# Capturing biological patterns from gene expression data of PCOS using unsupervised dimensionality reduction algorithms

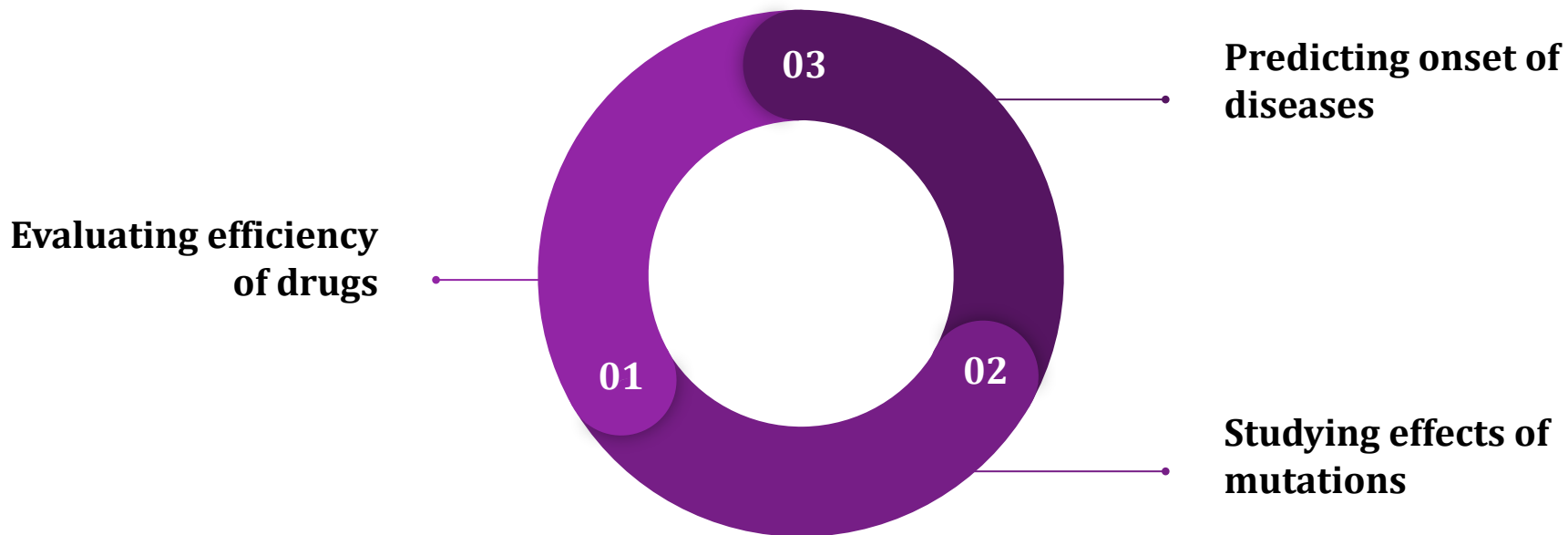N Sowmya Manojna (BE17B007)
A V Lakshmy (CS16B101)

