# MLDL

SAND2018-8213C

# Machine Learning and Deep Learning Conference 2018

# XPCA: Extending PCA for Combinations of Discrete and Continuous Data
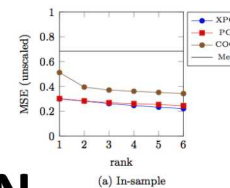
- Kina Kincher-Winoto/8762
- Tammy Kolda (8759), Cliff Anderson-Bergman (LLNL)

# The Project

MXD is a project that aims to provide alternatives to principal component analysis (PCA) to MiXeD data types

# Our Problem

## Principal Component Analysis (PCA)

- Standard Statistical Tool
- *Not all column marginals are gaussian*
- *Sensitive to scaling*
- *Sensitive to outliers*

Is this data Gaussian?

> ## Problem:
>
> **How can we relax the assumptions of PCA, ultimately to handle data that is binary, ordinal or continuous?**
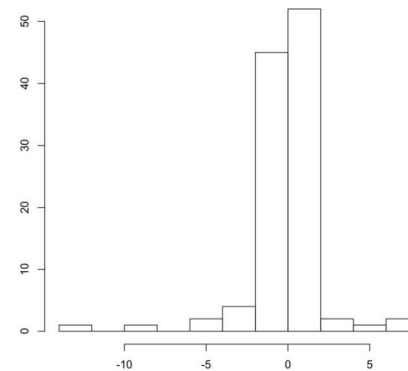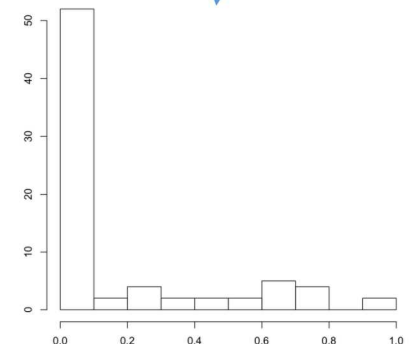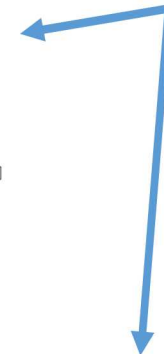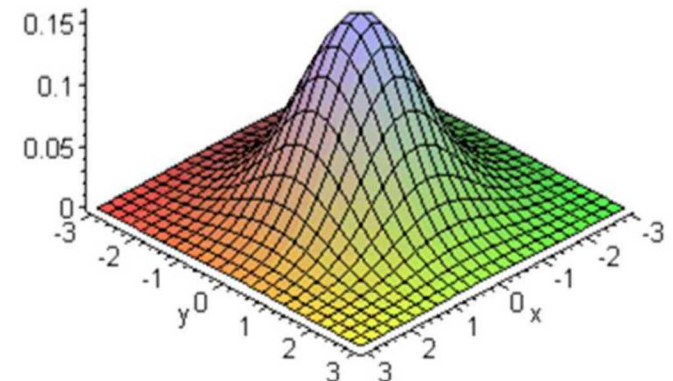
# Enter copulas

Copula (kɑ·pu·lə): A connecting word, in particular a form of the verb *be* connecting a subject and complement.

Copula (koʊ·pu·lə): A function that joins univariate distribution functions to form multivariate distribution functions. (Wolfram Mathworld)

