

## Correlation analysis 3: Measures of correlation (continued)

Ryan Tibshirani  
Data Mining: 36-462/36-662

February 21 2013

## Reminder: correlation, rank correlation

Last time we learned about **correlation**. In the **population**: for random variables  $X, Y \in \mathbb{R}$ ,

$$\text{Cor}(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)}\sqrt{\text{Var}(Y)}}$$

In the **sample**: for vectors  $x, y \in \mathbb{R}^n$ ,

$$\text{cor}(x, y) = \frac{\text{cov}(x, y)}{\sqrt{\text{var}(x)}\sqrt{\text{var}(y)}} = x^T y$$

the second equality holding if  $x, y$  have been centered and scaled

**Rank correlation** is only defined in the **sample**: given  $x, y \in \mathbb{R}^n$ , let  $r_x, r_y \in \mathbb{R}^n$  denote the ranks of  $x, y$ , respectively, and then

$$\text{rcor}(x, y) = \text{cor}(r_x, r_y)$$

## Reminder: maximal correlation

**Maximal correlation** is only defined in the **population**: for random variables  $X, Y \in \mathbb{R}$ ,

$$\text{mCor}(X, Y) = \max_{f, g} \text{Cor}(f(X), g(Y))$$

the maximum taken over all functions  $f, g$ . We were able to show that the optimal  $f, g$  are also optimal for the problem

$$\min_{\substack{E[f(X)] = E[g(Y)] = 0 \\ \|f(X)\| = \|g(Y)\| = 1}} E[(f(X) - g(Y))^2]$$

where  $\|Z\| = \sqrt{E[Z^2]}$ . The fixed points of this problem,

$$\begin{aligned} g(Y) &= E[f(X)|Y] / \|E[f(X)|Y]\| \\ f(X) &= E[g(Y)|X] / \|E[g(Y)|X]\| \end{aligned}$$

suggested an **alternating algorithm** for finding  $\text{mCor}$

## Reminder: alternating conditional expectations algorithm

The **alternating conditional expectations** (ACE) algorithm can be used to compute  $\text{mCor}$ . Algorithm:

- ▶ Set  $f_0(X) = (X - \mathbb{E}[X])/\|X - \mathbb{E}[X]\|$
- ▶ For  $k = 1, 2, 3 \dots$ 
  1. Let  $g_k(Y) = \mathbb{E}[f_{k-1}(X)|Y]/\|\mathbb{E}[f_{k-1}(X)|Y]\|$
  2. Let  $f_k(X) = \mathbb{E}[g_k(Y)|X]/\|\mathbb{E}[g_k(Y)|X]\|$
  3. Stop if  $\mathbb{E}[f_k(X)g_k(Y)] = \mathbb{E}[f_{k-1}(X)g_{k-1}(Y)]$
- ▶ Upon convergence,  $\text{mCor}(X, Y) = \mathbb{E}[f_k(X)g_k(Y)]$

Maximal correlation isn't well-defined in the **sample**, because for any vectors  $x, y \in \mathbb{R}^n$ ,

$$\max_{f,g} \text{cor}(f(x), g(y)) = 1$$

But the ACE algorithm can be adapted to the sample and we **define its output** to be the sample maximal correlation  $\text{mcor}$

## Reminder: ACE algorithm in the sample

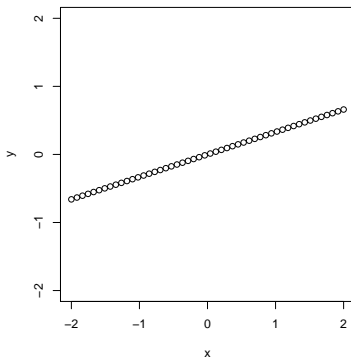
The **sample** version of the ACE algorithm requires us to pick a **smoother**  $\mathcal{S}$  to approximate the conditional expectations (e.g., kernel regression or local linear regression). Algorithm:

