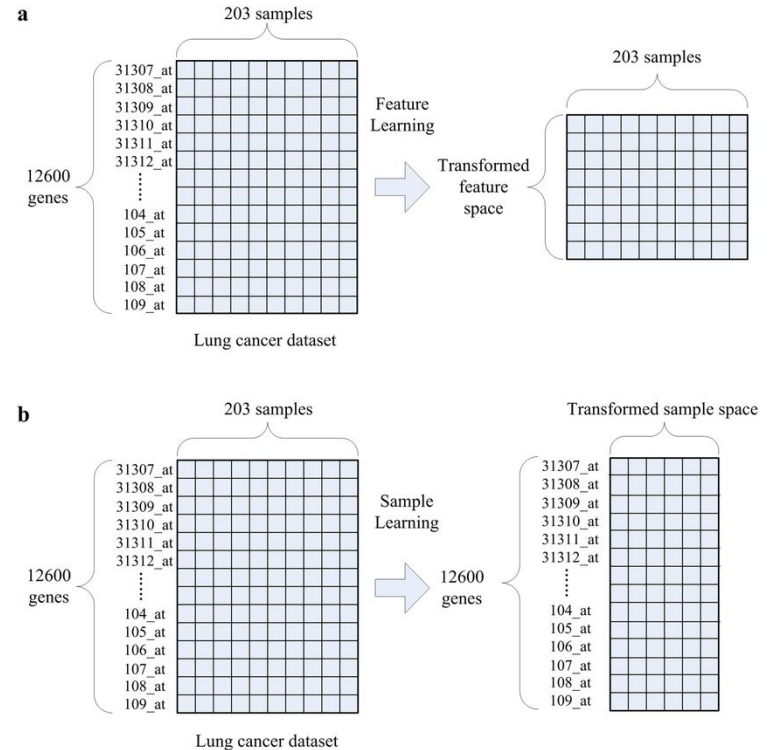# High Dimensional Data Classification

Mentee: Elvin Liu
Mentor: Zhaoxing Wu

# Overview

1. **Introduction**
2. **Background** (Fisher's Linear Discriminant Analysis)
3. **Solution** (Penalized Discriminant Analysis Projection Pursuit)
4. **Analysis**

# Premise and Motivation

- Classification is predicting labels from a dataset's features
- Prediction algorithms rely on sufficiently *large* **sample sizes** (n) to train such **features** (p)
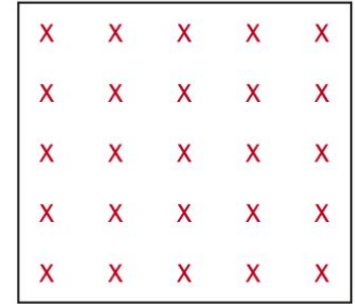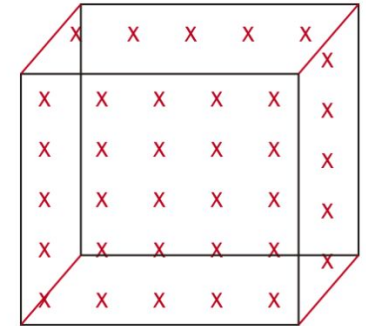
# Curse of Dimensionality

Accuracies of classification algorithms tend to dip in high dimensions due to the **curse of dimensionality**



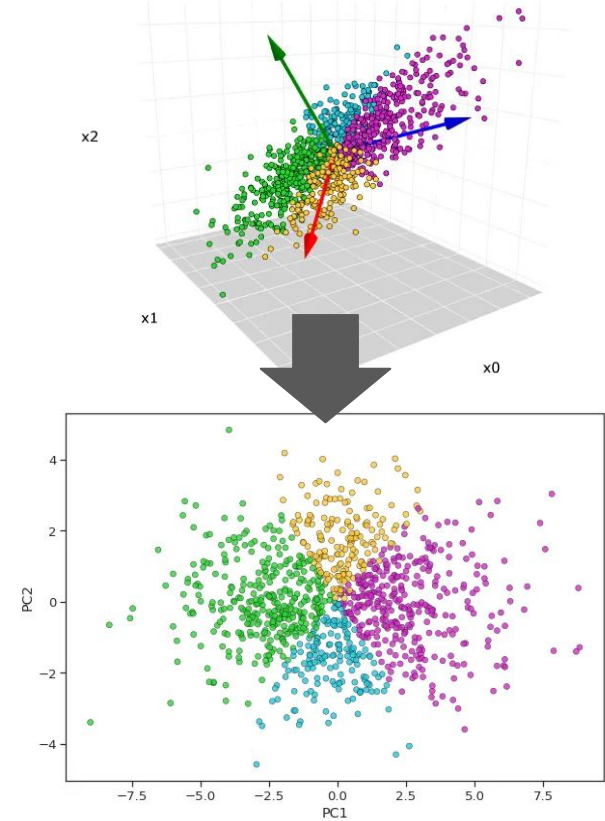*A one-dimensional features space with five data points*



