

# Assignment Project Exam Help

COMP3223 Foundations of ML: Dimensionality reduction

by

Principal Components Analysis (PCA)

<https://powcoder.com>

Srinandan Dasmahapatra

Add WeChat powcoder

Are more complex models better?

# Assignment Project Exam Help

- With large number of features, n, accidental correlations can obscure genuine patterns.

<https://powcoder.com>

Add WeChat powcoder

# Assignment Project Exam Help

- With large number of features,  $p$ , accidental correlations can obscure genuine patterns.
- In minimising (say) SSR you may be trying to reduce residuals to “noise” at the expense of “signal”. Eg., linear model  $y^n = f(x^n; w) := \sum_{i=1}^p w_i x_i^n$

<https://powcoder.com>

$$\sum_{n=1}^N (r^n)^2 = \sum (y^n - f(x^n; w))^2$$

Add WeChat powcoder

# Assignment Project Exam Help

- With large number of features  $n$ , accidental correlations can obscure genuine patterns.
- In minimising (say) SSR you may be trying to reduce residuals to “noise” at the expense of “signal”. Eg., linear model  $y^n = f(x^n; w) := \sum_{i=1}^p w_i x_i^n$

<https://powcoder.com>

$$\sum_{n=1}^N (r^n)^2 = \sum (y^n - f(x^n; w))^2$$

Add WeChat powcoder

- In trying to minimise the residuals, you may be adjusting coefficients  $w_k$  that correspond to features irrelevant to the actual signal.

# Assignment Project Exam Help

- With large number of features  $p$ , accidental correlations can obscure genuine patterns.
- In minimising (say) SSR you may be trying to reduce residuals to “noise” at the expense of “signal”. Eg., linear model  $y^n = f(x^n; w) := \sum_{i=1}^p w_i x_i^n$

<https://powcoder.com>

$$\sum_{n=1}^N (r^n)^2 = \sum (y^n - f(x^n; w))^2$$

Add WeChat powcoder

- In trying to minimise the residuals, you may be adjusting coefficients  $w_k$  that correspond to features irrelevant to the actual signal.
- This is the  $p > N$  problem

## Revision: linear regression

- Linear model:  $y = Xw$ ,  $X$  design matrix,  $\hat{w} = (X^T X)^{-1} X^T y$ ; also  $w_0$  is the difference between averages of inputs and outputs

Assignment Project Exam Help

$$0 = \sum_{n=1}^N r^n = \sum_n y^n - Nw_0 - \sum_n (w_1 x_1^n + \dots + w_p x_p^n)$$
$$\Rightarrow w_0 = \langle y \rangle - \sum_{i=1}^p w_i \langle x_i \rangle$$

Add WeChat powcoder

## Revision: linear regression

- Linear model:  $\mathbf{y} = \mathbf{X}\mathbf{w}$ ,  $\mathbf{X}$  design matrix,  $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ ; also  $w_0$  is the difference between averages of inputs and outputs.

# Assignment Project Exam Help

$$0 = \sum_{n=1}^N r^n = \sum_n y^n - Nw_0 - \sum_n (w_1 x_1^n + \dots + w_p x_p^n)$$
$$\Rightarrow w_0 = \langle y \rangle - \sum_{i=1}^p w_i \langle x_i \rangle$$

- Centering:**  $\tilde{\mathbf{y}}^T = \mathbf{y}^T - \langle \mathbf{y} \rangle$ ,  $\tilde{x}_j^T = \mathbf{x}_j^T - \langle \mathbf{x}_j \rangle$  and drop the column of ones from design matrix to form  $\tilde{\mathbf{X}}$  with matrix elements  $(\tilde{\mathbf{X}})_{nj} = x_j^n - \langle x_j \rangle$ .

# Add WeChat powcoder

## Revision: linear regression

- Linear model:  $\mathbf{y} = \mathbf{X}\mathbf{w}$ ,  $\mathbf{X}$  design matrix,  $\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$ ; also  $w_0$  is the difference between averages of inputs and outputs.

# Assignment Project Exam Help

$$0 = \sum_{n=1}^N r^n = \sum_n y^n - Nw_0 - \sum_n (w_1 x_1^n + \dots + w_p x_p^n)$$
$$\Rightarrow w_0 = \langle y \rangle - \sum_{i=1}^p w_i \langle x_i \rangle$$

- Centering:**  $\tilde{\mathbf{y}}^\top = \mathbf{y}^\top - \langle \mathbf{y} \rangle$ ,  $\tilde{x}_j^\top = \mathbf{x}_j^\top - \langle \mathbf{x}_j \rangle$  and drop the column of ones from design matrix to form  $\tilde{\mathbf{X}}$  with matrix elements  $(\tilde{\mathbf{X}})_{nj} = x_j^n - \langle x_j \rangle$ .

- $\frac{1}{N} \tilde{\mathbf{X}}^\top (\tilde{\mathbf{y}} - \tilde{\mathbf{X}} \hat{\mathbf{w}}) = 0$  implies  $\hat{\mathbf{w}} = \left( \frac{1}{N} (\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}}) \right)^{-1} \left( \frac{1}{N} \tilde{\mathbf{X}}^\top \mathbf{y} \right)$ , so

$$\hat{\mathbf{w}} = [\text{cov}(\mathbf{X})^{-1}] [\text{cov}(\mathbf{X}, \mathbf{y})].$$

Covariance as trace suggests reduction of matrix size

- $\tilde{y}^n = y^n - \langle y \rangle$ ,  $\tilde{x}_i^n = x_i^n - \langle x_i \rangle$  and drop the column of ones from design matrix to form  $\tilde{X}$  with matrix elements  $(\tilde{X})_{nj} = x_j^n - \langle x_j \rangle$ .

# Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

## Covariance as trace suggests reduction of matrix size

- $\tilde{y}^n = y^n - \langle y \rangle$ ,  $\tilde{x}_i^n = x_i^n - \langle x_i \rangle$  and drop the column of ones from design matrix to form  $\tilde{X}$  with matrix elements  $(\tilde{X})_{nj} = x_j^n - \langle x_j \rangle$ .
- $(n \times n)$  and  $(p \times p)$  matrices formed out of  $X$ .

$$(\tilde{X} \tilde{X}^T)_{mn} = \sum_{j=1}^p (\tilde{X})_{mj} (\tilde{X}^T)_{jn} = \sum_{i=1}^p (x_j^m - \langle x_j \rangle)(x_j^n - \langle x_j \rangle)$$

$$(\tilde{X}^T \tilde{X})_{ij} = \sum_{n=1}^N (\tilde{X}^T)_{in} (\tilde{X})_{nj} = \sum_{n=1}^N (x_i^n - \langle x_i \rangle)(x_j^n - \langle x_j \rangle)$$

Add WeChat powcoder

## Covariance as trace suggests reduction of matrix size

- $\tilde{y}^n = y^n - \langle y \rangle$ ,  $\tilde{x}_i^n = x_i^n - \langle x_i \rangle$  and drop the column of ones from design matrix to form  $\tilde{X}$  with matrix elements  $(\tilde{X})_{nj} = x_j^n - \langle x_j \rangle$ .
- $(n \times n)$  and  $(p \times p)$  matrices formed out of  $X$ .

# Assignment Project Exam Help

$$(\tilde{X} \tilde{X}^T)_{mn} = \sum_{j=1}^p (\tilde{X})_{mj} (\tilde{X}^T)_{jn} = \sum_{i=1}^p (x_j^m - \langle x_j \rangle)(x_j^n - \langle x_j \rangle)$$

$$(\tilde{X}^T \tilde{X})_{ij} = \sum_{n=1}^N (\tilde{X}^T)_{in} (\tilde{X})_{nj} = \sum_{n=1}^N (x_i^n - \langle x_i \rangle)(x_j^n - \langle x_j \rangle)$$

Add WeChat powcoder

- Trace of  $(\tilde{X} \tilde{X}^T)$  equals trace( $\tilde{X}^T \tilde{X}$ ):

$$\sum_{n=1}^N \sum_{j=1}^p (\tilde{X})_{mj} (\tilde{X}^T)_{jn} \delta_{mn} = \sum_{j=1}^p \sum_{n=1}^N (\tilde{X}^T)_{in} (\tilde{X})_{nj} \delta_{ij}$$

## Covariance as trace suggests reduction of matrix size

- $\tilde{y}^n = y^n - \langle y \rangle$ ,  $\tilde{x}_i^n = x_i^n - \langle x_i \rangle$  and drop the column of ones from design matrix to form  $\tilde{X}$  with matrix elements  $(\tilde{X})_{nj} = x_j^n - \langle x_j \rangle$ .
- $(n \times n)$  and  $(p \times p)$  matrices formed out of  $X$ .

# Assignment Project Exam Help

$$(\tilde{X}\tilde{X}^T)_{mn} = \sum_{j=1}^p (\tilde{X})_{mj}(\tilde{X}^T)_{jn} = \sum_{i=1}^p (x_j^m - \langle x_j \rangle)(x_j^n - \langle x_j \rangle)$$

$$(\tilde{X}^T\tilde{X})_{ij} = \sum_{n=1}^N (\tilde{X}^T)_{in}(\tilde{X})_{nj} = \sum_{n=1}^N (x_i^n - \langle x_i \rangle)(x_j^n - \langle x_j \rangle)$$

Add WeChat powcoder

- Trace of  $(\tilde{X}\tilde{X}^T)$  equals trace( $\tilde{X}^T\tilde{X}$ ):

$$\sum_{n=1}^N \sum_{j=1}^p (\tilde{X})_{mj}(\tilde{X}^T)_{jn} \delta_{mn} = \sum_{j=1}^p \sum_{n=1}^N (\tilde{X}^T)_{in}(\tilde{X})_{nj} \delta_{ij}$$

- Work with  $p$ -dim matrix not  $N$ -dim.

## Use of SVD in low rank approximation gives PCA

- Total variance of data  $X$  is the sum of eigenvalues of  $\frac{1}{N} \text{tr}(\tilde{X}^T \tilde{X})$

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

## Use of SVD in low rank approximation gives PCA

- Total variance of data  $X$  is the sum of eigenvalues of  $\frac{1}{N} \text{tr}(\tilde{X}^T \tilde{X})$
- If  $\tilde{X} = USV^T$ ,  $y = USV^T w$  and  $\hat{w} = V(S)^{-1}U^T y$ .  
 $S \in \text{diag}(\sigma_1, \dots, \sigma_{\min(p, N)})$ .  $U, V$  contain singular vectors of size  $N$  and  $p$  respectively.

$$\hat{w} = \sum_{k=1}^{\min(p, N)} \frac{u_k^T y}{\sigma_k} v_k$$

Add WeChat powcoder

## Use of SVD in low rank approximation gives PCA

- Total variance of data  $\tilde{X}$  is the sum of eigenvalues of  $\frac{1}{N} \text{tr}(\tilde{X}^\top \tilde{X})$
- If  $\tilde{X} = USV^\top$ ,  $y = USV^\top w$  and  $\hat{w} = V(S)^{-1}U^\top y$ .  
 $S \in \text{diag}(\sigma_1, \dots, \sigma_{\min(p, N)})$ .  $U, V$  contain singular vectors of size  $N$  and  $p$  respectively.

