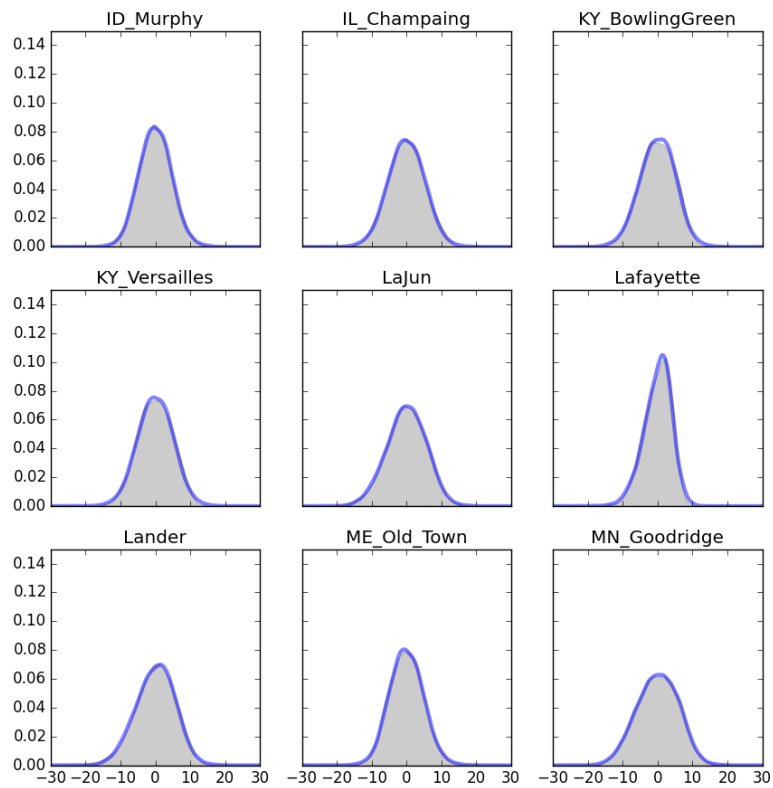


Figure 12: Reconstruction of distribution at each city by 5 archetypes. The original distributions are depicted as shadows; their approximation by mixtures of archetypes as solid curves.

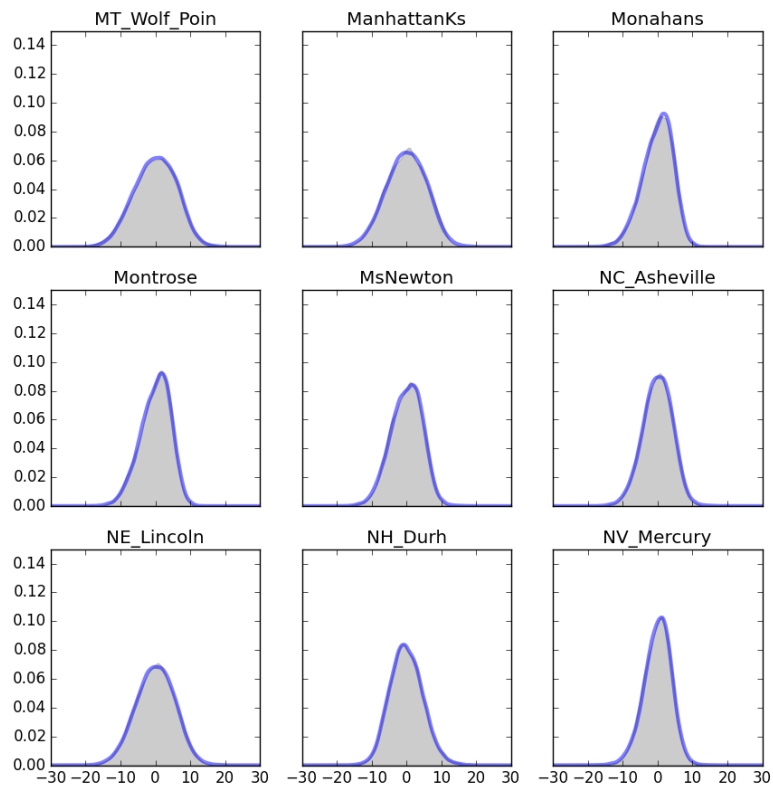


Figure 13: Reconstruction of distribution at each city by 5 archetypes. The original distributions are depicted as shadows; their approximation by mixtures of archetypes as solid curves.

