

Lecture 3: Nonlinear Canonical Component Analysis*Lecturer: Pramod Viswanath**Scribe: Bryan Clifford, Sept. 19, 2017*

3.1 Problem Setup

The setup is the same as that of CCA.

- We're given 2 data sets each consisting of N samples

$$\{\mathbf{x}_n\}_{n=1}^N \text{ and } \{\mathbf{y}_n\}_{n=1}^N.$$

- The data are two totally different measurement modalities of the same hidden (latent) variable. For example, \mathbf{x}_n could be a video (vectorized) of a person saying a word and \mathbf{y}_n could be the audio from the video, and the latent variable is the word the person is saying. Hence \mathbf{x}_n and \mathbf{y}_n live in different dimensional spaces:

$$\mathbf{x}_n \in \mathbb{R}^{d_x}, \quad \mathbf{y}_n \in \mathbb{R}^{d_y} \quad \forall n.$$

- We will organize the data into row matrices

$$\begin{aligned} \mathbf{X} &= (\mathbf{x}_1, \dots, \mathbf{x}_N)^T \in \mathbb{R}^{N \times d_x} \\ \mathbf{Y} &= (\mathbf{y}_1, \dots, \mathbf{y}_N)^T \in \mathbb{R}^{N \times d_y} \end{aligned} \tag{3.1}$$

3.1.1 Goal:

Our goal is to find two transforms for \mathbf{x}_n and \mathbf{y}_n that we can use to reduce the dimensionality of each data set to the dimension that we think the latent variable is. However, because both data sets are generated from the same latent variable, we want the transforms to be such that the transformed \mathbf{x} s and \mathbf{y} s are maximally covariant. Furthermore we want the transformed \mathbf{x} s to be uncorrelated from themselves (and the same for the transformed \mathbf{y} s).

This is almost the same goal as for CCA, however in CCA we restricted ourselves to only linear transforms of our variables. In this case we want to remove this restriction. This approach is called **Nonlinear Canonical Component/Correlation Analysis (CCA)**.

More formally we want to find two sets of functions $\{f_m\}_{m=1}^{d_z}$ such that $f_m : \mathbb{R}^{d_x} \rightarrow \mathbb{R}$ and $\{g_m\}_{m=1}^{d_z}$ such that $g : \mathbb{R}^{d_y} \rightarrow \mathbb{R}$ which solve the problem

$$\begin{aligned} & \max_{\substack{\{f_m\}_{m=1}^{d_z} \\ \{g_m\}_{m=1}^{d_z}}} \frac{\mathbb{E}[f_m(\mathbf{x})g_m(\mathbf{y})]}{\sqrt{\text{Var}[f_m(\mathbf{x})]\text{Var}[g_m(\mathbf{y})]}} \\ \text{s.t. } & \mathbb{E}[f_m(\mathbf{x})f_k(\mathbf{x})] = \mathbb{E}[g_m(\mathbf{y})g_k(\mathbf{y})] = 0 \quad \text{if } m \neq k. \end{aligned} \tag{3.2}$$

Note, this is no longer the Pearson correlation. Instead this is called the **maximum correlation** or also the **Hirschfeld-Gebelein-Renyi maximal correlation** and it is always greater than or equal to the Pearson-correlation since

$$\max_{f,g} \frac{\mathbb{E}[f(\mathbf{x})g(\mathbf{y})]}{\sqrt{\text{Var}[f(\mathbf{x})] \text{Var}[g(\mathbf{y})]}} \geq \max_{\boldsymbol{\alpha}, \boldsymbol{\beta}} \frac{\mathbb{E}[(\mathbf{x}^T \boldsymbol{\alpha})(\mathbf{y}^T \boldsymbol{\beta})]}{\sqrt{\text{Var}[\mathbf{x}^T \boldsymbol{\alpha}] \text{Var}[\mathbf{y}^T \boldsymbol{\beta}]}} \quad (3.3)$$

owing to the fact that $(\mathbf{x}^T \boldsymbol{\alpha}, \mathbf{y}^T \boldsymbol{\beta})$ is a special case of $(f(\mathbf{x}), g(\mathbf{y}))$.