- ▶ Set  $f_0(x) = (x - \bar{x}\mathbf{1})/\|x - \bar{x}\mathbf{1}\|_2$
- ▶ For  $k = 1, 2, 3, \dots$ 
  1. Let  $G(y) = \mathcal{S}(f_{k-1}(x)|y)$ , and center and scale,  
 $g_k(y) = (G(y) - \overline{G(y)}\mathbf{1})/\|G(y) - \overline{G(y)}\mathbf{1}\|_2$
  2. Let  $F(x) = \mathcal{S}(g_k(y)|x)$ , and center and scale,  
 $f_k(x) = (F(x) - \overline{F(x)}\mathbf{1})/\|F(x) - \overline{F(x)}\mathbf{1}\|_2$
  3. Stop if  $|f_k(x)^T g_k(y) - f_{k-1}(x)^T g_{k-1}(y)|$  is small
- ▶ Upon convergence, define  $\text{mcor}(x, y) = f_k(x)^T g_k(y)$

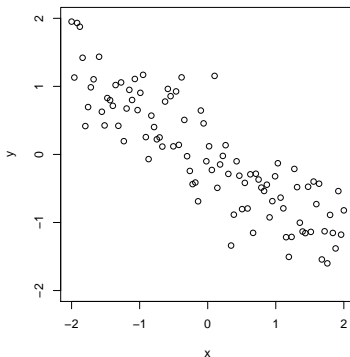
This not always guaranteed to converge, and the answer depends on the smoother. But most of the time it **works well** in practice

## Example: maximal correlation

Perfect linear



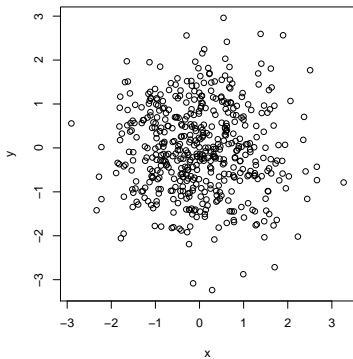
Noisy linear



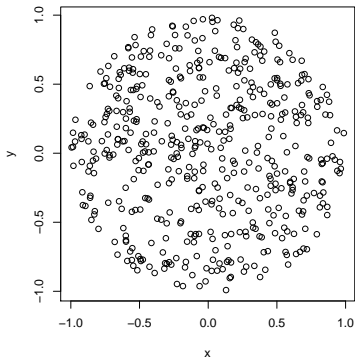
mcor =	1.000
rcor =	1.000
cor =	1.000

0.896
-0.872
-0.866

Independent



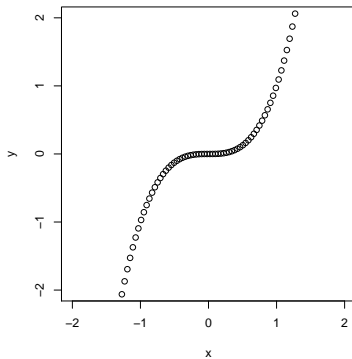
Ball



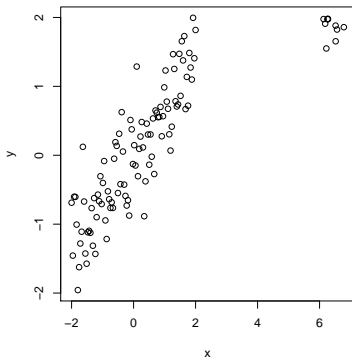
mcor = 0.124  
rcor = -0.021  
cor = -0.023