December 14, 2020

# Introduction

**Gene expression analysis** is a powerful technique with many applications:



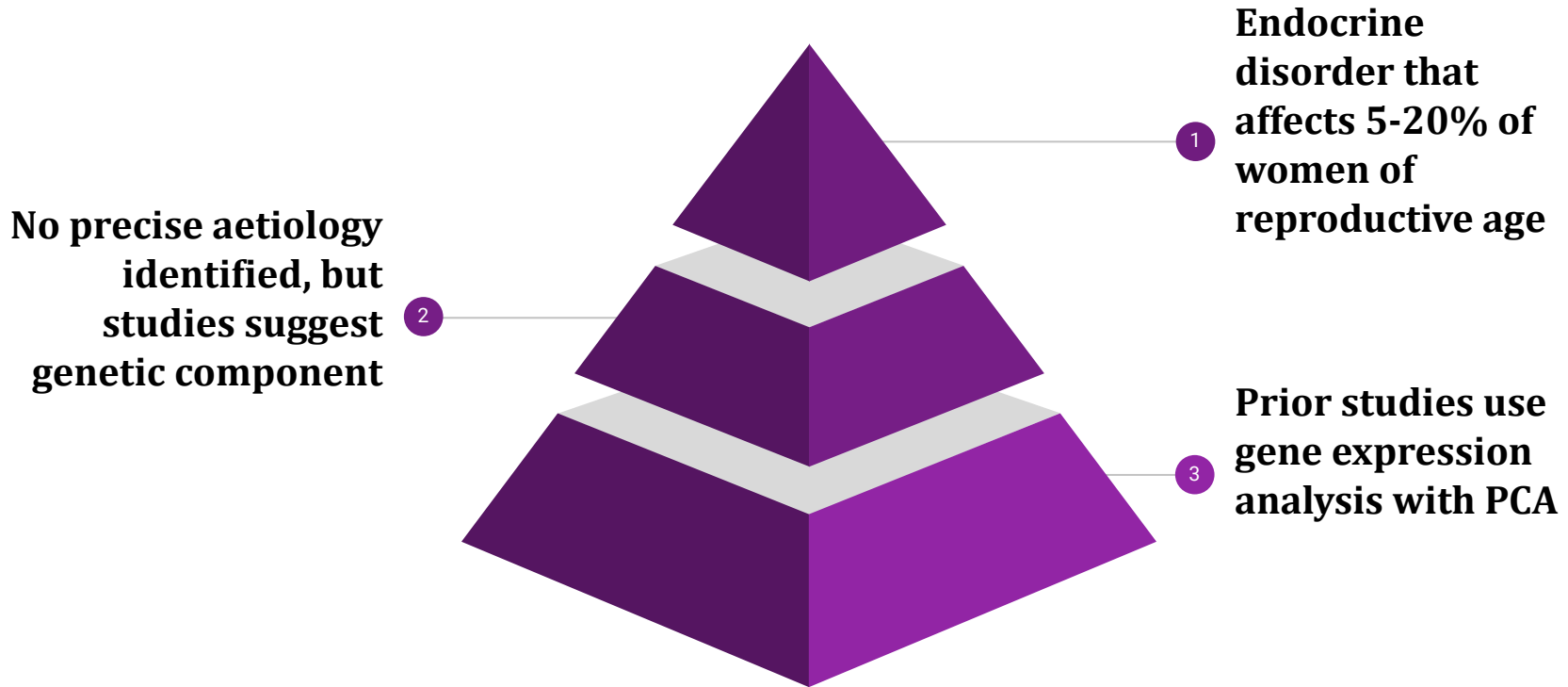Predicting onset of diseases

Evaluating efficiency of drugs

Studying effects of mutations

01 02 03

**Gene expression data analysis is computationally expensive!!**

Unsupervised Dimensionality Reduction Methods

- PCA, ICA, NMF
- Autoencoders: DAE, VAE

# Problem Statement

*"Capturing biological patterns from gene expression data of Polycystic Ovarian Syndrome (PCOS) using unsupervised dimensionality reduction algorithms, particularly autoencoders"*
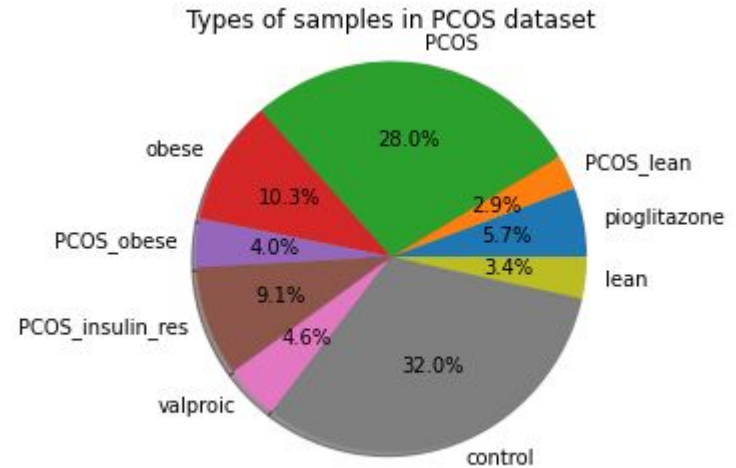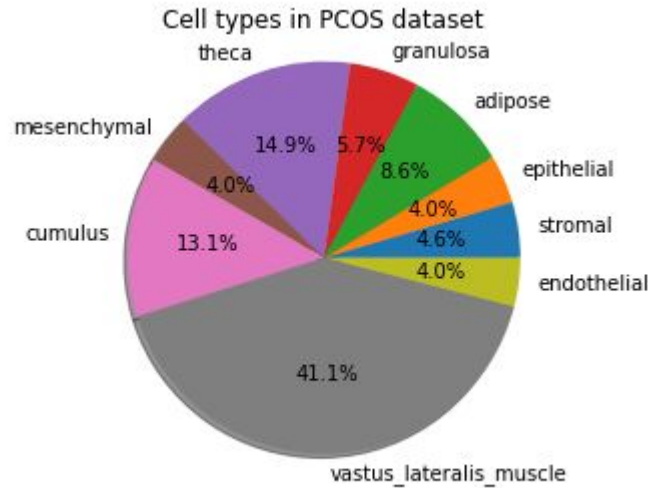
# What is PCOS? Why PCOS?