The maximum correlation ρ has many important properties. For scalars (i.e. $d_x = d_y = 1$) these properties are

1. $\rho : P_{XY} \rightarrow [0, 1]$, where P_{XY} is the joint probability distribution of the random variables x and y .
2. $\rho = 0$ if and only if x and y are independent.
3. $\rho = 1$ if and only if there exist either or both smooth functions h_1 and h_2 such that $x = h_1(y)$ and $y = h_2(x)$.
4. $\rho(x, y) = \rho(y, x)$
5. If x and y are jointly Gaussian then $\rho(x, y) = \rho_{\text{Pearson}}(x, y)$.

It is interesting to note that these properties are also used to define the maximum correlation. Most of the work goes in to showing the “only if”s and property 5.

3.2 Solution

This section shows a method for solving (3.2). It is very similar to CCA in that it involves a change of variables and the SVD; however as will be seen the solution in this section is not very practical for high-dimensional data (and hence in general, since part of the goal is to be able to reduce dimensionality of high-dimensional data).

We begin by further restricting ourselves to only centered functions with unit variance. This can be done without loss of generality since any uncentered function $f_{\text{uncentered}}$ can be centered by defining a new function $f_{\text{centered}} = f_{\text{uncentered}} - \mathbb{E}[f_{\text{uncentered}}]$. Likewise, the variance can be set to one by dividing by the standard deviation. Our new optimization problem is thus

$$\begin{aligned} & \max_{\substack{\{f_m\}_{m=1}^{d_z} \\ \{g_m\}_{m=1}^{d_z}}} \mathbb{E}[f_m(\mathbf{x})g_m(\mathbf{y})] \\ \text{s.t.} \quad & \mathbb{E}[f_m(\mathbf{x})f_k(\mathbf{x})] = \mathbb{E}[g_m(\mathbf{y})g_k(\mathbf{y})] = \delta_{m,k} \\ & \mathbb{E}[f_m(\mathbf{x})] = \mathbb{E}[g_m(\mathbf{y})] = 0. \end{aligned} \quad (3.4)$$

It is instructive to consider the case when \mathbf{x} and \mathbf{y} can only take on a finite number of values. This case was first proven in using a “geometric and operator theoretic interpretation” in [Wit75]. In

this case, denoting the probability mass function of the joint distribution as $p_{\mathbf{x},\mathbf{y}}$, we have

$$\begin{aligned} \mathbb{E}[f_m(\mathbf{x})g_m(\mathbf{y})] &= \sum_{\mathbf{x},\mathbf{y}} p_{\mathbf{x},\mathbf{y}} f_m(\mathbf{x})g_m(\mathbf{y}) \\ &= \sum_{\mathbf{x}} f_m(\mathbf{x}) \sum_{\mathbf{y}} p_{\mathbf{x},\mathbf{y}} g_m(\mathbf{y}) \\ &= \mathbf{f}_m^T \mathbf{P} \mathbf{g}_m, \end{aligned} \quad (3.5)$$

where \mathbf{f}_m^T , \mathbf{g}_m are the vectors of values of $f_m(\mathbf{x})$ and $g_m(\mathbf{y})$, respectively and \mathbf{P} is the matrix of probability mass values. We also have

$$\begin{aligned} \mathbb{E}[f_m(\mathbf{x})] &= \mathbb{E}[f_m(\mathbf{x})\mathbf{1}] \\ &= \mathbf{f}_m^T \mathbf{P} \mathbf{1}, \end{aligned} \quad (3.6)$$

and

$$\begin{aligned} \mathbb{E}[g_m(\mathbf{y})] &= \mathbb{E}[\mathbf{1}g_m(\mathbf{y})] \\ &= \mathbf{1}^T \mathbf{P} \mathbf{g}_m, \end{aligned} \quad (3.7)$$

where $\mathbf{1}$ is a vector with 1 for every element, so that Eq. (3.4) becomes

$$\begin{aligned} &\max_{\substack{\{\mathbf{f}_m\}_{m=1}^{d_z} \\ \{\mathbf{g}_m\}_{m=1}^{d_z}}} \mathbf{f}_m^T \mathbf{P} \mathbf{g}_m \\ \text{s.t.} \quad &\mathbf{f}_m^T \mathbf{P} \mathbf{1} = \mathbf{1}^T \mathbf{P} \mathbf{g}_m = 0 \\ &\mathbf{f}_m^T \mathbf{P}_x \mathbf{f}_k = \mathbf{g}_m^T \mathbf{P}_y \mathbf{g}_k = \delta_{m,k}. \end{aligned} \quad (3.8)$$

where $\mathbf{P}_x = \text{diag}\{\mathbf{p}_x\}$, $\mathbf{P}_y = \text{diag}\{\mathbf{p}_y\}$, and $\mathbf{p}_x = \mathbf{P}\mathbf{1}$ and $\mathbf{p}_y = \mathbf{P}^T\mathbf{1}$ are the marginal mass function vectors.

