# Canonical Correlation Analysis (CCA)

**SETTING:** * Unlike PCA we're going to focus on population (i.e. means & expectation) rather than sample averages (sample means & sample covariances)

**GIVEN:** A random vector $\vec{z}$ in $(p+q)$-dimensional space which we partition into a $p$-dimensional $\vec{x}$

$q$-dimensional $\vec{y}$

LET $\quad \vec{\mu}_x = E[\vec{x}] \quad \Sigma_x = E\left[(\vec{x}-\mu_x)(\vec{x}-\mu_x)^T\right] \in \mathbb{R}^{p \times p}$

$\vec{\mu}_y = E[\vec{y}] \quad \Sigma_y = E\left[(\vec{y}-\vec{\mu}_y)(\vec{y}-\mu_y)^T\right] \in \mathbb{R}^{q \times q}$

DEFINE $\quad \Sigma_{xy} = E\left[(\vec{x}-\mu_x)(\vec{y}-\mu_y)^T\right] \in \mathbb{R}^{p \times q}$

$\Sigma_{yx} = E\left[(\vec{y}-\vec{\mu}_y)(\vec{x}-\vec{\mu}_x)^T\right] \in \mathbb{R}^{q \times p}$

$\qquad \Sigma_{xy} = \Sigma_{yx}^T$

CROSS-COVARIANCE MATRIX!

IF $\quad Z = \begin{bmatrix} \vec{x} \\ \vec{y} \end{bmatrix}$ THEN

$$\mu_Z = \begin{bmatrix} \vec{\mu}_x \\ \vec{\mu}_y \end{bmatrix} \qquad \Sigma_Z = \begin{bmatrix} \Sigma_x & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_y \end{bmatrix} \in \mathbb{R}^{(p+q) \times (p+q)}$$

**GOAL:** We look for $\vec{a} \in \mathbb{R}^p$, $\vec{b} \in \mathbb{R}^q$ so that

$$\text{Corr}\left(\vec{a}^T \vec{x}, \vec{b}^T \vec{y}\right)$$

is MAXIMIZED!

**NOTE:** 1) If we think of $\vec{x}$ as causing $\vec{y}$ then $\vec{a}^T\vec{x}$ may be called the "best predictor" and $\vec{b}^T\vec{y}$ as "most predictable criterion". HOWEVER, there is no assumption of causal asymmetry; $\vec{x}$ and $\vec{y}$ treated symmetrically.

2) Connections to multiple regression (see Zelterman Ch. 7) and reduced rank regression.

# Canonical Correlation Analysis

The correlation between $\eta = \vec{a}^T \vec{x}$ and $\xi = \vec{b}^T \vec{y}$ is

$$\rho(\vec{a}, \vec{b}) = \text{Corr}(\eta, \xi) = \frac{\text{Cov}(\eta, \xi)}{\sqrt{\text{Var}(\eta)\,\text{Var}(\xi)}}$$

DEPENDS ON $\vec{a}, \vec{b}$

$$= \frac{E\left[\vec{a}^T \vec{x}\, \vec{y}^T \vec{b}\right] - \vec{a}^T E[\vec{x}] E[\vec{y}^T]\vec{b}}{\sqrt{\text{Var}(\vec{a}^T \vec{x})\,\text{Var}(\vec{b}^T \vec{y})}}$$

$$= \frac{\vec{a}\left(E[\vec{x}\vec{y}^T] - E[\vec{x}]E[\vec{y}^T]\right)\vec{b}}{\sqrt{\vec{a}^T \Sigma_x \vec{a}\ \vec{b}^T \Sigma_y \vec{b}}} = \frac{\vec{a}^T \Sigma_{xy} \vec{b}}{\sqrt{\vec{a}^T \Sigma_x \vec{a}\ \vec{b}^T \Sigma_y \vec{b}}}$$

Note:

1) NEED $\quad \vec{a}^T \Sigma_x \vec{a} \neq 0$
   $\qquad\qquad \vec{b}^T \Sigma_y \vec{b} \neq 0$ $\Big\}$ ASSUME $\Sigma_x, \Sigma_y$ HAVE FULL RANK

