# On Robustness of Principal Component Regression

**Anish Agarwal**

**Devavrat Shah, Dennis Shen, Dogyoon Song**

**MIT**

1

# 1 What is PCR?
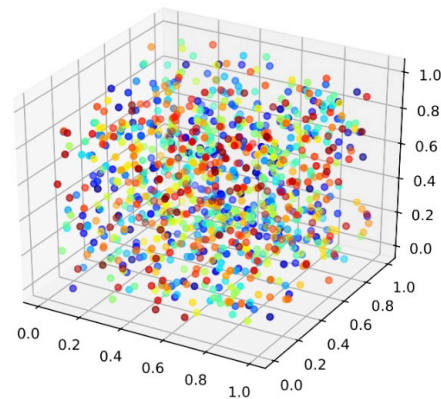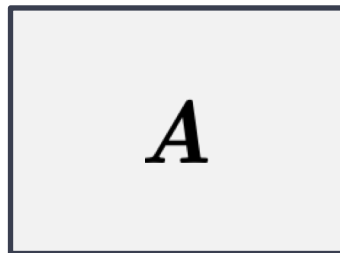
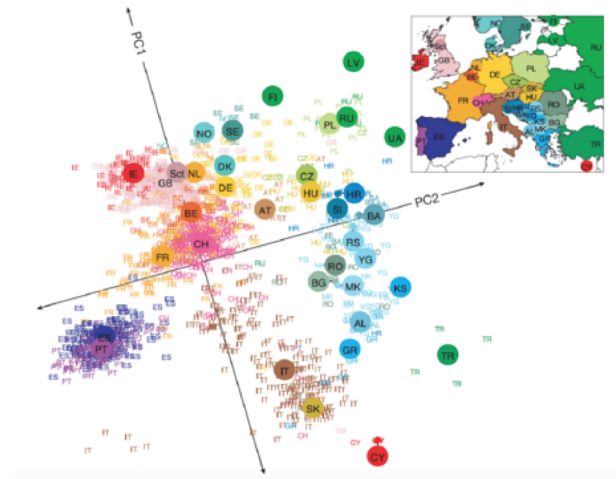$$Y = A\beta + \epsilon$$

$$Y = A \beta + \epsilon$$

Step 1: PCA

$$A$$

$$Y = A\, \beta + \epsilon$$

## Step 1: PCA
($k$-components)

$$A_k$$

# What is PCR?

$$Y = A_k \beta + \epsilon$$

## Step 2: Regression

$$\beta_k = \underset{\theta}{\text{minimize}} \; \left\| Y - A_k \theta \right\|_2^2$$

# What is PCR?

$$Y = A_k \, \beta_k + \epsilon$$

## Step 3: Prediction

$$\widehat{Y} = A_k \, \beta_k$$

# 2 When & Why Use PCR

*"IF DATA IS (APPROXIMATELY) LOW-DIMENSIONAL, USE PCR!"*

*-- Anonymous Data Scientists*

When exactly should we be using PCR?

## 2 Key Questions We Answer

Theoretical properties of PCR?

Is dimension-reduction only benefit to PCR?

Our Theoretical Analysis of PCR helps answer following questions..

How low-rank do covariates need to be?

How many principal components to pick?

How well does PCR perform on a test data (i.e. generalization properties)?
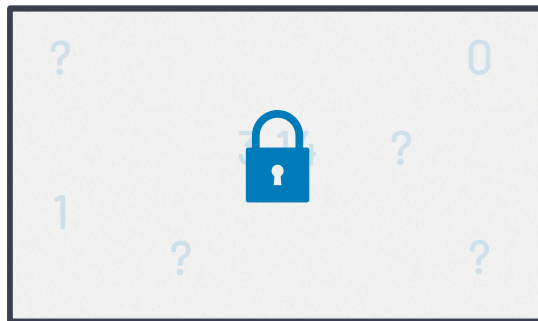
# Is Dimension-Reduction Only Benefit?