To simplify the constrains a bit we define

$$\begin{aligned} \tilde{\mathbf{f}} &= \mathbf{P}_x^{1/2} \mathbf{f} \\ \tilde{\mathbf{g}} &= \mathbf{P}_y^{1/2} \mathbf{g}. \end{aligned} \quad (3.9)$$

which gives us

$$\begin{aligned} &\max_{\substack{\{\tilde{\mathbf{f}}_m\}_{m=1}^{d_z} \\ \{\tilde{\mathbf{g}}_m\}_{m=1}^{d_z}}} \tilde{\mathbf{f}}_m^T \mathbf{P}_x^{-1/2} \mathbf{P} \mathbf{P}_y^{-1/2} \tilde{\mathbf{g}}_m \\ \text{s.t.} \quad &\tilde{\mathbf{f}}_m^T \mathbf{p}_x^{1/2} = \tilde{\mathbf{g}}_m^T \mathbf{p}_y^{1/2} = 0 \\ &\tilde{\mathbf{f}}_m^T \tilde{\mathbf{f}}_k = \tilde{\mathbf{g}}_m^T \tilde{\mathbf{g}}_k = \delta_{m,k}. \end{aligned} \quad (3.10)$$

where we have made use of the fact that $\mathbf{P}_x^{-1/2} \mathbf{P} \mathbf{1} = \mathbf{P}_x^{-1/2} \mathbf{p}_x = \mathbf{p}_x^{1/2}$ (the vector of square-roots of marginal mass function values). Similarly, $\mathbf{1}^T \mathbf{P} \mathbf{P}_y^{-1/2} = (\mathbf{P}_y^{-1/2} \mathbf{p}_y)^T = (\mathbf{p}_y^{1/2})^T$.

Like with CCA we are very close to the SVD problem. The only difference seems to be the extra constraint

$$\tilde{\mathbf{f}}_m^T \mathbf{p}_x^{1/2} = \tilde{\mathbf{g}}_m^T \mathbf{p}_y^{1/2} = 0 \quad (3.11)$$

which says that \mathbf{f}_m and \mathbf{g}_m must be orthogonal to the vectors $\mathbf{p}_x^{1/2}$ and $\mathbf{p}_y^{1/2}$ respectively. However, this is actually not a problem.

Note that if we let $\tilde{\mathbf{f}}_m = \mathbf{p}_x^{1/2}$ and $\tilde{\mathbf{g}}_m = \mathbf{p}_y^{1/2}$ then the objective function equals 1 which is the **maximum** it can be! Hence from this observation we see that the first left and right singular vectors of the matrix

$$\mathbf{Q} := \mathbf{P}_x^{-1/2} \mathbf{P} \mathbf{P}_y^{-1/2} \quad (3.12)$$

are in fact $\tilde{\mathbf{f}}_1 = \mathbf{p}_x^{1/2}$ and $\tilde{\mathbf{g}}_1 = \mathbf{p}_y^{1/2}$ respectively and the corresponding singular value is $\sigma_1 = 1$.

Thus we see that we are really trying to compute the SVD of \mathbf{Q} , but we want to ignore the first singular vectors and focus on the remaining set. Thus **the maximum correlation is given by the second singular value of \mathbf{Q} and the corresponding left and right singular vectors are the optimal values of \mathbf{f} and \mathbf{g} .**

It is worth noting that this can be generalized using tensor products to independent pairs of random variables $((\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2))$ coming from independent distributions $p_{\mathbf{x}_1, \mathbf{y}_1}$ and $p_{\mathbf{x}_2, \mathbf{y}_2}$. Again, this was first shown in [Wit75].

