

CANONICAL CORRELATION ANALYSIS (CCA)

Ke Chen

Department of Computer Science, The University of Manchester

Ke.Chen@manchester.ac.uk

OUTLINE

BACKGROUND

History, motivation, task and scope of applications

PRINCIPLE

CCA Derivation: learning maximum correlation between two random vectors

ALGORITHM

CCA algorithms and Proportion of Variance (PoV)

ILLUSTRATIVE EXAMPLE

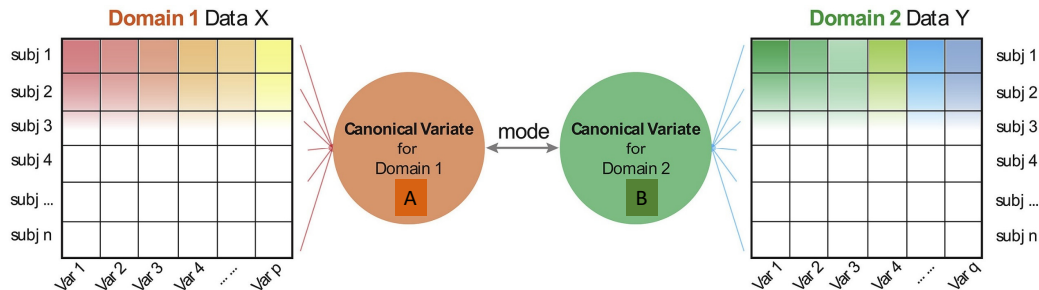
CCA on synthetic and real datasets

EXTENSION

Sparse CCA, kernel CCA, Deep CCA and Probabilistic CCA

- Concept originated by Camille Jordan in 1875, independently developed and name coined by Harold Hotelling in 1930s
- An approach to exploratory data analysis especially for **two relevant random vectors**, e.g., visual and audio channels in video
- An important linear model for contemporary representation learning
- Applied in economics, medicine, meteorology, neuroscience, pattern recognition, psychology and many other fields
- In context of ML and statistics, widely used for **feature extraction** and **data interpretation**

- Find **optimal linear projections** for two relevant random vectors to **maximise their correlation** in the new spaces
- New **low-dimensional representations** of two random vectors achieved by projecting data points to new CCA spaces



BACKGROUND

- Find **optimal linear projections** for two relevant random vectors to **maximise their correlation** in the new spaces
- New **low-dimensional representations** of two random vectors achieved by projecting data points to new CCA spaces

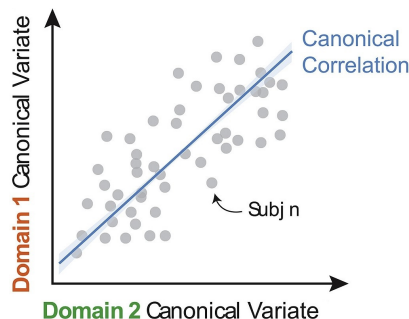
For $X_{p \times N} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ and

$Y_{q \times N} = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N\}$, find $\mathbf{a}_{p \times 1}$ and $\mathbf{b}_{q \times 1}$:

Domain 1: $z_{X,n} = \mathbf{x}_n^T \mathbf{a} \quad n = 1, 2, \dots, N,$

Domain 2: $z_{Y,n} = \mathbf{y}_n^T \mathbf{b} \quad n = 1, 2, \dots, N.$

i.e., $(\mathbf{z}_X, \mathbf{z}_Y) = (X^T \mathbf{a}, Y^T \mathbf{b})$ such that
maximise correlation $(\mathbf{z}_X, \mathbf{z}_Y)$



PRINCIPLE: CCA DERIVATION

Find out the 1st pair of canonical vectors, \mathbf{a} and \mathbf{b}

Given a dataset of N data points, (X, Y) , $X_{p \times N} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ and $Y_{q \times N} = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N\}$ are **centralised**, respectively.