2) IF SCALARS $c, d \neq 0$ THEN

$$\rho(\vec{a}, \vec{b}) = \rho(c\vec{a}, d\vec{b})$$

SO WE'LL FOCUS ON CONSTRAINTS $\vec{a}^T \Sigma_x \vec{a} = \vec{b}^T \Sigma_y \vec{b} = 1$

$$\ast \quad \boxed{(\vec{a}, \vec{b})^* = \underset{\vec{a}^T \Sigma_x \vec{a} = \vec{b}^T \Sigma_y \vec{b} = 1}{\arg\max} \rho(a, b)}$$

3) CAN MINIMIZE CORRELATION $\vec{a} \longmapsto -a \qquad$ OR $\qquad \vec{b} \longmapsto -\vec{b}$

$$\rho(-\vec{a}, \vec{b}) = -\rho(\vec{a}, \vec{b}) \qquad\qquad \rho(\vec{a}, -\vec{b}) = -\rho(\vec{a}, \vec{b})$$

BUT NOT BOTH $\vec{a} \to -\vec{a}$
$\qquad\qquad\qquad\quad \vec{b} \to -\vec{b}$

4) WE CAN SOLVE $\ast$ USING i) LAGRANGE MULTIPLIERS LIKE PCA
   TO GET GENERALIZED EIGENVALUE PROBLEM (MESSY & VERY TECHNICAL)
   OR USING ii) SVD (JUST MESSY)

## Aside: Fractional Powers of a Matrix

Let $A$ be a diagonalizable $n \times n$ matrix so that there exists invertible $P \in \mathbb{R}^{n \times n}$ and diagonal $D \in \mathbb{R}^{n \times n}$ such that

$$A = PDP^{-1} \quad , \quad D = \begin{bmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_n \end{bmatrix}$$

For $n = 0, 1, 2, \dots$ we define

$$A^n = PD^n P^{-1} \qquad D^\alpha = \begin{bmatrix} \lambda_1^n & & 0 \\ & \ddots & \\ 0 & & \lambda_n^n \end{bmatrix}$$

If $\lambda_1, \dots \lambda_n$ are $> 0$ we can extend this to negative and fractional powers, e.g.

$$A^{-1} = P \begin{bmatrix} 1/\lambda_1 & & 0 \\ & \ddots & \\ 0 & & 1/\lambda_n \end{bmatrix} P^{-1}$$

$$A^{1/2} = P \begin{bmatrix} \sqrt{\lambda_1} & & 0 \\ & \ddots & \\ 0 & & \sqrt{\lambda_n} \end{bmatrix} P^{-1}$$

---

Since $\Sigma_x, \Sigma_y$ are full rank covariance matrices (positive definite)

$$\Sigma_x = \Sigma_x^{1/2} \Sigma_x^{1/2} \qquad \Sigma_x = P_x \Lambda_x P_x^T \qquad \Sigma_x^{1/2} = P_x \Lambda_x^{1/2} P_x^T$$

$$\Sigma_y = \Sigma_y^{1/2} \Sigma_y^{1/2} \qquad \Sigma_y = P_y \Lambda_y P_y^T \qquad \Sigma_y^{1/2} = P_y \Lambda_y^{1/2} P_y^T$$

$$\Lambda_x, \Lambda_y \text{ diagonal}$$

$$\frac{\vec{a}^T \Sigma_{xy} \vec{b}}{\sqrt{\vec{a}^T \Sigma_x \vec{a} \; \vec{b}^T \Sigma_y \vec{b}}}$$

Change of variable
$$\vec{\alpha} = \Sigma_x^{1/2} \vec{a}$$
$$\vec{\beta} = \Sigma_y^{1/2} \vec{b}$$
$$\vec{a} = \Sigma_x^{-1/2} \vec{\alpha}$$
$$\vec{b} = \Sigma_y^{-1/2} \vec{\beta}$$

$$\longrightarrow \frac{\vec{\alpha}^T \Sigma_x^{-1/2} \Sigma_{xy} \Sigma_y^{-1/2} \vec{\beta}}{\sqrt{\vec{\alpha}^T \Sigma_x^{-1/2} \Sigma_x \Sigma_x^{-1/2} \vec{\alpha} \; \vec{\beta}^T \Sigma_y^{-1/2} \Sigma_y \Sigma_y^{-1/2} \vec{\beta}}}$$

$$\parallel$$

$$\frac{\vec{\alpha}^T \Sigma_x^{-1/2} \Sigma_{xy} \Sigma_y^{-1/2} \vec{\beta}}{\sqrt{\vec{\alpha}^T \vec{\alpha} \; \vec{\beta}^T \vec{\beta}}}$$

