

Recap

- (Ch 10) Data in high dimensions
 - Visualizing data
 - Summarizing data

Today

- (Ch 10) Data in high dimensions
 - Dimensionality reduction
 - Principal components analysis

Covariance matrix of multidimensional data

- Given a dataset $\{\mathbf{x}\}$ of N d -dimensional vectors \mathbf{x}_i , the covariance matrix is a $d \times d$ matrix

$$\text{Covmat}(\{\mathbf{x}\}) = \frac{\sum_i (\mathbf{x}_i - \text{mean}(\{\mathbf{x}\})) (\mathbf{x}_i - \text{mean}(\{\mathbf{x}\}))^T}{N}$$

- Properties
 - The (j, k) entry of $\text{Covmat}(\{\mathbf{x}\})$ is $\text{cov}(\{\mathbf{x}^{(j)}\}, \{\mathbf{x}^{(k)}\})$
 - The (j, j) entry of $\text{Covmat}(\{\mathbf{x}\})$ is $\text{var}(\{\mathbf{x}^{(j)}\})$
 - $\text{Covmat}(\{\mathbf{x}\})$ is symmetric

Translation properties

- Translating the data translates the mean

$$\text{mean}(\{\mathbf{x} + \mathbf{c}\}) = \text{mean}(\{\mathbf{x}\}) + \mathbf{c}$$

- Translating the data leaves the covariance matrix unchanged

$$\text{Covmat}(\{\mathbf{x} + \mathbf{c}\}) = \text{Covmat}(\{\mathbf{x}\})$$

Linear transformation properties

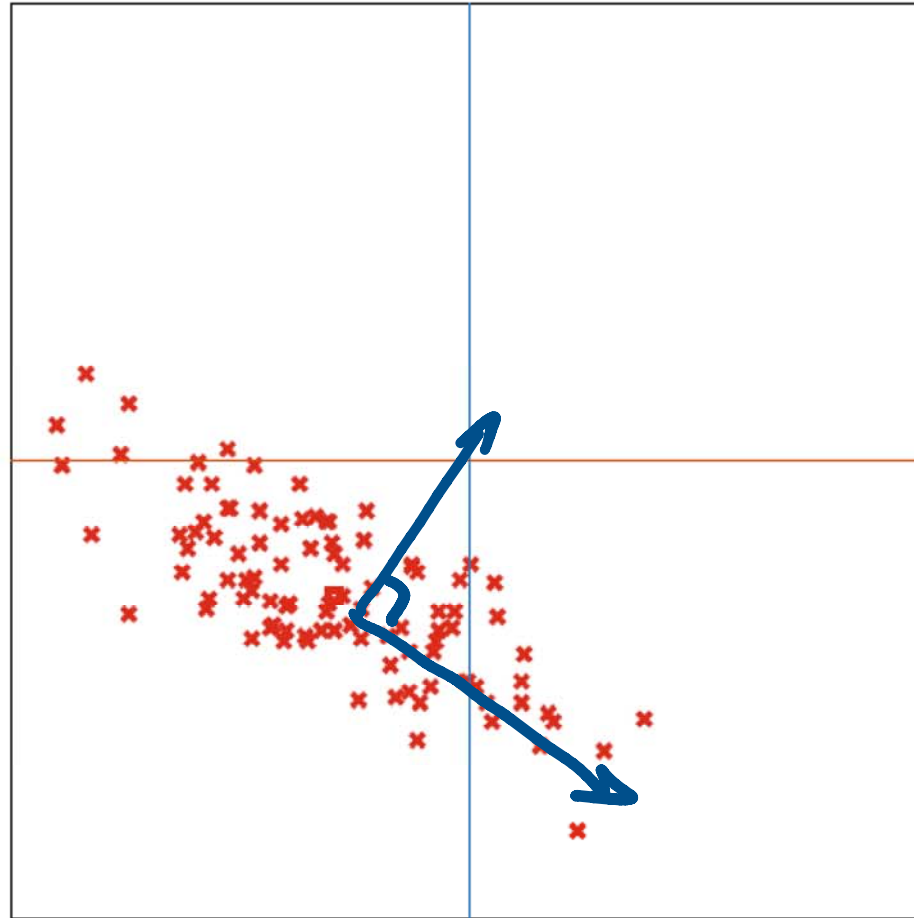
- Linearly transforming the data linearly transforms the mean