3.2.1 Generalization to continuous variables

The above solution easily generalizes to the case of continuous variables, albeit with a slight caveat. Following the same argument, we see that with continuous variables Eq. (3.5) becomes

$$\mathbb{E}[f_m(\mathbf{x})g_m(\mathbf{y})] = \int \int p_{xy}(\mathbf{x}, \mathbf{y}) f_m(\mathbf{x}) g_m(\mathbf{y}) d\mathbf{x} d\mathbf{y}. \quad (3.13)$$

Now, if we restrict f_m and g_m to be members of a finite dimensional Hilbert space (see [NS82] or ECE 513 notes) then we can write them as a linear combinations of a set of basis functions

$$\begin{aligned} f_m(\mathbf{x}) &= \sum_{i=1}^{L_x} f_{m,i} \phi_i(\mathbf{x}) \\ g_m(\mathbf{y}) &= \sum_{i=1}^{L_y} g_{m,i} \psi_i(\mathbf{y}) \end{aligned} \quad (3.14)$$

An example of this could be the polynomial basis but with a domain restriction on \mathbf{x} and \mathbf{y} (i.e. they both have finite support) in which case the basis functions could be given by (for $d_x = d_y = 1$)

$$\phi_i(x) = x^{i-1}, \quad \psi_i(y) = y^{i-1}. \quad (3.15)$$

With the Hilbert space restriction Eq. (3.13) becomes

$$\begin{aligned} \mathbb{E}[f_m(\mathbf{x})g_m(\mathbf{y})] &= \sum_{i=1}^{L_x} \sum_{j=1}^{L_y} f_{m,i} g_{m,j} \int \int p_{xy}(\mathbf{x}, \mathbf{y}) \phi_i(\mathbf{x}) \psi_j(\mathbf{y}) d\mathbf{x} d\mathbf{y} \\ &= \mathbf{f}_m^T \mathbf{P} \mathbf{g}_m, \end{aligned} \quad (3.16)$$

where now we have

$$P_{i,j} = \int \int p_{xy}(\mathbf{x}, \mathbf{y}) \phi_i(\mathbf{x}) \psi_j(\mathbf{y}) d\mathbf{x} d\mathbf{y}. \quad (3.17)$$

Note that the only difference is in the meaning our matrix and vector elements which are now the basis coefficients and basis cross-correlation values.

Continuing along the same steps and without much loss of generality (just restricting the set of Hilbert spaces we can have) assuming that $\phi_1(\mathbf{x}) = \psi_1(\mathbf{y}) = 1$ we have

$$\begin{aligned} \mathbb{E}[f_m(\mathbf{x})] &= \mathbb{E}[f_m(\mathbf{x})1] \\ &= \mathbf{f}_m^T \mathbf{P} \mathbf{e}_1 \\ &= \mathbf{f}_m^T \bar{\boldsymbol{\phi}}, \end{aligned} \tag{3.18}$$

and

$$\begin{aligned} \mathbb{E}[g_m(\mathbf{y})] &= \mathbb{E}[1g_m(\mathbf{y})] \\ &= \mathbf{e}_1^T \mathbf{P} \mathbf{g}_m \\ &= \bar{\boldsymbol{\psi}}^T \mathbf{g}_m \end{aligned} \tag{3.19}$$

where \mathbf{e}_1 is the standard basis vector with all zero elements except for the first one which is 1, $\bar{\boldsymbol{\phi}} = \mathbf{P} \mathbf{e}_1$, $\bar{\boldsymbol{\psi}} = \mathbf{P}^T \mathbf{e}_1$ are the vectors of expected values of the basis functions. Now replacing \mathbf{P} with $\Sigma_{\phi, \psi} \mathbf{1}$ with \mathbf{e}_1 and \mathbf{p}_x and \mathbf{p}_y with $\bar{\boldsymbol{\phi}}_x$ and $\bar{\boldsymbol{\psi}}_y$ respectively, in the equations for the discrete case, the results are the same.

Instead of requiring that $\phi_1(\mathbf{x}) = \psi_1(\mathbf{y}) = 1$ we could also require that our basis functions be defined such that $\bar{\boldsymbol{\phi}} = \mathbf{0}$ and similarly for $\bar{\boldsymbol{\psi}}$. With this setup, the zero-mean constraint vanishes, and the maximum correlation becomes the first singular value (with respective f and g given by the singular vectors). In this case the N-CCA problem is solved by replacing \mathbf{P}_x and \mathbf{P}_y with the cross-correlation functions of the original basis functions, i.e.

$$\{\mathbf{P}_x\}_{i,j} = \int p_x(\mathbf{x}) \phi_i(\mathbf{x}) \phi_j(\mathbf{x}) d\mathbf{x} \tag{3.20}$$

and respectively for \mathbf{P}_y .