0.316  
-0.033  
-0.029

Perfect cubic



Outliers

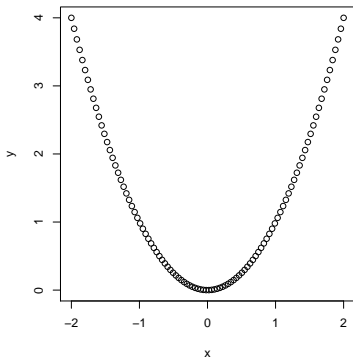


$mcor = 1.000$   
 $rcor = 1.000$   
 $cor = 0.920$

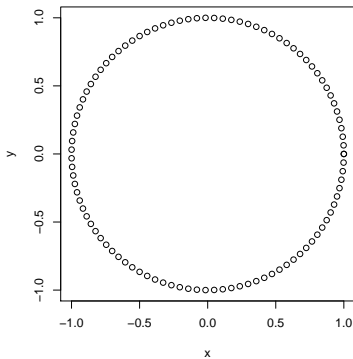
$0.913$   
 $0.905$   
 $0.834$



Perfect quadratic



Perfect circle

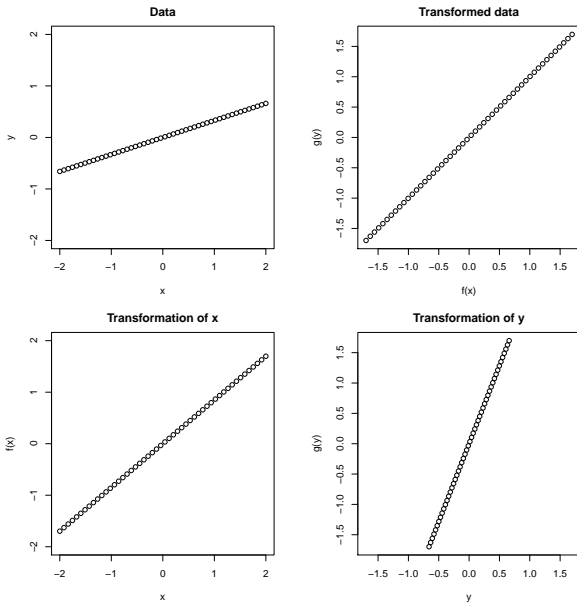


mcor =	1.000
rcor =	0.013
cor =	0.000

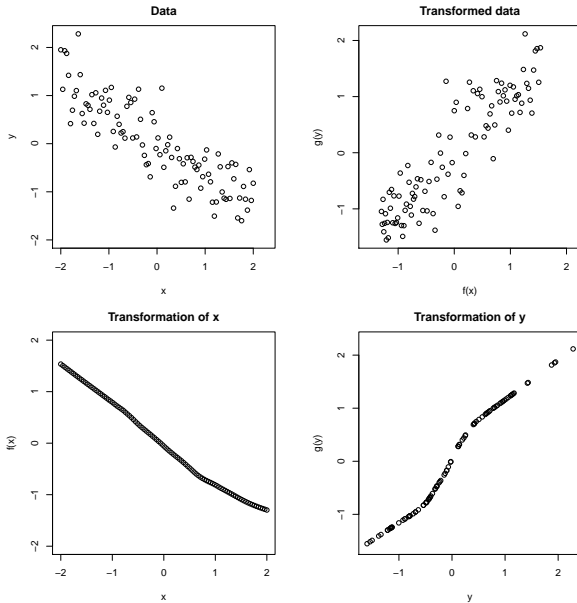
1.000
-0.001
0.000

# Example: ACE algorithm

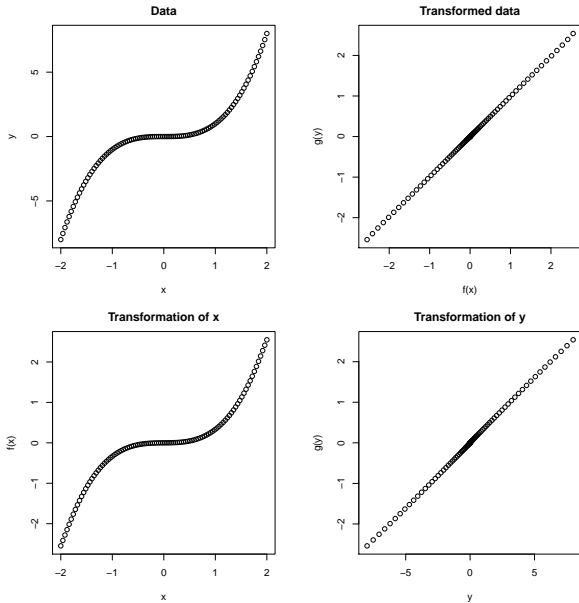
Perfect linear



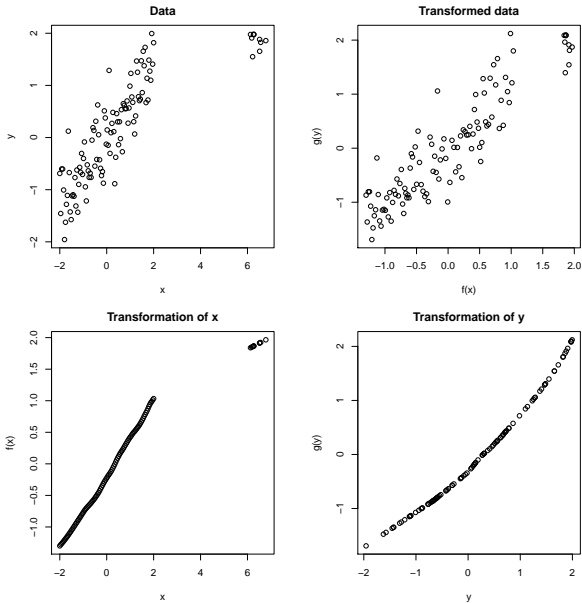
# Noisy linear



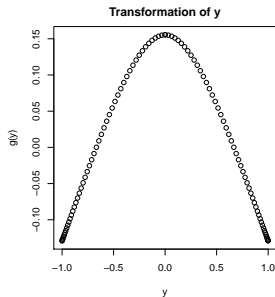
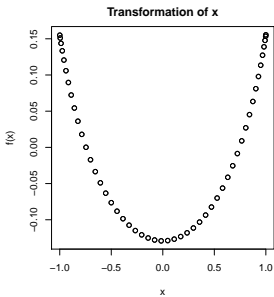
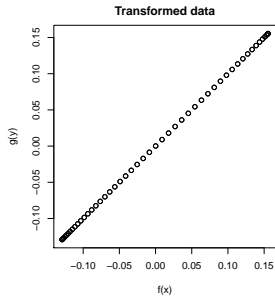
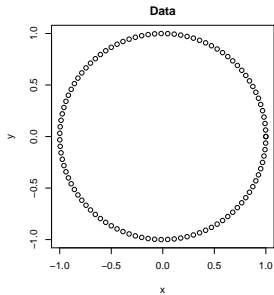
# Perfect cubic



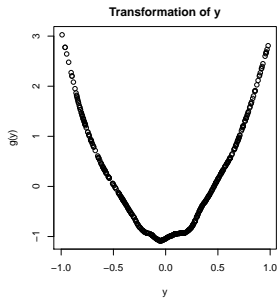
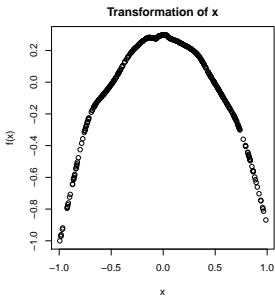
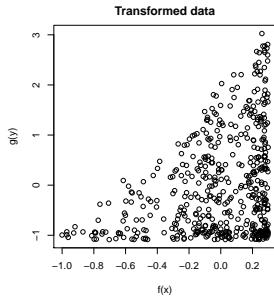
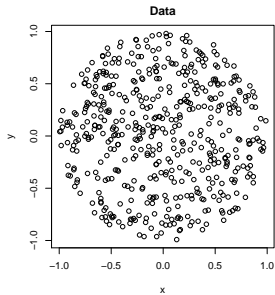
# Outliers



# Perfect circle



# Ball



## Problems with sample maximal correlation and ACE

As mentioned before, a big problem with maximal correlation is that it is **not well-defined** in the sample

The sample ACE algorithm is a nice tool, but defining the sample maximal correlation in terms of its output is not ideal, because the answer depends on the **choice of smoother**. It's also hard to compute, compared to the usual correlation (Homework 3)

Aside from these problems, maximal correlation and the ACE algorithm also have some undesirable properties. E.g., there are cases in which  $\text{mCor}(X, Y) = 1$  for random variables  $X, Y$  that **don't exhibit perfect determinism**, and similarly, cases in which the sample ACE algorithm returns  $\approx 1$  for vectors  $x, y$  whose components are not completely determined from one another (Homework 3)



