

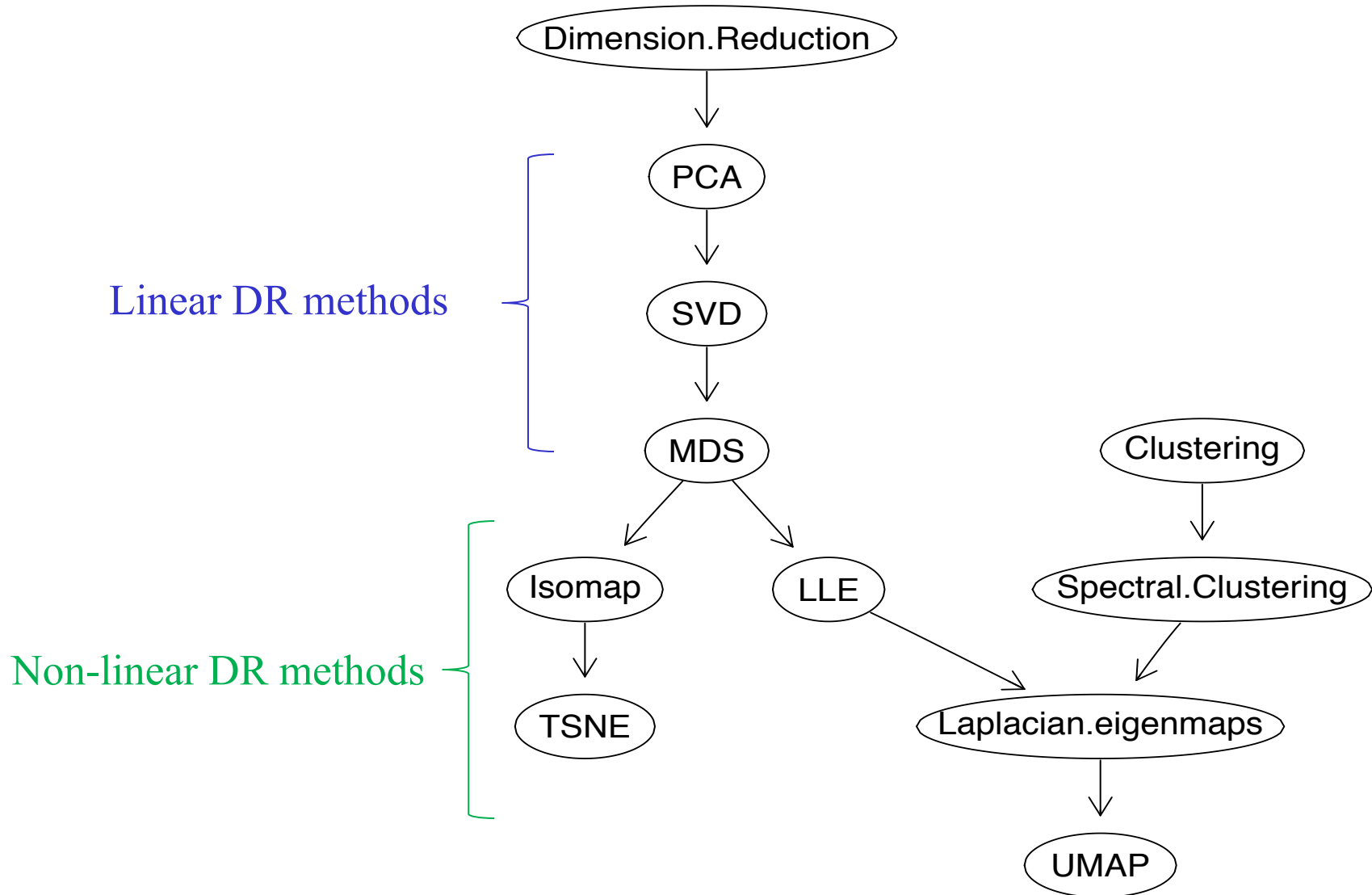
# **Dimension Reduction Methods: From PCA to TSNE and UMAP**

Maxwell Lee

High-dimension Data Analysis Group  
Laboratory of Cancer Biology and Genetics  
Center for Cancer Research  
National Cancer Institute

April 23, 2020

# Outline for Dimension Reduction Methods



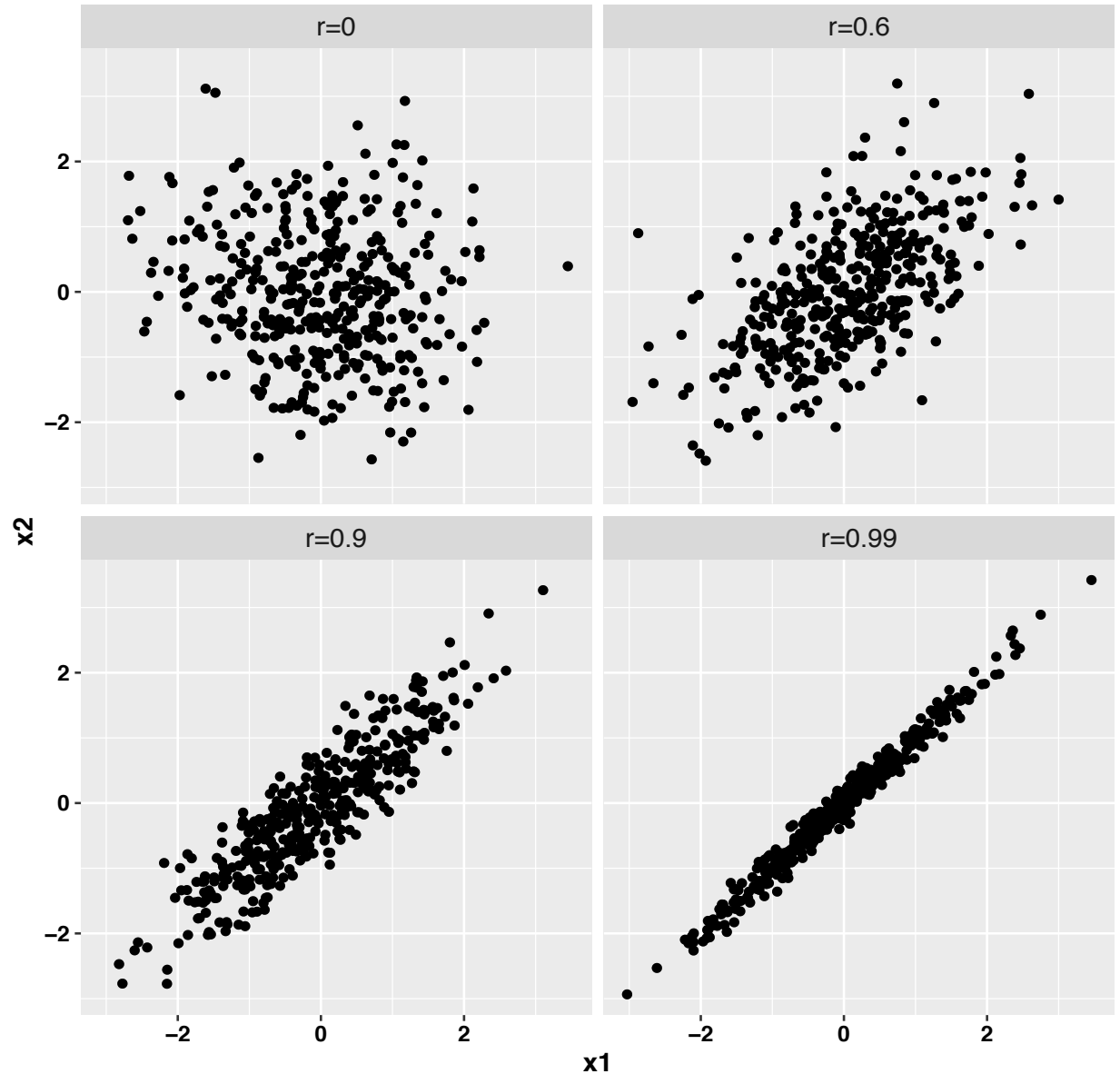
# The Presence of Correlation Between Variables Is the Reason Why We Can Reduce Dimension by PCA

$$\begin{pmatrix} X_1 \\ X_2 \end{pmatrix} \sim \mathcal{N}\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}\right)$$