3.2.2 Computational challenges

While the above SVD solutions to non-linear CCA can be rather practical when either the number of discrete values of \mathbf{x} and \mathbf{y} or the dimension of the spaces of f_m and g_m (L_x and L_y respectively) are small. In these cases the SVD can be computed efficiently and, more importantly, it is feasible to estimate the probability mass function or basis cross-correlation values (\mathbf{P}) from the data (via the sample estimates). Unfortunately these methods don't scale well with the size of the problem. The number of terms in \mathbf{P} scales quadratically with the size of the spaces. Hence when \mathbf{x} and \mathbf{y} take on many values or L_x and L_y are large, vast amounts of training data would be needed to obtain reasonable estimates of \mathbf{P} . There are two methods to do this in practice Kernel-CCA (K-CCA) [BJ02] and the more prominent method, the Alternating Conditional Expectation (ACE) algorithm [BF85] which are detailed in the following sections.

3.3 K-CCA

K-CCA was proposed in [BJ02] as a biproduct of the development of a form of independent component analysis that used the maximum correlation as its contrast function. This method is very similar to the continuous case discussed above. The difference is the type of Hilbert space used. In K-CCA the Hilbert space is required to have a very special structure (called a **Reproducing**

Kernel Hilbert Space) (RKHS). Specifically a RKHS is a Hilbert space of functions defined by a Mercer kernel function $K(\mathbf{x}_1, \mathbf{x}_2)$ which is defined such that the ‘‘Gram matrix’’ \mathbf{K} with elements

$$K_{i,j} = K(\mathbf{x}_i, \mathbf{x}_j) \quad (3.21)$$

is semi-positive definite for any set of \mathbf{x} s, i.e. for \mathbf{K} defined from $K(\mathbf{x}_1, \mathbf{x}_2)$ and $\{\mathbf{x}_n\}_{n=1}^N$

$$\boldsymbol{\alpha}^T \mathbf{K} \boldsymbol{\alpha} \geq 0 \quad \forall \boldsymbol{\alpha} \neq \mathbf{0} \in \mathbb{R}^N \quad (3.22)$$

and such that any function $f(\mathbf{x})$ in the RKHS can be written as

$$f(\mathbf{x}) = \sum_{n=1}^N f_n K(\mathbf{x}, \mathbf{x}_n). \quad (3.23)$$

and with an inner product between elements (functions) in the space defined as

$$\langle f_1(\mathbf{x}), f_2(\mathbf{x}) \rangle := \sum_{i,j} f_{1,i} f_{2,j} K(\mathbf{x}_i, \mathbf{x}_j). \quad (3.24)$$

It should be noted that this definition of the RKHS is based on the Moor-Aronszajn theorem. In our previous notation we would say that the set of basis functions $\{\phi_l(\mathbf{x})\}$ is given by $\{K(\mathbf{x}, \mathbf{x}_n)\}_{n=1}^N$, and L_x is given by N .

Given the above definition every function f in an RKHS has the **reproducing property** which is

$$f(\mathbf{x}') = \langle K(\mathbf{x}, \mathbf{x}'), f(\mathbf{x}) \rangle. \quad (3.25)$$

This property leads to an important interpretation of K which is as a ‘‘feature map’’, i.e. a map from the \mathbf{x} space to a function. Denoting the function defined by a particular \mathbf{x}' as $\Phi_{\mathbf{x}'}(\mathbf{x}) = K(\mathbf{x}, \mathbf{x}')$ we have

$$\langle \Phi_{\mathbf{x}'}(\mathbf{x}), \Phi_{\mathbf{x}''}(\mathbf{x}) \rangle = K(\mathbf{x}', \mathbf{x}'') \quad (3.26)$$

In K-CCA the RKHS for f and g are the ones defined by the sets of samples $\{\mathbf{x}_n\}_{n=1}^N$ and $\{\mathbf{y}_n\}_{n=1}^N$. Let the feature maps for the spaces be denoted as $\Phi_{\mathbf{x}'}(\mathbf{x})$ and $\Psi_{\mathbf{y}'}(\mathbf{y})$, and the corresponding kernels by K_x and K_y , respectively. Given these definitions we can see that Eq. (3.27) becomes

$$\begin{aligned} P_{i,j} &= \int \int p_{xy}(\mathbf{x}, \mathbf{y}) \Phi_{\mathbf{x}_i}(\mathbf{x}) \Psi_{\mathbf{y}_j}(\mathbf{y}) d\mathbf{x} d\mathbf{y} \\ &= \int \int p_{xy}(\mathbf{x}, \mathbf{y}) K_x(\mathbf{x}, \mathbf{x}_i) K_y(\mathbf{y}, \mathbf{y}_i) d\mathbf{x} d\mathbf{y}. \end{aligned} \quad (3.27)$$