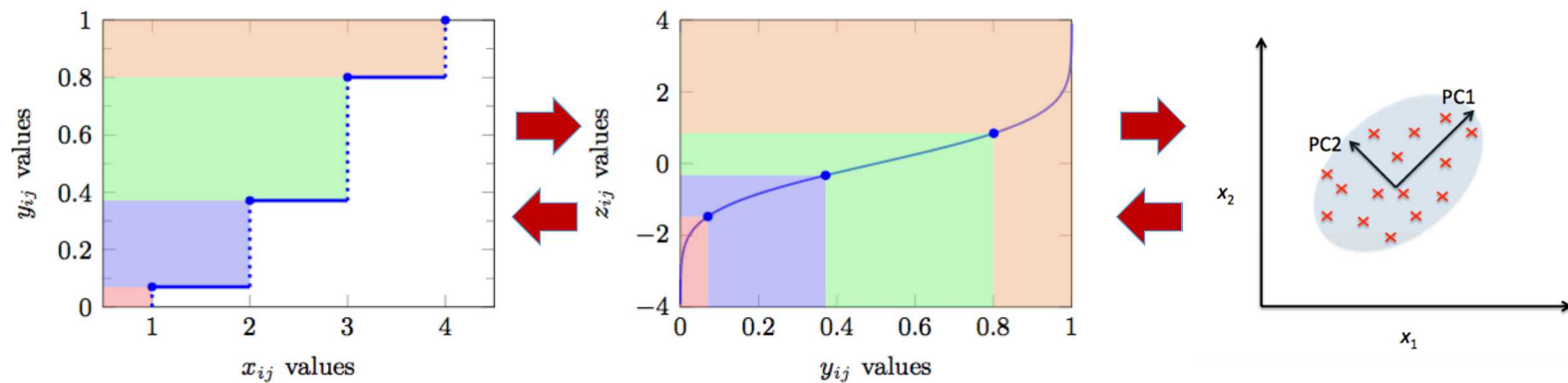A copula model provides the decomposition of the dependency of the marginal distributions such that the copula contains the dependency structure only…We uncouple variance and dependency structure such that PCA is only influenced by the dependency in the data. (Egger et al., 2016)

Gaussian copula: Assumes multivariate normal dependency.

# XPCA: eXtension of PCA

XPCA estimates the marginal distributions of each column and accounts for discrete variables in the likelihood calculation by integrating over appropriate intervals.
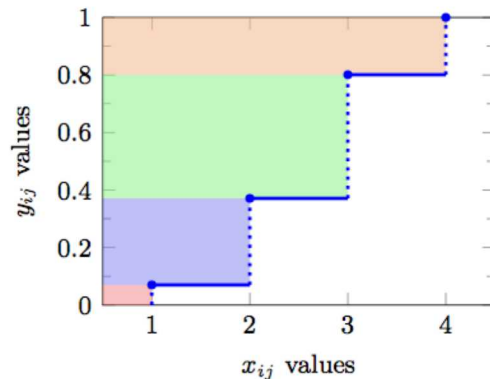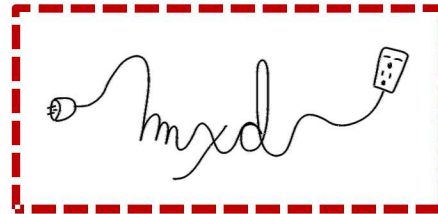


- XPCA uses copulas to better handle discrete values and outliers – both of which PCA struggles with
- XPCA is an improvement of academic research Copula Component Analysis (COCA)[1,2]

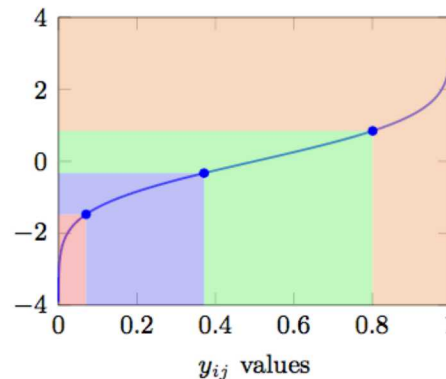[1] Han & Liu. Semiparametric principal component analysis. 2012
[2] Egger et al. Copula eigenfaces - semiparametric principal component analysis for facial appearance modeling. 2016
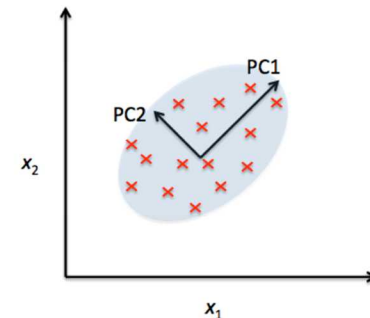
# XPCA: eXtension of PCA