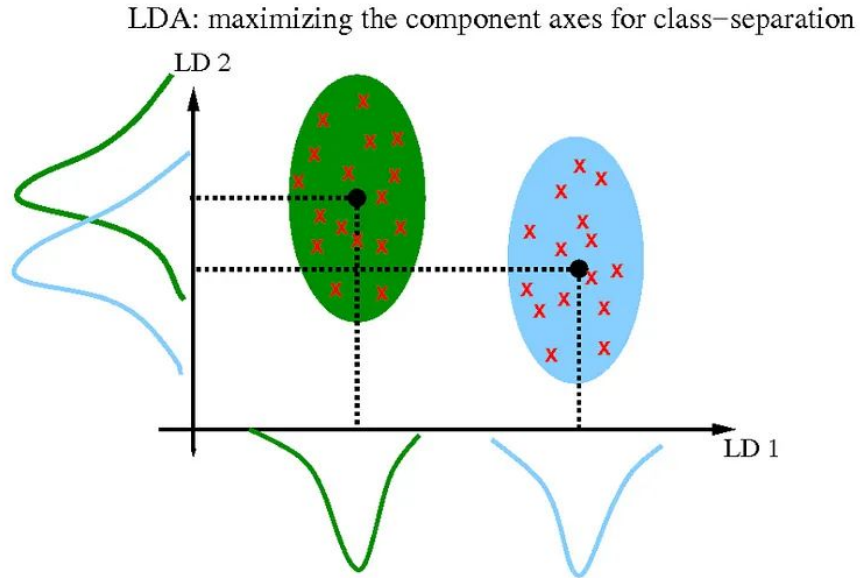*A two-dimensional features space with 25 data points*



*A three-dimensional features space with 125 data points*

V. Pappu, & P.M. Pardalos. High Dimensional Data Classification. Clusters, Orders and Trees: Methods and Applications. 4

# Projection Pursuit

● Seek interesting *low-dimensional* projections



Lee, E.-K., & Cook, D. (2009). A projection pursuit index for large P Small N Data. Statistics and Computing, 20(3), 381–392. 5

# Fisher's Linear Discriminant Analysis (LDA)



LDA: maximizing the component axes for class−separation

- Supervised dimensionality reduction
- Reduce dimensions but preserve features relevant for class discrimination

Lee, E.-K., & Cook, D. (2009). A projection pursuit index for large P Small N Data. Statistics and Computing, 20(3), 381–392. 6

# Between-Class and Within-Class Scatter Matrices

1. The **between-class** scatter matrix $S_B$ is computed by the following equation,

$$S_B = \sum_{i=1}^{c} N_i (m_i - m)(m_i - m)^T$$

where **m** is the overall mean, and $m_i$ and $N_i$ are the sample mean and sizes of the respective classes.

Lee, E.-K., & Cook, D. (2009). A projection pursuit index for large P Small N Data. Statistics and Computing, 20(3), 381–392.

# Fisher's LDA

Intuitively, makes sense to *maximize* **between-class** scatter matrix and *minimize* **within-class** scatter matrix.



Lee, E.-K., & Cook, D. (2009). A projection pursuit index for large P Small N Data. Statistics and Computing, 20(3), 381–392. 8

# Fisher's LDA

$$\mathbf{Q}\mathbf{v} = \lambda\mathbf{v}$$

where…

$$\mathbf{Q} = \mathbf{S_w}^{-1}\mathbf{S_B}$$

$\mathbf{v}$ = eigenvector

$\lambda$ = eigenvalue

*Eigenvectors* of largest eigenvalues *maximize* **between-class/within-class variance** most

*k* most important eigenvectors of $\mathbf{S_w}^{-1}\mathbf{S_B}$ are **linear discriminant function**