- To find out the 1st pair of **canonical vectors**, \mathbf{a} and \mathbf{b} , of **maximum correlation**, we need to establish a utility function by encoding this learning goal.
- Projecting X and Y onto \mathbf{a} and \mathbf{b} leads to $\mathbf{z}_X = X^T \mathbf{a}$ and $\mathbf{z}_Y = Y^T \mathbf{b}$.
- The **correlation** between two vectors, \mathbf{z}_X and \mathbf{z}_Y , is

$$\text{correlation}(\mathbf{z}_X, \mathbf{z}_Y) = \frac{\mathbf{z}_X^T \mathbf{z}_Y}{\|\mathbf{z}_X\| \|\mathbf{z}_Y\|} = \frac{\mathbf{a}^T X Y^T \mathbf{b}}{\sqrt{\mathbf{a}^T X X^T \mathbf{a}} \sqrt{\mathbf{b}^T Y Y^T \mathbf{b}}} = \frac{\mathbf{a}^T S_{XY} \mathbf{b}}{\sqrt{\mathbf{a}^T S_{XX} \mathbf{a}} \sqrt{\mathbf{b}^T S_{YY} \mathbf{b}}},$$

where S_{XX} , S_{YY} and S_{XY} are **covariance** matrices on X , Y and XY (defined on both X and Y with size $p \times q$), respectively; i.e.

$$S_{XX} \equiv \frac{1}{N-1} X X^T, \quad S_{YY} \equiv \frac{1}{N-1} Y Y^T, \quad S_{XY} \equiv \frac{1}{N-1} X Y^T.$$

Find out the 1st pair of canonical vectors, \mathbf{a} and \mathbf{b}

- The denominators in $\text{correlation}(\mathbf{z}_X, \mathbf{z}_Y)$ simply play a role of normalisation to make the correlation **invariant** with respect to rescaling \mathbf{a} and \mathbf{b} .

$$\text{correlation}(\mathbf{z}_X, \mathbf{z}_Y) = \frac{\mathbf{z}_X^T \mathbf{z}_Y}{\|\mathbf{z}_X\| \|\mathbf{z}_Y\|} = \frac{\mathbf{a}^T S_{XY} \mathbf{b}}{\sqrt{\mathbf{a}^T S_{XX} \mathbf{a}} \sqrt{\mathbf{b}^T S_{YY} \mathbf{b}}}$$

- Thus, the **learning objective** to find out the 1st pair of \mathbf{a} and \mathbf{b} for maximum correlation is formulated along with two **normalisation constraints**:

$$\mathbf{z}_X^T \mathbf{z}_Y \text{ s.t. } \|\mathbf{z}_X\|^2 = 1, \|\mathbf{z}_Y\|^2 = 1 \implies \mathbf{a}^T S_{XY} \mathbf{b} \text{ s.t. } \mathbf{a}^T S_{XX} \mathbf{a} = 1, \mathbf{b}^T S_{YY} \mathbf{b} = 1.$$

- Apply the **Lagrangian multipliers** to form the unconstrained utility function:

$$\mathcal{L}(\mathbf{a}, \mathbf{b}, \lambda_a, \lambda_b; X, Y) = \mathbf{a}^T S_{XY} \mathbf{b} + \frac{\lambda_a}{2}(1 - \mathbf{a}^T S_{XX} \mathbf{a}) + \frac{\lambda_b}{2}(1 - \mathbf{b}^T S_{YY} \mathbf{b}).$$

Find out the 1st pair of canonical vectors, \mathbf{a} and \mathbf{b}

- To find out the optimal \mathbf{a} and \mathbf{b} , apply the optimality conditions:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{a}} = 0 \implies S_{XY}\mathbf{b} - \lambda_a S_{XX}\mathbf{a} = 0 \implies S_{XY}\mathbf{b} = \lambda_a S_{XX}\mathbf{a}, \quad (1)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{b}} = 0 \implies S_{YX}\mathbf{a} - \lambda_b S_{YY}\mathbf{b} = 0 \implies S_{YX}\mathbf{a} = \lambda_b S_{YY}\mathbf{b}, \quad (2)$$