# NO!

Noisy

Missing

Mixed valued

Sensitive

# We show PCR is surprisingly robust to problems that plague large-scale modern datasets

## Main Contribution of this Work

# Error-In Variable Regression

## (Setting We Consider)

$$Y = A\,\beta + \epsilon$$

$$Y = Z\beta + \epsilon$$

Representative of modern datasets

Noisy
Missing
Mixed valued
Sensitive

Time Series Analysis (measurement noise)

Causal Inference (Synthetic Control) (measurement noise)

Differentially-private Regression (noise by design)

Mixed Valued Regression (structural noise)

Time Series Analysis (measurement noise)

Causal Inference (Synthetic Control) (measurement noise)

**Differentially-private Regression (noise by design)**

Mixed Valued Regression (structural noise)

# Formal Results

# Theorem (Informal): Training Error

If principal components chosen correctly ($k = r$)

PCR implicitly denoises covariates!

number of covariates

$$\frac{1}{n}\mathbb{E}\left[\|\boldsymbol{\widehat{A}_r}\beta_r - \boldsymbol{A}\beta\|_2^2\right] \sim \frac{\sigma^2 r}{n}\frac{\mathrm{poly}(\log p)}{\rho^4}$$

fraction of observations

OLS minmax error rate
(low-dimensional, noiseless, fully observed covariates)

If principal components **not** chosen correctly ($k \neq r$)

$$\text{Test Error} \sim \text{Train Error with PCR}(k) + C\frac{k^{3/2}}{\sqrt{n}}$$

PCR implicitly performs $l_0$-regularization

PCR implicitly de-noises covariates
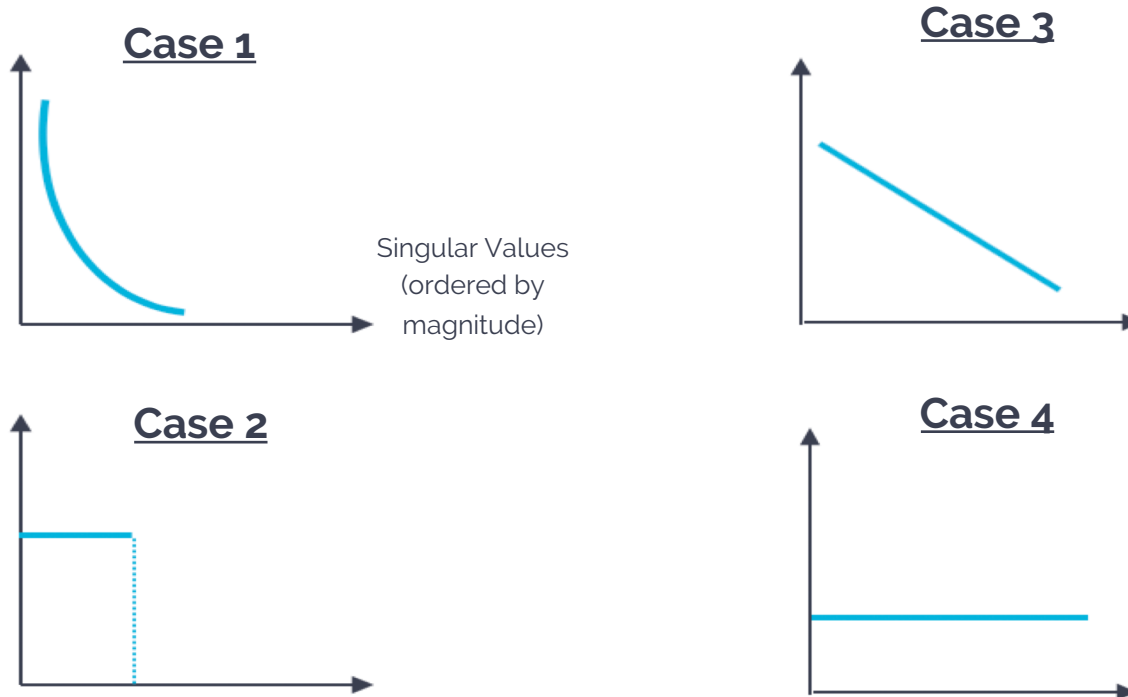
Choose $k$ that minimizes above