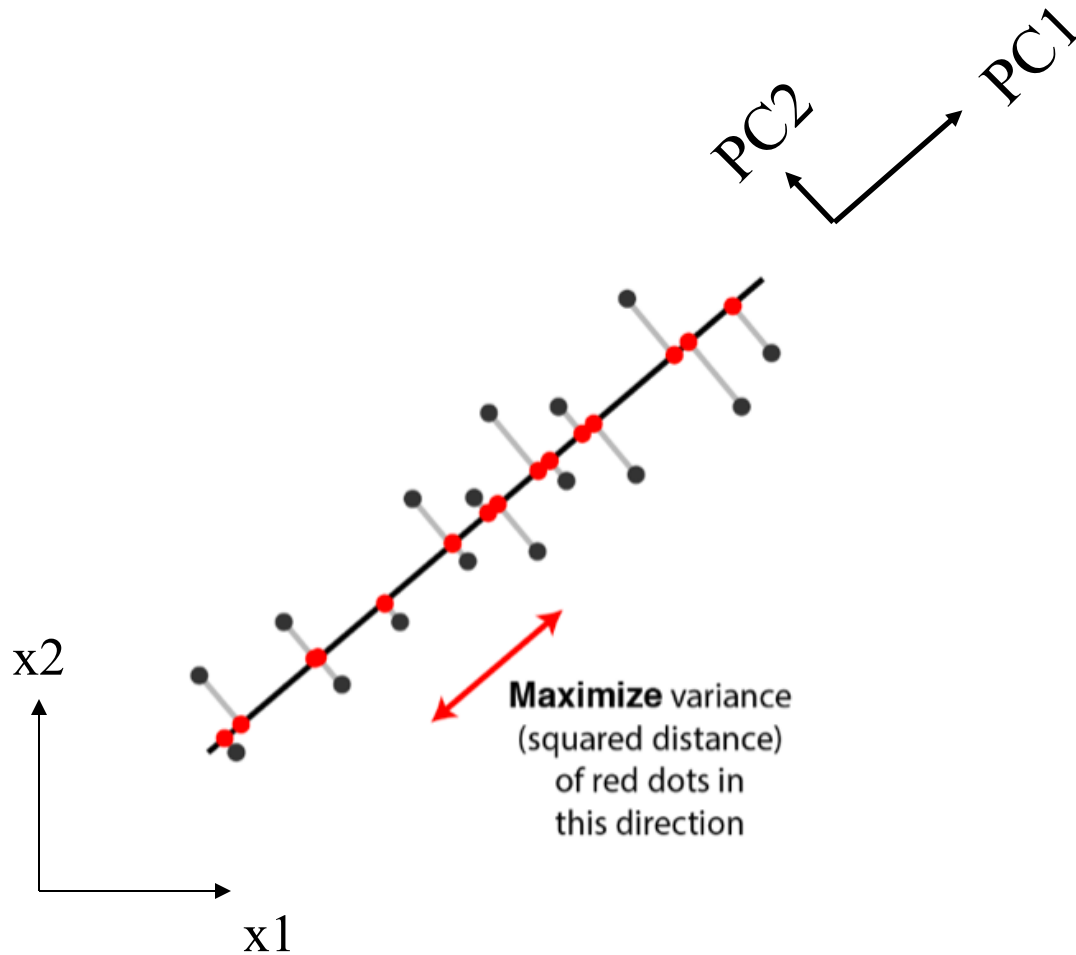
$$r = \rho$$

	$\mu_1$	$\mu_2$	$\sigma_1$	$\sigma_2$	$\rho$
d1	0	0	1	1	0
d2	0	0	1	1	0.6
d3	0	0	1	1	0.9
d4	0	0	1	1	0.99

$$n = 400$$



# Principal Component Analysis (PCA)



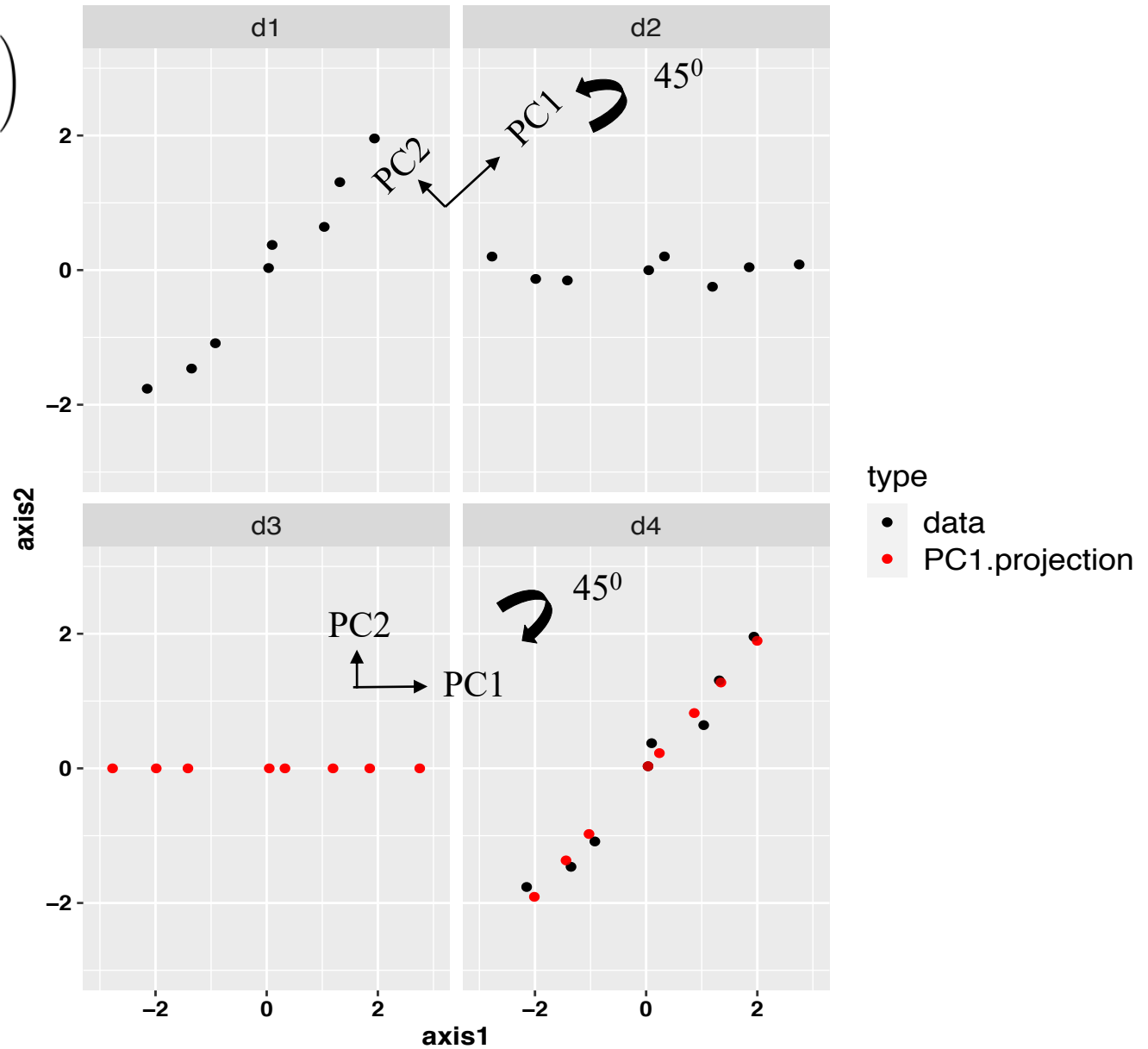
Karl Pearson 1901; Harold Hotelling 1933-1936

# Geometric View of PCA: Rotation of Coordinates

$$\begin{pmatrix} X_1 \\ X_2 \end{pmatrix} \sim \mathcal{N} \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix} \right)$$

