
DATA, MSML, BIOI 602

Principles of Data Science

SVD and PCA

Heng Huang
Department of Computer Science

Orthogonal Matrix

- Suppose \mathbf{A} is a square matrix. \mathbf{A} is called orthogonal matrix if

$$\mathbf{A}^T \mathbf{A} = \mathbf{A} \mathbf{A}^T = \mathbf{I}$$

where \mathbf{I} is an identity matrix, \mathbf{A}^T is the transpose of \mathbf{A} .

- For an orthogonal matrix, we have

$$\mathbf{A}^{-1} = \mathbf{A}^T$$

Eigenvalue & Eigenvector

- Suppose \mathbf{A} is a square matrix. If having a number, λ , and a non-zero vector, \mathbf{X} , satisfy

$$\mathbf{AX} = \lambda\mathbf{X}$$

- We called λ the eigenvalue of \mathbf{A} , and \mathbf{X} is the eigenvector of \mathbf{A}
- If we know the eigenvalues of \mathbf{A} , the eigenvectors can be determined by substituting the eigenvalues into above equation.

Calculation of Eigenvalues

- We can determine the eigenvalues of \mathbf{A} by solving the following equation:

$$|\mathbf{A} - \lambda \mathbf{I}| = 0$$

where \mathbf{I} is an identity matrix

Singular Values

- Suppose \mathbf{A} is a $m \times n$ matrix and its rank is r ($r \leq n$).
We can calculate the non-zero eigenvalues of $\mathbf{A}^T \mathbf{A}$,
e.g.,
$$\lambda_1 \geq \lambda_2 \dots \geq \lambda_r$$
- We call $\mu_i = \sqrt{\lambda_i}$ ($i = 1, 2, \dots, r$) as the singular values of \mathbf{A}

What is SVD?

Any $m \times n$ matrix \mathbf{A} with rank of r , can be decomposed into

$$\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{V}^T$$

where \mathbf{U} and \mathbf{V} are orthogonal matrices and \mathbf{D} is a diagonal matrix containing singular values, $\{\mu_i, i = 1, 2, \dots, r\}$. This factored matrix representation is known as the SVD.

SVD More Formally

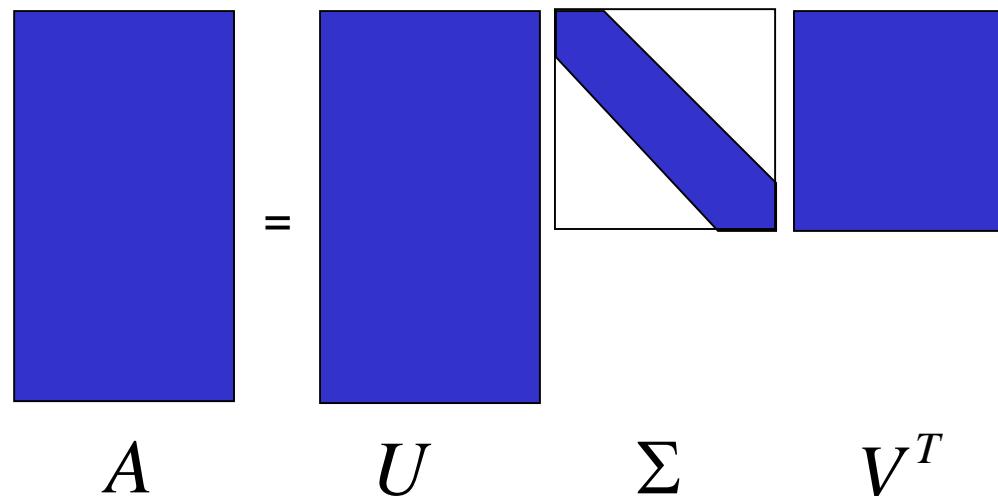
- The diagonal values of Σ (μ_1, \dots, μ_n) are called the **singular values**. It is accustomed to sort them: $\mu_1 \geq \mu_2 \geq \dots \geq \mu_n$
- The columns of U ($\mathbf{u}_1, \dots, \mathbf{u}_n$) are called the **left singular vectors**. They are the axes of the ellipsoid.
- The columns of V ($\mathbf{v}_1, \dots, \mathbf{v}_n$) are called the **right singular vectors**. They are the preimages of the axes of the ellipsoid.

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

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

Reduced SVD

- For rectangular matrices, we have two forms of SVD. The reduced SVD looks like this:
 - The columns of U are orthonormal
 - Cheaper form for computation and storage

$$A = U \Sigma V^T$$
The diagram illustrates the Reduced Singular Value Decomposition (SVD) of a matrix A. It shows a blue square matrix A on the left, followed by an equals sign. To the right of the equals sign is another blue square matrix U. Next is a white square matrix Σ, which contains a blue diagonal band from top-left to bottom-right. Finally, there is a blue square matrix V^T on the far right.

Full SVD

- We can complete U to a full orthogonal matrix and pad Σ by zeros accordingly

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

Example of SVD

- Suppose

$$\mathbf{A} = \begin{pmatrix} 1.2 & 0.9 & -4 \\ 1.6 & 1.2 & 3 \end{pmatrix}$$

- 1) Calculate the eigenvalues of $\mathbf{A}^T \mathbf{A}$

$$\lambda_1 = 25, \lambda_2 = 6.25$$

- 2) The non-zero singular values, $\mu_i = \sqrt{\lambda_i}$, e.g.,

$$\mu_1 = 5, \mu_2 = 2.5$$

- 3) \mathbf{V} formed from the orthonormal eigenvectors as columns,

$$\mathbf{A}^T \mathbf{A} = \mathbf{V} \boldsymbol{\Sigma}^T \mathbf{U}^T \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^T = \mathbf{V} \boldsymbol{\Sigma}^T \boldsymbol{\Sigma} \mathbf{V}^T.$$

Example of SVD

$$\mathbf{V} = \begin{pmatrix} v_1 \\ v_2 \\ v_3 \end{pmatrix}^T = \begin{pmatrix} 0 & 0.8 & -0.6 \\ 0 & 0.6 & 0.8 \\ 1 & 0 & 0 \end{pmatrix}$$

- Calculate $\mathbf{U} = \{\mathbf{A}v_j / \mu_j, j=1,2\}$
Where v_j is the non-zero eigenvector of $\mathbf{A}^T \mathbf{A}$