$$\hat{w} = \sum_{k=1}^{\min(p, N)} \frac{u_k^\top y}{\sigma_k} v_k$$

- Total variance is  $\frac{1}{N} \text{tr}(\tilde{X}^\top \tilde{X}) = \frac{1}{N} \sum_{k=1}^{\min(p, N)} \sigma_k^2$

Add WeChat powcoder

## Use of SVD in low rank approximation gives PCA

- Total variance of data  $X$  is the sum of eigenvalues of  $\frac{1}{N} \text{tr}(\tilde{X}^T \tilde{X})$
- If  $\tilde{X} = USV^T$ ,  $y = USV^T w$  and  $\hat{w} = V(S)^{-1}U^T y$ .  
 $S \in \text{diag}(\sigma_1, \dots, \sigma_{\min(p, N)})$ .  $U, V$  contain singular vectors of size  $N$  and  $p$  respectively.

$$\hat{w} = \sum_{k=1}^{\min(p, N)} \frac{u_k^T y}{\sigma_k} v_k$$

- Total variance is  $\frac{1}{N} \text{tr}(\tilde{X}^T \tilde{X}) = \frac{1}{N} \sum_{k=1}^{\min(p, N)} \sigma_k^2$
- **Discard small values of  $\sigma_k$** , keep largest  $r$  singular vectors. Fraction of variation in data explained by  $r$  components is

$$\left( \sum_{k=1}^r \sigma_k^2 \right) / \left( \sum_{k=1}^{\min(p, N)} \sigma_k^2 \right)$$

## Use of SVD in low rank approximation gives PCA

- Total variance of data  $X$  is the sum of eigenvalues of  $\frac{1}{N} \text{tr}(\tilde{X}^T \tilde{X})$
- If  $\tilde{X} = USV^T$ ,  $y = USV^T w$  and  $\hat{w} = V(S)^{-1}U^T y$ .  
 $S \in \text{diag}(\sigma_1, \dots, \sigma_{\min(p, N)})$ .  $U, V$  contain singular vectors of size  $N$  and  $p$  respectively.

<https://powcoder.com>

- Total variance is  $\frac{1}{N} \text{tr}(\tilde{X}^T \tilde{X}) = \frac{1}{N} \sum_{k=1}^{\min(p, N)} \sigma_k^2$
- **Discard small values of  $\sigma_k$** , keep largest  $r$  singular vectors. Fraction of variation in data explained by  $r$  components is

$$\left( \sum_{k=1}^r \sigma_k^2 \right) / \left( \sum_{k=1}^{\min(p, N)} \sigma_k^2 \right)$$

- $v_1, \dots, v_r$  are the **principal components** (of variation).

SVD is a low rank approximation to any matrix



# Assignment Project Exam Help

<https://powcoder.com>

- Single image as matrix  $M$

Add WeChat powcoder

Reconstruction for  $r = 1, \dots, 10$

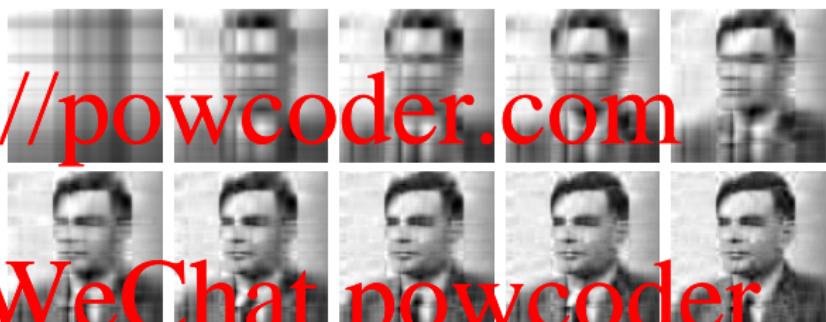


SVD is a low rank approximation to any matrix



# Assignment Project Exam Help

<https://powcoder.com>



- Single image as matrix  $M$
- $M_{1200 \times 1200}$

Add WeChat powcoder

Reconstruction for  $r = 1, \dots, 10$

SVD is a low rank approximation to any matrix



# Assignment Project Exam Help

<https://powcoder.com>



- Single image as matrix  $M$

Add WeChat powcoder

- $M_{1200 \times 1200}$

Reconstruction for  $r = 1, \dots, 10$

- SVD  $M = U\Sigma V^T$

SVD is a low rank approximation to any matrix



# Assignment Project Exam Help

<https://powcoder.com>



- Single image as matrix  $M$
- $M_{1200 \times 1200}$
- SVD  $M = U\Sigma V^T$
- Reconstruction:  
 $\widetilde{M} = \sum_i^r \sigma_i u_i v_i^T$

Add WeChat powcoder

Reconstruction for  $r = 1, \dots, 10$

First principal component captures greatest variation

- Data (mean subtracted)  $x^n \in \mathbb{R}^p$   $i = n, \dots, N$  arranged in data matrix  $X$ .

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

First principal component captures greatest variation

- Data (mean subtracted)  $x^n \in \mathbb{R}^p$   $i = n, \dots, N$  arranged in data matrix  $X$ .
- Linear combinations of data vectors:  $c^n = \sum_{j=1}^p w_{ij} x_j^n$  written as  $\mathbf{z}^n w$ .

# Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

First principal component captures greatest variation

- Data (mean subtracted)  $x^n \in \mathbb{R}^p$   $i = n, \dots, N$  arranged in data matrix  $X$ .
- Linear combinations of data vectors  $c^T = \sum_{j=1}^p w_j x_j^n$  written as  $\mathbf{X}\mathbf{w}$ .
- $\text{var}(X\mathbf{w}) = \langle \mathbf{w}^T X^T X \mathbf{w} \rangle = \mathbf{w}^T S \mathbf{w}$  where  $S = \text{var}(X)$ .

<https://powcoder.com>

Add WeChat powcoder

## First principal component captures greatest variation

- Data (mean subtracted)  $x^n \in \mathbb{R}^p$   $i = n, \dots, N$  arranged in data matrix  $X$ .
- Linear combinations of data vectors  $z^{(i)} = \sum_{j=1}^p w_j x_j^{(i)}$  written as  $\mathbf{z}^T \mathbf{w}$ .
- $\text{var}(X\mathbf{w}) = \langle \mathbf{w}^T X^T X \mathbf{w} \rangle = \mathbf{w}^T S \mathbf{w}$  where  $S = \text{var}(X)$ .
- Seek components of linear combination  $\mathbf{w}$  that maximise the variance  $\mathbf{w}^T S \mathbf{w}$  of the combined data points  $z^{(i)}$ ,  $i = 1, \dots, N$ . Normalise  $\mathbf{w}$ . Seek

$$\underset{\mathbf{w}}{\operatorname{argmax}} \left\{ \mathbf{w}^T S \mathbf{w} - \lambda(\mathbf{w}^T \mathbf{w} - 1) \right\} \implies S \mathbf{w} = \lambda \mathbf{w}.$$

Add WeChat powcoder

## First principal component captures greatest variation

- Data (mean subtracted)  $x^n \in \mathbb{R}^p$   $i = n, \dots, N$  arranged in data matrix  $X$ .
- Linear combinations of data vectors  $z^{(i)} = \sum_{j=1}^p w_j x_j^{(i)}$  written as  $\mathbf{z}'\mathbf{w}$ .
- $\text{var}(X\mathbf{w}) = \langle \mathbf{w}^T X^T X \mathbf{w} \rangle = \mathbf{w}^T S \mathbf{w}$  where  $S = \text{var}(X)$ .
- Seek components of linear combination  $\mathbf{w}$  that maximise the variance  $\mathbf{w}^T S \mathbf{w}$  of the combined data points  $z^{(i)}$ ,  $i = 1, \dots, N$ . Normalise  $\mathbf{w}$ . Seek

$$\underset{\mathbf{w}}{\operatorname{argmax}} \left\{ \mathbf{w}^T S \mathbf{w} - \lambda(\mathbf{w}^T \mathbf{w} - 1) \right\} \implies S \mathbf{w} = \lambda \mathbf{w}.$$

- Order eigenvalues  $(\lambda_1, \dots, \lambda_p)$ .  $\mathbf{w}_1$  first principal component.

## First principal component captures greatest variation

- Data (mean subtracted)  $x^n \in \mathbb{R}^p$   $i = n, \dots, N$  arranged in data matrix  $X$ .
- Linear combinations of data vectors  $z^{(n)} = \sum_{j=1}^p w_j x_j^n$  written as  $\mathbf{z}^n \mathbf{w}$ .
- $\text{var}(X\mathbf{w}) = \langle \mathbf{w}^T X^T X \mathbf{w} \rangle = \mathbf{w}^T S \mathbf{w}$  where  $S = \text{var}(X)$ .
- Seek components of linear combination  $\mathbf{w}$  that maximise the variance  $\mathbf{w}^T S \mathbf{w}$  of the combined data points  $z^{(n)}$ ,  $n = 1, \dots, N$ . Normalise  $\mathbf{w}$ . Seek

<https://powcoder.com>

$$\underset{\mathbf{w}}{\operatorname{argmax}} \{ \mathbf{w}^T S \mathbf{w} - \lambda(\mathbf{w}^T \mathbf{w} - 1) \} \implies S \mathbf{w} = \lambda \mathbf{w}.$$

- Order eigenvalues  $(\lambda_1, \dots, \lambda_p)$ .  $\lambda_1$  first principal component.
- Linear combinations  $z_1^n = w_{1,1} x_1^n + w_{1,2} x_2^n + \dots + w_{1,p} x_p^n$  constitute representation of data in terms of first principal component. Instead of  $p$  components  $(x_1^n, \dots, x_p^n)$ , **one** component  $z_1^n$  represents  $x^n$ .

## First principal component captures greatest variation

- Data (mean subtracted)  $x^n \in \mathbb{R}^p$   $i = n, \dots, N$  arranged in data matrix  $X$ .
- Linear combinations of data vectors  $z^{(n)} = \sum_{j=1}^p w_j x_j^n$  written as  $\mathbf{z}^n \mathbf{w}$ .
- $\text{var}(X\mathbf{w}) = \langle \mathbf{w}^T X^T X \mathbf{w} \rangle = \mathbf{w}^T S \mathbf{w}$  where  $S = \text{var}(X)$ .
- Seek components of linear combination  $\mathbf{w}$  that maximise the variance  $\mathbf{w}^T S \mathbf{w}$  of the combined data points  $z^{(n)}$ ,  $n = 1, \dots, N$ . Normalise  $\mathbf{w}$ . Seek

<https://powcoder.com>

$$\underset{\mathbf{w}}{\operatorname{argmax}} \{ \mathbf{w}^T S \mathbf{w} - \lambda(\mathbf{w}^T \mathbf{w} - 1) \} \implies S \mathbf{w} = \lambda \mathbf{w}.$$

- Add WeChat powcoder
- Order eigenvalues  $(\lambda_1, \dots, \lambda_p)$ .  $\lambda_1$  first principal component.
- Linear combinations  $z_1^n = w_{1,1} x_1^n + w_{1,2} x_2^n + \dots + w_{1,p} x_p^n$  constitute representation of data in terms of first principal component. Instead of  $p$  components  $(x_1^n, \dots, x_p^n)$ , **one** component  $z_1^n$  represents  $x^n$ .
- Variance of  $\{z_1^n\}$  is  $\lambda_1$ .