$$\rho = 0.9$$

$$n = 8$$



# Correlation Between Variables Can Result from Heterogeneity in Sample

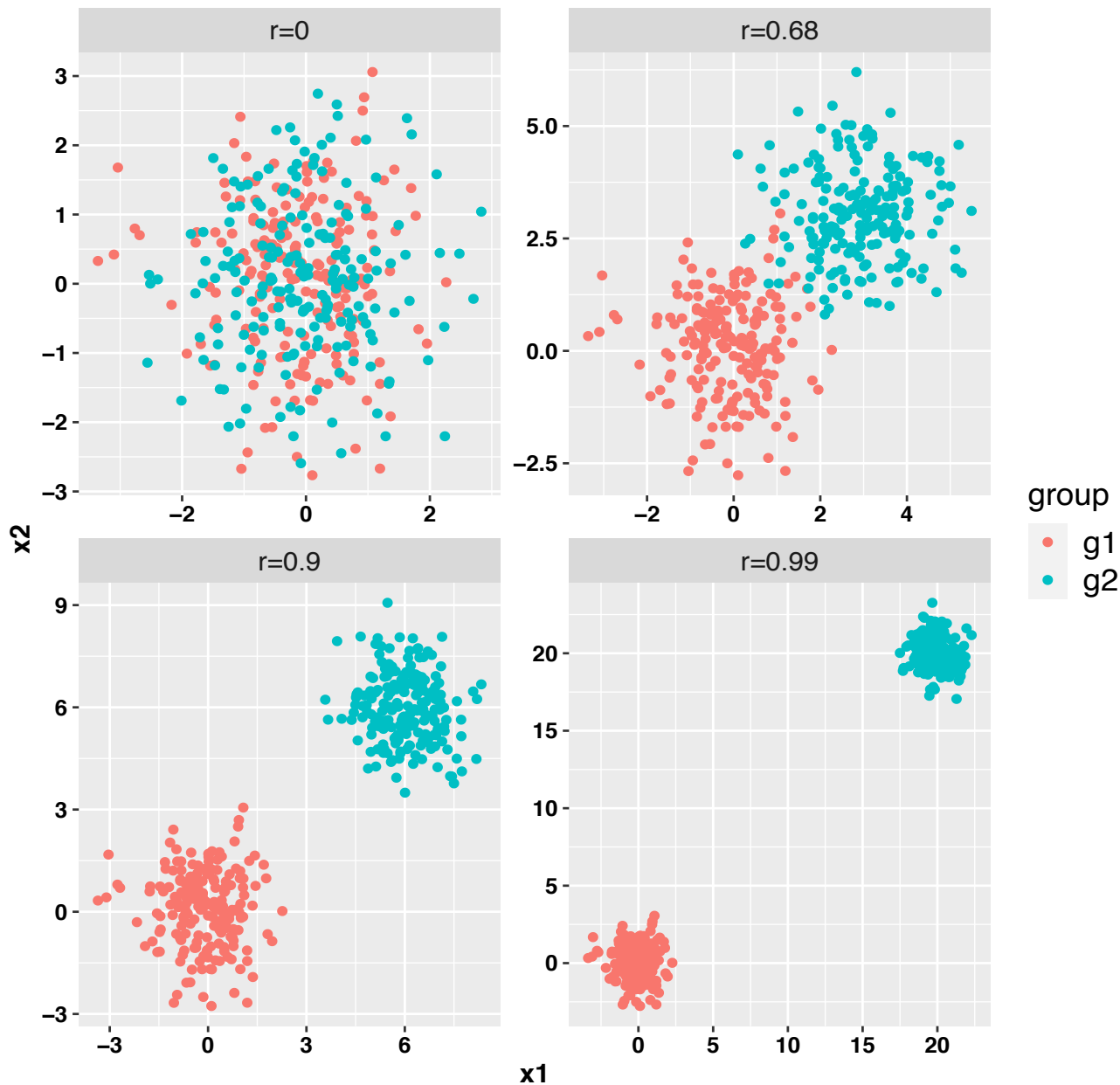
$$\begin{pmatrix} X_1 \\ X_2 \end{pmatrix} \sim \mathcal{N} \left( \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix} \right)$$