We have

$$\mathbf{U} = \begin{pmatrix} -0.8 & 0.6 \\ 0.6 & 0.8 \end{pmatrix}$$

Example of SVD

- Now the SVD of \mathbf{A} is

$$\mathbf{A} = \mathbf{UDV}^T = \begin{pmatrix} -0.8 & 0.6 \\ 0.6 & 0.8 \end{pmatrix} \begin{pmatrix} 5 & 0 & 0 \\ 0 & 2.5 & 0 \end{pmatrix} \begin{pmatrix} 0 & 0.8 & -0.6 \\ 0 & 0.6 & 0.8 \\ 1 & 0 & 0 \end{pmatrix}^T$$

$$\mathbf{A}^+ = \mathbf{VD}^+ \mathbf{U}^T = \begin{pmatrix} 0 & 0.8 & -0.6 \\ 0 & 0.6 & 0.8 \\ 1 & 0 & 0 \end{pmatrix} \begin{pmatrix} \frac{1}{5} & 0 \\ 0 & \frac{1}{2.5} \\ 0 & 0 \end{pmatrix} \begin{pmatrix} -0.8 & 0.6 \\ 0.6 & 0.8 \end{pmatrix}^T$$

Matrix Inverse and Solving Linear Systems

- Matrix inverse: $A = U \Sigma V^T$

$$A^{-1} = (U \Sigma V^T)^{-1} = (V^T)^{-1} \Sigma^{-1} U^{-1} =$$

$$= V \begin{bmatrix} \frac{1}{\sigma_1} & & \\ & \ddots & \\ & & \frac{1}{\sigma_n} \end{bmatrix} U^T$$

- So, to solve $A\mathbf{x} = \mathbf{b}$

$$\mathbf{x} = V \Sigma^{-1} U^T \mathbf{b}$$

Application: Image Compression

- Uncompressed m by n pixel image: $m \times n$ numbers
- Rank q approximation of image:
 - q singular values
 - The first q columns of \mathbf{U} (m -vectors)
 - The first q columns of \mathbf{V} (n -vectors)
 - Total: $q \times (m + n + 1)$ numbers

Example: Yogi (Uncompressed)

- Source: [Will]
- Yogi: Rock photographed by Sojourner Mars mission.
- 256×264 grayscale bitmap $\rightarrow 256 \times 264$ matrix M
- Pixel values $\in [0,1]$
- ~ 67584 numbers

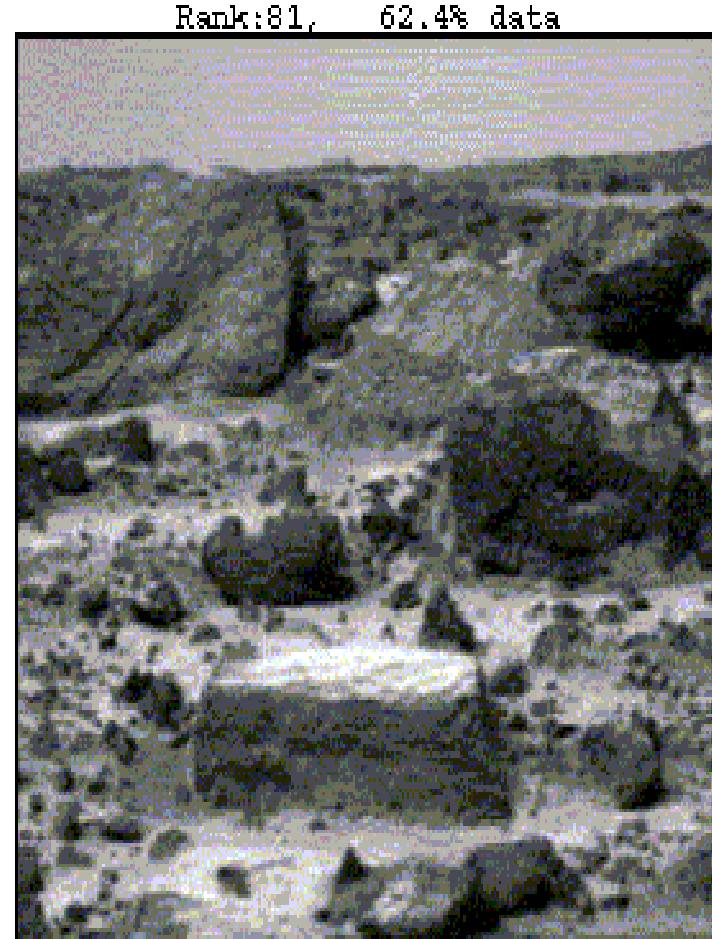


Example: Yogi (Compressed)

- M has 256 singular values
- Rank 81 approximation of M :
- $81 \times (256 + 264 + 1) =$
 ~ 42201 numbers



Example: Yogi (Both)



Application: Noise Filtering

- Data compression: Image degraded to reduce size
- Noise Filtering: Lower-rank approximation used to improve data.
 - Noise effects primarily manifest in terms corresponding to smaller singular values.
 - Setting these singular values to zero removes noise effects.

Principal Components Analysis (PCA)