# Assignment Project Exam Help

- Other eigenvectors  $w_k$  (from  $S w_k = \lambda_k w_k$ ) represent combinations that create  $z_k^n = w_{k,1}x_1^n + w_{k,2}x_2^n + \dots + w_{k,p}x_p^n$  decorrelated from  $z_1$ .  
 $\text{cov}(z_i^n, z_j^n) = 0$  for  $i \neq j$ .

<https://powcoder.com>

Add WeChat powcoder

# Assignment Project Exam Help

- Other eigenvectors  $w_k$  (from  $S w_k = \lambda_k w_k$ ) represent combinations that create  $z_k^n = w_{k,1}x_1^n + w_{k,2}x_2^n + \dots + w_{k,p}x_p^n$  decorrelated from  $z_1$ .  
 $\text{cov}(z_i^n, z_j^n) = 0$  for  $i \neq j$ .
- Largest  $n$  eigenvalues give new ‘co-ordinates,’  $(z_1^n, \dots, z_r^n)$  of the data and account for fraction of variation

Add WeChat  $\frac{\left(\sum\limits_{k=1}^r \lambda_k\right)}{\left(\sum\limits_{k=1}^p \lambda_k\right)}$  powcoder

# Assignment Project Exam Help

- Other eigenvectors  $w_k$  (from  $S w_k = \lambda_k w_k$ ) represent combinations that create  $z_k^n = w_{k,1}x_1^n + w_{k,2}x_2^n + \dots + w_{k,p}x_p^n$  decorrelated from  $z_1$ .  
 $\text{cov}(z_i^n, z_j^n) = 0$  for  $i \neq j$ .
- Largest  $n$  eigenvalues give new ‘co-ordinates,’  $(z_1^n, \dots, z_r^n)$  of the data and account for fraction of variation

Add WeChat  $\frac{\left(\sum\limits_{k=1}^r \lambda_k\right)}{\left(\sum\limits_{k=1}^p \lambda_k\right)}$  powcoder

- $w_1, \dots, w_r$  are the **principal components** (of variation).

# Assignment Project Exam Help

To make PCA independent of the scale of the data  $x^{n_i}$ , they might have to be first standardised such that all data entries  $x_k^{n_i}$  have equal range, eg by transforming them to unit variance.

<https://powcoder.com>

Add WeChat powcoder

## Eigenfaces as features

We will express arbitrary data vectors as linear combinations  $v^n = \sum_{i=1}^r \alpha_i^n w_i$ .  
of a set of eigenvectors  $\{w_i\}$ .

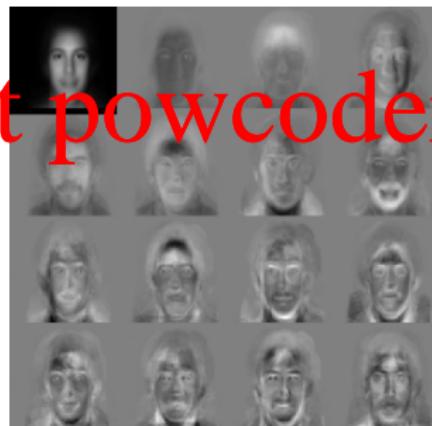
# Assignment Project Exam Help

These are obtained from a certain data dependent matrix. Often most of the  $w_i$  are small (except for a small set of 'important' eigenvectors). In such a case we can represent a high dimensional data vector  $v^n$  effectively by a much smaller (and more robust) set of new coordinates  $\{\alpha_i^n\}$ .

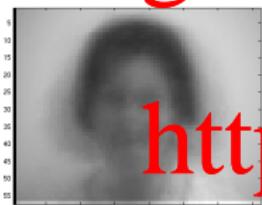
<https://powcoder.com>

## Add WeChat powcoder

For a face recognition problem  
the vector  $v^n$  and the dominant  
**eigenfaces** may look like this:



# Assignment Project Exam Help



<https://powcoder.com>

Mean face      eigenface 1      eigenface 2      eigenface 30  
Higher eigenfaces show only some random structure.

Add WeChat powcoder

- In regression  $y = Xw$

# Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

- In regression  $y = Xw$
- Take output  $y$  is matrix of  $x^n$ ,  $n = 1, \dots, N$ . Instead of vector  $\hat{w}$  we have matrix  $\hat{W}$  (one column for each data point):

$$\widehat{W} = (X^T X)^{-1} X^T x = \mathbb{I}$$

and the “predicted” output is  $X\widehat{W} = X$ .

Add WeChat powcoder

- In regression  $y = Xw$
- Take output  $y$  is matrix of  $x^n, n = 1, \dots, N$ . Instead of vector  $\hat{w}$  we have matrix  $\hat{W}$  (one column for each data point):

$$\hat{W} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{x} = \mathbb{I}$$

and the “predicted” output is  $X\hat{W} = X$ .