$$\text{mean}(\{A\mathbf{x}\}) = A \text{mean}(\{\mathbf{x}\})$$

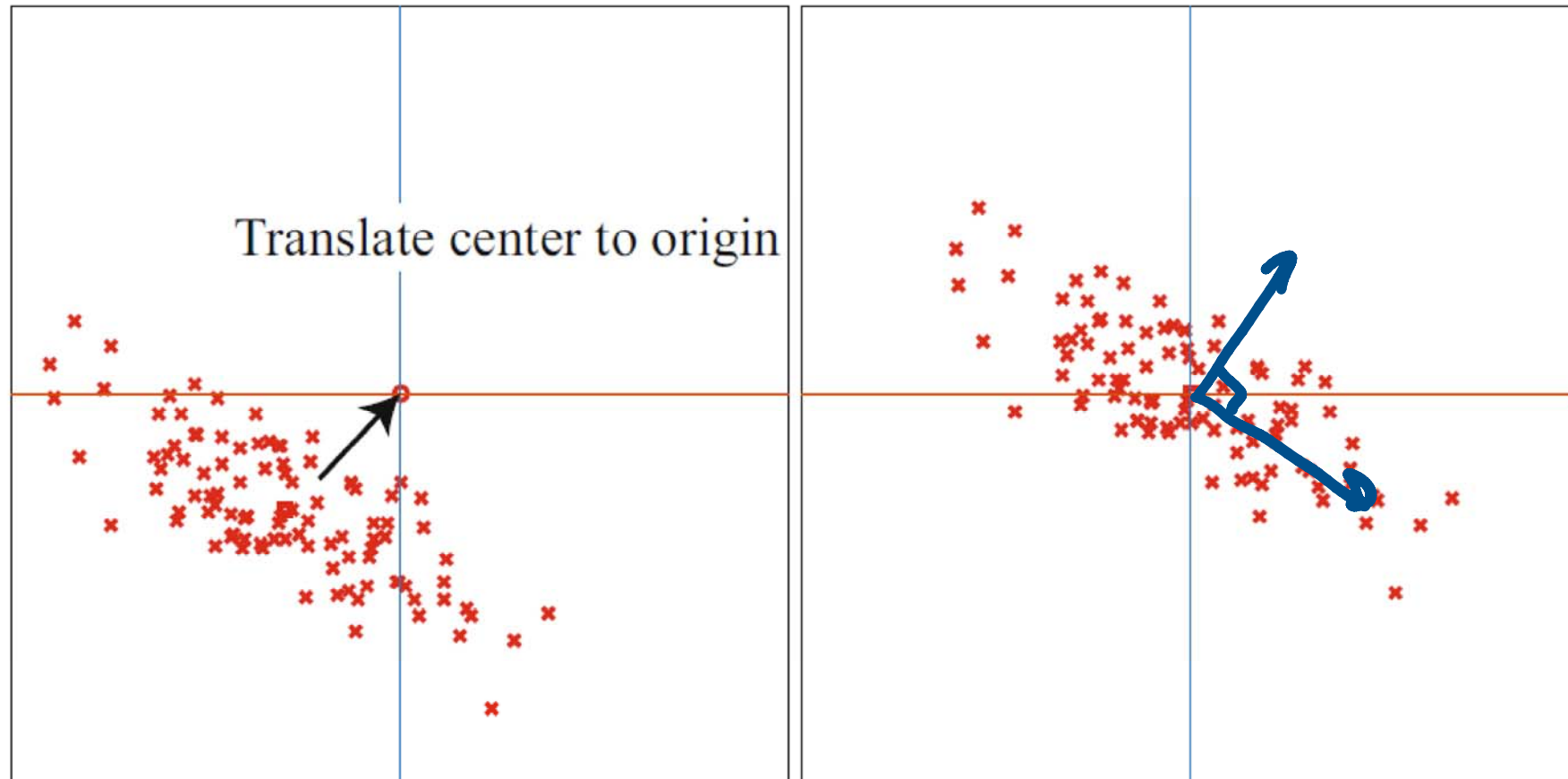
- Linearly transforming the data changes the covariance matrix