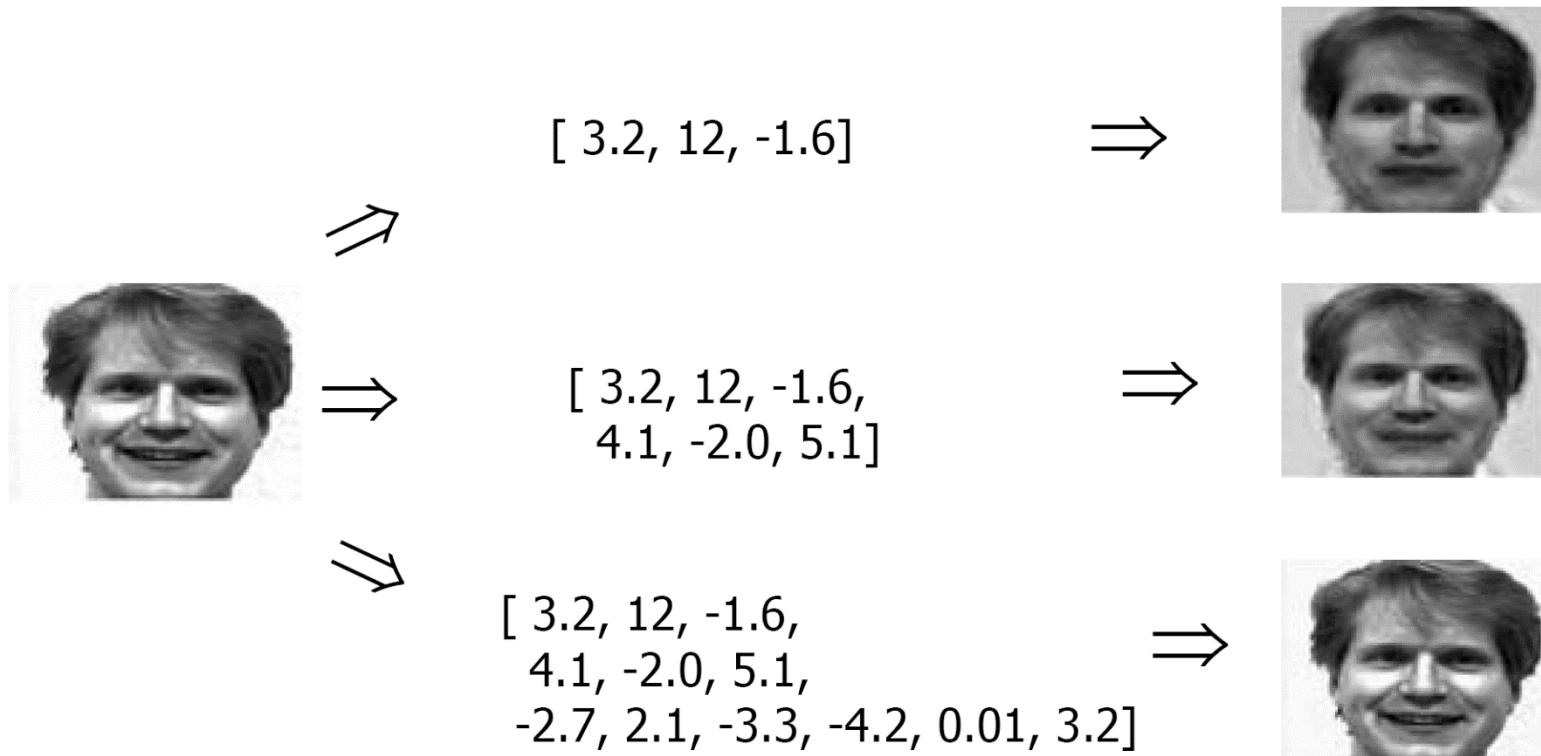
- Idea:
 - Given data points in d -dimensional space, project into *lower dimensional* space while *preserving as much information* as possible
 - Eg, find best planar approximation to 3D data
 - Eg, find best 12-D approximation to 10^4 -D data
 - In particular, choose projection that *minimizes squared error* in reconstructing original data

An Vision Application: Facial Recognition

- Want to identify specific person, based on facial image
- Robust to ...
 - Facial hair, glasses, ...
 - Different lighting
- ⇒ Can't just use given 256 x 256 pixels
- Need another option!



An Vision Application: Facial Recognition



Why Do We Care

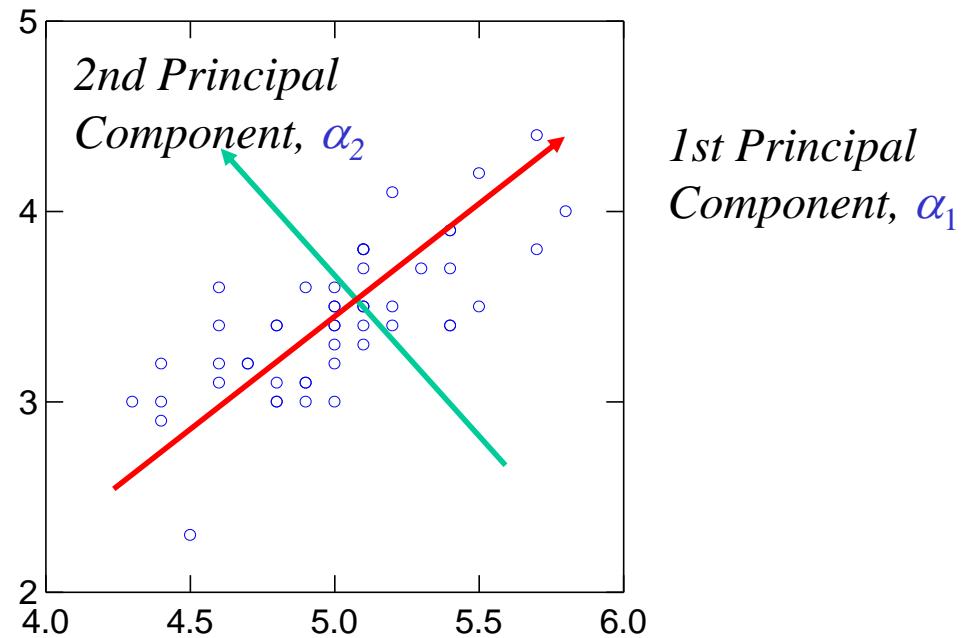
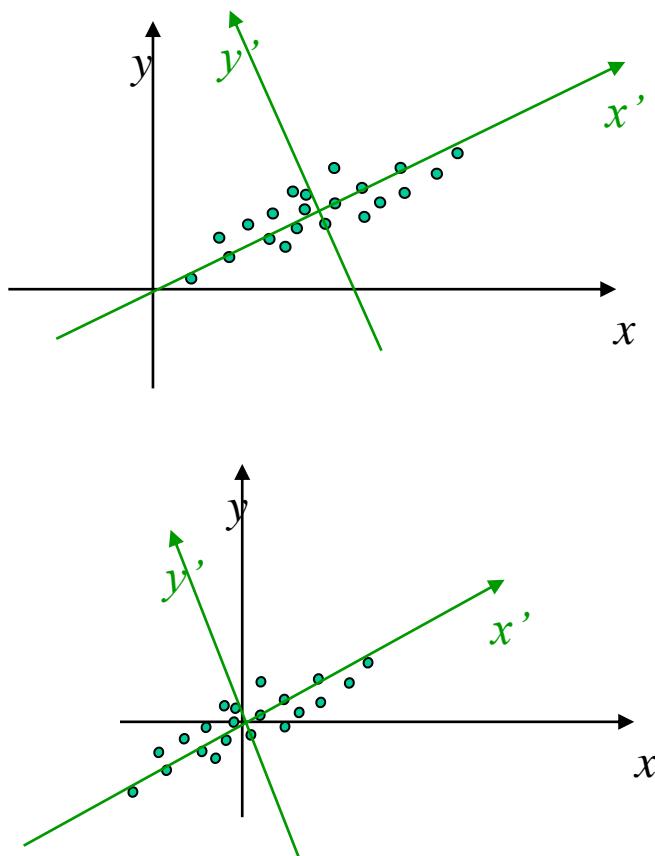
- Lower dimensional representations permit
 - Compression
 - Noise filtering
- As preprocessing for classification:
 - Reduces feature space dimension
 - Simpler Classifiers
 - Possibly better generalization
 - May facilitate simple (nearest neighbor) methods