- Choose a set of  $q$  orthonormal vectors  $\{v_i\}$  whose linear combinations  $z^n = \sum_{i=1}^q \alpha_i^n v_i$  minimise  $\|x^n - z^n\|^2$ .

Add WeChat powcoder

- In regression  $y = Xw$
- Take output  $y$  is matrix of  $x^n, n = 1, \dots, N$ . Instead of vector  $\hat{w}$  we have matrix  $\hat{W}$  (one column for each data point):

$$\hat{W} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{x} = \mathbb{I}$$

and the “predicted” output is  $X\hat{W} = X$ .

- Choose a set of  $q$  orthonormal vectors  $\{v_i\}$  whose linear combinations  $z^n = \sum_{i=1}^q \alpha_i^n v_i$  minimise  $\|x^n - z^n\|^2$ .
- If  $q < p$ , dimensional reduction. Instead of basis vectors  $e_k = (0, \dots, 0, \underbrace{1}_k, 0, \dots, 0)^T, k = 1, \dots, p$ , use  $v_i, i = 1, \dots, q$ . The co-ordinates are  $\alpha_i^n$

## Exercise

- Minimising with respect to  $\begin{pmatrix} v_1 \\ v_2 \end{pmatrix}_{(1)}$  and  $\begin{pmatrix} v_1 \\ v_2 \end{pmatrix}_{(2)}$ ,

$$\left\| \begin{pmatrix} x_1^{(n)} \\ x_2^{(n)} \end{pmatrix} - \begin{pmatrix} v_{1;1} & v_{2;1} \\ v_{1;2} & v_{2;2} \end{pmatrix} \begin{pmatrix} \alpha_1^{(n)} \\ \alpha_2^{(n)} \end{pmatrix} \right\|^2$$

(in  $v_j|_i$ ,  $j$  the label for the vector,  $i$  the row index),

<https://powcoder.com>

$$\begin{pmatrix} \alpha_1^{(n)} \\ \alpha_2^{(n)} \end{pmatrix} = \begin{pmatrix} v_{1;1} & v_{2;1} \\ v_{1;2} & v_{2;2} \end{pmatrix}^{-1} \begin{pmatrix} x_1^{(n)} \\ x_2^{(n)} \end{pmatrix}.$$

Add WeChat powcoder

## Exercise

- Minimising with respect to  $\begin{pmatrix} v_1 \\ v_2 \end{pmatrix}_{(1)}$  and  $\begin{pmatrix} v_1 \\ v_2 \end{pmatrix}_{(2)}$ ,

$$\left\| \begin{pmatrix} x_1^{(n)} \\ x_2^{(n)} \end{pmatrix} - \begin{pmatrix} v_{1;1} & v_{2;1} \\ v_{1;2} & v_{2;2} \end{pmatrix} \begin{pmatrix} \alpha_1^{(n)} \\ \alpha_2^{(n)} \end{pmatrix} \right\|^2$$

(in  $v_j|_i$ ,  $j$  the label for the vector,  $i$  the row index),

<https://powcoder.com>

$$\begin{pmatrix} \alpha_1^{(n)} \\ \alpha_2^{(n)} \end{pmatrix} = \begin{pmatrix} v_{1;1} & v_{2;1} \\ v_{1;2} & v_{2;2} \end{pmatrix}^{-1} \begin{pmatrix} x_1^{(n)} \\ x_2^{(n)} \end{pmatrix}.$$

- Extend to  $q$ -dimensional representation of  $x^{(n)} \in \mathbb{R}^p$ .

$$\begin{pmatrix} \alpha_1^{(n)} \\ \vdots \\ \alpha_q^{(n)} \end{pmatrix} = \underbrace{\begin{pmatrix} v_{1;1} & v_{2;1} & \cdots & v_{q;1} \\ \vdots & \vdots & \ddots & \vdots \\ v_{1;p} & v_{2;p} & \cdots & v_{q;p} \end{pmatrix}}_{\text{columns of matrix } V_q \text{ are orthogonal vectors}}^T \begin{pmatrix} x_1^{(n)} \\ \vdots \\ x_p^{(n)} \end{pmatrix}.$$

columns of matrix  $V_q$  are orthogonal vectors

## PCA: Low rank projection

- q-dimensional representation of  $x^n \in \mathbb{R}^p$ , with orthonormal vectors as columns of  $V_q^T$ :

$$\begin{pmatrix} \alpha_1^{(n)} \\ \vdots \\ \alpha_q^{(n)} \end{pmatrix} = V_q^T \begin{pmatrix} x_1^{(n)} \\ \vdots \\ x_p^{(n)} \end{pmatrix}.$$

<https://powcoder.com>

Add WeChat powcoder

## PCA: Low rank projection

- q-dimensional representation of  $x^n \in \mathbb{R}^p$ , with orthonormal vectors as columns of  $V_q$ :

$$\begin{pmatrix} \alpha_1^{(n)} \\ \vdots \\ \alpha_q^{(n)} \end{pmatrix} = V_q^T \begin{pmatrix} x_1^{(n)} \\ \vdots \\ x_p^{(n)} \end{pmatrix}.$$

<https://powcoder.com>

- Since  $(V_q V_q^T)(V_q V_q^T) = (V_q V_q^T)$ ,  $P_q \triangleq V_q V_q^T$  is a projection onto q-dim subspace spanned by columns of  $V_q$ .