$$\text{Covmat}(\{A\mathbf{x}\}) = A \text{Covmat}(\{\mathbf{x}\}) A^T$$

Dimensionality reduction: 2D to 1D example

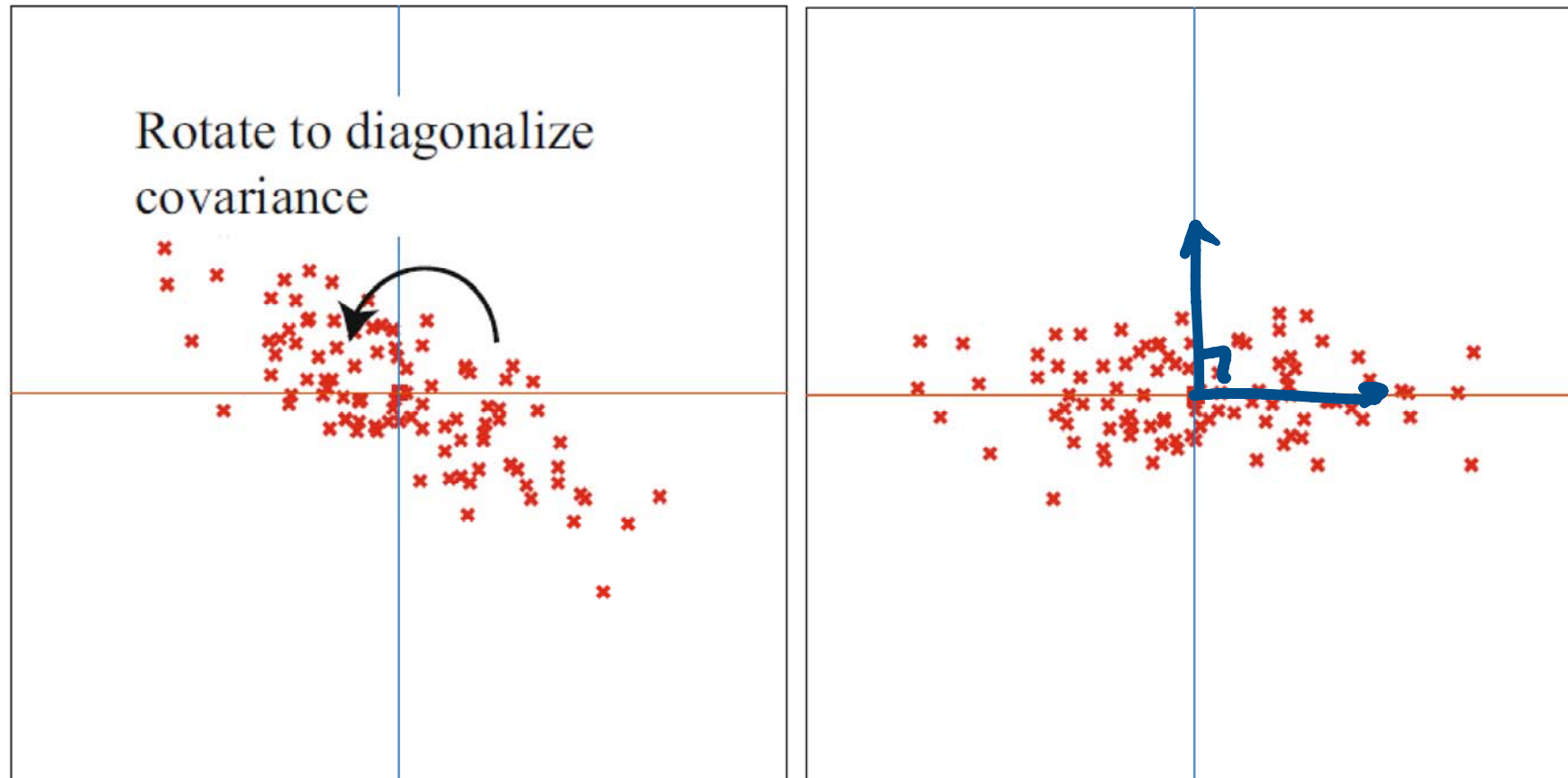


Step 1: subtract mean



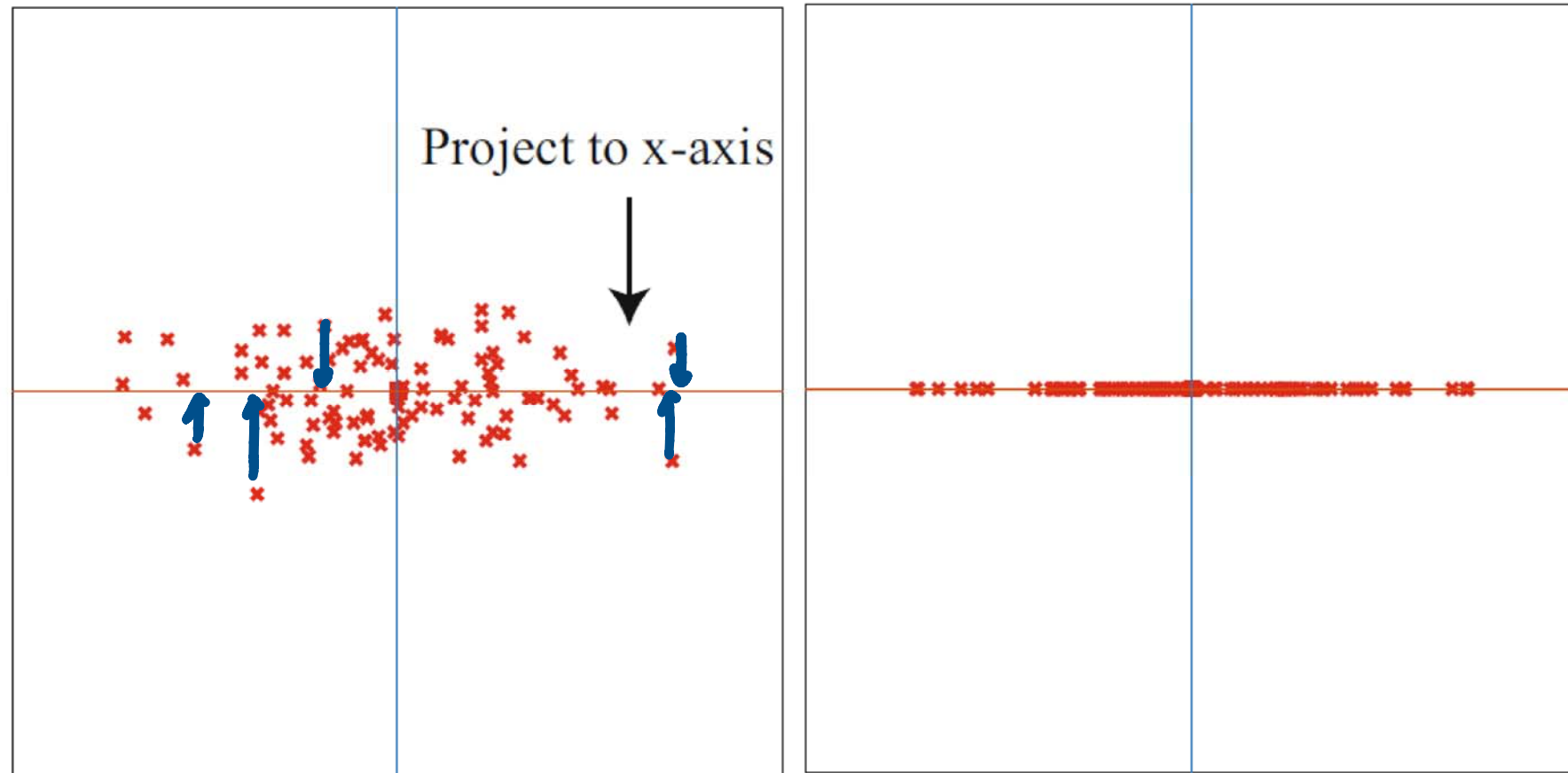
Step 2: apply linear transformation

$$\begin{bmatrix} \cdot \\ \cdot \end{bmatrix}$$



$$\begin{bmatrix} \cdot & 0 \\ 0 & \cdot \end{bmatrix}$$

Step 3: drop component(s)



Principal components analysis (PCA)