Projection

- Orthonormal basis \rightarrow trivial projection
- Given basis $\mathbf{U} = \{\mathbf{u}_1, \dots, \mathbf{u}_k\}$
can project any d -dim \mathbf{x} to k values
 - $\alpha_1 = \mathbf{u}_1^\top \mathbf{x}$ $\alpha_2 = \mathbf{u}_2^\top \mathbf{x}$... $\alpha_k = \mathbf{u}_k^\top \mathbf{x}$
 - $\alpha = \mathbf{U}^\top \mathbf{x}$
 - $\mathbf{x} \approx \sum_i \alpha_i \mathbf{u}_i = \sum_i (\mathbf{u}_i^\top \mathbf{x}) \mathbf{u}_i$ [“=” if all d values]
- We will use “centered” vectors:
 $\mathbf{x}' = \mathbf{x} - \underline{\mathbf{x}}$ where $\underline{\mathbf{x}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}^n$

$$\alpha_i = \mathbf{u}_i^\top (\mathbf{x} - \underline{\mathbf{x}})$$

Principal Components Analysis



Minimize Reconstruction Error

- Assume data is set of N d -dimensional vectors, $\mathbf{x}^n = \langle x_1^n \dots x_d^n \rangle$
- Represent each in terms of any d orthogonal basis vectors

$$\mathbf{x}^n = \sum_{i=1}^d z_i^n \mathbf{u}_i; \quad \mathbf{u}_i^T \mathbf{u}_j = \delta_{ij}$$

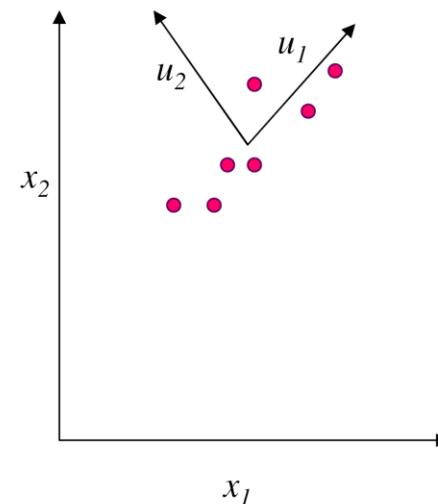
PCA: given $k < d$. Find $\{ \mathbf{u}_1, \dots, \mathbf{u}_k \}$

that minimizes $E_k = \sum_{n=1}^N \| \mathbf{x}^n - \hat{\mathbf{x}}_k^n \|_2^2$

where $\hat{\mathbf{x}}_k^n = \underline{\mathbf{x}} + \sum_{i=1}^k \alpha_i^n \mathbf{u}_i$

Mean

$$\underline{\mathbf{x}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}^n$$



PCA

- Note $\hat{\mathbf{x}}_d^n = \underline{\mathbf{x}} + \sum_{i=1}^d \alpha_i^n \mathbf{u}_i \equiv \mathbf{x}^n$
- So... $\mathbf{x}^n - \hat{\mathbf{x}}_k^n = \sum_{i=k+1}^d \alpha_i^n \mathbf{u}_i = \sum_{i=k+1}^d ((\mathbf{x}^n - \underline{\mathbf{x}})^T \mathbf{u}_i) \mathbf{u}_i$
- So... $E_k = \sum_{n=1}^N \left\| \sum_{i=k+1}^d ((\mathbf{x}^n - \underline{\mathbf{x}})^T \mathbf{u}_i) \mathbf{u}_i \right\|^2 = \sum_{n=1}^N \sum_{i=k+1}^d [(\mathbf{x}^n - \underline{\mathbf{x}})^T \mathbf{u}_i]^2$
 $= \sum_{i=k+1}^d \sum_{n=1}^N [\mathbf{u}_i^T (\mathbf{x}^n - \underline{\mathbf{x}})] [(\mathbf{x}^n - \underline{\mathbf{x}})^T \mathbf{u}_i]$
 $= \sum_{i=k+1}^d \mathbf{u}_i^T \Sigma \mathbf{u}_i$

Covariance matrix:

$$\Sigma = \sum_n (\mathbf{x}^n - \bar{\mathbf{x}})(\mathbf{x}^n - \bar{\mathbf{x}})^T$$

PCA: given $k < d$. Find $\{\mathbf{u}_1, \dots, \mathbf{u}_k\}$

that minimizes $E_k = \sum_{n=1}^N \| \mathbf{x}^n - \hat{\mathbf{x}}_k^n \|_2^2$

where $\hat{\mathbf{x}}_k^n = \underline{\mathbf{x}} + \sum_{i=1}^k \alpha_i^n \mathbf{u}_i$

ang

Justifying Use of Eigenvectors

- Goal
 - minimize: $\mathbf{u}^T \Sigma \mathbf{u}$
 - subject to: $\mathbf{u}^T \mathbf{u} = 1$
- Use Lagrange Multipliers... minimize:
$$f(\mathbf{u}) = \mathbf{u}^T \Sigma \mathbf{u} - \lambda [\mathbf{u}^T \mathbf{u} - 1]$$

- Set derivative to 0:

$$\Sigma \mathbf{u} - \lambda \mathbf{u} = 0$$

- Def'n of eigenvalue λ , eigenvector \mathbf{u} !
- If multiple vectors \mathbf{u}_i :
 - Minimize sum of independent terms...
 - Each is eigen value/vector

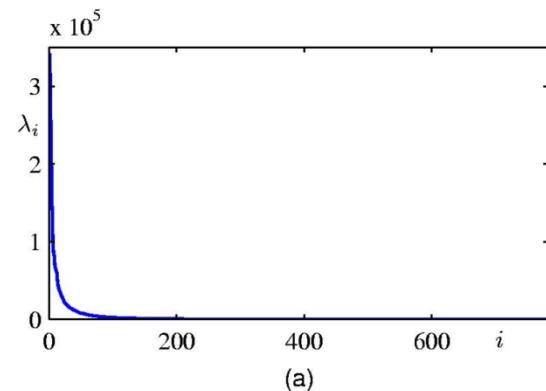
PCA

$$\text{Minimize } E_k = \sum_{i=k+1}^d \mathbf{u}_i^\top \Sigma \mathbf{u}_i$$

$$\rightarrow \sum \mathbf{u}_i = \lambda_i \mathbf{u}_i$$

Eigenvalue Eigenvector

$$\begin{aligned} \Rightarrow E_k &= \sum_{i=k+1}^d \mathbf{u}_i^\top \Sigma \mathbf{u}_i = \sum_{i=k+1}^d \mathbf{u}_i^\top \lambda_i \mathbf{u}_i \\ &= \sum_{i=k+1}^d \lambda_i \mathbf{u}_i^\top \mathbf{u}_i = \sum_{i=k+1}^d \lambda_i \end{aligned}$$