Add WeChat powcoder

## PCA: Low rank projection

- q-dimensional representation of  $x^n \in \mathbb{R}^p$ , with orthonormal vectors as columns of  $V_q$ :

$$\begin{pmatrix} \alpha_1^{(n)} \\ \vdots \\ \alpha_q^{(n)} \end{pmatrix} = V_q^T \begin{pmatrix} x_1^{(n)} \\ \vdots \\ x_p^{(n)} \end{pmatrix}.$$

<https://powcoder.com>

- Since  $(V_q V_q^T)(V_q V_q^T) = (V_q V_q^T)$ ,  $P_q \triangleq V_q V_q^T$  is a projection onto q-dim subspace spanned by columns of  $V_q$ .
- $Y = (x_i^{(n)})_{i=1}^N \in V_q^T \mathbb{R}^p$  is the matrix of N q-1Dm vectors  $\mathbf{z}^i \in \mathbb{R}^q$

## PCA: Low rank projection

- q-dimensional representation of  $x^n \in \mathbb{R}^p$ , with orthonormal vectors as columns of  $V_q$ :

$$\begin{pmatrix} \alpha_1^{(n)} \\ \vdots \\ \alpha_q^{(n)} \end{pmatrix} = V_q^T \begin{pmatrix} x_1^{(n)} \\ \vdots \\ x_p^{(n)} \end{pmatrix}.$$

<https://powcoder.com>

- Since  $(V_q V_q^T)(V_q V_q^T) = (V_q V_q^T)$ ,  $P_q \triangleq V_q V_q^T$  is a projection onto q-dim subspace spanned by columns of  $V_q$ .
- $Y = \sum_{q=1}^N V_q z_q^T$  is the matrix of N q-1m vectors  $z_q^T \in \mathbb{R}^m$
- Among all rank q matrices  $A$ , quantity  $\|X - A\|^2$  is minimised is  $A = V_q V_q^T X$ , q leading vectors (in order of singular values  $\sigma_i$  in  $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_{\min(p,N)})$  in  $V$  in  $X = U \Sigma V^T$ )

## PCA: Low rank projection

- q-dimensional representation of  $x^n \in \mathbb{R}^p$ , with orthonormal vectors as columns of  $V_q$ :

$$\begin{pmatrix} \alpha_1^{(n)} \\ \vdots \\ \alpha_q^{(n)} \end{pmatrix} = V_q^T \begin{pmatrix} x_1^{(n)} \\ \vdots \\ x_p^{(n)} \end{pmatrix}.$$

<https://powcoder.com>

- Since  $(V_q V_q^T)(V_q V_q^T) = (V_q V_q^T)$ ,  $P_q \triangleq V_q V_q^T$  is a projection onto q-dim subspace spanned by columns of  $V_q$ .
- $Y = \sum_{i=1}^q \sigma_i v_i v_i^T$  is the matrix of N (1-1)m vectors  $v_i \in \mathbb{R}^m$
- Among all rank q matrices  $A$ , quantity  $\|X - A\|^2$  is minimised is  $A = V_q V_q^T X$ , q leading vectors (in order of singular values  $\sigma_i$  in  $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_{\min(p, N)})$  in  $V$  in  $X = U \Sigma V^T$ )
- $\|X - V_q V_q^T X\|^2$  is the minimum residual.

## Compare PCA pre-processed regression with regularisation

- Regularisation using the 2-norm of  $w$  involves minimising the loss function

$$\begin{aligned}\hat{w}_{\text{ridge}} &= \underset{w}{\operatorname{argmin}} \|y - Xw\|^2 + \lambda \|w\|^2 \\ \implies \hat{w}_{\text{ridge}} &= (X^T X + \lambda \mathbb{I})^{-1} X^T y \\ &= V(S^2 + \lambda \mathbb{I})^{-1} S U^T y \\ &= \sum_{k=1}^{\min(p, N)} (u_k^T y) \frac{\sigma_k}{\sigma_k^2 + \lambda} v_k\end{aligned}$$

Add WeChat powcoder

## Compare PCA pre-processed regression with regularisation

- Regularisation using the 2-norm of  $w$  involves minimising the loss function

$$\begin{aligned}\hat{w}_{\text{ridge}} &= \underset{w}{\operatorname{argmin}} \|y - Xw\|^2 + \lambda \|w\|^2 \\ \implies \hat{w}_{\text{ridge}} &= (X^T X + \lambda \mathbb{I})^{-1} X^T y \\ &= V(S^2 + \lambda \mathbb{I})^{-1} S U^T y \\ &= \sum_{k=1}^{\min(p, N)} (u_k^T y) \frac{\sigma_k}{\sigma_k^2 + \lambda} v_k\end{aligned}$$

- PCA drops all terms with small variance contributions, ridge regression performs a “soft” re-weighting.

Add WeChat powcoder

## Compare PCA pre-processed regression with regularisation

- Regularisation using the 2-norm of  $w$  involves minimising the loss function

$$\begin{aligned}\widehat{w}_{\text{ridge}} &= \operatorname{argmin}_w \|y - Xw\|^2 + \lambda\|w\|^2 \\ \implies \widehat{w}_{\text{ridge}} &= (X^T X + \lambda \mathbb{I})^{-1} X^T y \\ &= V(S^2 + \lambda \mathbb{I})^{-1} S U^T y \\ &= \sum_{k=1}^{\min(p, N)} (u_k^T y) \frac{\sigma_k}{\sigma_k^2 + \lambda} v_k\end{aligned}$$

- PCA drops all terms with small variance contributions, ridge regression performs a “soft” re-weighting.