$$\rho = 0$$

Group1

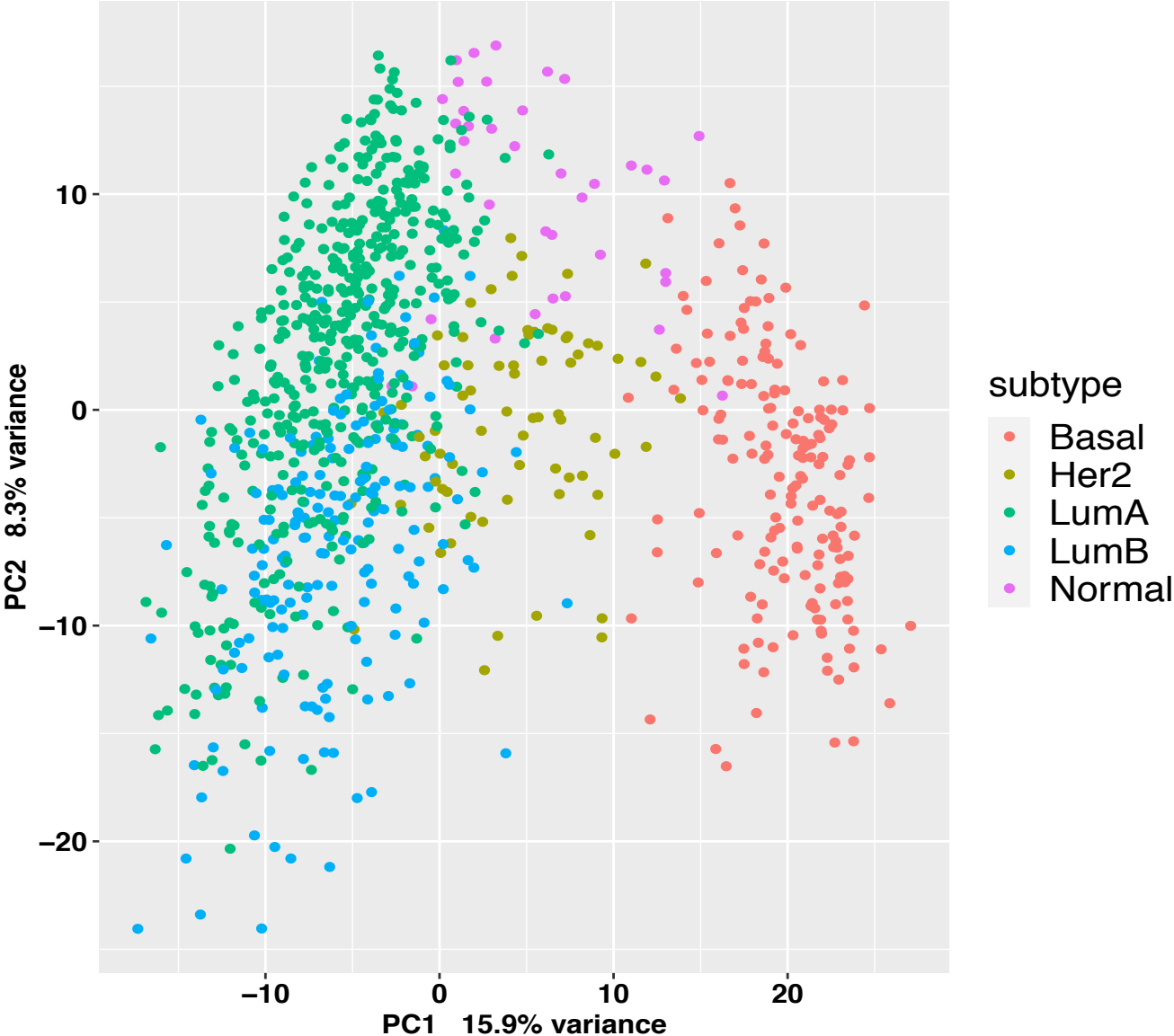
Group2

	$\mu_1$	$\mu_2$	$\mu_1$	$\mu_2$
d1	0	0	0	0
d2	0	0	3	3
d3	0	0	6	6
d4	0	0	20	20

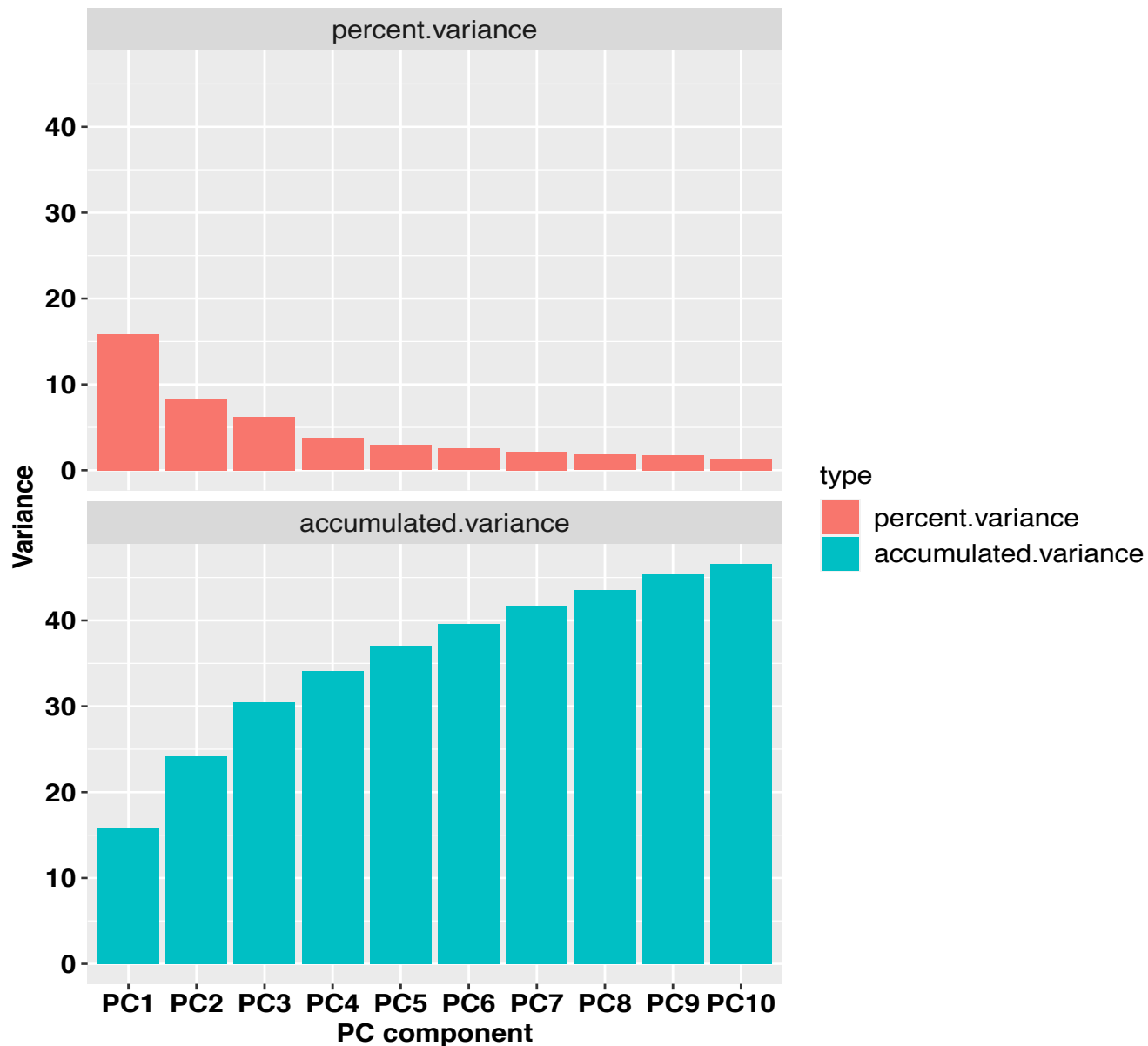


# PCA Analysis of TCGA Breast Cancer RNAseq Data

TCGA BRCA samples: n=977, top 5k most variable genes



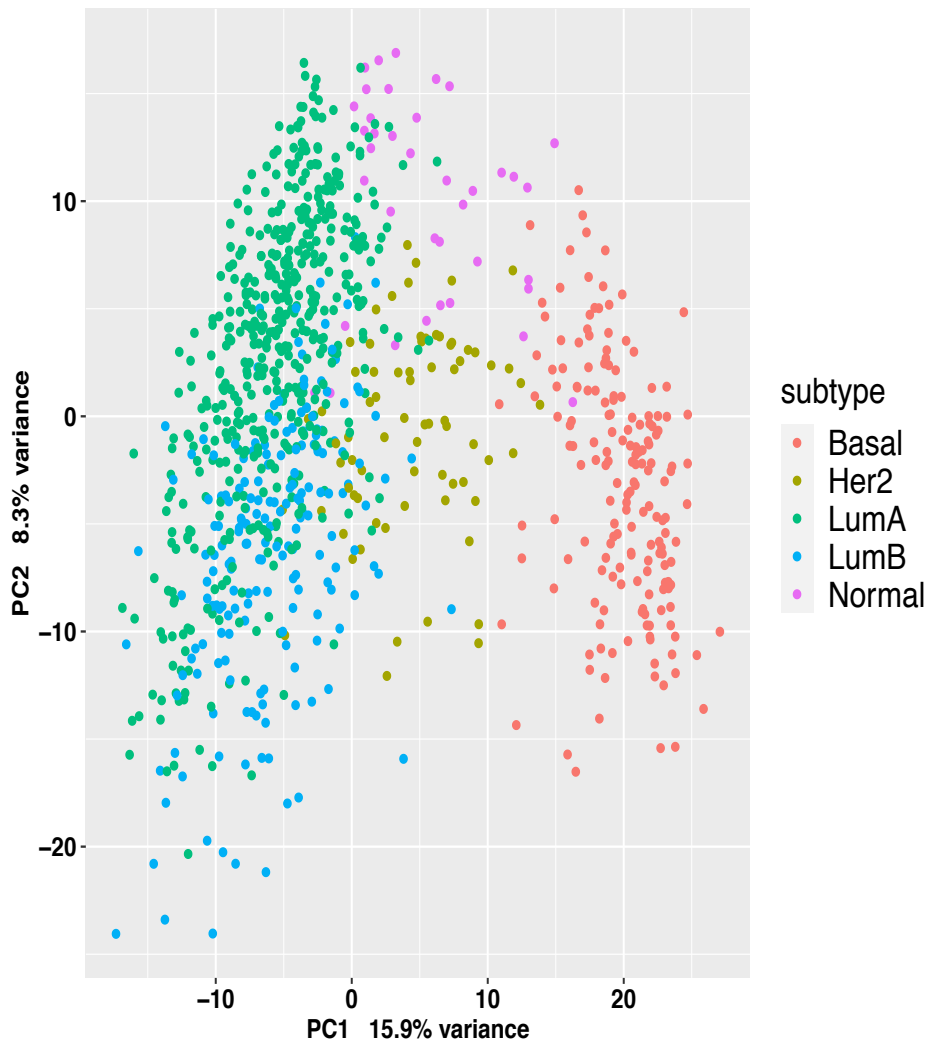
# Variance of Principal Components Are Ranked from the Highest to the Lowest