So... to minimize E_k , take **SIMILAR** eigenvalues $\{ \lambda_i \}$

PCA Algorithm

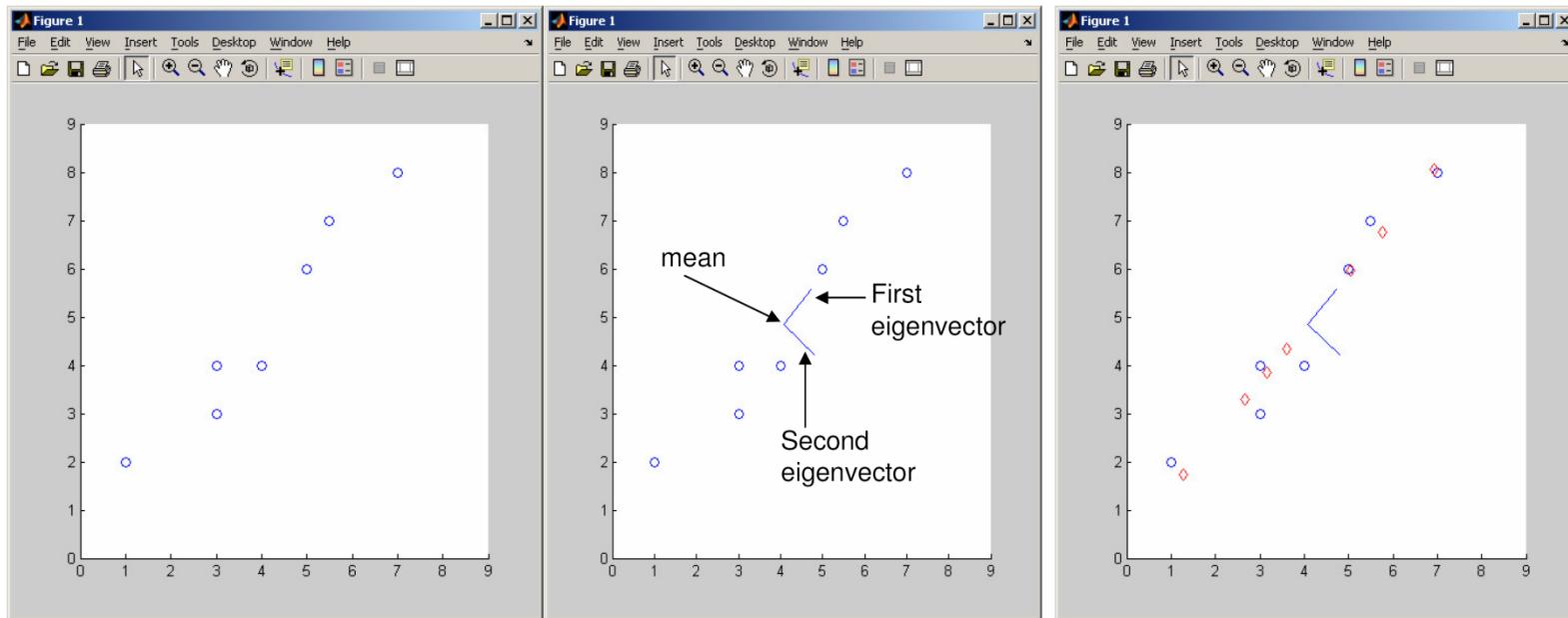
PCA algorithm(\mathbf{X} , k): top k eigenvalues/eigenvectors

% $\mathbf{X} = d \times N$ data matrix,

% ... each data point \mathbf{x}^n = column vector

- $\underline{\mathbf{x}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}^n$
- $\mathbf{A} \leftarrow$ subtract mean $\underline{\mathbf{x}}$ from each column vector \mathbf{x}^n in \mathbf{X}
- $\Sigma \leftarrow \mathbf{A} \mathbf{A}^T$... covariance matrix of \mathbf{A}
- $\{ \lambda_i, \mathbf{u}_i \}_{i=1..d} =$ eigenvectors/eigenvalues of Σ
... $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$
- Return $\{ \lambda_i, \mathbf{u}_i \}_{i=1..k}$
% top k principle components

PCA Example



Reconstructed data using
only first eigenvector ($k=1$)

PCA and SVD

- We can compute the principal components by SVD of X :

$$\begin{aligned} X &= U\Sigma V^T \\ XX^T &= U\Sigma V^T(U\Sigma V^T)^T = \\ &= U\Sigma V^T V \Sigma^T U^T = \underline{U \tilde{\Sigma}^2 U^T} \end{aligned}$$

- Thus, the **left singular vectors** of X are the principal components! We sort them by the size of the singular values of X .

PCA for Image Compression



d=1



d=2



d=4



d=8



d=16



d=32



d=64



d=100



**Original
Image**

Eigenfaces



- Example data set: Images of faces
 - Famous Eigenface approach [Turk & Pentland], [Sirovich & Kirby]
- Each face \mathbf{a} is ...
 - 256×256 values (luminance at location)
 - \mathbf{a} in $\mathbb{R}^{256 \times 256}$ (view as 1D vector)
- Form $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_m]$
- Compute $\Sigma = \mathbf{A}\mathbf{A}^T$
- Problem: Σ is $64K \times 64K$... HUGE!!!

$$\mathbf{A} = \begin{bmatrix} \mathbf{a}_1, \dots, \mathbf{a}_m \end{bmatrix}$$

$\brace{m \text{ faces}}$ $\brace{256 \times 256 \text{ real values}}$

18

Computational Complexity

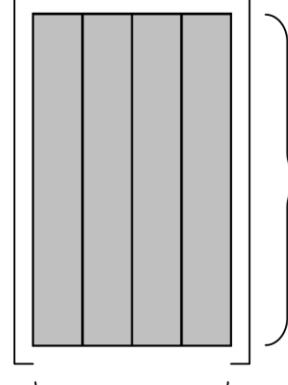
- Suppose m instances, each of size d
 - Eigenfaces: $m=500$ faces, each of size $d=64K$
- Given $d \times d$ covariance matrix Σ , can compute
 - all d eigenvectors/eigenvalues in $O(d^3)$
 - first k eigenvectors/eigenvalues in $O(k d^2)$
- But if $d=64K$, EXPENSIVE!