**1** Empirical Distribution Function estimates marginal distribution.

**2** A *Gaussian Copula* relates non-Gaussian variables to Gaussian latent variables via Gaussian Cumulative Distribution Function
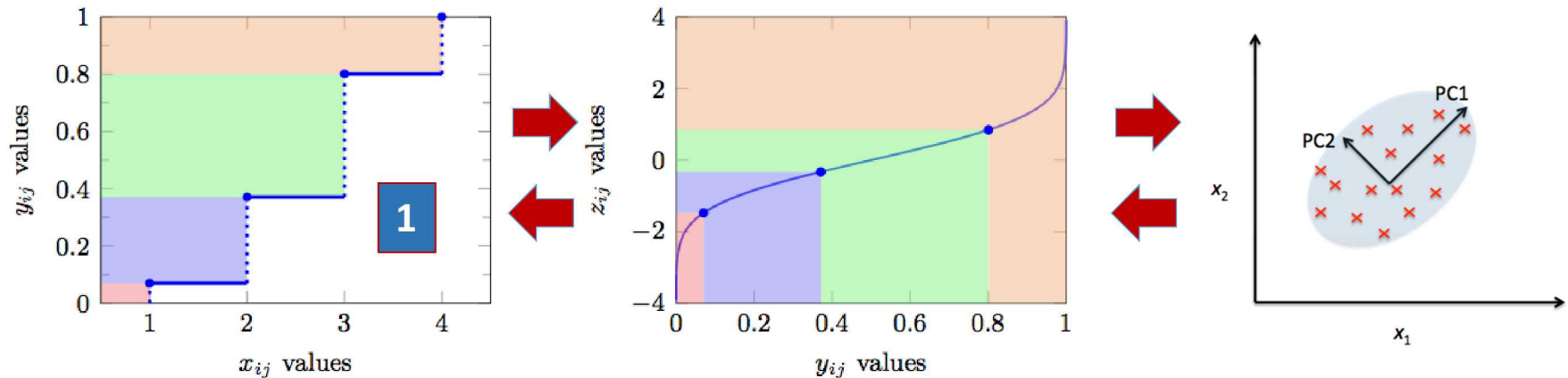
**3** *Latent variables* follow same model as probabilistic PCA BUT subject to a different loss

**4** Reverse each of those steps by taking inverse of each step

# XPCA: eXtension of PCA



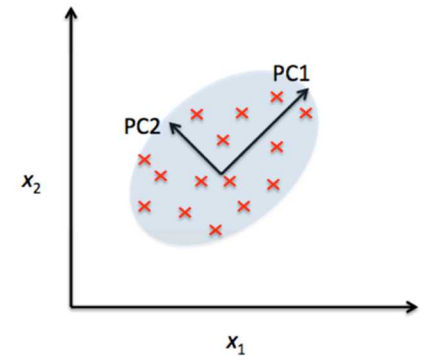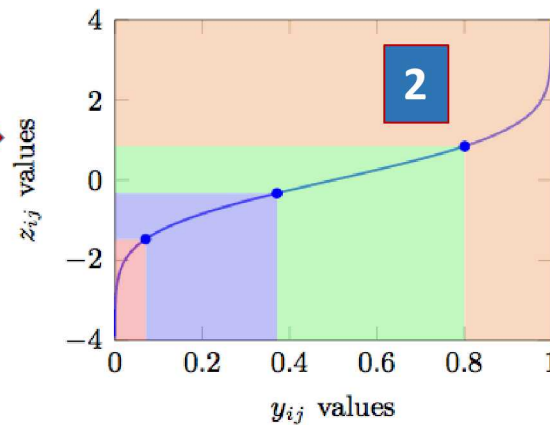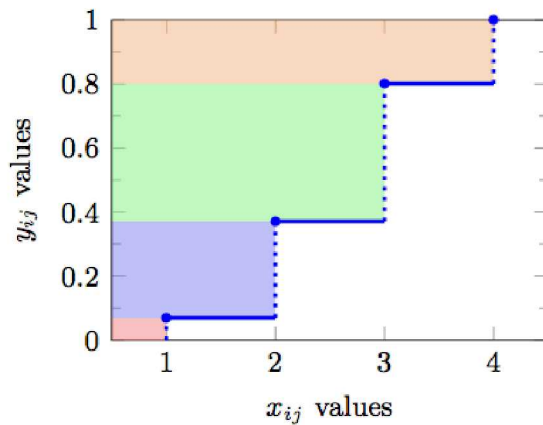**1**    Empirical Distribution Function estimates marginal distribution.