We will reduce the dimensionality of dataset $\{\mathbf{x}\}$ from d to s

- Step 1: define $\{\mathbf{m}\}$ such that $\mathbf{m}_i = \mathbf{x}_i - \text{mean}(\{\mathbf{x}\})$
- Step 2: define $\{\mathbf{r}\}$ such that $\mathbf{r}_i = U^T \mathbf{m}_i$
where $\Lambda = U^T \text{Covmat}(\{\mathbf{x}\})U$ is the diagonalization of $\text{Covmat}(\{\mathbf{x}\})$ with the eigenvalues sorted in decreasing order
- Step 3: define $\{\mathbf{p}\}$ such that each \mathbf{p}_i is \mathbf{r}_i with the last $d - s$ components zeroed out (or discarded)

What happens to the mean?

- Step 1:

$$\text{mean}(\{\mathbf{m}\}) = \text{mean}(\{\mathbf{x}\}) - \text{mean}(\{\mathbf{x}\}) = \mathbf{0}$$

- Step 2:

$$\text{mean}(\{\mathbf{r}\}) = U^T \text{mean}(\{\mathbf{m}\}) = U^T \mathbf{0} = \mathbf{0}$$

- Step 3:

$$\text{mean}(\{\mathbf{p}\}) = \text{mean}(\{\mathbf{r}\}) = \mathbf{0}$$

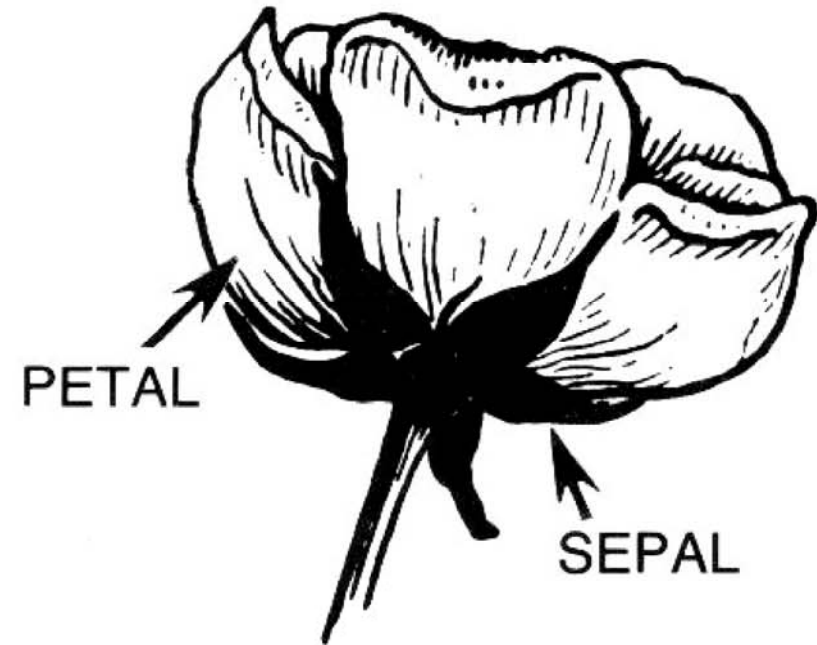
Mean square error of the projection (step 3)

$$\begin{aligned}
 \frac{1}{N} \sum_i \|\mathbf{r}_i - \mathbf{p}_i\|^2 &= \frac{1}{N} \sum_i \sum_{j=s+1}^d (r_i^{(j)})^2 = \sum_{j=s+1}^d \underbrace{\frac{1}{N} \sum_i (r_i^{(j)})^2}_{\text{var}(\{r^{(j)}\})} \\
 &= \sum_{j=s+1}^d \text{var}(\{r^{(j)}\}) \quad \text{since mean}(\{r\}) = 0 \\
 &= \sum_{j=s+1}^d \lambda_j \quad \text{since the variances are diagonal} \\
 &\quad \text{entries of Covmat}(\{r\}) = \Lambda
 \end{aligned}$$

The mean square error is the sum of the smallest $d - s$ eigenvalues in Λ

PCA: iris dataset example

- The Iris dataset is a famous dataset consisting of measurements of three different varieties of iris flowers
 - Iris-setosa
 - Iris-versicolor
 - Iris-virginica
- There are 4 measurements per item
 - Sepal length (cm)
 - Sepal width (cm)
 - Petal length (cm)
 - Petal width (cm)
- See today's Jupyter notebook