# Use PCR!

# Don't Use PCR!

**Case 1**

Magnitude of
Singular Values

Singular Values
(ordered by
magnitude)

**Case 3**

**Case 2**

**Case 4**

Magnitude of
Singular Values

**Case 1**

Singular Values
(ordered by
magnitude)

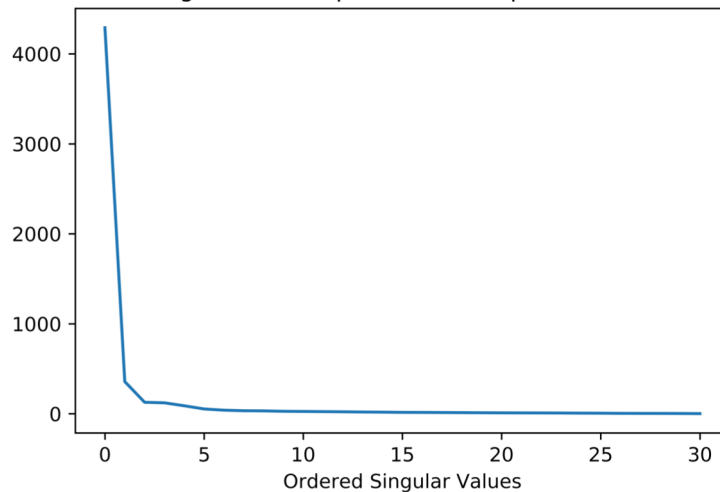## GDP Trajectories (Macroeconomics)



Singular Value Spectrum of Basque Country Dataset



Singular Value Spectrum of Prop 99 Dataset

Avito Ad-Click Dataset (E-Commerce)

Magnitude of
Singular Values

**Case 1**

Singular Values
(ordered by
magnitude)



Eigenspectrum of Avito Context Ads Dataset (190 M Clicks)