## Distance correlation

**Distance correlation**<sup>1</sup> is a recent notion of correlation that also characterizes independence completely. It is well-defined in both the population and in the sample, and is easy to compute

In the **population**: given random variables  $X, Y \in \mathbb{R}$ , let  $X', Y'$  and  $X'', Y''$  be independent pairs of random variables taken from the same joint distribution as that of  $X, Y$ . The **distance covariance** of  $X, Y$  is then defined as the square root of

$$\begin{aligned} \text{dCov}^2(X, Y) &= E[|X - X'| |Y - Y'|] + E[|X - X'|]E[|Y - Y'|] \\ &\quad - E[|X - X'| |Y - Y''] - E[|X - X''| |Y - Y'|] \\ &= E[|X - X'| |Y - Y'|] + E[|X - X'|]E[|Y - Y'|] \\ &\quad - 2E[|X - X'| |Y - Y''] \end{aligned}$$

---

<sup>1</sup>Szekely et al. (2007), “Measuring and Testing Dependence By Correlation of Distances”; Szekely et al. (2009) “Brownian Distance Covariance”

Now we can define a notion of distance variance and distance correlation in an analogous manner to their usual definitions. I.e., the **distance variance** of  $X$  is defined by

$$\text{dVar}^2(X) = \text{dCov}^2(X, X)$$

and the **distance correlation** of  $X, Y$  is defined by

$$\text{dCor}^2(X, Y) = \frac{\text{dCov}^2(X, Y)}{\sqrt{\text{dVar}^2(X)}\sqrt{\text{dVar}^2(Y)}}$$

**Properties:**

- ▶  $\text{dCor}(aX + b, Y) = \text{dCor}(X, Y)$  for any  $a, b \in \mathbb{R}, a \neq 0$
- ▶  $0 \leq \text{dCor}(X, Y) \leq 1$
- ▶  $\text{dCor}(X, Y) = 0$  if and only if  $X, Y$  are independent

(There are many more)

## Distance correlation in the sample

Now in the **sample**: given vectors  $x, y \in \mathbb{R}^n$ , we first define the distance matrices  $A, B \in \mathbb{R}^{n \times n}$  by

$$A_{ij} = |x_i - x_j| \quad \text{and} \quad B_{ij} = |y_i - y_j|$$

for all  $i, j = 1, \dots, n$ . Then we **double center**  $A, B$  to get  $\tilde{A}, \tilde{B}$ , respectively, i.e., we center both the rows and columns of each matrix. Recall that, letting  $M = \mathbb{1}\mathbb{1}^T/n$ , this is

$$\tilde{A} = (I - M)A(I - M) \quad \text{and} \quad \tilde{B} = (I - M)B(I - M)$$

The **distance covariance** of  $x, y$  is then defined as the square root of

$$\text{dcov}^2(x, y) = \frac{1}{n^2} \sum_{i,j=1}^n \tilde{A}_{ij} \tilde{B}_{ij}$$

Now we follow the same steps as in the population, first defining the **distance variance** or  $x$  by

$$\text{dvar}^2(x) = \text{dcov}^2(x, x)$$

and then defining the **distance correlation** of  $x, y$  by

$$\text{dcor}^2(x, y) = \frac{\text{dcov}^2(x, y)}{\sqrt{\text{dvar}^2(x)}\sqrt{\text{dvar}^2(y)}}$$

**Properties:**

- ▶  $\text{dcor}(ax + b, y) = \text{dcor}(x, y)$  for any  $a, b \in \mathbb{R}, a \neq 0$
- ▶  $0 \leq \text{dcor}(x, y) \leq 1$
- ▶  $\text{dcor}(x, y) = 1$  if and only if  $y = ax + b$  for some  $a, b \in \mathbb{R}, a \neq 0$

(There are many more)

## Why this sample definition?

What's the **connection** between this sample definition and the population definition?