A Clever Workaround

- Note that $m \ll 64K$
- Use $L = A^T A$ instead of $\Sigma = AA^T$
- If v is eigenvector of L
then Av is eigenvector of Σ

$$A = \begin{bmatrix} a_1, \dots, a_m \end{bmatrix}$$

256×256
real values
 m faces



Proof: $L v = \gamma v$

$$A^T A v = \gamma v$$

$$A (A^T A v) = A(\gamma v) = \gamma Av$$

$$(A A^T) A v = \gamma (Av)$$

$$\Sigma (Av) = \gamma (Av)$$

Principle Components



DATA, MSML, BIOI 602

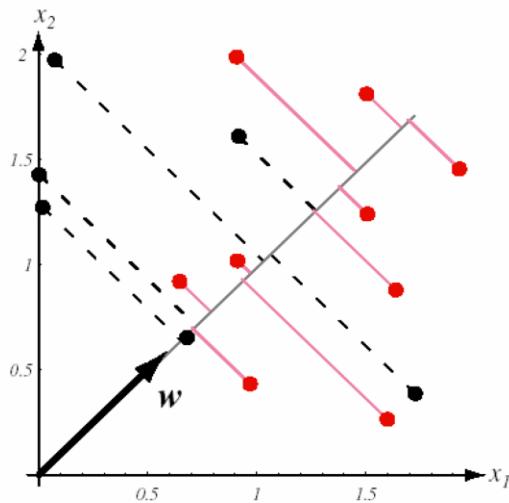
Principles of Data Science

Fisher Linear Discriminant Analysis

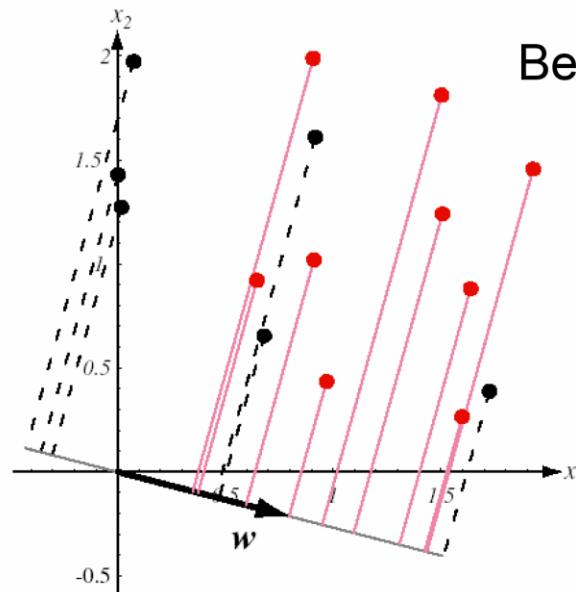
Heng Huang
Department of Computer Science

Fisher Linear Discriminant

Classes mixed



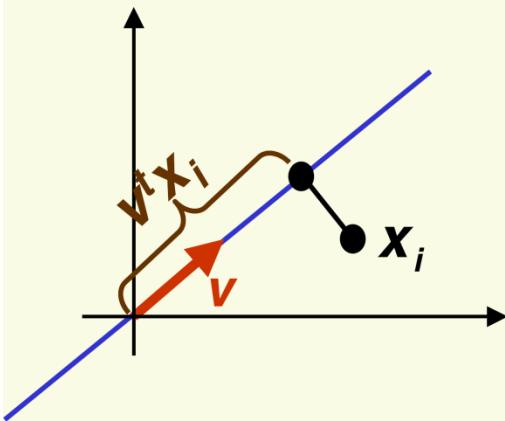
Better Separation



The figure on the right shows greater separation between subsets, one set of the points with dashed line, another with solid line.

Fisher Linear Discriminant

- Suppose we have 2 classes and d -dimensional samples $\mathbf{x}_1, \dots, \mathbf{x}_n$ where
 - n_1 samples come from the first class
 - n_2 samples come from the second class
- consider projection on a line
- Let the line direction be given by unit vector \mathbf{v}



- Scalar $\mathbf{v}^t \mathbf{x}_i$ is the distance of projection of \mathbf{x}_i from the origin
- Thus it $\mathbf{v}^t \mathbf{x}_i$ is the projection of \mathbf{x}_i into a one dimensional subspace

Fisher Linear Discriminant

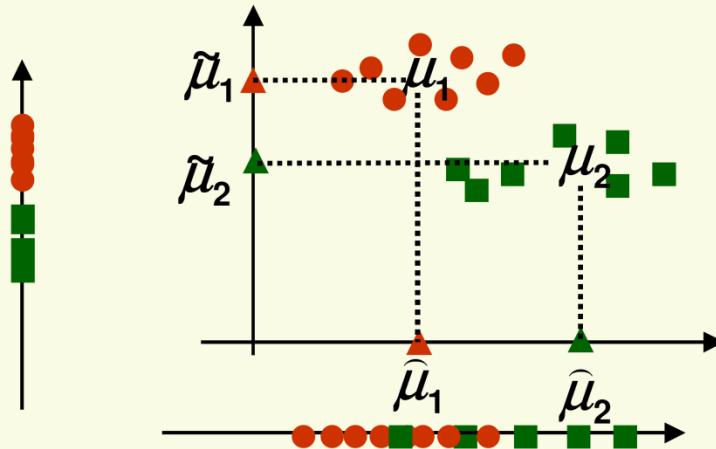
- Thus the projection of sample \mathbf{x}_i onto a line in direction \mathbf{v} is given by $\mathbf{v}^t \mathbf{x}_i$,
- How to measure separation between projections of different classes?
- Let $\tilde{\mu}_1$ and $\tilde{\mu}_2$ be the means of projections of classes 1 and 2
- Let μ_1 and μ_2 be the means of classes 1 and 2
- $|\tilde{\mu}_1 - \tilde{\mu}_2|$ seems like a good measure

$$\tilde{\mu}_1 = \frac{1}{n_1} \sum_{x_i \in C1}^{n_1} \mathbf{v}^t \mathbf{x}_i = \mathbf{v}^t \left(\frac{1}{n_1} \sum_{x_i \in C1}^{n_1} \mathbf{x}_i \right) = \mathbf{v}^t \mu_1$$