Reconstructing the data

- Given the projected dataset $\{\mathbf{p}\}$ and $\text{mean}(\{\mathbf{x}\})$, we can approximately reconstruct the original dataset as $\{\hat{\mathbf{x}}\}$

$$\hat{\mathbf{x}}_i = U\mathbf{p}_i + \text{mean}(\{\mathbf{x}\})$$

- Notice that each reconstructed data item $\hat{\mathbf{x}}_i$ is $\text{mean}(\{\mathbf{x}\})$ plus a linear combination of the columns of U weighted by the entries in \mathbf{p}_i
- The columns of U are the normalized eigenvectors of $\text{Covmat}(\{\mathbf{x}\})$ and are called the **principal components** of the data $\{\mathbf{x}\}$

End-to-end mean square error

- Each \mathbf{x}_i becomes \mathbf{r}_i by translation and rotation
- Each \mathbf{p}_i becomes $\hat{\mathbf{x}}_i$ by the opposite rotation and translation
- Therefore, the end-to-end mean square error is

$$\frac{1}{N} \sum_i \|\hat{\mathbf{x}}_i - \mathbf{x}_i\|^2 = \frac{1}{N} \sum_i \|\mathbf{r}_i - \mathbf{p}_i\|^2 = \sum_{j=s+1}^d \lambda_j$$