Let  $A, B$  contain the pairwise **absolute differences** of  $x, y$ , and let  $\tilde{A}, \tilde{B}$  be their double centered versions. It turns out (Homework 3, bonus):

$$\begin{aligned} \text{dcov}^2(x, y) &= \frac{1}{n^2} \sum_{i,j=1}^n \tilde{A}_{ij} \tilde{B}_{ij} \\ &= \frac{1}{n^2} \sum_{i,j=1}^n A_{ij} B_{ij} - \frac{1}{n} \sum_{j=1}^n A_{\cdot j} B_{\cdot j} - \frac{1}{n} \sum_{i=1}^n A_{i \cdot} B_{i \cdot} + A_{\cdot \cdot} B_{\cdot \cdot} \end{aligned}$$

(where  $\cdot$  means sum over that component). Now compare the **population** version:

$$\begin{aligned} \text{dCov}^2(X, Y) &= \text{E}[|X - X'| | Y - Y'|] + \text{E}[|X - X'|] \text{E}[|Y - Y'|] \\ &\quad - \text{E}[|X - X'| | Y - Y''] - \text{E}[|X - X''| | Y - Y'|] \end{aligned}$$

Do these terms match (which ones)?

## Distance correlation and characteristic functions

There is an alternate definition in terms of **characteristic functions** (this was actually the original motivation)

In the **population**: given random variables  $X, Y \in \mathbb{R}$ , let  $h_X(t) = \mathbb{E}[\exp(itX)]$ ,  $h_Y(t) = \mathbb{E}[\exp(itY)]$  denote their characteristic functions, and let  $h_{X,Y}(s, t) = \mathbb{E}[\exp(i(sX + tY))]$  denote their joint characteristic function. Then

$$\text{dCov}(X, Y) = \|h_{X,Y} - h_X h_Y\|$$

where the above  $\|\cdot\|$  is a specially defined **norm on functions** (actually a double integral)

From this, one can show that  $\text{dCov}(X, Y) = \|h_{X,Y} - h_X h_Y\| = 0$  if and only if  $h_{X,Y}(s, t) = h_X(s)h_Y(t)$  for all  $s, t \in \mathbb{R}$ , which is true if and only if  $X, Y$  are independent

In the **sample**: given vectors  $x, y \in \mathbb{R}^n$ , define the marginal and joint **empirical characteristic functions**,

$$h_x(t) = \frac{1}{n} \sum_{k=1}^n \exp(itx_k) \quad \text{and} \quad h_y(t) = \frac{1}{n} \sum_{k=1}^n \exp(ity_k)$$
$$h_{x,y}(s, t) = \frac{1}{n} \sum_{k=1}^n \exp(i(sx_k + ty_k))$$

Then it turns out that

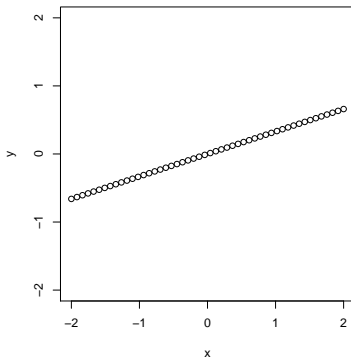
$$\text{dcov}(x, y) = \|h_{x,y} - h_x h_y\|,$$

where  $\|\cdot\|$  is the same special norm on functions as in the previous slide

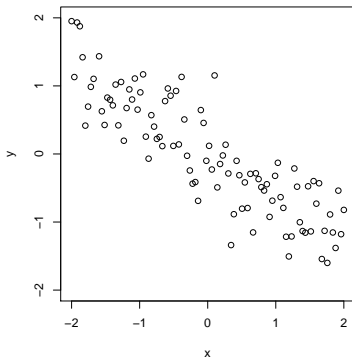
(It's a lot easier to compute  $\text{dcor}$  using the first definition!)