similarly, $\tilde{\mu}_2 = \mathbf{v}^t \mu_2$

Fisher Linear Discriminant

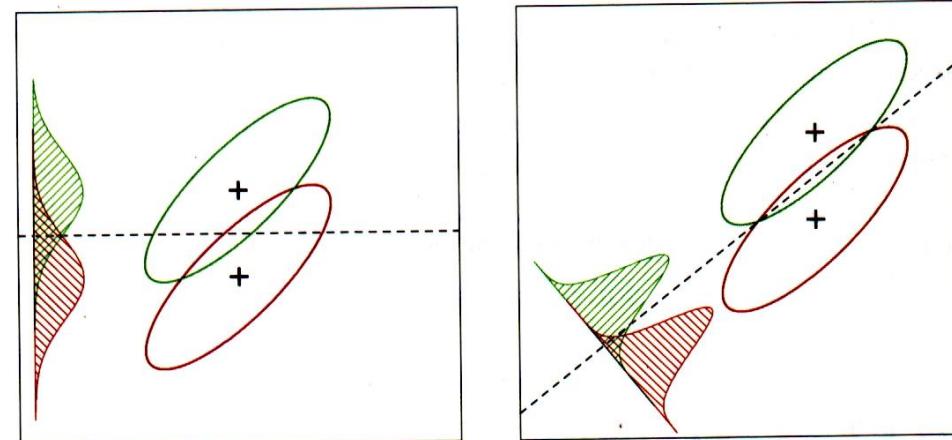
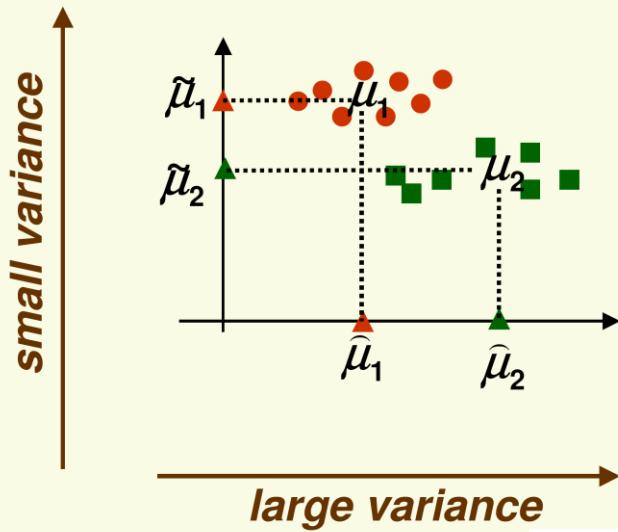
- How good is $|\tilde{\mu}_1 - \tilde{\mu}_2|$ as a measure of separation?
 - The larger $|\tilde{\mu}_1 - \tilde{\mu}_2|$, the better is the expected separation



- the vertical axes is a better line than the horizontal axes to project to for class separability
- however $|\hat{\mu}_1 - \hat{\mu}_2| > |\tilde{\mu}_1 - \tilde{\mu}_2|$

Fisher Linear Discriminant

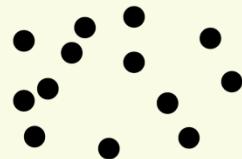
- The problem with $|\hat{\mu}_1 - \hat{\mu}_2|$ is that it does not consider the variance of the classes



Fisher Linear Discriminant

- We need to normalize $|\tilde{\mu}_1 - \tilde{\mu}_2|$ by a factor which is proportional to variance
- Have samples $\mathbf{z}_1, \dots, \mathbf{z}_n$. Sample mean is $\mu_z = \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i$
- Define their **scatter** as
$$\mathbf{s} = \sum_{i=1}^n (\mathbf{z}_i - \mu_z)^2$$
- Thus scatter is just sample variance multiplied by n
 - scatter measures the same thing as variance, the spread of data around the mean
 - scatter is just on different scale than variance

larger scatter:



smaller scatter:



Fisher Linear Discriminant

- Fisher Solution: normalize $|\tilde{\mu}_1 - \tilde{\mu}_2|$ by scatter
- Let $\mathbf{y}_i = \mathbf{v}^t \mathbf{x}_i$, i.e. \mathbf{y}_i 's are the projected samples
- Scatter for projected samples of class 1 is

$$\tilde{s}_1^2 = \sum_{\mathbf{y}_i \in \text{Class 1}} (\mathbf{y}_i - \tilde{\mu}_1)^2$$

- Scatter for projected samples of class 2 is

$$\tilde{s}_2^2 = \sum_{\mathbf{y}_i \in \text{Class 2}} (\mathbf{y}_i - \tilde{\mu}_2)^2$$

Fisher Linear Discriminant

- We need to normalize by both scatter of class 1 and scatter of class 2
- Thus Fisher linear discriminant is to project on line in the direction v which maximizes

want projected means are far from each other

$$J(v) = \frac{(\tilde{\mu}_1 - \tilde{\mu}_2)^2}{\tilde{s}_1^2 + \tilde{s}_2^2}$$

want scatter in class 1 is as small as possible, i.e. samples of class 1 cluster around the projected mean $\tilde{\mu}_1$

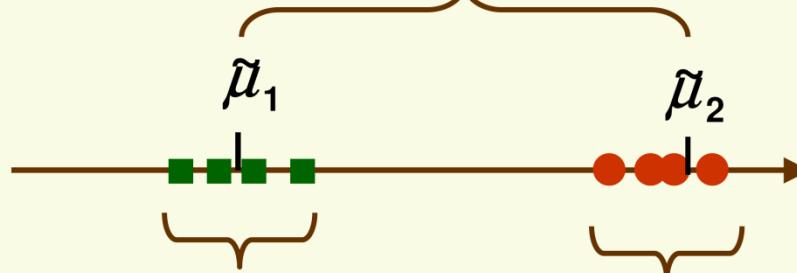
want scatter in class 2 is as small as possible, i.e. samples of class 2 cluster around the projected mean $\tilde{\mu}_2$

Fisher Linear Discriminant

$$J(\mathbf{v}) = \frac{(\tilde{\mu}_1 - \tilde{\mu}_2)^2}{\tilde{s}_1^2 + \tilde{s}_2^2}$$

- If we find \mathbf{v} which makes $J(\mathbf{v})$ large, we are guaranteed that the classes are well separated

projected means are far from each other



small \tilde{s}_1 implies that projected samples of class 1 are clustered around projected mean

small \tilde{s}_2 implies that projected samples of class 2 are clustered around projected mean

Fisher Linear Discriminant

$$J(v) = \frac{(\tilde{\mu}_1 - \tilde{\mu}_2)^2}{\tilde{s}_1^2 + \tilde{s}_2^2}$$

- All we need to do now is to express J explicitly as a function of v and maximize it
 - straightforward but need linear algebra and Calculus
- Define the separate class scatter matrices S_1 , and S_2 for classes 1 and 2. These measure the scatter of original samples x_i (before projection)

$$S_1 = \sum_{x_i \in \text{Class 1}} (x_i - \mu_1)(x_i - \mu_1)^t$$

$$S_2 = \sum_{x_i \in \text{Class 2}} (x_i - \mu_2)(x_i - \mu_2)^t$$

Fisher Linear Discriminant

- Now define the **within** the class scatter matrix

$$\mathbf{S}_W = \mathbf{S}_1 + \mathbf{S}_2$$

- Recall that $\tilde{\mathbf{s}}_1^2 = \sum_{y_i \in \text{Class 1}} (\mathbf{y}_i - \tilde{\boldsymbol{\mu}}_1)^2$