rank of $s$ in column $j$      unique values in column $j$

$$\hat{F}_j(x) = \frac{\max\left\{ r_j(s) \mid s \leq x \text{ and } s \in \mathcal{C}_j \right\}}{m_j},$$

total number of entries in column $j$

# XPCA: eXtension of PCA



A *Gaussian Copula* relates non-Gaussian variables to Gaussian latent variables via Gaussian Cumulative Distribution Function
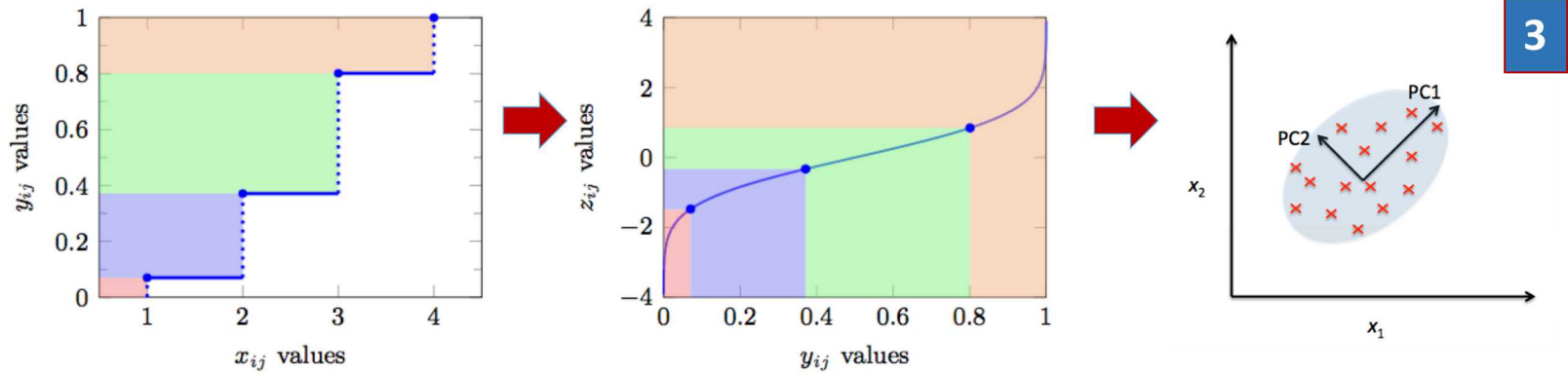
$$z_{ij} \in \left( \Phi^{-1}(\hat{F}_j(x_{ij} - \epsilon)), \Phi^{-1}(\hat{F}_j(x_{ij})) \right)$$

left side of interval          right side of interval
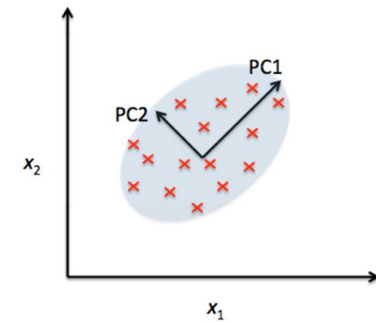
# XPCA: eXtension of PCA



**3**

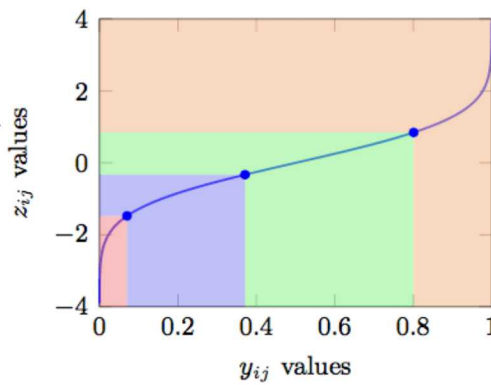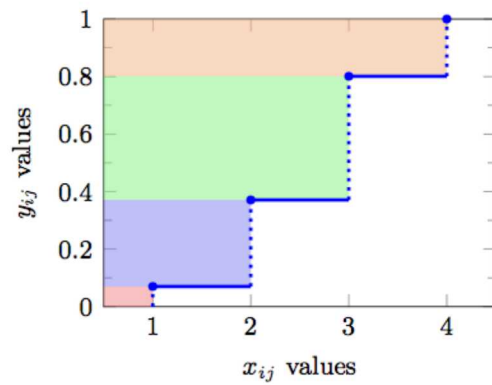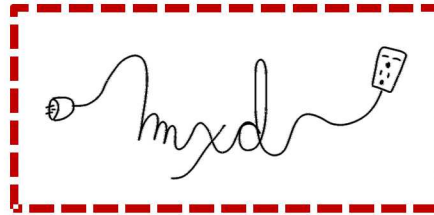*Latent variables* follow same model as probabilistic PCA BUT subject to a different loss.