In K-CCA we then replace the integral by the sample expectation so that

$$P_{i,j} = \frac{1}{N} \sum_{n=1}^N K_x(\mathbf{x}_n, \mathbf{x}_i) K_y(\mathbf{y}_n, \mathbf{y}_i) \quad (3.28)$$

and hence

$$\mathbf{P} = \frac{1}{N} \mathbf{K}_x \mathbf{K}_y. \quad (3.29)$$

By further restricting the kernels such that

$$\sum_{n=1}^N \Phi_{\mathbf{x}_n}(\mathbf{x}) = 0 \quad (3.30)$$

and respectively for $\{\Psi_{\mathbf{y}_n}\}_{n=1}^N$ then when we replace the population expectations with sample expectations we have the same problem as discussed in the previous section.

Equation (3.29) is particularly important since it shows that \mathbf{P} will be written as the product of two semi-positive definite matrices (and similarly for the \mathbf{P}_x and \mathbf{P}_y). This fact enables the K-CCA algorithm to be computed efficiently using incomplete Cholesky decompositions, which leads to linear scaling with the number of data points [BJ02].

3.4 ACE

The ACE algorithm is a very widely used algorithm that enables the maximum correlation and respective mappings f and g to be found through a simple alternating procedure which under certain conditions is proven to converge globally [BF85].

The method can be derived as follows. First note that

$$\mathbb{E}[(f(\mathbf{x}) - g(\mathbf{y}))^2] = \mathbb{E}[f^2(\mathbf{x})] + \mathbb{E}[g^2(\mathbf{y})] - 2\mathbb{E}[f(\mathbf{x})g(\mathbf{y})]. \quad (3.31)$$

Thus if we include the mean-zero and unit variance constraints on f and g we have

$$\mathbb{E}[(f(\mathbf{x}) - g(\mathbf{y}))^2] = 2(1 - \mathbb{E}[f(\mathbf{x})g(\mathbf{y})]). \quad (3.32)$$

we will be able to solve the N-CCA problem (for the first pair of functions) with the solution to

$$\begin{aligned} & \min_{f,g} \mathbb{E}[(f(\mathbf{x}) - g(\mathbf{y}))^2] \\ & \text{s.t.} \quad \mathbb{E}[f(\mathbf{x})] = \mathbb{E}[g(\mathbf{y})] = 0 \\ & \quad \mathbb{E}[f^2(\mathbf{x})] = \mathbb{E}[g^2(\mathbf{y})] = 1. \end{aligned} \quad (3.33)$$

The ACE algorithm does this alternatively by iterating between the follow problems

$$\begin{aligned} 1. \quad & f_k(\mathbf{x}) = \arg \min_f \mathbb{E}[(f(\mathbf{x}) - g_{k-1}(\mathbf{y}))^2] \quad \text{s.t.} \quad \mathbb{E}[f^2(\mathbf{x})] = 1 \\ & = \frac{\mathbb{E}[g_{k-1}(\mathbf{y}) \mid \mathbf{x}]}{\sqrt{\text{Var}[\mathbb{E}[g_{k-1}(\mathbf{y}) \mid \mathbf{x}]]}} \end{aligned} \quad (3.34)$$

$$\begin{aligned} 2. \quad & g_k(\mathbf{y}) = \arg \min_g \mathbb{E}[(f_k(\mathbf{x}) - g(\mathbf{y}))^2] \quad \text{s.t.} \quad \mathbb{E}[g^2(\mathbf{y})] = 1 \\ & = \frac{\mathbb{E}[f_k(\mathbf{x}) \mid \mathbf{y}]}{\sqrt{\text{Var}[\mathbb{E}[f_k(\mathbf{x}) \mid \mathbf{y}]]}} \end{aligned} \quad (3.35)$$

where k is the iteration number. The solutions to the above optimization problems are proven as follows. In practice, dividing by the variance in the first problem can be neglected. The algorithm can be shown to converge to an optimal solution [BF85].