- Using $\mathbf{y}_i = \mathbf{v}^t \mathbf{x}_i$ and $\tilde{\boldsymbol{\mu}}_1 = \mathbf{v}^t \boldsymbol{\mu}_1$

$$\begin{aligned}\tilde{\mathbf{s}}_1^2 &= \sum_{y_i \in \text{Class 1}} (\mathbf{v}^t \mathbf{x}_i - \mathbf{v}^t \boldsymbol{\mu}_1)^2 \\ &= \sum_{y_i \in \text{Class 1}} (\mathbf{v}^t (\mathbf{x}_i - \boldsymbol{\mu}_1))^t (\mathbf{v}^t (\mathbf{x}_i - \boldsymbol{\mu}_1)) \\ &= \sum_{y_i \in \text{Class 1}} ((\mathbf{x}_i - \boldsymbol{\mu}_1)^t \mathbf{v})^t ((\mathbf{x}_i - \boldsymbol{\mu}_1)^t \mathbf{v}) \\ &= \sum_{y_i \in \text{Class 1}} \mathbf{v}^t (\mathbf{x}_i - \boldsymbol{\mu}_1)(\mathbf{x}_i - \boldsymbol{\mu}_1)^t \mathbf{v} = \mathbf{v}^t \mathbf{S}_1 \mathbf{v}\end{aligned}$$

Fisher Linear Discriminant

- Similarly $\tilde{\mathbf{S}}_2^2 = \mathbf{v}^t \mathbf{S}_2 \mathbf{v}$
- Therefore $\tilde{\mathbf{S}}_1^2 + \tilde{\mathbf{S}}_2^2 = \mathbf{v}^t \mathbf{S}_1 \mathbf{v} + \mathbf{v}^t \mathbf{S}_2 \mathbf{v} = \mathbf{v}^t \mathbf{S}_w \mathbf{v}$
- Define between class scatter matrix
 - $\mathbf{S}_B = (\mu_1 - \mu_2)(\mu_1 - \mu_2)^t$
- \mathbf{S}_B measures separation between the means of two classes (before projection)
- Let's rewrite the separations of the projected means

$$\begin{aligned}(\tilde{\mu}_1 - \tilde{\mu}_2)^2 &= (\mathbf{v}^t \mu_1 - \mathbf{v}^t \mu_2)^2 \\&= \mathbf{v}^t (\mu_1 - \mu_2)(\mu_1 - \mu_2)^t \mathbf{v} \\&= \mathbf{v}^t \mathbf{S}_B \mathbf{v}\end{aligned}$$

Fisher Linear Discriminant

- Thus our objective function can be written:

$$J(\mathbf{v}) = \frac{(\tilde{\mu}_1 - \tilde{\mu}_2)^2}{\tilde{\mathbf{s}}_1^2 + \tilde{\mathbf{s}}_2^2} = \frac{\mathbf{v}^t \mathbf{S}_B \mathbf{v}}{\mathbf{v}^t \mathbf{S}_W \mathbf{v}}$$

- Maximize** $J(\mathbf{v})$ by taking the derivative w.r.t. \mathbf{v} and setting it to 0

$$\begin{aligned}\frac{d}{d\mathbf{v}} J(\mathbf{v}) &= \frac{\left(\frac{d}{d\mathbf{v}} \mathbf{v}^t \mathbf{S}_B \mathbf{v} \right) \mathbf{v}^t \mathbf{S}_W \mathbf{v} - \left(\frac{d}{d\mathbf{v}} \mathbf{v}^t \mathbf{S}_W \mathbf{v} \right) \mathbf{v}^t \mathbf{S}_B \mathbf{v}}{(\mathbf{v}^t \mathbf{S}_W \mathbf{v})^2} \\ &= \frac{(2 \mathbf{S}_B \mathbf{v}) \mathbf{v}^t \mathbf{S}_W \mathbf{v} - (2 \mathbf{S}_W \mathbf{v}) \mathbf{v}^t \mathbf{S}_B \mathbf{v}}{(\mathbf{v}^t \mathbf{S}_W \mathbf{v})^2} = \mathbf{0}\end{aligned}$$

Fisher Linear Discriminant

- Need to solve $\mathbf{v}^t \mathbf{S}_W \mathbf{v} (\mathbf{S}_B \mathbf{v}) - \mathbf{v}^t \mathbf{S}_B \mathbf{v} (\mathbf{S}_W \mathbf{v}) = 0$

$$\Rightarrow \frac{\mathbf{v}^t \mathbf{S}_W \mathbf{v} (\mathbf{S}_B \mathbf{v})}{\mathbf{v}^t \mathbf{S}_W \mathbf{v}} - \frac{\mathbf{v}^t \mathbf{S}_B \mathbf{v} (\mathbf{S}_W \mathbf{v})}{\mathbf{v}^t \mathbf{S}_W \mathbf{v}} = 0$$

$$\Rightarrow \mathbf{S}_B \mathbf{v} - \frac{\mathbf{v}^t \mathbf{S}_B \mathbf{v} (\mathbf{S}_W \mathbf{v})}{\mathbf{v}^t \mathbf{S}_W \mathbf{v}} = 0$$

$$\Rightarrow \underbrace{\mathbf{S}_B \mathbf{v}}_{\lambda} = \underbrace{\mathbf{S}_W \mathbf{v}}_{\lambda}$$

generalized eigenvalue problem

Multiple Discriminant Analysis

- We want to find a transform matrix \mathbf{W} that maximizes the ratio of the determinants of the between-class scatter to the within-class scatter:

$$J(\mathbf{W}) = \frac{|\tilde{\mathbf{S}}_B|}{|\tilde{\mathbf{S}}_W|} = \frac{|\mathbf{W}^t \mathbf{S}_B \mathbf{W}|}{|\mathbf{W}^t \mathbf{S}_W \mathbf{W}|}$$

- The columns of an optimal \mathbf{W} are the generalized eigenvectors corresponding to the largest eigenvalues

$$\mathbf{S}_B \mathbf{w}_i = \lambda_i \mathbf{S}_W \mathbf{w}_i$$

- If \mathbf{S}_W is nonsingular, $\mathbf{S}_W^{-1} \mathbf{S}_B \mathbf{w}_i = \lambda_i \mathbf{w}_i$