Projected data trims least important features

$$I_{\mathrm{LDA}}(\mathbf{A}) = \begin{cases} 1 - \frac{|\mathbf{A}^T\mathbf{S_w}\mathbf{A}|}{|\mathbf{A}^T(\mathbf{S_w}+\mathbf{S_B})\mathbf{A}|}, & \text{for } |\mathbf{A}^T(\mathbf{S_w}+\mathbf{S_B})\,\mathbf{A}| \neq 0, \\ 0, & \text{for } |\mathbf{A}^T(\mathbf{S_w}+\mathbf{S_B})\,\mathbf{A}| = 0, \end{cases}$$

Lee, E.-K., & Cook, D. (2009). A projection pursuit index for large P Small N Data. Statistics and Computing, 20(3), 381–392. 9

# Penalized Discriminant Analysis

n < p leads to data piling issues in LDA

Lee, E.-K., & Cook, D. (2009). A projection pursuit index for large P Small N Data. Statistics and Computing, 20(3), 381–392. 10

# PDA PP-Index

Use $\tilde{\Sigma}(\lambda) = (1 - \lambda)\hat{\Sigma} + \lambda \cdot \mathbf{diag}(\hat{\bar{\Sigma}})$ to estimate the variance-covariance matrix

$\mathbf{\Sigma}$ (**hat**) = maximum likelihood estimator (MLE) of $\mathbf{\Sigma}$

Using standardized data vectors, the above equation turns to

$\tilde{\mathbf{R}}(\lambda) = (1 - \lambda)\hat{\mathbf{R}} + \lambda\mathbf{I},$ where

$\mathbf{R}$ (**hat**) = MLE of correlation matrix

$\mathbf{I}$ = identity matrix

Lee, E.-K., & Cook, D. (2009). A projection pursuit index for large P Small N Data. Statistics and Computing, 20(3), 381–392. 11

# PDA PP-Index

$$\mathbf{\Sigma^B} = \sum_{i=1}^{g} n_i (\bar{\mathbf{X}}^*_{i.} - \bar{\mathbf{X}}^*_{..})(\bar{\mathbf{X}}^*_{i.} - \bar{\mathbf{X}}^*_{..})^T,$$

$$\mathbf{\Sigma^S} = \sum_{i=1}^{g}\sum_{j=1}^{n_i} (\mathbf{X}^*_{ij} - \bar{\mathbf{X}}^*_{i.})(\mathbf{X}^*_{ij} - \bar{\mathbf{X}}^*_{i.})^T$$

$\mathbf{X}^*_i$ = $i$-th group mean of the standardized data and

$\mathbf{X^s}_{**}$ = 0 = total mean of the standardized data

$$I_{\text{PDA}}(\mathbf{A}, \lambda) = 1 - \frac{|\mathbf{A}^T\{(1-\lambda)\,\mathbf{\Sigma^S} + n\lambda\mathbf{I}_p\}\mathbf{A}|}{|\mathbf{A}^T\{(1-\lambda)\,(\mathbf{\Sigma^B} + \mathbf{\Sigma^S}) + n\lambda\mathbf{I}_p\}\mathbf{A}|}$$

Lee, E.-K., & Cook, D. (2009). A projection pursuit index for large P Small N Data. Statistics and Computing, 20(3), 381–392. 12

# Example Projections



Histogram on the left is 1-D projection
Scatter plot on the right is 2-D projection

# Sample Data Background

**Sample**:

- 18 malignant cancer tumors paired
- 18 normal tissue from the *same patient*

**Features**:

- 3200 full-length human cDNA
- 3400 expressed sequence tags

**Classifier**: Malignant Tumor | Normal Tissue

**Features (≈6600)**



**Samples (36)**

Notterman, et al, Cancer Research, vol. 61: 2001 14

# Training/Testing Split

Let's split *training*/*testing* in a 70:30 ratio

**Training dataset**:

**25** samples

**13** normal tissue/**12** tumors

**Original dataset**:

**36** samples

**18** normal tissue/**18** tumors

70%

30%

**Testing dataset**:

**11** samples

**5** normal tissue/**6** tumors

# Predictions

| | Feature 1 | Feature 2 | Feature 3 | … |
|---|---|---|---|---|
| Tumor μ | 2 | 3 | 4 | … |
| Normal μ | 7 | 8 | 9 | … |



| | Feature 1 | Feature 2 | Feature 3 | … |
|---|---|---|---|---|
| Tumor 1… | 1 | 6 | 10 | … |
| Tumor 2… | 2 | 7 | 11 | … |
| Tumor *x*… | … | … | … | … |
| Normal 1… | 5 | 5 | 7 | … |
| Normal 2… | 6 | 6 | 8 | … |
| Normal *y*… | … | … | … | … |