- Multiplying \mathbf{a}^T and \mathbf{b}^T on both sides in Eqs. (1) and (2) and further applying two constraints, $\mathbf{a}^T S_{XX}\mathbf{a} = 1$, $\mathbf{b}^T S_{YY}\mathbf{b} = 1$, we have

$$\mathbf{a}^T (S_{XY}\mathbf{b}) = \lambda_a \mathbf{a}^T S_{XX}\mathbf{a} = \lambda_a, \quad \mathbf{b}^T (S_{YX}\mathbf{a}) = \lambda_b \mathbf{b}^T S_{YY}\mathbf{b} = \lambda_b.$$

- $\mathbf{a}^T S_{XY}\mathbf{b}$ and $\mathbf{b}^T S_{YX}\mathbf{a}$ are scalar, and $(\mathbf{a}^T S_{XY}\mathbf{b})^T = \mathbf{b}^T S_{YX}\mathbf{a}$. Therefore, we must have $\lambda_a = \lambda_b$ at the optimum. $\lambda_a = \lambda_b = \lambda$ is the canonical correlation coefficient.

Find out the 1st pair of canonical vectors, \mathbf{a} and \mathbf{b}

- By using the common λ , we rewrite Eqs. (1) and (2) as follows:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{a}} = 0 \implies S_{XY} \mathbf{b} = \lambda S_{XX} \mathbf{a}, \quad (3)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{b}} = 0 \implies S_{YX} \mathbf{a} = \lambda S_{YY} \mathbf{b}, \quad (4)$$

- Assume S_{XX} and S_{YY} are invertible, we rewrite Eq.(4) as $\mathbf{b} = \frac{1}{\lambda} S_{YY}^{-1} S_{YX} \mathbf{a}$ and insert it into Eq.(3):

$$S_{XY} S_{YY}^{-1} S_{YX} \mathbf{a} = \lambda^2 S_{XX} \mathbf{a} \implies \boxed{S_{XX}^{-1} S_{XY} S_{YY}^{-1} S_{YX} \mathbf{a} = \lambda^2 \mathbf{a}}. \quad (5)$$

- With eigen-analysis on Eq.(5), choose \mathbf{a}^* with the largest λ^* and $\mathbf{b}^* = \frac{1}{\lambda^*} S_{YY}^{-1} S_{YX} \mathbf{a}^*$.
- The choice of \mathbf{a}^* and \mathbf{b}^* leads to the **the 1st pair of canonical vectors, \mathbf{a} and \mathbf{b}** .

CCA Formulation

- Similarly, we can find out more pairs of canonical vectors with the same idea, e.g., the learning objectives of the **2nd pair** of \mathbf{a} and \mathbf{b} as follows:

$$\mathbf{a}^T S_{XY} \mathbf{b} \text{ s.t. } \mathbf{a}^T S_{XX} \mathbf{a} = 1, \mathbf{b}^T S_{YY} \mathbf{b} = 1, \mathbf{a}^T \mathbf{a}^* = 0, \mathbf{b}^T \mathbf{b}^* = 0.$$

- In general, M ($M \leq \min(p, q)$) pairs of canonical vectors, $A = \{\mathbf{a}_1, \dots, \mathbf{a}_M\}$ and $B = \{\mathbf{b}_1, \dots, \mathbf{b}_M\}$, are achieved by doing the **eigen analysis** on the composite matrix, $S_{XX}^{-1} S_{XY} S_{YY}^{-1} S_{YX}$, to form $A_{p \times M}$ with the eigen vectors, $\mathbf{a}_1, \dots, \mathbf{a}_M$, of top M eigenvalues and $B_{q \times M}$ with $\mathbf{b}_i = \frac{1}{\lambda_i} S_{YY}^{-1} S_{YX} \mathbf{a}_i$, $i = 1, \dots, M$
- **Property**: all M canonical vectors in A and B are **uncorrelational**, i.e., $\mathbf{a}_i^T \mathbf{a}_j \neq 0$ if $i = j$ and 0 otherwise, and $\mathbf{b}_i^T \mathbf{b}_j \neq 0$ if $i = j$ and 0 otherwise.
- For dataset $(X_{p \times N}, Y_{q \times N})$, we can use CCA projection matrices, $A_{p \times M} = \{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_M\}$ and $B_{q \times M} = \{\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_M\}$, to achieve the **M -dimensional representations**:

$$(\mathbf{z}_{X,i}, \mathbf{z}_{Y,i}) = (X^T \mathbf{a}_i, Y^T \mathbf{b}_i) \quad i = 1, 2, \dots, M.$$