## Example: distance correlation

Perfect linear



Noisy linear

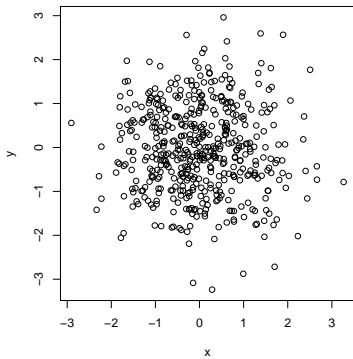


dcor =	1.000
mcor =	1.000
rcor =	1.000
cor =	1.000

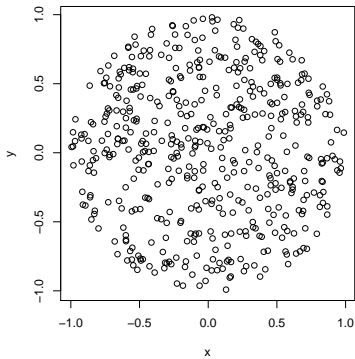
dcor =	0.867
mcor =	0.896
rcor =	-0.872
cor =	-0.866



# Independent



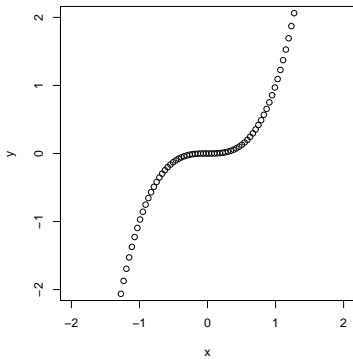
# Ball



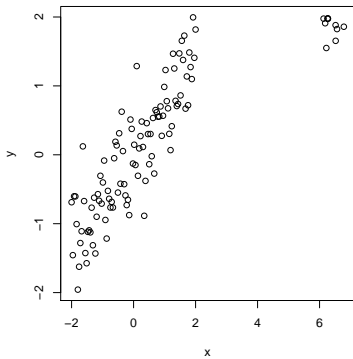
dcor =	0.078
mcor =	0.124
rcor =	-0.021
cor =	-0.023

0.099
0.316
-0.033
-0.029

Perfect cubic



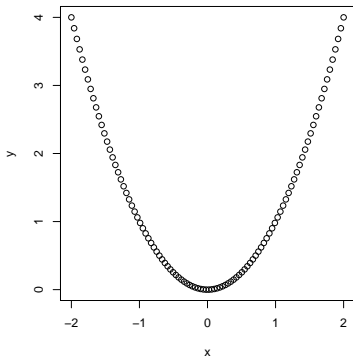
Outliers



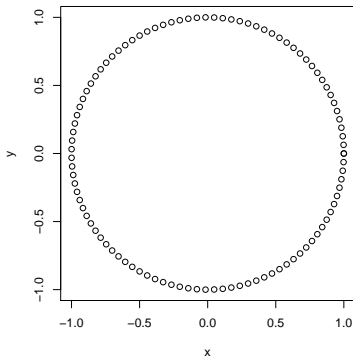
dcor = 0.920  
 mcor = 1.000  
 rcor = 1.000  
 cor = 0.920

0.854  
 0.913  
 0.905  
 0.834

Perfect quadratic



Perfect circle



dcor = 0.492  
 mcor = 1.000  
 rcor = 0.013  
 cor = 0.000

0.200  
 1.000  
 -0.001  
 0.000

## Distance correlation in R

The `dcor` function in the `energy` package can be used to compute distance correlation. E.g.,

```
dcor(x, y)
```

## Maximal information coefficient

Reshef et al. (2011), "Detecting Novel Associations in Large Data Sets"

## RESEARCH ARTICLES

### Detecting Novel Associations in Large Data Sets

David N. Beebe,<sup>1,2,3\*</sup> Yahir A. Beebe,<sup>1,2,3</sup> Hilary E. Grogan,<sup>1</sup> Sharon R. Grogan,<sup>2</sup>  
Glaan McVean,<sup>1,2</sup> Peter J. Tambough,<sup>2</sup> Eric J. Lander,<sup>1,2,3</sup>

Identifying interrelated relationships between pairs of variables in large data sets is increasingly important, thus, we present a measure of dependence for two-variable relationships: the mutual information coefficient (MIC). MIC captures a wide range of associations both functional and not, and for functional relationships provides a score that roughly equals the coefficient of determination ( $R^2$ ) of the data relative to the regression function. MIC belongs to a larger class of mutual information-based nonparametric information (NPMI) statistics for identifying and classifying relationships. We apply MIC and NPMI to data sets in global health, gene expression, map-locus bandwidth, and the human gut microbiome and identify known and new relationships.