# Exponential-decaying spectrum is ubiquitous in real-world data

## Cricket Trajectories (Sports)
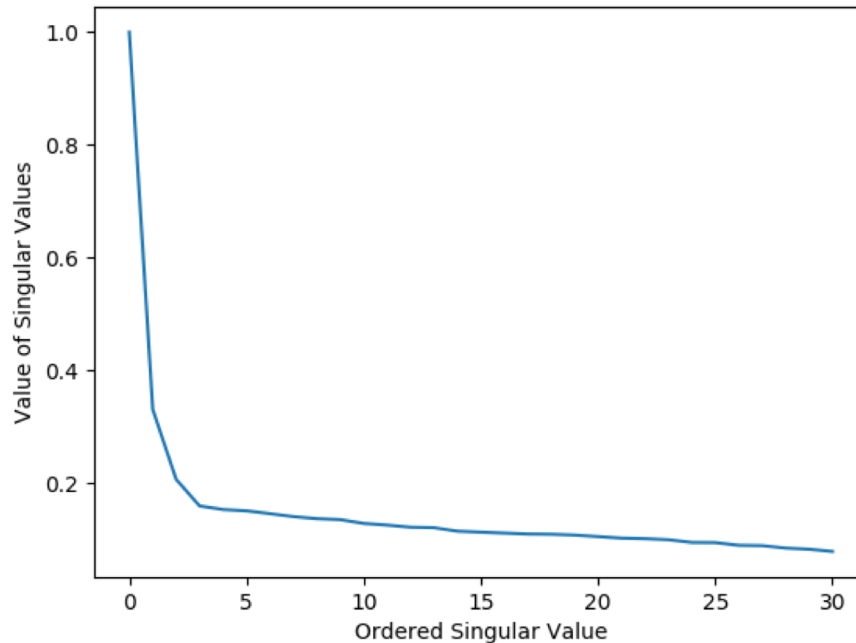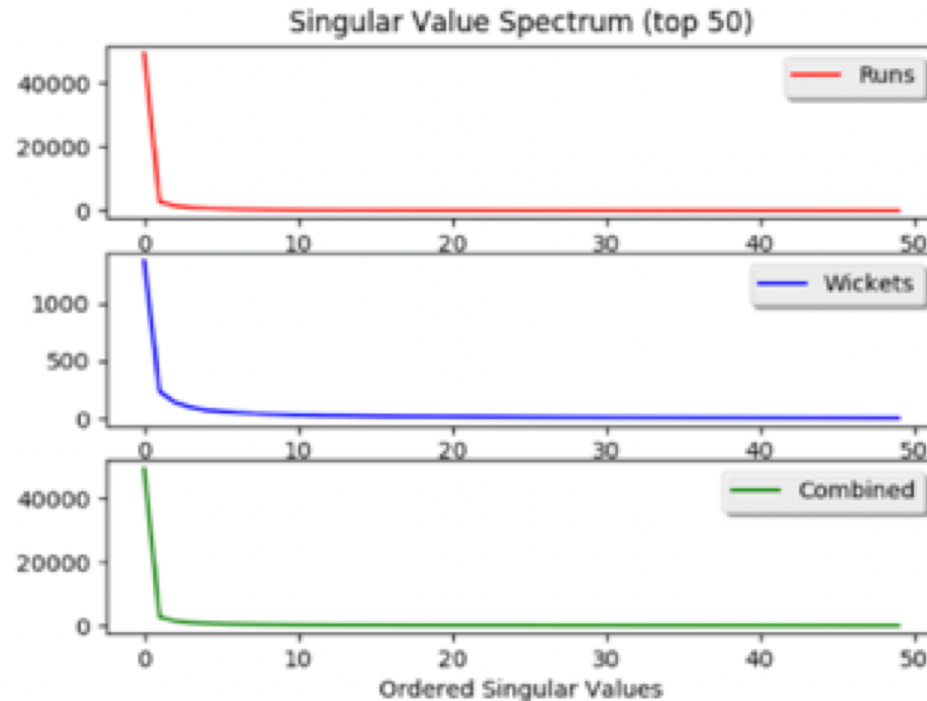


Singular Value Spectrum (top 50)



Case 1

Magnitude of Singular Values

Singular Values (ordered by magnitude)

# 3 Surprising Applications of PCR

Time Series Analysis   (measurement noise)

Causal Inference (Synthetic Control) (measurement noise)

**Differentially-private Regression** **(noise by design)**

Mixed Valued Regression (structural noise)

Data privacy is top-of-mind as we increasingly apply ML on sensitive user data (genetic data, purchase history etc.)