$$\mathbf{Z} \sim \mathrm{MVN}\left(\boldsymbol{\Theta}, \sigma^2 \mathbf{I}\right) \quad \text{where} \quad \boldsymbol{\Theta} = \mathbf{U}\mathbf{V}^T$$

σ captures the magnitude of variance on top of the low rank structure

# XPCA: eXtension of PCA

**4** Reverse each of those steps by taking inverse of each step

# Application: the data

**1** Senator Voting Data from January 1989 – 2017
- 271 senators and their votes over $101^{st}$ – $114^{th}$ congressional sessions
  - 1 : Yay | -1: Nay | 0: abstain
- 271 rows (senators) x 9044 columns (bills)
- 63% of data is missing

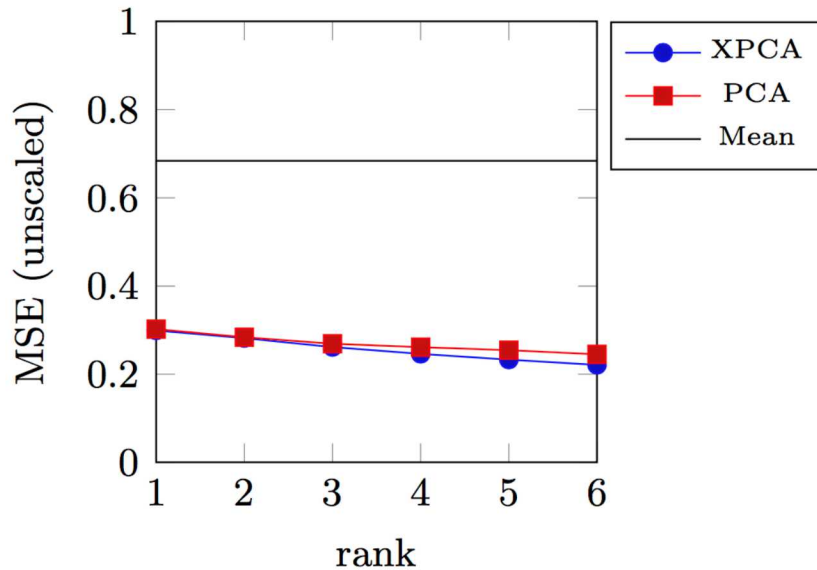www.senate.gov/legislative/votes.htm

**2** NBA Basketball Statistics from 2015-16 season
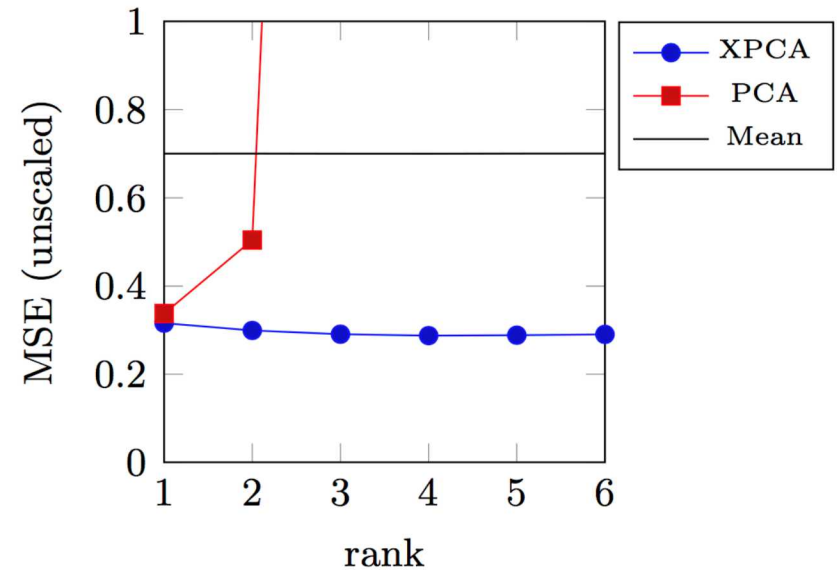- 476 players and their various statistics