**I**magine a data set with hundreds of variables, many of which may contain important, unexplored relationships. There are tens of thousands of possible variable pairs—but how many to examine manually? If you do not already know what kinds of relationships to search for, how do you efficiently identify the important ones? Data sets of this size are increasingly common in fields as varied as genomics, physics, political science, and economics, so, making this question an important and growing challenge (1, 2).

One way to begin exploring a large data set is to search for pairs of variables that are closely associated. To do this, we could calculate some measure of dependency for each pair, rank the pairs by their scores, and examine the top-scoring pairs. For this strategy to work, the statistic used to measure dependency should have two basic properties: generality and variability.

By generality, we mean that with sufficient sample size the statistic should capture a wide range of interesting associations, not limited to specific functional types (such as linear, exponential or periodic), or even to all functional relationships (1). The latter condition is desirable because

[illegible]

By suitability, we mean that the same  $\sigma$  should give similar scores to equally noisy relationships of different types. For example, we want noisy linear relationships to differ structurally from noisy relationships of any other type. It is difficult to formulate for ourselves a general test, but we have a clear impression of the basic use of functional relationships: *an* suitable model should give similar scores to functionally related relationships with similar  $\sigma$  values.

Here, we describe an exploratory data analysis tool, the maximal information coefficient (MIC), that unifies these two heuristic practices. We establish MIC's generality through proofs of equitability on functional relations through simulation, and observe that this results into intuitively equitable behavior on empirical associations. Furthermore, we show that MIC gives rise to a larger family of statistics, which we refer to as MINE, or maximal

MINE statistics can be used not only to identify interesting associations, but also to characterize according to properties such as novelty and economicity. We demonstrate application of MDC and MINE to data sets: health, baseball, genomics, and the human microbiome.

The maximal information coefficient, briefly, MIC, is based on the idea that if a relationship exists between two variables, the grid can be drawn on the scatterplot of the variables that partitions the data to maximize that relationship. Thus, to calculate the MIC out of two-variable data, we explore all grids to a maximal grid resolution, dependent on sample size (Fig. 3A), computing the scores

of integers  $\{x, y\}$  the largest possible mutual information achievable by any  $x$ -by- $y$  grid applied to the data. We then normalize these mutual information values to ensure a fair comparison between grids of different dimensions and to obtain modified values between 0 and 1. We define the characteristic matrix  $M = (m_{x,y})$ , where  $m_{x,y}$  is the highest normalized mutual information achieved by any  $x$ -by- $y$  grid, and the matrix MRC to be the maximum value in

More formally, for a grid  $G$ , let  $I_G$  denote the mutual information of the probability dis-

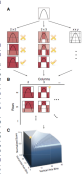


Fig. 3. Computing MIC (for each pair  $(x, y)$ ), the MIC algorithm finds the  $xy$ -grid with the highest induced mutual information. (a) The algorithm normalizes the mutual information scores and computes a matrix that stores, for each resource, the best grid at that resolution and its normalized score. (b) The normalized scores form the characteristic matrix, which can be visualized as a surface. MIC corresponds to the highest point on this surface; in this example, there are many grids that achieve the highest score. The star in (b) marks a sample grid achieving this score, and the star in (c) marks that grid's corresponding location on the surface.

This has been called a “correlation for the 21st century”, and it has also been heavily criticized ... read it yourself and make your own judgements!

## No free lunch

Generally speaking, we can't design a measure of correlation that performs well in every situation



This is because methods that are designed to detect very broad notions of associations will tend to give rise to false positives

## Recap: measures of correlation (continued)

In this lecture we saw examples of the **alternating conditional expectations** (ACE) algorithm in practice. The answers depend on the choice of smoother

In the population, **distance covariance** is defined in terms of independent copies of the pair of random variables. From this we can also define **distance variance** and **distance correlation**. Importantly, the latter is zero if and only if the random variables are independent

In the sample, distance covariance is defined by forming two absolute distance matrices, double centering them, and summing the elementwise products. Distance variance and distance correlation follow as before

Both population and sample versions have alternative definitions in terms of **characteristic functions**

Next time: midterm 1

Good luck!