# Predictions

|  | Feature 1 | Feature 2 | Feature 3 | … |
|---|---|---|---|---|
| Test Datapoint | 3 | 4 | 5 | … |

|  | Feature 1 | Feature 2 | Feature 3 | … |
|---|---|---|---|---|
| Tumor μ | 2 | 3 | 4 | … |
| Normal μ | 7 | 8 | 9 | … |

|  | Feature 1 | Feature 2 | Feature 3 | … |
|---|---|---|---|---|
| (Test Datapoint - Normal μ)$^2$ | $(2 - 7)^2$ | $(3 - 8)^2$ | $(4 - 9)^2$ | … |

|  | Feature 1 | Feature 2 | Feature 3 | … |
|---|---|---|---|---|
| (Test Datapoint - Tumor μ)$^2$ | $(3 - 2)^2$ | $(4 - 3)^2$ | $(5 - 4)^2$ | … |

# Predictions



Testing datapoint is classified as **Tumor** since its distance to the **Tumor mean** is *smaller* than to **Normal Tissue's mean**

# Prediction Results

2-Dimensional LDA
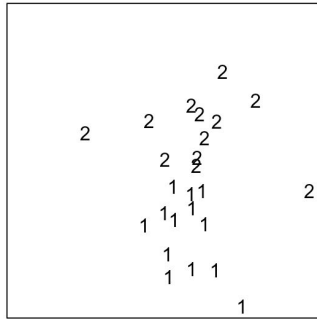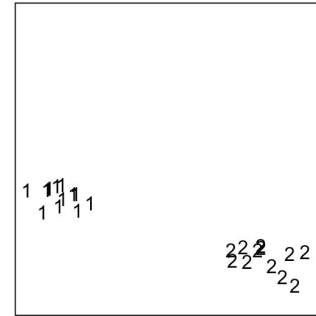projection of testing data



[1] 0.8181818

2-Dimensional PDA ($\lambda = 0.5$)
projection of testing data



[1] 0.9090909

# Support Vector Machines

```
Call:
svm(formula = classifier ~ ., data = co

Parameters:
   SVM-Type:  C-classification
 SVM-Kernel:  radial
       cost:  1

Number of Support Vectors:  25

        Predicted
Actual 0 1
      0 2 3
      1 0 6
```

*e1071* package

**2/5** normal tissues *correctly* classified

**3/5** normal tissues *incorrectly* classified as tumors

**6/6** tumors *correctly* classified

**Accuracy: (2 + 6)/(5 + 6) = 8/11 ≈ 72.73%**

# References (images)

1201904. "Dimensionality Reduction Techniques Skill Test for Data Scientists (Updated 2023)." Analytics Vidhya, 14 Mar. 2023, www.analyticsvidhya.com/blog/2017/03/questions-dimensionality-reduction-data-scientist/.

Liu, Jian, et al. "Cancer Characteristic Gene Selection via Sample Learning Based on Deep Sparse Filtering." The Differences between Sample Learning and Feature Learning., May 2018, www.researchgate.net/figure/The-differences-between-sample-learning-and-feature-learning-a-A-feature-learning_fig1_325426703.

Raschka, Sebastian. "Linear Discriminant Analysis." Sebastian Raschka, PhD, 3 Aug. 2014, sebastianraschka.com/Articles/2014_python_lda.html.

Shetty, Badreesh. "What Is the Curse of Dimensionality?" Built In, 19 Aug. 2022, builtin.com/data-science/curse-dimensionality.

Vungarala, Seshu Kumar. "PCA vs Lda-No More Confusion!" Medium, Medium, 30 Apr. 2023, medium.com/@seshu8hachi/pca-vs-lda-no-more-confusion-fc21fb8d06e9#:~:text=PCA%20is%20an%20unsupervised%20method%20that%20aims%20to%20find%20the,the%20characteristics%20of%20the%20dataset.

# References (papers)

Kabacoff, R. (2015). R in Action, 2nd edition, Chapter 17 (ISBN: 9781617291388). Manning Publications Co.

Lee, E.-K., & Cook, D. (2009). A projection pursuit index for large P Small N Data. Statistics and Computing, 20(3), 381–392.

Notterman, et al, Cancer Research vol. 61: 2001

V. Pappu, & P.M. Pardalos. High Dimensional Data Classification. Clusters, Orders and Trees: Methods and Applications.