**Endocrine disorder that affects 5-20% of women of reproductive age**

1

**No precise aetiology identified, but studies suggest genetic component**

2

**Prior studies use gene expression analysis with PCA**

3

# Datasets

Cell types in PCOS dataset

theca · granulosa · adipose · epithelial · stromal · endothelial · vastus_lateralis_muscle · cumulus · mesenchymal

14.9% · 5.7% · 8.6% · 4.0% · 4.6% · 4.0% · 41.1% · 13.1% · 4.0%

Types of samples in PCOS dataset

PCOS · obese · PCOS_obese · PCOS_insulin_res · valproic · control · lean · pioglitazone · PCOS_lean

28.0% · 10.3% · 4.0% · 9.1% · 4.6% · 32.0% · 3.4% · 5.7% · 2.9%

Total 175 samples from 9 different datasets; 1671 gene expressions

Heterogeneous network analysis on merged PCOS dataset*

*Not all IDs are mapped



Most genes that play a role in PCOS are largely related to **immunological** datasets!

# Techniques Used

**1**

PCA, ICA, NMF initial results

**2**

DAE, VAE parameters and reconstruction costs

**3**

Gene set coverage, stability and SVCCA

**4**

Biological interpretations

* Way, Gregory P., Michael Zietz, Vincent Rubinetti, Daniel S. Himmelstein, and Casey S. Greene. "Compressing gene expression data using multiple latent space dimensionalities learns complementary biological representations." Genome Biology 21, no. 1 (2020): 1-27.

# Techniques Used

## PCA



## ICA

- The dimensions used are "independent to each other"
- Like a rotation of PCA

## NMF

- Like PCA, but, except the coefficients in the linear combination must be non-negative
- Dimensions that don't contribute much have a zero coefficient.

DAE models in addition, set a random fraction of the input data to 0.

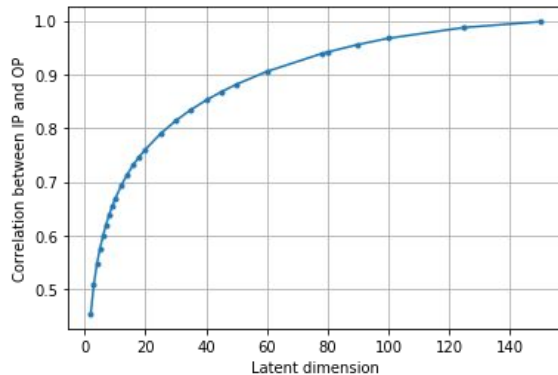# Key results and interpretations

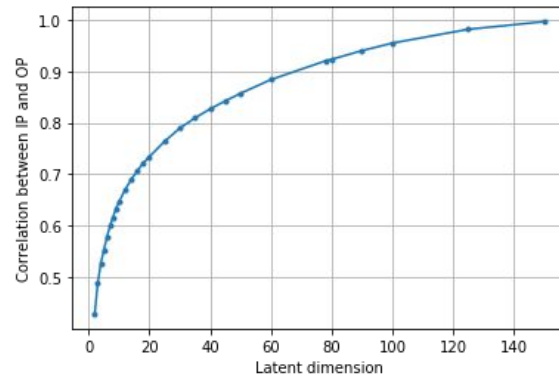# PCA, ICA, NMF - Reconstruction Cost

# PCA, ICA, NMF - Strongly Associated Dimensions

# PCA, ICA, NMF - SVCCA

# PCA, ICA, NMF Correlation and Stability

Reconstruction cost for DAE model

Reconstruction cost for VAE model

# Challenges and Future Work

## Challenges

- Small dataset
- Absence of all gene mappings

## Future Work

- Visualizing plots
- Modifications of autoencoder models

# Thank you!
# Any questions?