  - e.g. shots scored, number of assists, draft number
- 476 rows (players) x 40 columns (stats)
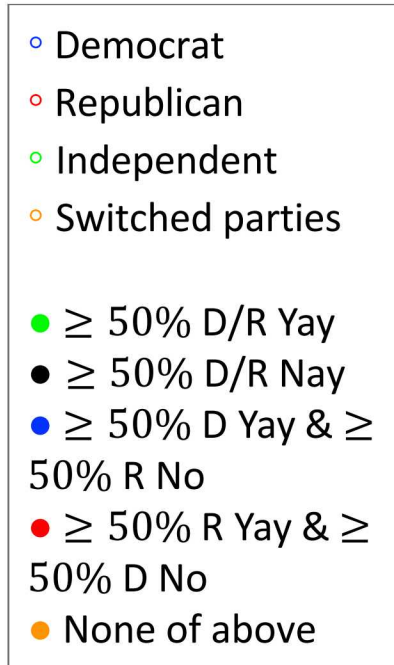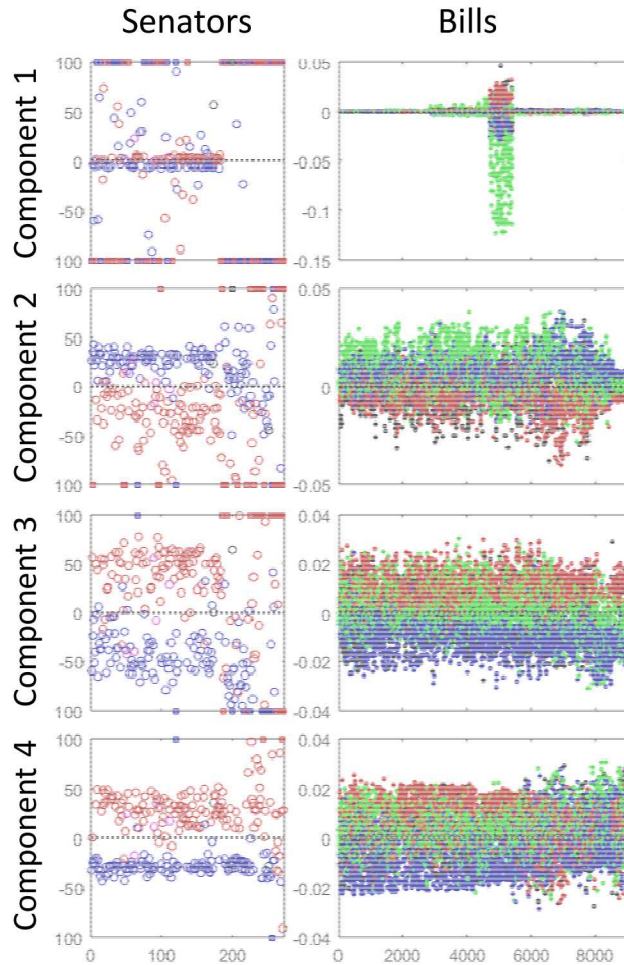
www.nba.com/stats

# Senator Data: model fit



**In-Sample Error**
the amount of error when fit is made across all available data

**4 Fold Cross-Validation Error**
the amount of error when ¼ of data is left out, calculated for each quarter

# Senator Data: components



PCA

Senators | Bills

Component 1
Component 2
Component 3
Component 4

- ○ Democrat
- ○ Republican
- ○ Independent
- ○ Switched parties

- ● ≥ 50% D/R Yay
- ● ≥ 50% D/R Nay
- ● ≥ 50% D Yay & ≥ 50% R No
- ● ≥ 50% R Yay & ≥ 50% D No
- ● None of above