- **Lasso:** regularisation method that drops components

$$\widehat{w}_{\text{lasso}} = \operatorname{argmin}_w \|y - Xw\|^2 \text{ s.t. } \lambda \sum_i |w_i| < t.$$

## Compare PCA pre-processed regression with regularisation

- Regularisation using the 2-norm of  $w$  involves minimising the loss function

$$\begin{aligned}\widehat{w}_{\text{ridge}} &= \operatorname{argmin}_w \|y - Xw\|^2 + \lambda\|w\|^2 \\ \implies \widehat{w}_{\text{ridge}} &= (X^T X + \lambda \mathbb{I})^{-1} X^T y \\ &= V(S^2 + \lambda \mathbb{I})^{-1} S U^T y \\ &= \sum_{k=1}^{\min(p, N)} (u_k^T y) \frac{\sigma_k}{\sigma_k^2 + \lambda} v_k\end{aligned}$$

- PCA drops all terms with small variance contributions, ridge regression performs a “soft” re-weighting.

- **Lasso:** regularisation method that drops components

$$\widehat{w}_{\text{lasso}} = \operatorname{argmin}_w \|y - Xw\|^2 \text{ s.t. } \lambda \sum_i |w_i| < t.$$

## Compare PCA pre-processed regression with regularisation

- Regularisation using the 2-norm of  $w$  involves minimising the loss function

$$\begin{aligned}\widehat{w}_{\text{ridge}} &= \operatorname{argmin}_w \|y - Xw\|^2 + \lambda\|w\|^2 \\ \implies \widehat{w}_{\text{ridge}} &= (X^T X + \lambda I)^{-1} X^T y \\ &= V(S^2 + \lambda I)^{-1} S U^T y \\ &= \sum_{k=1}^{\min(p, N)} (u_k^T y) \frac{\sigma_k}{\sigma_k^2 + \lambda} v_k\end{aligned}$$

- PCA drops all terms with small variance contributions, ridge regression performs a “soft” re-weighting.

- **Lasso:** regularisation method that drops components

$$\widehat{w}_{\text{lasso}} = \operatorname{argmin}_w \|y - Xw\|^2 \text{ s.t. } \lambda \sum_i |w_i| < t. \\ (\|w\|_1 \triangleq \sum_i |w_i|, \text{ l-norm.})$$

Refine representation: rotation in q-dim PC space for  
interpretable/sparse combinations; introducing non-linearity

- Using ideas from

# Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Refine representation: rotation in q-dim PC space for interpretable/sparse combinations; introducing non-linearity

- Using ideas from

$$\text{classic Lasso: } \hat{\boldsymbol{\omega}}_{\text{lasso}} = \arg \min_{\boldsymbol{\omega}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\omega}\|^2 \text{ s.t. } \sum_i |\omega_i| \leq t.$$

# Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Refine representation: rotation in q-dim PC space for interpretable/sparse combinations; introducing non-linearity

- Using ideas from

- **Lasso:**  $\hat{w}_{\text{lasso}} = \operatorname{argmin}_w \|y - Xw\|^2 \text{ s.t. } \lambda \sum_i |w_i| \leq t.$

- **Elastic net:**  $\operatorname{argmin}_w \|y - Xw\|^2 + \lambda_2 \|w\|^2 + \lambda_1 \|w\|_1.$

# Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Refine representation: rotation in q-dim PC space for interpretable/sparse combinations; introducing non-linearity

- Using ideas from

- **lasso:**  $\hat{w}_{\text{lasso}} = \operatorname{argmin}_w \|y - Xw\|^2 \text{ s.t. } \lambda \sum_i |w_i| \leq t.$

- **Elastic net:**  $\operatorname{argmin}_w \|y - Xw\|^2 + \lambda_2 \|w\|^2 + \lambda_1 \|w\|_1.$

- **Sparse PCA:**  $A_{p \times q} := (\alpha_1, \dots, \alpha_q)$ ,  $B_{p \times q} := (\beta_1, \dots, \beta_q)$ ,

$A^T A = I_q$ , so that  $A^T A^T = P_q$ , a q-dim projector:

$$\beta = \operatorname{argmin}_{A, B} \sum_{i=1}^N \|x_i - AB^T x_i\|^2 + \lambda_2 \|\beta\|^2 + \lambda_1 \|\beta\|_1$$

$$v_i = \hat{\beta}_i / \|\hat{\beta}\|$$

$v_i$  sparse PCA components

Add WeChat powcoder

Refine representation: rotation in q-dim PC space for interpretable/sparse combinations; introducing non-linearity

# Assignment Project Exam Help

- Using ideas from
  - **Lasso:**  $\hat{w}_{\text{lasso}} = \operatorname{argmin}_w \|y - Xw\|^2 \text{ s.t. } \lambda \sum_i |w_i| \leq t.$
  - **Elastic net:**  $\operatorname{argmin}_w \|y - Xw\|^2 + \lambda_2 \|w\|^2 + \lambda_1 \|w\|_1.$

- **Sparse PCA:**  $A_{p \times q} := (\alpha_1, \dots, \alpha_q)$ ,  $B_{p \times q} := (\beta_1, \dots, \beta_q)$ ,  
 $A^T A = I_q$ , so that  $A^T A^T = P_q$ , a q-dim projector:

$$\begin{aligned}\beta &= \operatorname{argmin}_{A, B} \sum_{i=1}^N \|x_i - AB^T x_i\|^2 + \lambda_2 \|\beta\|^2 + \lambda_1 \|\beta\|_1 \\ v_i &= \hat{\beta}_i / \|\hat{\beta}\|\end{aligned}$$

- $v_i$  sparse PCA components
- Implement PCA via neural network with linear activation function. If non-linear, we arrive at an **autoencoder**, trained via backprop.