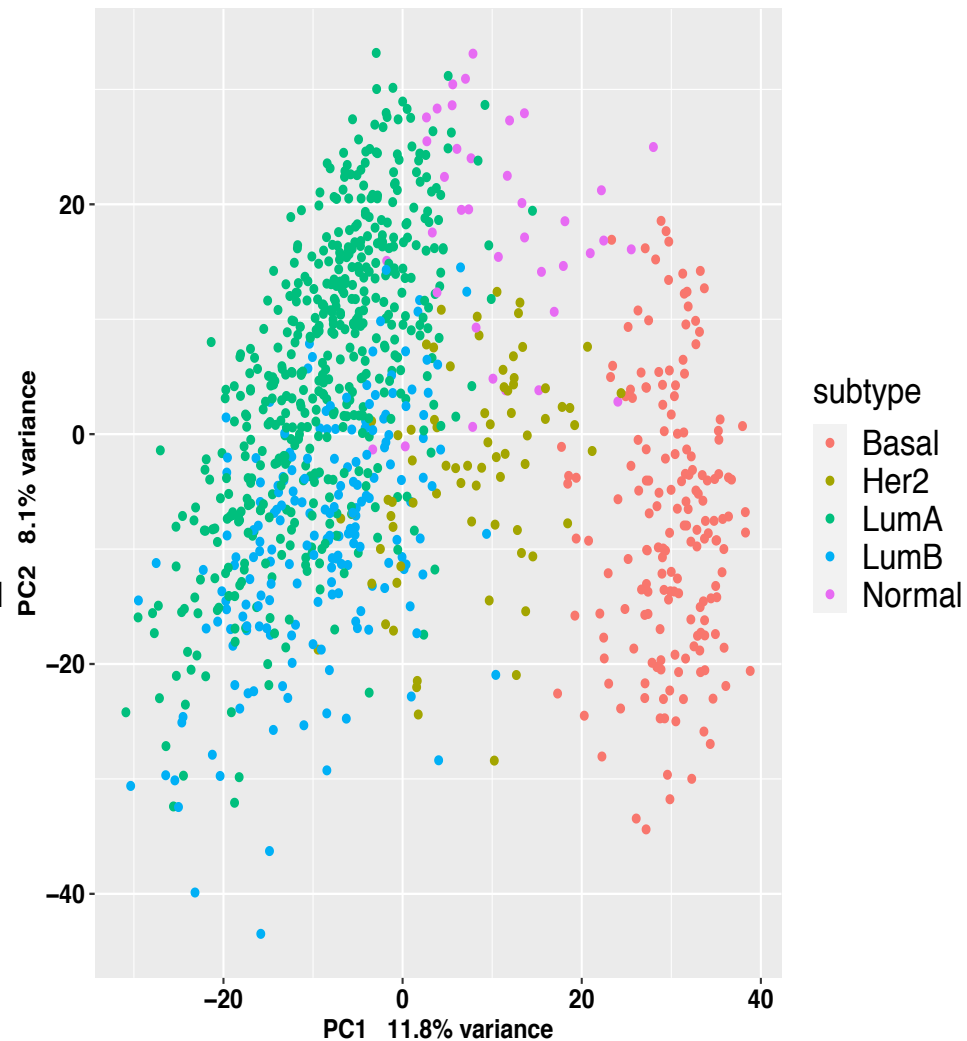


# Filtering Out Genes of Low Variance Increases Percent of Variance Accounted for by PC1

TCGA BRCA samples: n=977, top 5k most variable genes



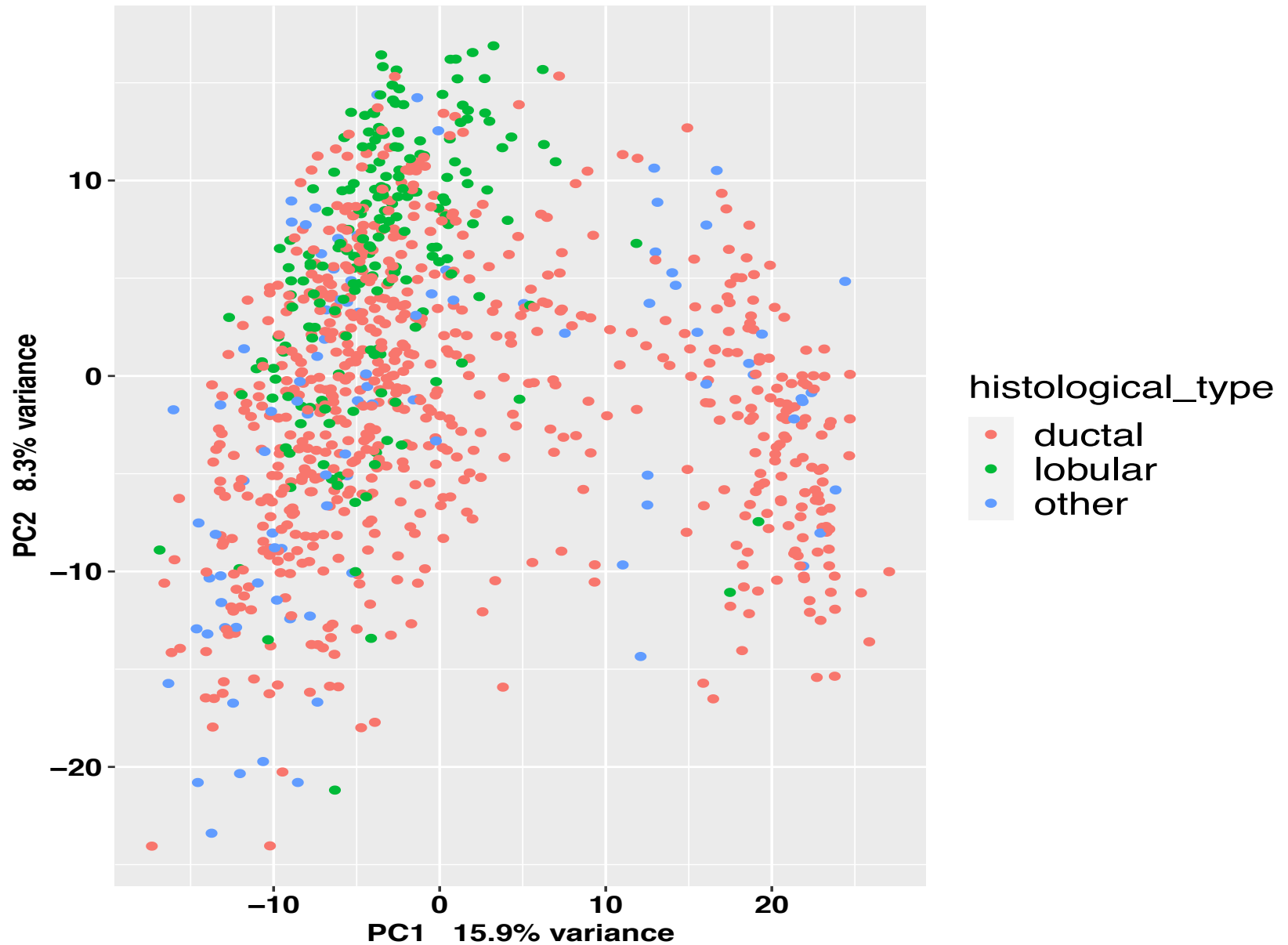
TCGA BRCA samples: n=977, all 20k genes





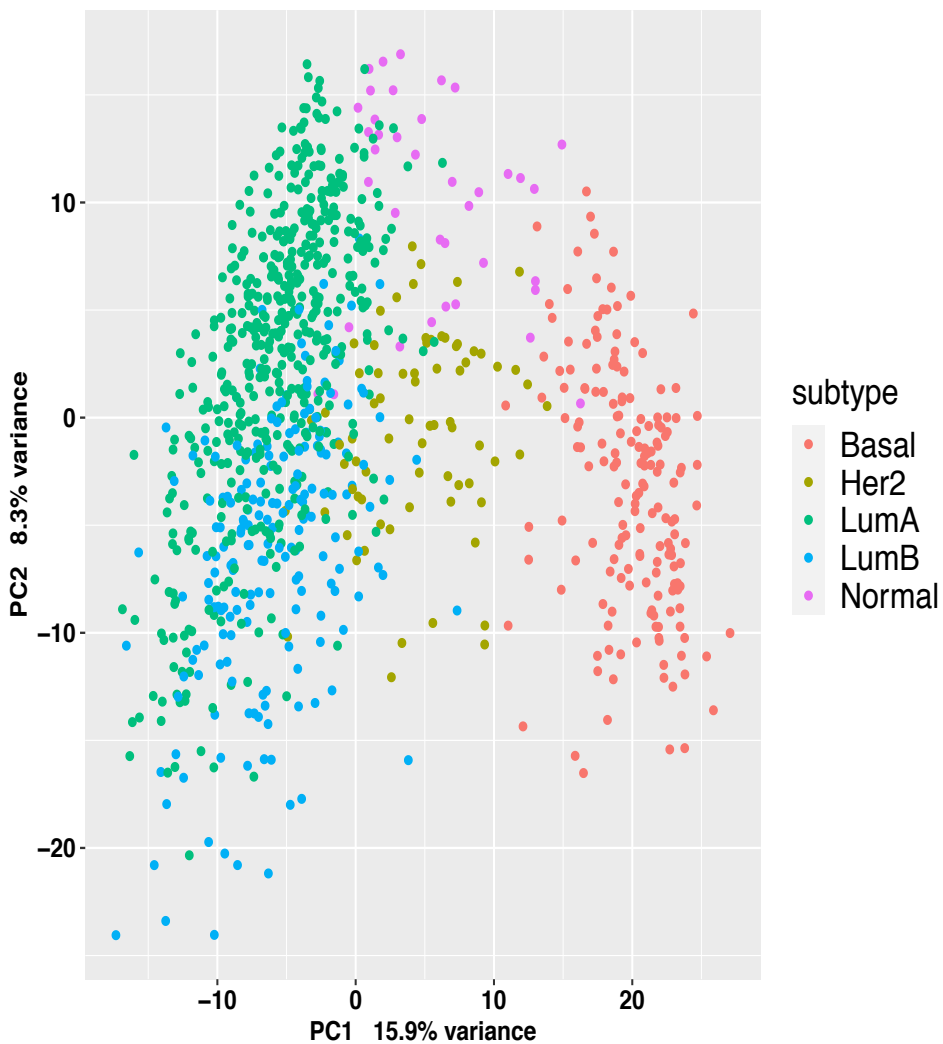
# Variation in Histological Type Is Associated with PC2