ALGORITHM: BASIC CCA

Training Phase: Find out CCA projection matrices

• Data centralisation

For a dataset, $X_{p \times N} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ ($p < N$) and $Y_{q \times N} = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N\}$ ($q < N$), calculate $\mathbf{m}_X = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$ and $\mathbf{m}_Y = \frac{1}{N} \sum_{n=1}^N \mathbf{y}_n$ to produce $\bar{X}_{p \times N} = \{\mathbf{m}_X, \mathbf{m}_X, \dots, \mathbf{m}_X\}$ and $\bar{Y}_{q \times N} = \{\mathbf{m}_Y, \mathbf{m}_Y, \dots, \mathbf{m}_Y\}$. Data centralisation: $\hat{X} = X - \bar{X}$ and $\hat{Y} = Y - \bar{Y}$.

• Eigen analysis

Calculate covariance matrices: $S_{XX} = \frac{1}{N-1} \hat{X} \hat{X}^T$, $S_{YY} = \frac{1}{N-1} \hat{Y} \hat{Y}^T$, $S_{XY} = \frac{1}{N-1} \hat{X} \hat{Y}^T$.

For $S_{XX}^{-1} S_{XY} S_{YY}^{-1} S_{YX}$, find out and rank all p eigenvalues so that $\lambda_1^2 \geq \lambda_2^2 \geq \dots \geq \lambda_p^2$ with their corresponding eigenvectors, $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_p$. Then, generate paired eigenvectors via $\mathbf{b}_i = \frac{1}{\lambda_i} S_{YY}^{-1} S_{YX} \mathbf{a}_i$, $i = 1, 2, \dots, \min(p, q)$.

• Construct projection matrices

Select top M ($M \leq \min(p, q)$) eigenvectors to form the paired projection matrices:

$A_{p \times M} = \{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_M\}$ and $B_{q \times M} = \{\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_M\}$

Deployment Phase: Generate low-dimensional representations in CCA space

$$(\mathbf{x}, \mathbf{y}) \implies (\mathbf{z}_x, \mathbf{z}_y), \quad \mathbf{z}_x = A^T (\mathbf{x} - \mathbf{m}_X) \text{ and } \mathbf{z}_y = B^T (\mathbf{y} - \mathbf{m}_Y).$$

ALGORITHM: SVD-BASED CCA

Training Phase: Find out CCA projection matrices

- **Data centralisation**

For a given dataset, $X_{p \times N} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ and $Y_{q \times N} = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N\}$, calculate $\mathbf{m}_X = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$ and $\mathbf{m}_Y = \frac{1}{N} \sum_{n=1}^N \mathbf{y}_n$ to produce $\bar{X}_{p \times N} = \{\mathbf{m}_X, \mathbf{m}_X, \dots, \mathbf{m}_X\}$ and $\bar{Y}_{q \times N} = \{\mathbf{m}_Y, \mathbf{m}_Y, \dots, \mathbf{m}_Y\}$. Data centralisation: $\hat{X} = X - \bar{X}$ and $\hat{Y} = Y - \bar{Y}$.

- **SVD solution to eigen analysis**

Calculate covariance matrices: $S_{XX} = \frac{1}{N-1} \hat{X} \hat{X}^T$, $S_{YY} = \frac{1}{N-1} \hat{Y} \hat{Y}^T$, $S_{XY} = \frac{1}{N-1} \hat{X} \hat{Y}^T$.

Apply SVD so as to $S_{XX}^{-\frac{1}{2}} S_{XY} S_{YY}^{-\frac{1}{2}} = U_{p \times p} \Sigma_{p \times q} V_{q \times q}^T$, where $U_{p \times p} = \{\mathbf{u}_1, \dots, \mathbf{u}_p\}$ and $V_{q \times q} = \{\mathbf{v}_1, \dots, \mathbf{v}_q\}$.