SUPPOSE RANK $\left(\Sigma_{xy}\right) = k \leq \min(p,q)$ SO THAT $\Sigma_x^{-1/2} \Sigma_{xy} \Sigma_y^{-1/2}$ HAS

RANK $k$ AS WELL.

$$\Rightarrow \quad \Sigma_x^{-1/2} \Sigma_{xy} \Sigma_y^{-1/2} = \overset{\cdot\cdot}{U} \overset{\cdot}{D} \overset{\cdot}{V}^T$$

$$= \left[\vec{u}_1 | \ldots | \vec{u}_k\right] \begin{bmatrix} \sigma_1 & \\ & \ddots \\ & & \sigma_k \end{bmatrix} \left[\vec{v}_1 | \ldots | \vec{v}_k\right]^T$$

$$\vec{u}_1, \ldots \vec{u}_k \quad \text{ORTHONORMAL VECTORS IN } \mathbb{R}^p$$

$$\vec{v}_1, \ldots \vec{v}_k \quad \text{'' \quad '' \quad , } \mathbb{R}^q$$

TO MAXIMIZE $\dfrac{\vec{\alpha}^T \Sigma_x^{-1/2} \Sigma_{xy} \Sigma_y^{-1/2}}{\sqrt{\vec{\alpha}^T \vec{\alpha} \, \vec{\beta}^T \vec{\beta}}} = \dfrac{\vec{\alpha}^T \left[\vec{u}_1, \ldots \vec{u}_k\right] \begin{bmatrix} \sigma_1 & 0 \\ 0 & \ddots & \sigma_k \end{bmatrix} \begin{bmatrix} \vec{v}_1^T \\ \vec{v}_k^T \end{bmatrix} \vec{\beta}}{\sqrt{\vec{\alpha}^T \vec{\alpha} \, \vec{\beta}^T \vec{\beta}}}$

WE CHOOSE $\vec{\alpha}_1 = \vec{u}_1$ , $\vec{\beta}_1 = \vec{v}_1$ TO ATTAIN (GLOBAL) MAXIMUM $\sigma_1$

SIMILARLY WE COULD CHOOSE $\vec{\alpha}_j = \vec{u}_j$ , $\vec{\beta}_j = \vec{v}_j$ TO ATTAIN (LOCAL) MAXIMUM $\sigma_j$ , $j=1, \ldots k$

<u>DEF</u>: LET $\vec{a}_j = \Sigma_x^{-1/2} \vec{\alpha}_j$ , $\vec{b}_j = \Sigma_y^{-1/2} \vec{\beta}_j$ AND $\sigma_1, \ldots \sigma_k$ BE AS ABOVE. FOR $j = 1, \ldots k$

(a) THE VECTORS $\vec{a}_j$ , $\vec{b}_j$ ARE CALLED THE $j^{th}$ CANONICAL

CORRELATION VECTORS !

(b) THE RANDOM VARIABLES $\eta_j = \vec{a}_j^T \vec{x}$ AND $\xi_j = \vec{b}_j^T \vec{y}_j$ ARE

CALLED THE $j^{th}$ CANONICAL CORRELATION VARIABLES

(c) $\sigma_j$ IS CALLED THE $j^{th}$ CANONICAL CORRELATION

<u>NOTE</u>: 1) THE NONZERO SINGULAR VALUES OF $\Sigma_x^{-1/2} \Sigma_{xy} \Sigma_y^{-1/2}$ ARE

THE SQUARE ROOTS OF THE NONZERO EIGENVALUES OF

$$\left(\Sigma_x^{-1/2} \Sigma_{xy} \Sigma_y^{-1/2}\right)\left(\Sigma_x^{-1/2} \Sigma_{xy} \Sigma_y^{-1/2}\right)^T = \Sigma_x^{-1/2} \Sigma_{xy} \Sigma_y^{-1} \Sigma_{yx} \Sigma_x^{-1/2} = N$$