TCGA BRCA samples: n=977, histological\_type

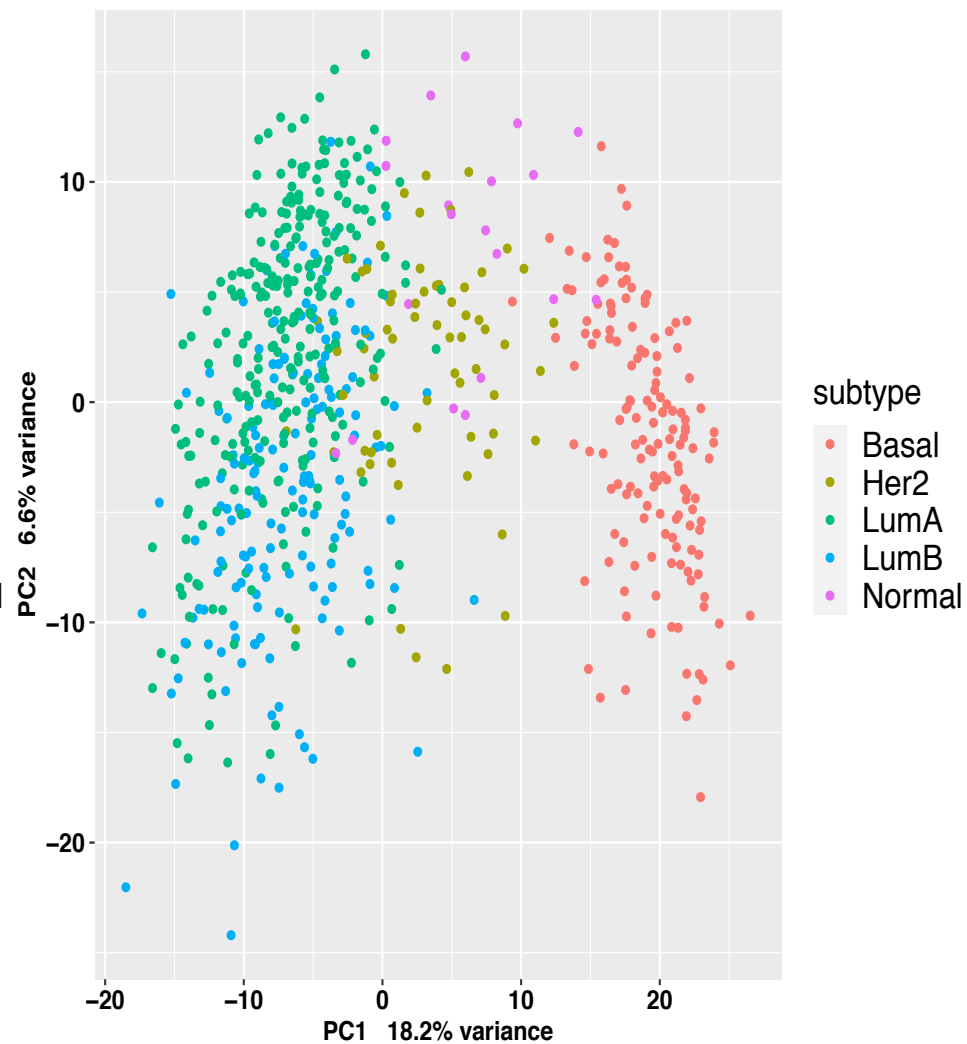


# Removing Heterogeneity in Histological Type Reduces PC2 Variance and Increases PC1 Variance

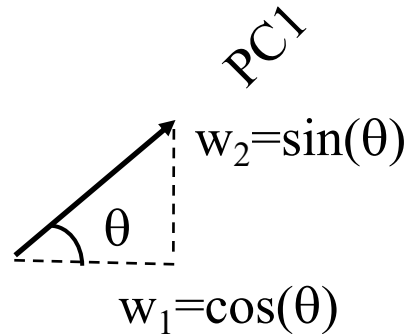
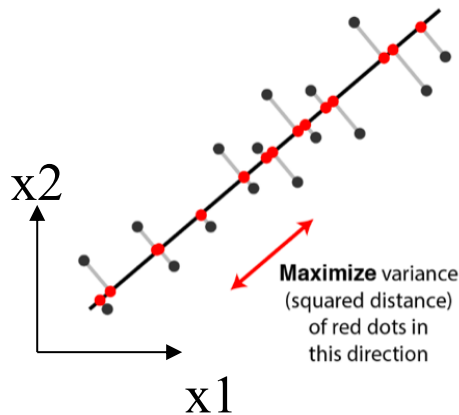
TCGA BRCA samples: n=977, top 5k most variable genes



TCGA BRCA samples: n=688, infiltrating ductal carcinoma



# Algorithm of PCA: How Does PCA Find the Direction of PC1?



$$\begin{matrix} & Xw & & z \\ \begin{bmatrix} X_{11} & X_{12} \\ X_{21} & X_{22} \\ \cdot & \cdot \\ \cdot & \cdot \\ X_{n1} & X_{n2} \end{bmatrix} & \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} & = & \begin{bmatrix} w_1 X_{11} + w_2 X_{12} \\ w_1 X_{21} + w_2 X_{22} \\ \cdot \\ \cdot \\ w_n X_{21} + w_2 X_{n2} \end{bmatrix}
 \end{matrix}$$

$$z = Xw$$

$$\text{var}(z) = (Xw)^T Xw$$

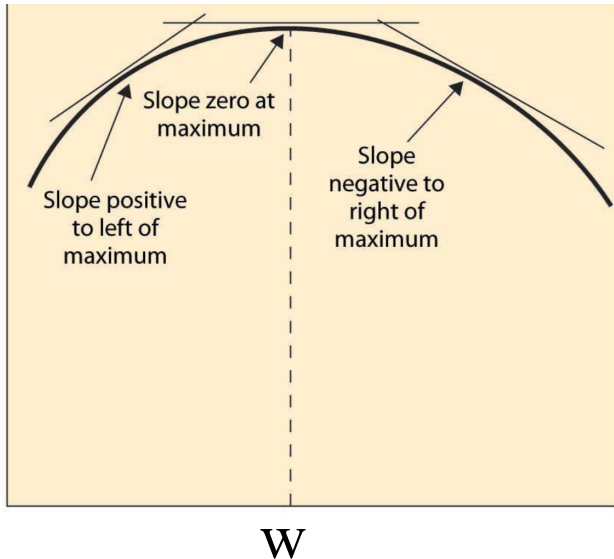
$$\text{var}(z) = w^T X^T X w = w^T S w$$

Choose  $w$  to maximize  $w^T S w$   
subject to  $w^T w = 1$

# The Direction of PC1 Is the Eigen Vector with the Highest Eigen Value

Choose  $w$  to maximize  $w^T S w$   
subject to  $w^T w = 1$

L



$$L(w, \lambda) = w^T S w - \lambda(w^T w - 1)$$

$$\frac{\partial L}{\partial w} = 2S w - 2\lambda w$$

$$S w = \lambda w$$

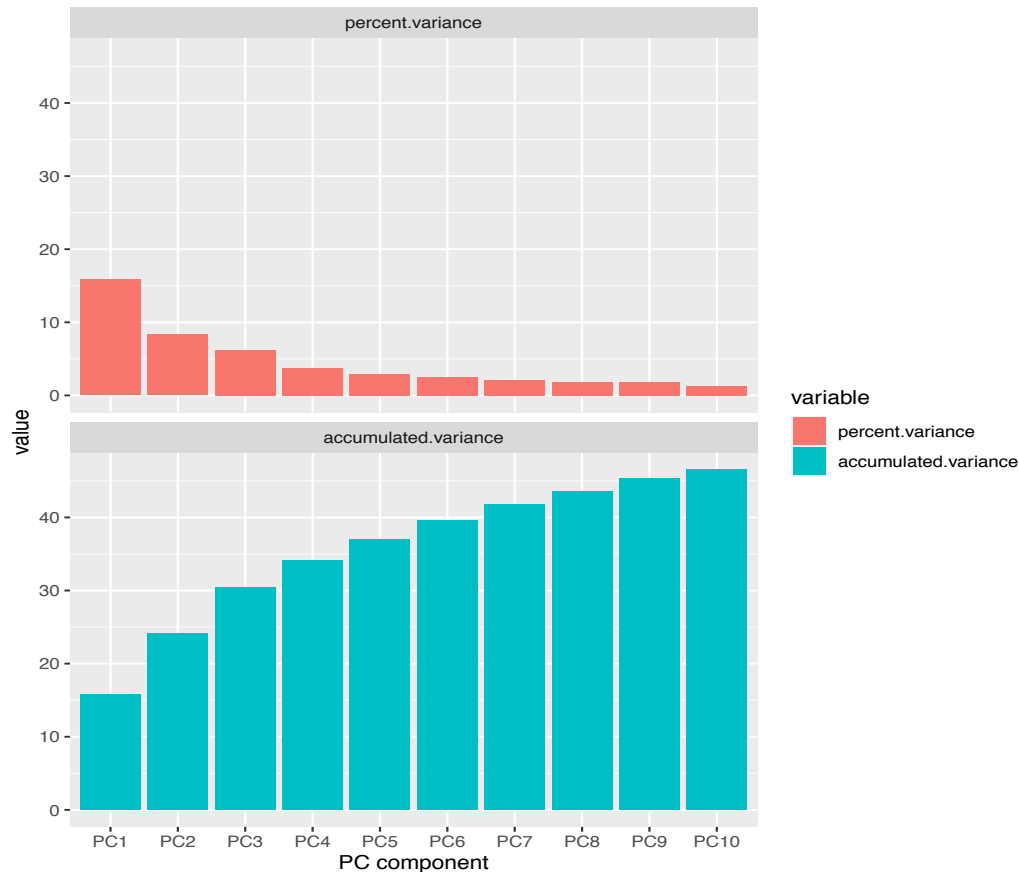
$w$  is the eigen vector and  $\lambda$  is eigen value

# Variance of PCs Are Eigen Value and Are Additive

$$\begin{aligned}\text{var}(z) &= \mathbf{w}^T \mathbf{S} \mathbf{w} \\ &= \mathbf{w}^T \boldsymbol{\lambda} \mathbf{w} \\ &= \lambda\end{aligned}$$

There are  $p$  pairs of eigen vectors and eigen values

$$\text{var}(Z) = \lambda_1 + \lambda_2 \dots + \lambda_p$$



# Singular Value Decomposition (SVD)

$$\mathbf{Z} = \mathbf{X}\mathbf{W}$$

$$\mathbf{Z}_s = \mathbf{X}\mathbf{W}\mathbf{\Lambda}^{-1/2}$$

$$\mathbf{Z}_s\mathbf{\Lambda}^{1/2}\mathbf{W}^T = \mathbf{X}$$

$$\mathbf{X} = \mathbf{Z}_s\mathbf{\Lambda}^{1/2}\mathbf{W}^T$$

$\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$  from standard SVD notation



# Right and Left Singular Vectors of SVD

p column vectors

n row vectors

$$\begin{bmatrix} X_{11} & X_{12} & \cdots & X_{1p} \\ X_{21} & X_{22} & \cdots & X_{2p} \\ \cdot & & & \cdot \\ \cdot & & & \cdot \\ X_{n1} & X_{n2} & \cdots & X_{np} \end{bmatrix}$$