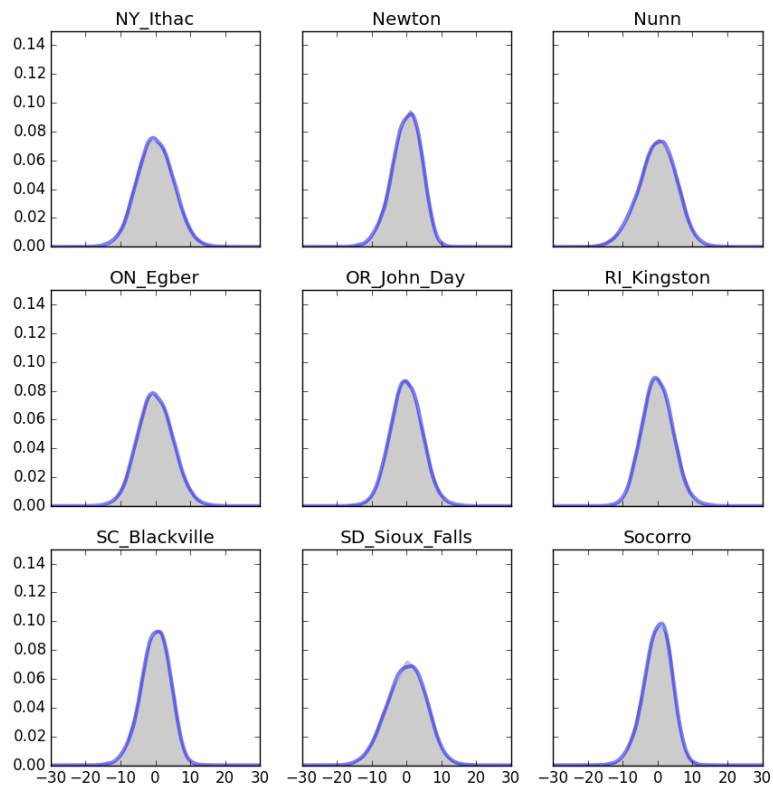


Figure 14: Reconstruction of distribution at each city by 5 archetypes. The original distributions are depicted as shadows; their approximation by mixtures of archetypes as solid curves.

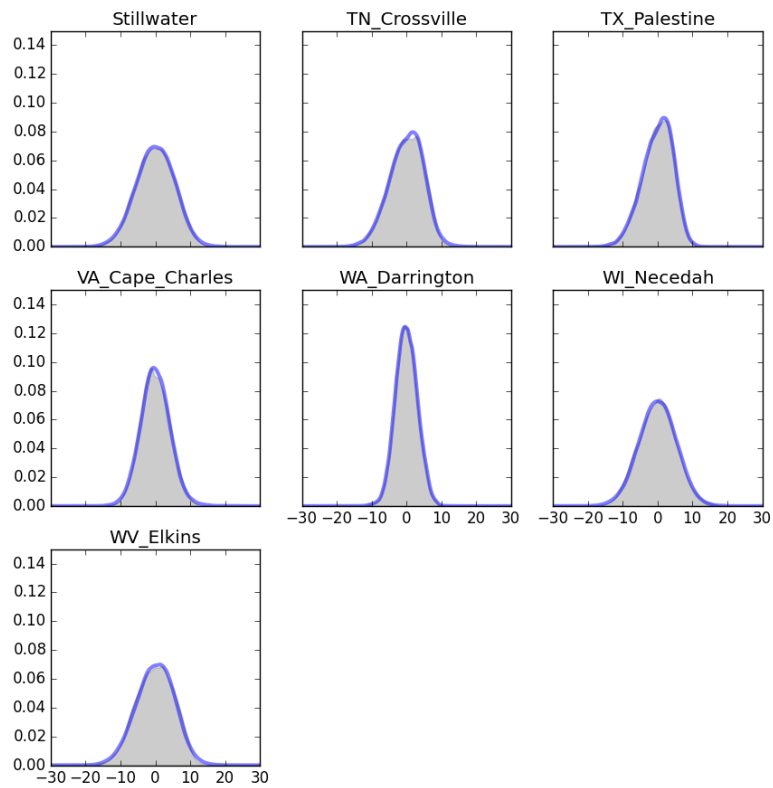


Figure 15: Reconstruction of distribution at each city by 5 archetypes. The original distributions are depicted as shadows; their approximation by mixtures of archetypes as solid curves.