XPCA

Senators | Bills

Component 1
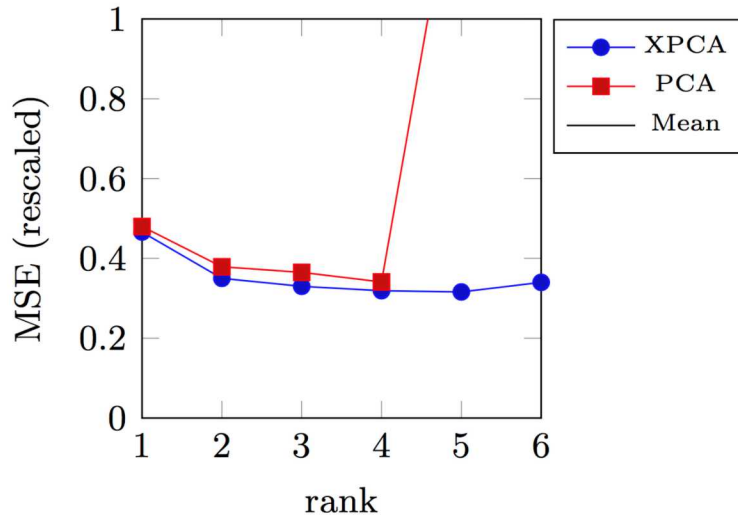Component 2
Component 3
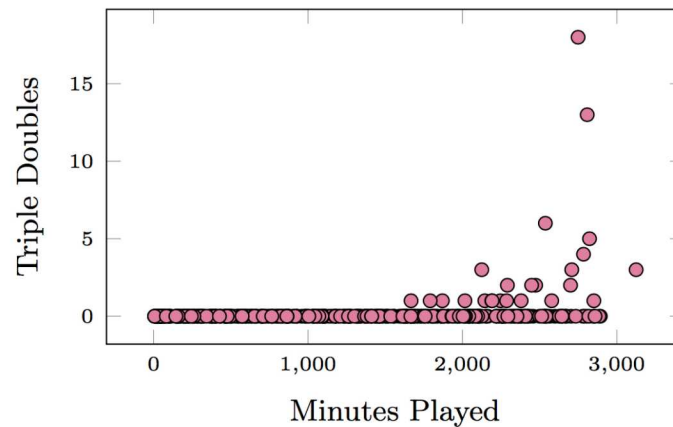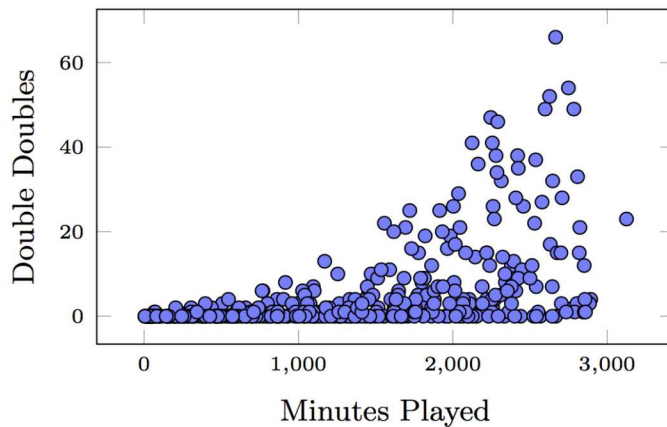Component 4

# NBA Basketball Data: Model fit



**20 Fold Cross-Validation Error**
the amount of error when ¼ of data is left out, calculated for each quarter

| | Kevin Durant: MVP? | Anthony Morrow: MVP? |
|---|---|---|
| True Value | 1 | 0 |
| XPCA estimate at best rank (r=5) | 0.3833959 | 0 |
| PCA estimate at best rank (r=4) | 0.1427923 | -0.05687824 |

PCA estimates values outside range

# NBA Basketball Data: Components

- The first component of both XPCA and PCA represent the minutes played and the count variables (e.g. shots made or times fouled or double doubles).
- However, XPCA also captures "triple doubles" influence
  - PCA cannot capture this influence possibly due to the heavy atom at 0



These is a clear correlation between minutes played and triple doubles. However, PCA struggles to find this correlation, while XPCA does not.

# Top 5 last thoughts

**Ongoing Threads:**

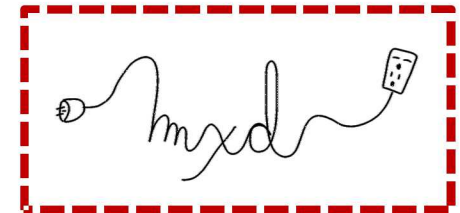    Application into mission problems

    Clustering post-XPCA in the analysis stage!

    Scaling up to tens of thousands of rows & columns

**Future Work:**

    Pipeline into other machine learning applications

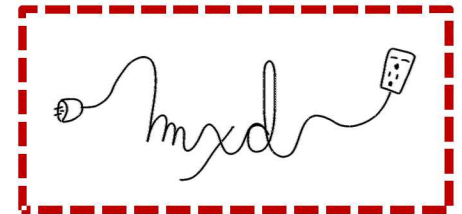    Scaling up to millions of rows and columns

# Learn more

**To read more:**

Paper is written (currently in R&A) and is available to read!

**To try it out:**

MXDLIB, the software package, is written in both R (stable) and Python (testing)

**POC:** Kina Kincher-Winoto
kwinoto@sandia.gov