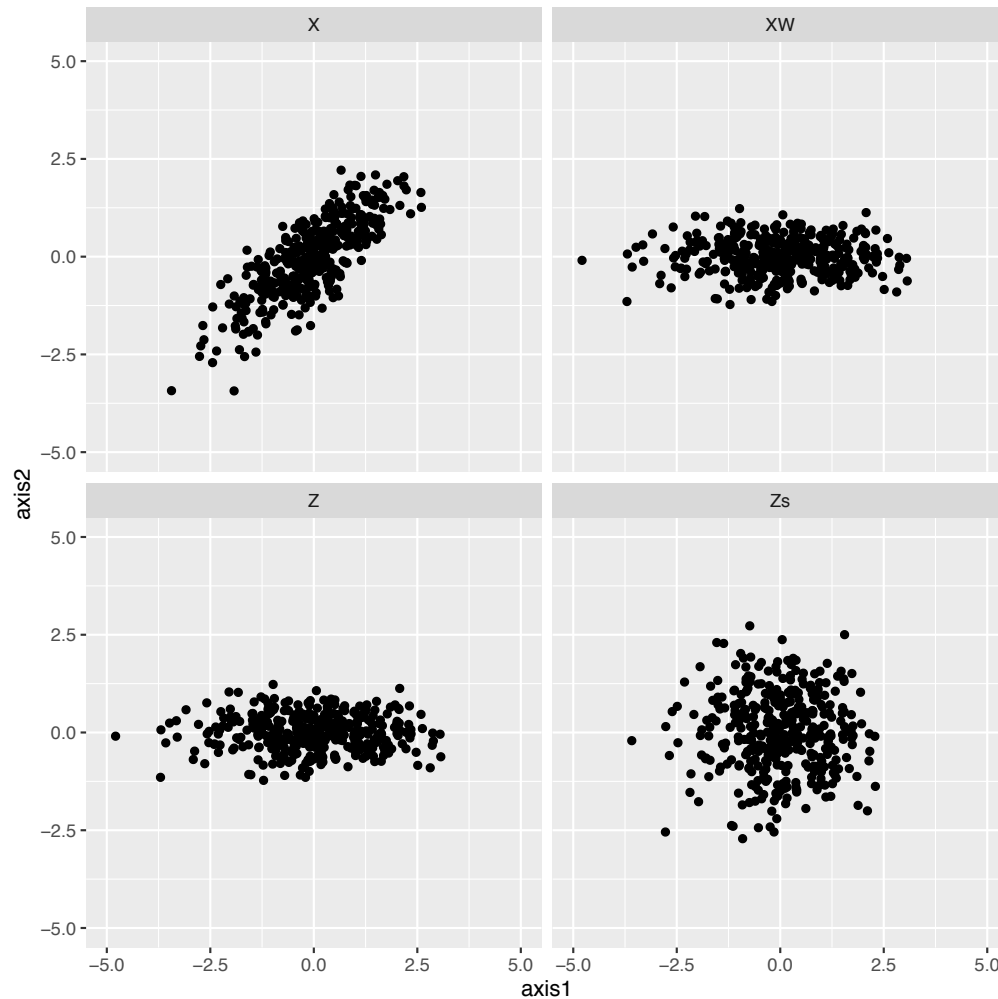
$$X = U\Sigma V^T$$

$$\begin{aligned} X^T X &= V\Sigma U^T U \Sigma V^T \\ &= V\Sigma^2 V^T \end{aligned}$$

$$\begin{aligned} X X^T &= U\Sigma V^T V \Sigma U^T \\ &= U\Sigma^2 U^T \end{aligned}$$

# Interconversion Between Principal Components and Standardized Principal Components

$$Z_s = XW D^{-1/2}$$



$$Z = Z_s D^{1/2}$$

# **Three Dimensional Object and Its Two dimensional Image: Ellipse As Shadow of Sphere**



# Europe Map on a Globe and on Google Map



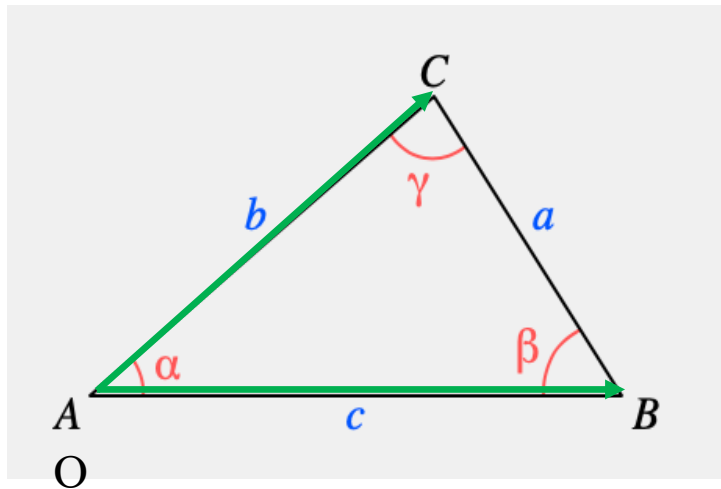
# Eight European Cities Pairwise Distance Matrix

	Athens	Berlin	Dublin	London	Madrid	Paris	Rome	Warsaw
Athens	0	1119	1777	1486	1475	1303	646	1013
Berlin	1119	0	817	577	1159	545	736	327
Dublin	1777	817	0	291	906	489	1182	1135
London	1486	577	291	0	783	213	897	904
Madrid	1475	1159	906	783	0	652	856	1483
Paris	1303	545	489	213	652	0	694	859
Rome	646	736	1182	897	856	694	0	839
Warsaw	1013	327	1135	904	1483	859	839	0

# Law of Cosines: Dot Product of Two Vectors Can Be Expressed By the Difference Between the Squared Distances

Law of cosine

$$a^2 = b^2 + c^2 - 2bc \cos(\alpha)$$



$$2bc \cos(\alpha) = b^2 + c^2 - a^2$$

$$bc \cos(\alpha) = -\frac{1}{2}(a^2 - b^2 + c^2)$$

$$\mathbf{b} \cdot \mathbf{c} = bc \cos(\alpha)$$

$$\mathbf{b} \cdot \mathbf{c} = -\frac{1}{2}(a^2 - b^2 - c^2)$$

# Eigen Decomposition Can Generate Projection Map From Kernel Matrix Derived From the Pairwise Distance Between Two Cities

$$\mathbf{b} \cdot \mathbf{c} = -1/2(a^2 - b^2 - c^2)$$

$$\begin{bmatrix} k_{11} & k_{12} & \dots & k_{1n} \\ k_{21} & k_{22} & \dots & k_{2n} \\ \cdot & & & \cdot \\ \cdot & & & \cdot \\ k_{n1} & k_{n2} & \dots & k_{nn} \end{bmatrix}$$