Following [BF85], it can be shown that the space of functions of random variables with the inner product between $f_1(\mathbf{x})$ and $f_2(\mathbf{x})$ given by $E[f_1(\mathbf{x})f_2(\mathbf{x})]$ and the induced norm is a Hilbert space and that the conditional expectation operator is the projection operator from the space of functions of random variables \mathbf{x} to the space of functions of random variables \mathbf{y} and hence

$$f^*(\mathbf{x}) = \arg \min_f E[(f(\mathbf{x}) - g(\mathbf{y}))^2] = E[g(\mathbf{y}) | \mathbf{x}]. \quad (3.36)$$

This can also be shown via the orthogonality principle.

Then since we have $E[g(\mathbf{y})] = 0$ by definition and by the tower property $E[E[g(\mathbf{y}) | \mathbf{x}]] = E[g(\mathbf{y})]$. Hence we automatically satisfy the mean-zero constraint on $f^*(\mathbf{x})$.

Now the set of functions of \mathbf{y} with unit variance is a subset of the functions of \mathbf{y} . By the tower property if we project a projection onto a subset of the first projections space its the same as projecting onto the subset. It can be shown that the projection of a vector \mathbf{f} in a Hilbert space onto the unit ball (subset of unit norm vectors) is given by dividing the vector by $\max\{1, \|\mathbf{f}\|\}$. Hence the solution to the optimization problem with the norm constraint is given by the conditional expectation divided by it's variance.

In practice, the ACE algorithm is applied to data samples and the population expectations are substituted for their sample equivalents. However, computing the conditional expectations from the data can be difficult. Consider the case of continuous variables, in which every value may be unique. Hence in the case of continuous variables the conditional expectation is often replaced with a **data smooth** which could be a regression of the data or a normalized histogram of bin frequencies [BF85]. An essential property of the smoother is its having a zero-mean (which can be done by subtracting the mean of the uncentered version).

One of the main problems with the ACE algorithm is that the performance often depends strongly on the choice of data smoother.

3.5 Example

It is instructive to consider a simple 1D example in which x and y are generated from a latent variable $z \sim \mathcal{N}(0, 1)$ with additive noise for $\epsilon_x \sim \mathcal{N}(0, \sigma)$ and $\epsilon_y \sim \mathcal{N}(0, \sigma)$ as

$$\begin{aligned} x &= w_x e^{-z} + \mu_x + \epsilon_x \\ y &= w_y e^{-z^2} + \mu_y + \epsilon_y \end{aligned} \quad (3.37)$$

Let us use as a smoother polynomial regression, i.e.

$$\begin{aligned} E[g(\mathbf{y}) | x] &\approx \sum_{m=0}^M a_{f,m} x^m \\ E[f(\mathbf{x}) | y] &\approx \sum_{m=0}^M a_{g,m} y^m \end{aligned} \quad (3.38)$$

where the least-squares regression solutions are given by

$$\hat{\mathbf{a}}_f = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{g} \quad (3.39)$$

where $\mathbf{X} = (\mathbf{1}, \mathbf{x}, \dots, \mathbf{x}^m)$, and $\hat{\mathbf{a}}_f$ is the vector of estimated coefficients of $f(x)$, \mathbf{x} is the vector of x samples, and \mathbf{g} is the vector of $g(y)$ values. The solution is similar for \mathbf{a}_g and b_g .

We initialize the coefficients to $f(x) = (x - E[()x]) / \text{Var}[()x]$ and similarly for $g(y)$ which is the optimal linear estimator. Then we estimate the coefficients of each estimator and normalize and center them repeatedly as per the ACE algorithm. The following MATLAB code simulates this for $M = 8$ and Fig. 3.1 shows the convergence of the estimated maximum correlation as well as the transformed variables. Note that since both x and y are generated from the same latent variable and an inverse transform exists to map x to z , the maximum correlation should be 1. We won't actually see this because our choice of data smoother precludes this inverse transform from the set of transforms we are projecting onto.