# Standard Notion of Privacy in ML
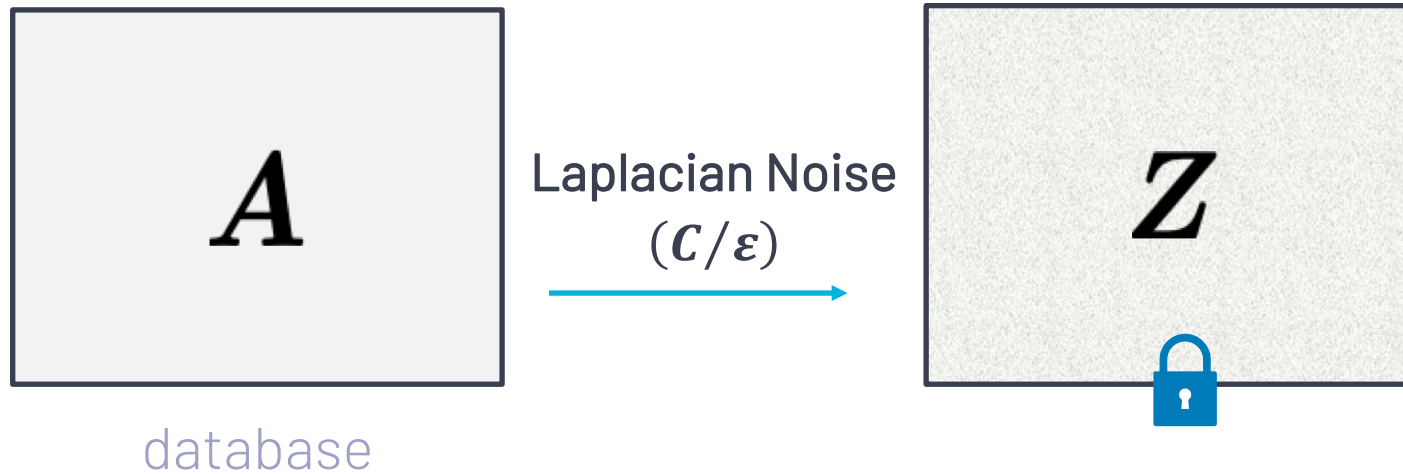
## ε-Differential Privacy

Intuitively, an algorithm is **ε**-differentially private if **outcome of a statistical query** on a database **cannot change by more than ε** due to **presence/absence of any user** data record

Example of Statistical Query:

*"Average Income of all users between ages 25 and 30"*

# How to achieve **ε**-differentially privacy?

## Laplace Mechanism



database

Laplacian Noise
$(C/\varepsilon)$

# Predictive Accuracy vs. Privacy Tradeoff

**Can we achieve good prediction error and still maintain privacy?**



$$Y = Z\beta + \epsilon$$

**Yes!**

# Predictive Accuracy vs. Privacy Tradeoff

**Can we achieve good prediction error and still maintain privacy?**

Step 1:

Data Owner adds Laplacian Noise
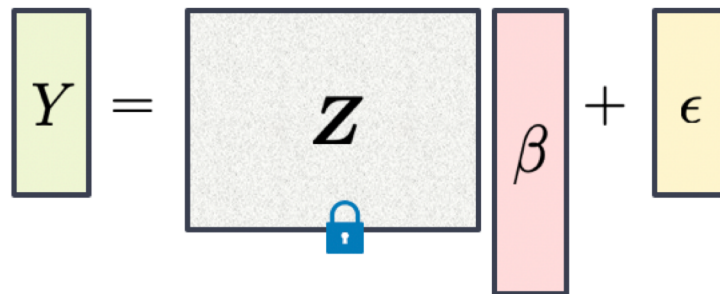


Step 2:

Analyst Performs PCR

**Done!**

# What is sample complexity cost for ε-differential privacy?

$$\text{Prediction Error} \quad \sim \quad \frac{\sigma^2 r}{n} \frac{\text{poly}(\log p)}{\rho^4} \left( \frac{1}{\varepsilon^2} \right)$$

## Does de-noising step (PCA) break privacy?

No, PCA only de-noises covariates *on average*

with respect to the $\| \cdot \|_{2,\infty}$ - norm

# 4 Conclusion

# Inspect spectrum of your covariate matrix



Magnitude of
Singular Values

**Case 1**

Singular Values
(ordered by
magnitude)

**Case 2**

Use PCR!

de-noises

regularizes

# Possible Implications for Modern ML

## Linear Case

Step 1: Dimension Reduction

$A$ → PCA → $A_k$

**Linear low-dimensional covariate pre-processing** has many implicit benefits (e.g. de-noising, regularizing)

## Non-Linear Case

$A$ → GANs?

Does non-linear covariate pre-processing (e.g. GANs) have similar benefits for unstructured data?

# Come Meet Us At Our Poster

**Poster #3 –** East Exhibition Hall B + C, 5-7pm, Thursday

## Shameless Plug :)

PCR for Time Series Analysis: **tspdb.mit.edu**

PCR for Causal Inference: **github.com/Romcos/SC_demo**