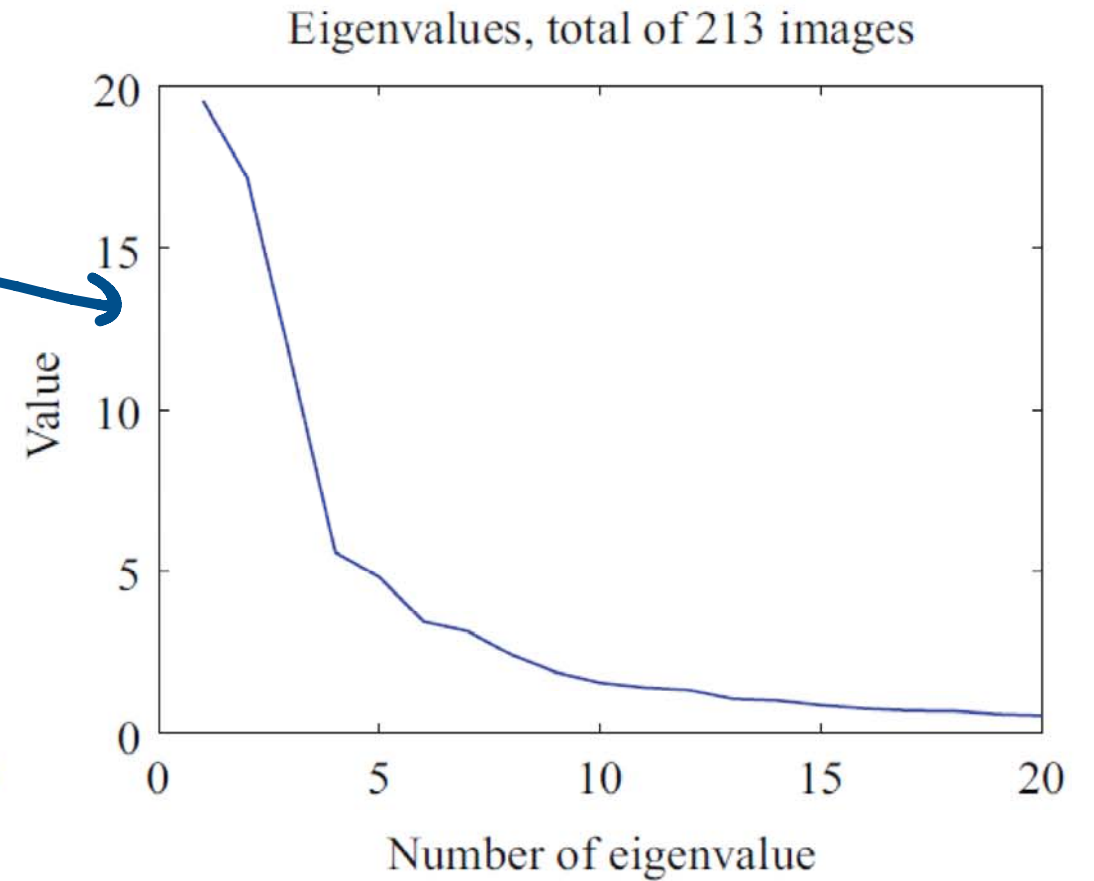
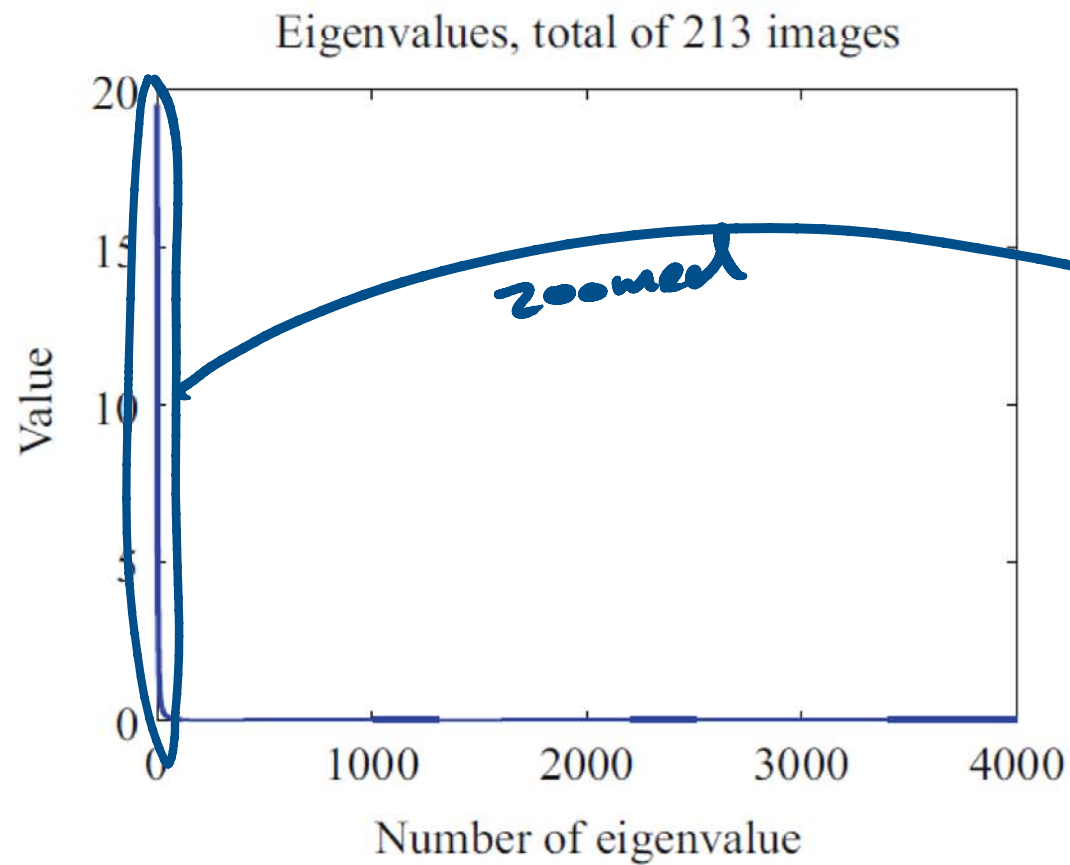
where $\lambda_{s+1}, \dots, \lambda_d$ are the smallest $d - s$ eigenvalues of $\text{Covmat}(\{\mathbf{x}\})$

PCA: Japanese face dataset example

- The dataset consists of 213 images of Japanese women
- Each image is grayscale and has 64×64 resolution
- We can treat each image as a vector of dimension $d = 4096$



How quickly do the eigenvalues drop off?



What do the principal components look like?

Mean image from Japanese Facial Expression dataset



First sixteen principal components



- The mean face is blurry
- The first few principal components capture
 - Shape of hair
 - Height of face
 - Height of eyebrows
 - Etc.

What do the reconstructions look like?



x_i

Number of principal components

Original image

mean

1

5

10

20

50

100

Reconstructions



Error