- **Construct projection matrices**

Select top M ($M \leq \min(p, q)$) eigenvectors to form the paired projection matrices:

$A_{p \times M} = \{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_M\}$ and $B_{q \times M} = \{\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_M\}$, where $\mathbf{a}_i = S_{XX}^{-\frac{1}{2}} \mathbf{u}_i$, $\mathbf{b}_i = S_{YY}^{-\frac{1}{2}} \mathbf{v}_i$

Deployment Phase: Generate low-dimensional representations in CCA space

$$(\mathbf{x}, \mathbf{y}) \implies (\mathbf{z}_x, \mathbf{z}_y), \quad \mathbf{z}_x = A^T (\mathbf{x} - \mathbf{m}_X) \text{ and } \mathbf{z}_y = B^T (\mathbf{y} - \mathbf{m}_Y).$$

Fact: `sklearn.cross_decomposition.CCA` in the scikit-learn library.

How to find out a proper dimension, M^* , to form CCA space?

The shared eigenvalues between A and B can still be used in PoV.

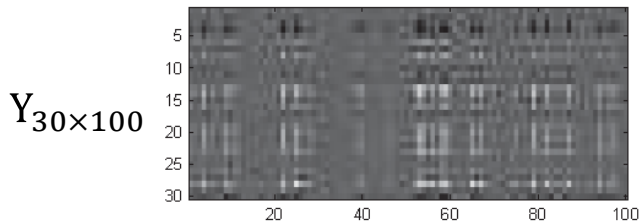
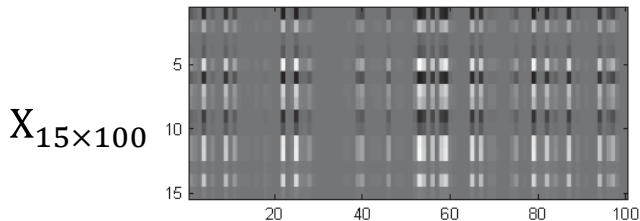
- **Proportion of Variance (PoV)**

$$\text{PoV}(k) = \frac{\sum_{i=1}^k \lambda_i^2}{\sum_{i=1}^d \lambda_i^2} = \frac{\lambda_1^2 + \cdots + \lambda_k^2}{\lambda_1^2 + \lambda_2^2 + \cdots + \lambda_{\min(p,q)}^2}, \quad \text{for } k < \min(p, q).$$

- In practice, find out k^* so that $\text{PoV}(k^*) \geq 90\%$. Then, set $M^* = k^*$.
- In CCA, $\text{PoV}(k)$ works up to $\min(p, q)$ non-zero eigenvalues achieved in the squared form with the basic or the SVD-based CCA algorithms.

Example 1: Synthetic Data

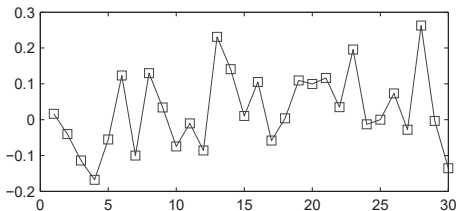
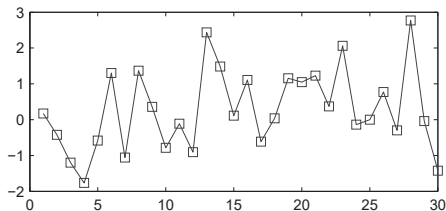
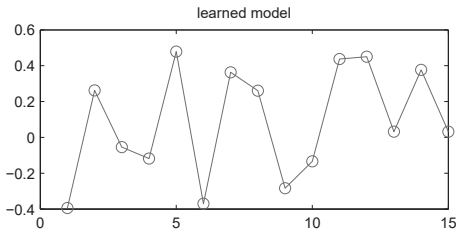
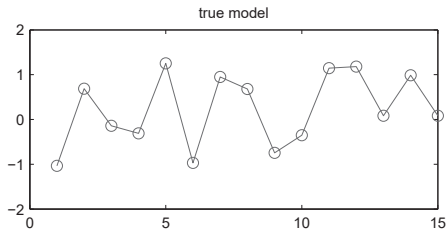
- Dataset: generated by $X_{15 \times 100} = \mathbf{a}_{15 \times 1} \mathbf{h}_{100 \times 1}^T$ and $Y_{30 \times 100} = \mathbf{b}_{30 \times 1} \mathbf{h}_{100 \times 1}^T$.
 \mathbf{a} , \mathbf{b} and \mathbf{h} are randomly chosen vectors.