$$\mathbf{K} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$$

$$\mathbf{Z} = \mathbf{U}\mathbf{\Lambda}^{1/2}$$

# MDS and PCA Are Equivalent

$$Z = U\Lambda^{1/2} \quad \text{from MDS}$$

$$X = U\Sigma V^T$$

$$XV = U\Sigma$$

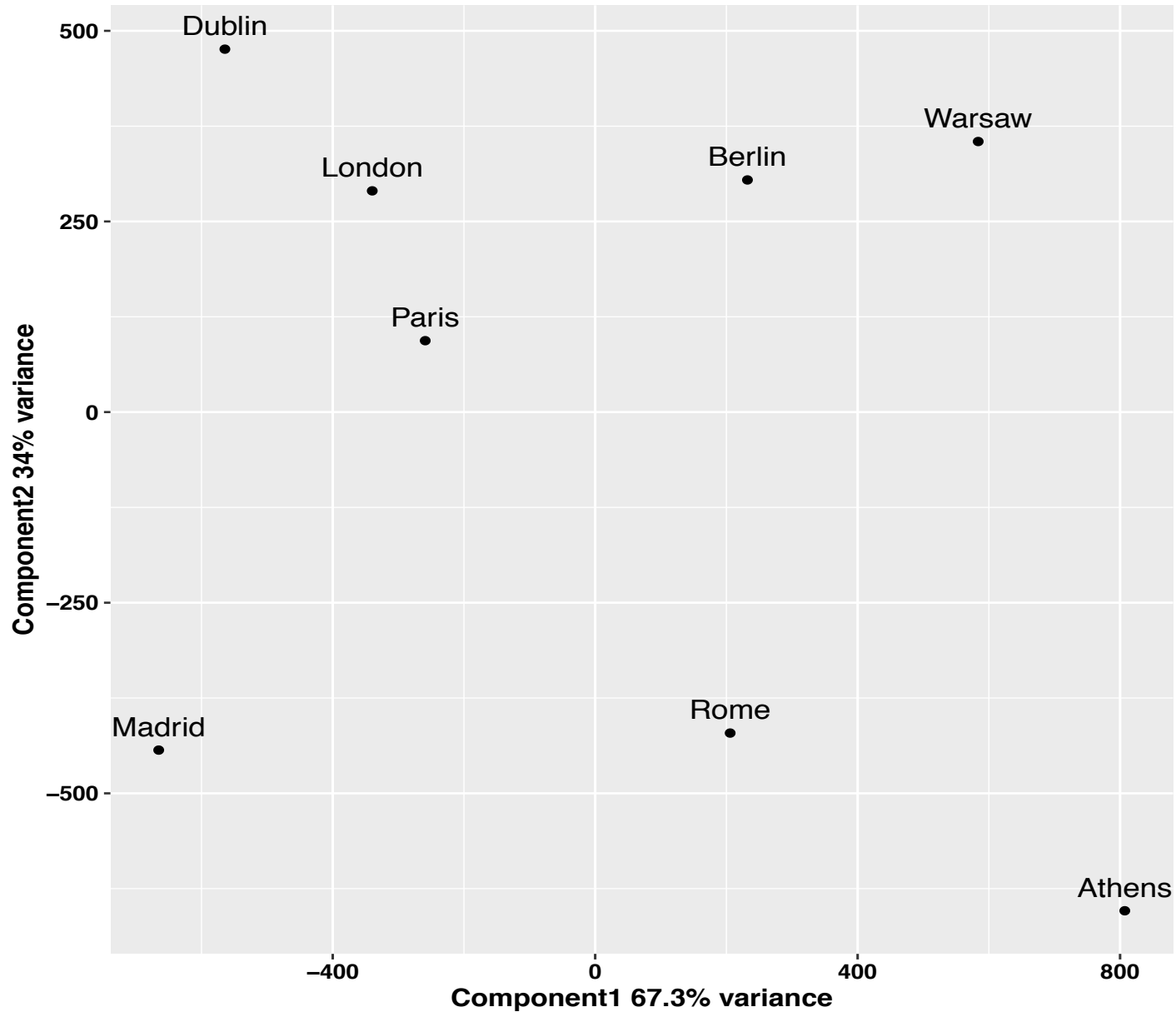
$$Z = U\Lambda^{1/2}$$

$$K = U\Lambda U^T$$

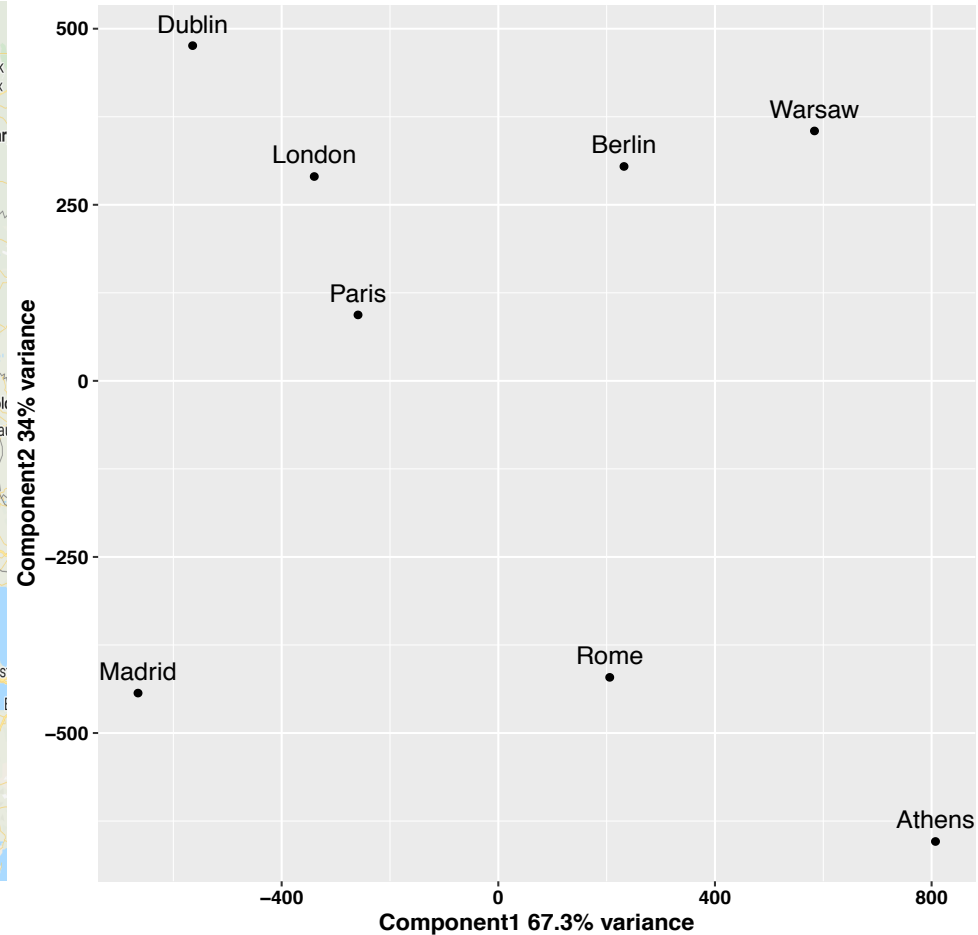
K was derived from the pairwise distance matrix without X



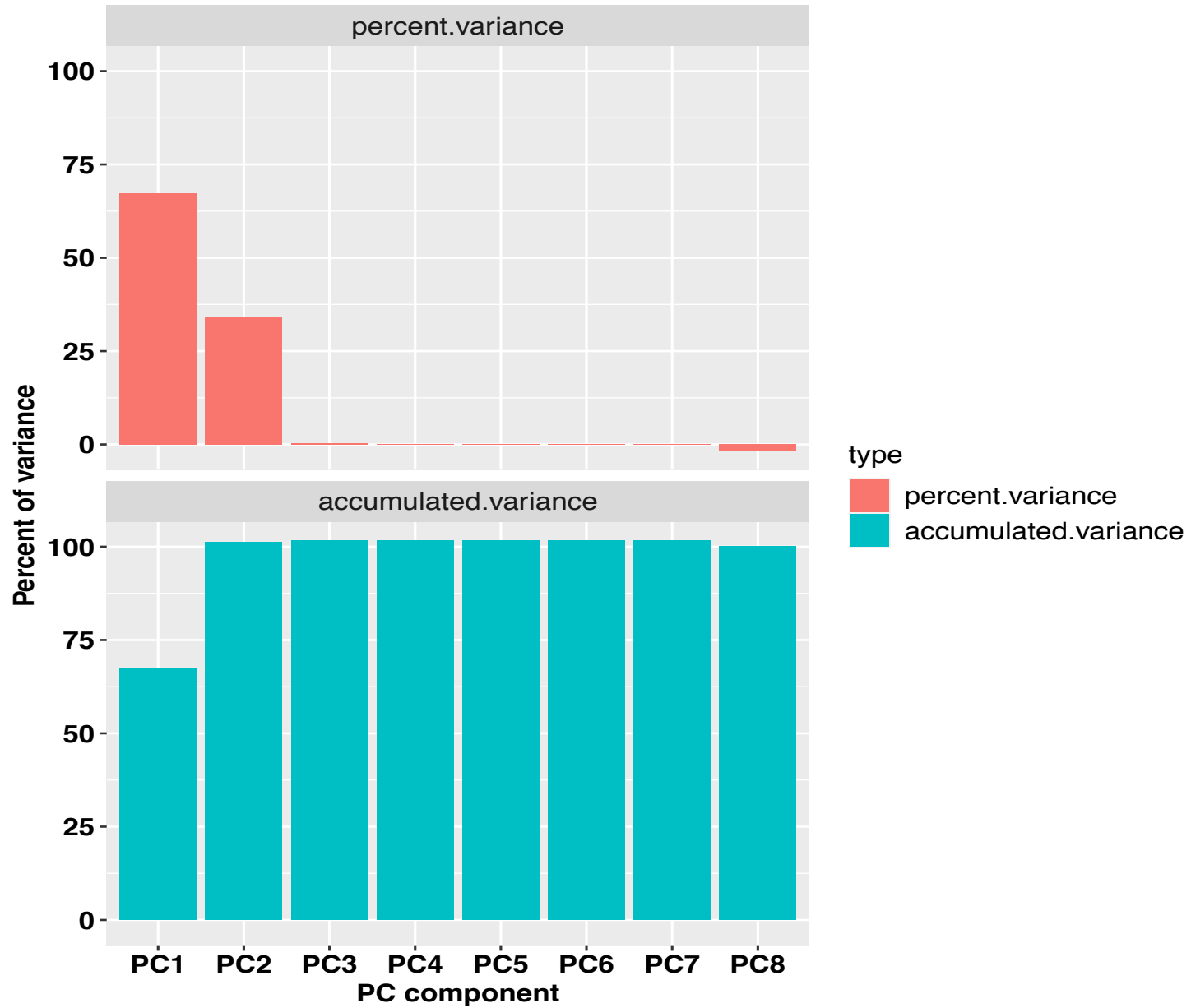
# Two Dimensional Map Generated with MDS



# Comparison Between Google Map and MDS Projection



# Variance of MDS Components



# Outline for Dimension Reduction Methods