BY SIMILARITY THE EIGENVALUES OF $N$ ARE THE SAME AS

$$M = \Sigma_x^{-1/2} N \Sigma_x^{1/2} = \boxed{\Sigma_x^{-1} \Sigma_{xy} \Sigma_y^{-1} \Sigma_{yx}}$$

3) NOTE THAT $\vec{\alpha}_1, \ldots, \vec{\alpha}_k$ ARE ORTHONORMAL, e.g. $\vec{\alpha}_j^T \vec{\alpha}_k = \delta_{jk}$

BUT $\vec{a}_1, \ldots, \vec{a}_k$ ARE IN GENERAL NEITHER PERPENDICULAR NOR

UNIT LENGTH, e.g. $\vec{a}_j^T \vec{a}_k \neq \delta_{jk}$. HOWEVER,

$$\vec{a}_j^T \Sigma_x \vec{a}_k = \vec{\alpha}_j^T \vec{\alpha}_k = \delta_{jk}$$

SIMILAR STATEMENT HOLDS FOR $\vec{\beta}_1, \ldots, \vec{\beta}_k, \vec{b}_1, \ldots, \vec{b}_k, \Sigma_y$

4) SINCE $\vec{\alpha}_j, \vec{\beta}_j$ WERE CHOSEN AS SINGULAR VECTORS OF $\Sigma_x^{-1/2} \Sigma_{xy} \Sigma_y^{-1/2}$

$$\Sigma_x^{-1/2} \Sigma_{xy} \Sigma_y^{-1/2} \vec{\beta}_j = \sigma_j \vec{\alpha}_j$$

$$\vec{\alpha}_j^T \Sigma_x^{-1/2} \Sigma_{xy} \Sigma_y^{-1/2} = \sigma_j \vec{\beta}_j^T$$

## CCA AND DATA

GIVEN INDEPENDENT SAMPLES OF $\begin{pmatrix} \vec{x}_i \\ \vec{y}_i \end{pmatrix}$ $i = 1, \ldots, N$ WE CAN

COMPUTE THE SAMPLE COVARIANCE MATRICES

$$\hat{\Sigma}_x = \left( \frac{1}{N} \sum_{i=1}^{N} \vec{x}_i \vec{x}_i^T \right) - \bar{x}\bar{x}^T \qquad \bar{x} = \frac{1}{N} \sum_{i=1}^{N} \vec{x}_i$$

$$\hat{\Sigma}_y = \left( \frac{1}{N} \sum_{i=1}^{N} \vec{y}_i \vec{y}_i^T \right) - \bar{y}\bar{y}^T \qquad \bar{y} = \frac{1}{N} \sum_{i=1}^{N} \vec{y}_i$$

AND THE SAMPLE CROSS-COVARIANCE MATRIX

$$\hat{\Sigma}_{xy} = \left( \frac{1}{N} \sum_{i=1}^{N} \vec{x}_i \vec{y}_i^T \right) - \bar{x}\bar{y}^T$$

$$= \frac{1}{N} \sum_{i=1}^{N} (\vec{x}_i - \bar{x})(\vec{y}_i - \bar{y})^T$$

AND REPEAT THE ANALYSIS

NOTE: NEED $\hat{\Sigma}_x, \hat{\Sigma}_y$ TO BE POSITIVE DEFINITE

$\Rightarrow$ NEED $N \geq \max\{p, q\}$

GIVEN $\hat{\Sigma}_x, \hat{\Sigma}_y, \hat{\Sigma}_{xy}$ WE CAN COMPUTE

a) SAMPLE CANONICAL CORRELATIONS $\hat{\sigma}_1, \ldots, \hat{\sigma}_k$ $\quad k = \text{RANK}\left(\hat{\Sigma}_{xy}\right) \leq \min\{p, q\}$

b) SAMPLE CANONICAL VECTORS $\vec{a}_1, \ldots \vec{a}_k, \vec{b}_1, \ldots \vec{b}_k$

THEN FOR EACH SAMPLE $(\vec{x}_i, \vec{y}_i)$ WE CAN COMPUTE THE CANONICAL VARIABLES (SCORES)

$$\eta_{ij} = \vec{a}_j^T \vec{x}_i = \vec{x}_i^T \vec{a}_j$$

$$\vec{\xi}_{ij} = \vec{b}_j^T \vec{y}_i = \vec{y}_i^T \vec{b}_j \qquad j = 1, \ldots k$$