ILLUSTRATIVE EXAMPLE

Example 1: Synthetic Data

- Dataset: generated by $X_{15 \times 100} = \mathbf{a}_{15 \times 1} \mathbf{h}_{100 \times 1}^T$ and $Y_{30 \times 100} = \mathbf{b}_{30 \times 1} \mathbf{h}_{100 \times 1}^T$.
- Apply CCA to the dataset to learn \mathbf{a} and \mathbf{b} in contrast to ground-truth



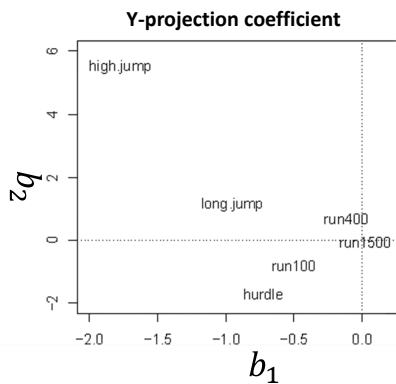
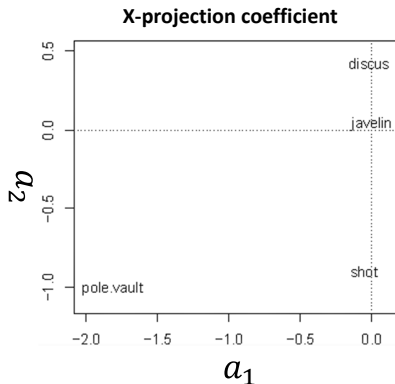
Example 2: Olympic Decathlon Data

- Dataset: men's 1988 Olympic decathlon (34 athletes, 10 events)
- Split the events into two subsets: **arm (X)** versus **leg (Y)** events:
 X: {shot, discus, javelin, pole.vault},
 Y: {run100, run400, run1500, hurdle, long.jump, high.jump}
- Apply CCA to achieve low-dimensional representations for data interpretation
- To make interpretation consistent for all events, **change sign** of all the running time (+ \rightarrow -) to indicate the larger the the better for all 10 events consistently

	run100	long.jump	shot	high.jump	run400	hurdle	discus	pole.vault	javelin	run1500	score
Schenk	-11.25	7.43	15.48	2.27	-48.90	-15.13	49.28	4.7	61.32	-268.95	8488
Voss	-10.87	7.45	14.97	1.97	-47.71	-14.46	44.36	5.1	61.76	-273.02	8399
Steen	-11.18	7.44	14.20	1.97	-48.29	-14.81	43.66	5.2	64.16	-263.20	8328
Thompson	-10.62	7.38	15.02	2.03	-49.06	-14.72	44.80	4.9	64.04	-285.11	8306
Blondel	-11.02	7.43	12.92	1.97	-47.44	-14.40	41.20	5.2	57.46	-256.64	8286
Plaziat	-10.83	7.72	13.58	2.12	-48.34	-14.18	43.06	4.9	52.18	-274.07	8272
Bright	-11.18	7.05	14.12	2.06	-49.34	-14.39	41.68	5.7	61.60	-291.20	8216
De.Wit	-11.05	6.95	15.34	2.00	-48.21	-14.36	41.32	4.8	63.00	-265.86	8189
Johnson	-11.15	7.12	14.52	2.03	-49.15	-14.66	42.36	4.9	66.46	-269.62	8180