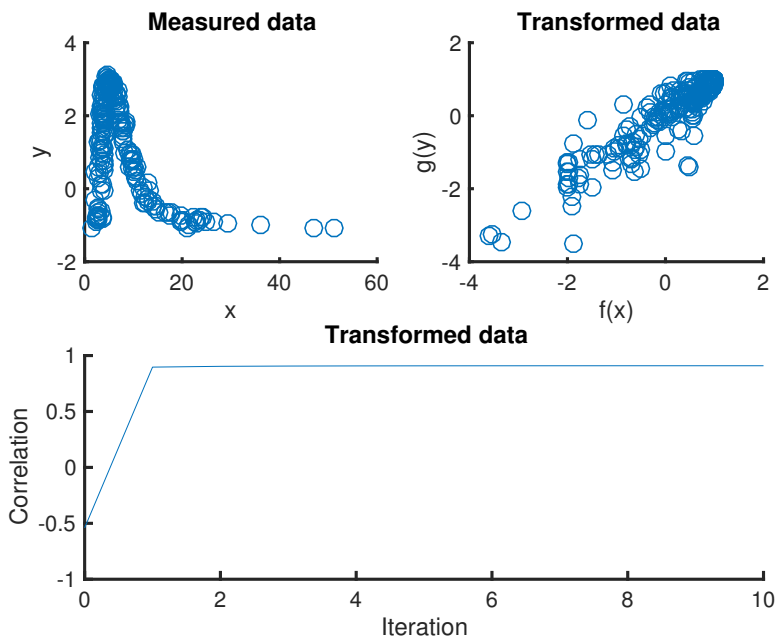


Figure 3.1: Simulation results for order 8 polynomial regression: The top left shows the original measured data and the top right shows the transformed data with the estimated optimal $f(x)$ and $g(y)$. The maximum correlation is shown on the bottom and converges very quickly.

```
%% ——— Perform ACE on random 1d example ———
```

```
N = 200; % number of samples
```

```
sigma_x = .5;
```

```
sigma_y = .1;
```

```
% generate random parameters and latent variable samples
```

```
z = randn(N,1);
```

```
% generate data
```

```
x = 3*exp(-z) + 2 + sigma_x*randn(N,1);
```

```
y = 4*exp(-z.^2) - 1 + sigma_y*randn(N,1);
```

```
% matrices for regression
```

```
M = 8; % polynomial order
```

```
X = ones(N,1);
```

```
Y = ones(N,1);
```

```
for m = 1:M
```

```
    X = [X,x.^m];
```

```
    Y = [Y,y.^m];
```

```
end
```

```
% perform the ACE algorithm
```

```
Nitr = 10; % number of iterations
```

```
f = (x - mean(x,1))/std(x);
```

```
g = (y - mean(y,1))/std(y);
```

```
a_f_vals = [];
```

```
b_f_vals = [];
```

```
a_g_vals = [];
```

```
b_g_vals = [];
```

```
cor_vals = mean(f.*g);
```

```
for k = 1:Nitr
```

```
    % estimate E[g(y)|x]
```

```
    coefs = (X'*X)\(X'*g);
```

```
    f = X*coefs;
```

```
    f = f - mean(f);
```

```
    f = f./std(f);
```

```
    % estimate E[f(x)|y]
```

```
    coefs = (Y'*Y)\(Y'*f);
```

```
    g = Y*coefs;
```

```
    g = g - mean(g);
```

```
    g = g./std(g);
```

```
    % compute correlation
```

```
    cor_vals(end+1) = mean(f.*g);
```

```
end
```

```
% display results
```

```
figure;
```

```
% —— measured ——
```

```
subplot(2,2,1);
```

```
scatter(x(:,1),y(:,1),128);
```

```
title('Measured_data');
```

```
xlabel('x');
```

```
ylabel('y');
```

```
% —— transformed ——
```

```
subplot(2,2,2);
```

```
scatter(f,g,128);
title('Transformed_data');
xlabel('f(x)');
ylabel('g(y)');
% —— convergence ——
subplot(2,1,2);
hold all;
plot(0:Nitr, cor_vals);
title('Transformed_data');
ylabel('Correlation');
xlabel('Iteration');
```

Bibliography

- [BF85] L. Breiman and J. Friedman. Estimating optimal transformations for multiple regression and correlation. *J. Am. Stat. Assoc.*, 80:580 – 598, 1985.
- [BJ02] F. Bach and M. Jordan. Kernel independent component analysis. *J. Machine Learn. Res.*, 3:1 – 48, 2002.
- [NS82] A. W. Naylor and G. R. Sell. *Linear Operator Theory in Engineering and Science*. Springer-Verlag, 2 edition, 1982.
- [Wit75] H. Witsenhausen. On sequences of pairs of dependent random variables. *SIAM J. Appl. Math.*, 28:100 – 113, 1975.