---

NOTE:

1) LET $\quad \vec{\eta} = \begin{bmatrix} \vec{a}_1^T \vec{x} \\ \vdots \\ \vec{a}_k^T \vec{x} \end{bmatrix} = \begin{bmatrix} \vec{a}_1^T \\ \vdots \\ \vec{a}_k^T \end{bmatrix} \vec{x}$ AND $\quad \hat{\xi} = \begin{bmatrix} \vec{b}_1^T \vec{y} \\ \vdots \\ \vec{b}_k^T \vec{y} \end{bmatrix} = \begin{bmatrix} \vec{b}_1^T \\ \vdots \\ \vec{b}_k^T \end{bmatrix} \vec{y}$

ONE CAN SHOW (HOMEWORK)

$$\Sigma_{\eta \xi} = \begin{bmatrix} \sigma_1 & & 0 \\ & \ddots & \\ 0 & & \sigma_k \end{bmatrix} \qquad \Sigma_\eta = \Sigma_\xi = I_k$$
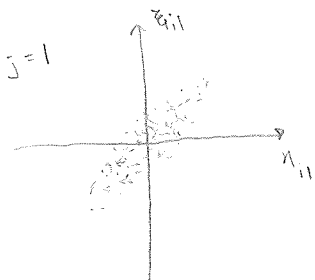
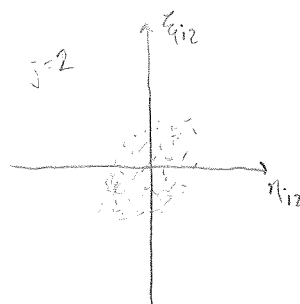SAME STATEMENT HOLDS FOR SAMPLE COVARIANCES AND CROSS-COVARIANCE

OF $\quad \vec{\eta}_i = \begin{bmatrix} \vec{a}_1^T \\ \vdots \\ \vec{a}_k \end{bmatrix} \vec{x}_i \qquad \hat{\xi}_i = \begin{bmatrix} \vec{b}_1^T \\ \vdots \\ \vec{b}_k \end{bmatrix} \vec{y}_i$

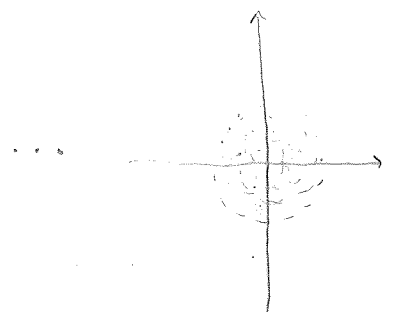2) IF WE PLOT $\left\{(\eta_{ij}, \xi_{ij})\right\}_{i=1}^N$ FOR $j = 1, \ldots, k$ THE

PLOTS WILL SHOW DECREASING POSITIVE CORRELATION



$j=1$ ($\xi_{i1}$ vs $\eta_{i1}$)   CORRELATION $= \sigma_1$

$j=2$ ($\xi_{i2}$ vs $\eta_{i2}$)   CORRELATION $= \sigma_2 < \sigma_1$

$\ldots$   CORRELATION $\sigma_k < \sigma_{k-1} \ldots$

3) In PCA, dimension reduction focused on minimizing the loss of 'information' (variance/correlation) within the data, $\bar{X}$.

CCA reduces dimensions with the goal of minimizing the loss of 'information' (correlation) between data $\bar{X}$ and $\bar{Y}$.

- If $\sigma_j \approx 0$ for $j \geq s$ and $s \ll k$, we can DISCARD $\eta_{ij}, \xi_{ij}$ for $j \geq s$ without losing much of the information in $\vec{x}_i$ about $\vec{y}_i$
  $\vec{y}_i$ about $\vec{x}_i$!

- $\vec{x}_i^{(k)} = \sum_{s=1}^{k} \eta_{ij} \vec{a}_s \in \text{SPAN}\{\vec{a}_1, \dots \vec{a}_k\} \subseteq \mathbb{R}^p$
  $\vec{y}_i^{(k)} = \sum_{s=1}^{k} \xi_{ij} \vec{b}_s \in \text{SPAN}\{\vec{b}_1, \dots, \vec{b}_k\} \subseteq \mathbb{R}^q$ } DIFFERENT SPACES!