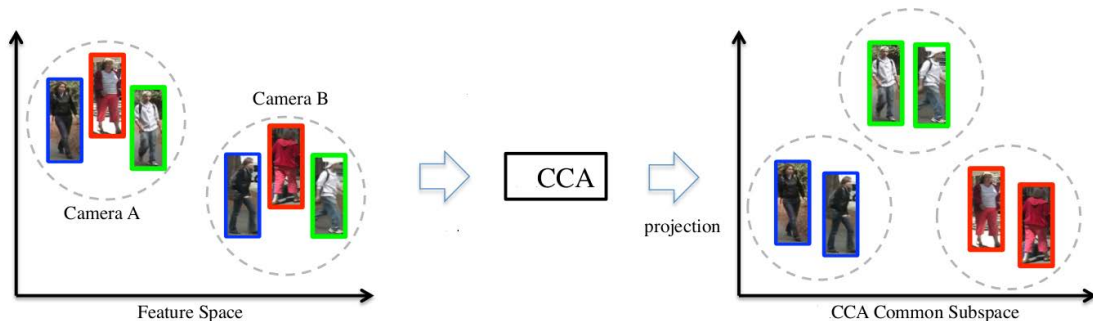
Example 2: Olympic Decathlon Data

- Dataset: men's 1988 Olympic decathlon (34 athletes, 10 events)
- Split the events into two subsets: arm (X) versus leg (Y) events:
 $X = \{\text{shot, discus, javelin, pole.vault}\},$
 $Y = \{\text{run100, run400, run1500, hurdle, long.jump, high.jump}\}$
- Projection matrices for $M = 2$, $A_{4 \times 2}$ and $B_{6 \times 2}$



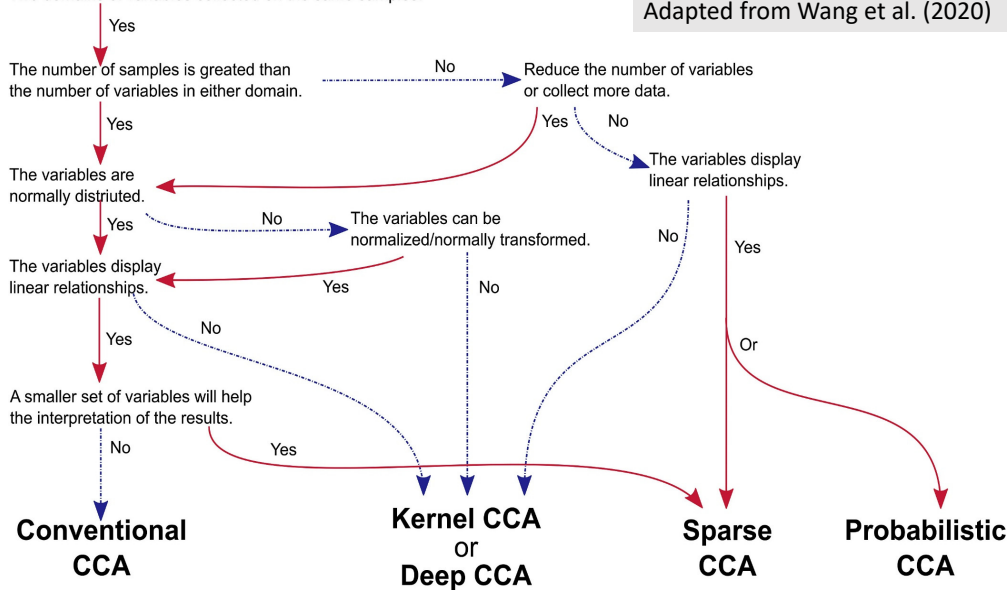
Example 3: Multi-view Alignment in Computer Vision

- Dataset: images captured by two cameras from different angles
- CCA leads to the common latent space that “groups” the correlated “objects”.



Adapted from Wang et al. (2020)

Two domains of variables collected on the same samples.



If you want to deepen your understanding and learn something beyond this lecture, you can self-study the optional references below.

[Alpaydin, 2014] Alpaydin E. (2014): *Introduction to Machine Learning* (3rd Ed.), MIT Press. (Section 6.9)

[Barber, 2012] Barber D. (2012): *Bayesian Reasoning and Machine Learning*, Cambridge University Press. (Sections 15.8)

[Uurtio et al., 2017] Uurtio V. et al. (2017): A tutorial on canonical correlation methods. *ACM Computing Survey*, Vol. 50, No. 6, Article 95.