  1) $\vec{a}_1, \dots \vec{a}_k$ NOT AN ORTHONORMAL BASIS FOR SPAN $\{\vec{a}_1, \dots \vec{a}_k\}$
     $\vec{b}_1, \dots \vec{b}_k$       SPAN $\{\vec{b}_1, \dots \vec{b}_k\}$

  2) CCA DOES NOT FIND A $k$-DIMENSIONAL SUBSPACE OF $\mathbb{R}^{(p+q)}$!

LIMITATIONS OF PCA

1) LIKE PCA, THE CANONICAL CORRELATION VECTORS AND VARIABLES (SCORES) REPRESENT LINEAR COMBINATIONS OF COMPONENTS OF $\vec{x}, \vec{y}$
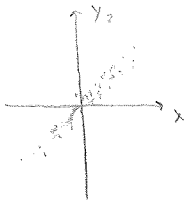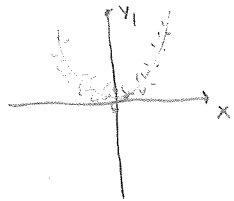
   $\Rightarrow$ INTERPRETABILITY IS A ISSUE

2) CORRELATION, CANONICAL OF OTHERWISE, IS A MEASUREMENT OF LINEAR DEPENDENCE

$\underline{E_x}:$    $x \sim N(0,1) \subset \mathbb{R}^1$

$$\vec{y}|x \sim \begin{bmatrix} x^2 \\ x \end{bmatrix} + \vec{w} \quad, \quad \vec{w} \sim N\left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \sigma^2 I \right)$$

$\underline{\text{DATA DRAWN FROM}} \uparrow$



FIND CANONICAL CORRELATIONS AND VECTORS

- $E_x = 0$ , $\Sigma_x = 1$

- $E y = E\left[ \begin{bmatrix} x^2 \\ x \end{bmatrix} + \vec{w} \right] = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$ , $\Sigma_y = E y \vec{y}^T - \mu_y \mu_y^T$

$$= E\left[ \begin{bmatrix} x^2 \\ x \end{bmatrix} [x^2 \ x] \right] + \sigma^2 I - \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$= E\left[ \begin{matrix} x^4 & x^3 \\ x^3 & x^2 \end{matrix} \right] + \sigma^2 I - \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$$

$$= \begin{bmatrix} 2+\sigma^2 & 0 \\ 0 & 1+\sigma^2 \end{bmatrix}$$

- $\Sigma_{xy} = E[x \vec{y}] - \cancel{(E_x)(E\vec{y})}^{0} = E\left[ [x^3, x^2] + x \vec{w} \right] = [0, 1]$

$$\Rightarrow \Sigma_x^{-1/2} \Sigma_{xy} \Sigma_y^{-1/2} = [1][0,1]\begin{bmatrix} \frac{1}{\sqrt{2+\sigma^2}} & 0 \\ 0 & \frac{1}{\sqrt{1+\sigma^2}} \end{bmatrix} = \left[ 0, \frac{1}{\sqrt{1+\sigma^2}} \right]$$

$$= [1]\left[ \frac{1}{\sqrt{1+\sigma^2}}, 0 \right] \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

$$\Rightarrow \sigma_1 = \frac{1}{\sqrt{1+\sigma^2}}$$

$\alpha_1 = 1 \quad \Rightarrow a_1 = \Sigma_x^{-1/2} \alpha_1 = 1$     $\overset{\text{NO WEIGHT TO } y_1 = x^2 + \sigma_{w_1}!}{\nearrow}$

$\vec{\beta}_1 = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \Rightarrow b_1 = \Sigma_y^{-1/2} \vec{\beta}_1 = \begin{bmatrix} 0 \\ \frac{1}{\sqrt{1+\sigma^2}} \end{bmatrix}$    * CORRELATION BETWEEN $x^2$ AND $x$ IS 0!

* CANONICAL CORRELATIONS MAY ALL BE NEAR BUT THERE MAY BE A STRONG DEPENDENCE BETWEEN (COMPONENTS OF) $\vec{x}$ AND (COMPONENTS OF) $\vec{y}$ !