

Assignment Project Exam Help

COMP3223 Foundations of ML: Dimensionality reduction

by

<https://eduassistpro.github.io>

Srinandan Dasmah

Add WeChat edu_assist_pro

Are more complex models better?

Assignment Project Exam Help

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

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

Are more complex models better?

Assignment Project Exam Help

- With large number of features, p , accidental correlations can obscure genuine patterns.
- In min the e

$$\sum_{i=1}^p w_i x_i^n$$

<https://eduassistpro.github.io>

$$(r^n)^2 = (y^n - f(x^n; \boldsymbol{w}))^2$$

Add WeChat edu_assist_pro

Assignment Project Exam Help

- With large number of features, p , accidental correlations can obscure genuine patterns.
- In minimising the error function, $\sum_{i=1}^p w_i x_i^n$

<https://eduassistpro.github.io/>

$$(r^n)^2 = (y^n - f(x^n; \mathbf{w}))^2$$

Add WeChat `edu_assist_pro`

- In trying to minimise the residuals, you may be adjusting weights 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 min the e

$$\sum_{i=1}^p w_i x_i^n$$

<https://eduassistpro.github.io>

$$(r^n)^2 = (y^n - f(x^n; \mathbf{w}))^2$$

Add WeChat `edu_assist_pro`

- In trying to minimise the residuals, you may be adjusting 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

N

$w_p x_p^n$)

<https://eduassistpro.github.io>

i i

i=1

Add WeChat edu_assist_pro

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

N

$w_p x_p^n$)

<https://eduassistpro.github.io>

i i

i=1

- Centering:** $\tilde{y}_i = y_i - \bar{y}$, $\tilde{x}_{ij} = x_{ij} - \langle x_j \rangle$ from design matrix to form \tilde{X} with matrix

Add WeChat edu_assist_pro

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

N

$w_p x_p^n)$

<https://eduassistpro.github.io>

i i

i=1

- Centering:** $\tilde{\mathbf{y}}^i = \mathbf{y}^i - \langle \mathbf{y} \rangle$, $\tilde{x}_{ij}^i = x_{ij}^i - \langle x_j \rangle$ from design matrix to form $\tilde{\mathbf{X}}$ with matrix

- $\frac{1}{N} \tilde{\mathbf{X}}^\top (\tilde{\mathbf{y}} - \tilde{\mathbf{X}} \hat{\mathbf{w}}) = 0$ implies $\hat{\mathbf{w}} = (\frac{1}{N} (\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}}))^{-1} (\frac{1}{N} \tilde{\mathbf{X}}^\top \tilde{\mathbf{y}})$, 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://eduassistpro.github.io>

Add WeChat edu_assist_pro

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 .

($x_j^n - \langle x_j \rangle$)

$$(\tilde{X}^\top \tilde{X})_{ij} = (\tilde{X}^\top)_{in} (\tilde{X})_{jn} = \langle x_i^n - \langle x_i \rangle, x_j^n - \langle x_j \rangle \rangle$$

Add WeChat edu_assist_pro

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 \tilde{X} .

($x_j^n - \langle x_j \rangle$)

$$(\tilde{X}^T \tilde{X})_{ij} = (\tilde{X}^T)_{in} (\tilde{X})_{jn} \quad n \quad n \quad \langle x_j \rangle)$$

- Add WeChat edu_assist_pr

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

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 \tilde{X} .

($x_j^n - \langle x_j \rangle$)

$$(\tilde{X}^\top \tilde{X})_{ij} = (\tilde{X}^\top)_{in} (\tilde{X})_{jn} \quad n \quad n \quad \langle x_j \rangle)$$

- Add WeChat edu_assist_pr

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

- 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://eduassistpro.github.io>

Add WeChat edu_assist_pro

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 $\tilde{X} = USV^T$, $y = USV^Tw$ and $\hat{w} = V(S)^{-1}U^Ty$.

$S \in \text{diag}(\sigma_1, \dots, \sigma_{\min(p, N)})$. U, V contain singular vectors of size N and p respectively.

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

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 $\tilde{X} = USV^T$, $y = USV^Tw$ and $\hat{w} = V(S)^{-1}U^Ty$.

$S \in \text{diag}(\sigma_1, \dots, \sigma_{\min(p, N)})$. U, V contain singular vectors of size N and p respectively.

<https://eduassistpro.github.io>

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

Add WeChat `edu_assist_pro`

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 $\tilde{X} = USV^T$, $y = USV^Tw$ and $\hat{w} = V(S)^{-1}U^Ty$.

$S \in \text{diag}(\sigma_1, \dots, \sigma_{\min(p, N)})$. U, V contain singular vectors of size N and p respectively.

<https://eduassistpro.github.io>

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

- Add WeChat** `edu_assist_pro`
Discard small values of σ_k , keep largest variation in data explained by r components

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

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 $\tilde{X} = USV^T$, $y = USV^Tw$ and $\hat{w} = V(S)^{-1}U^Ty$.

$S \in \text{diag}(\sigma_1, \dots, \sigma_{\min(p, N)})$. U, V contain singular vectors of size N and p respectively.

<https://eduassistpro.github.io>

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

- Add WeChat** `edu_assist_pro`
Discard small values of σ_k , keep largest variation in data explained by r components

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

- 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://eduassistpro.github.io>

- Single image as
matrix M

Add WeChat edu_assist_pro

Reconstruction for r

SVD is a low rank approximation to any matrix

Assignment Project Exam Help

<https://eduassistpro.github.io>

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

Add WeChat edu_assist_pro
Reconstruction for r

SVD is a low rank approximation to any matrix

Assignment Project Exam Help

<https://eduassistpro.github.io>

- Single image as

matrix M

Add WeChat edu_assist_pro

- $M_{1200 \times 1200}$

Reconstruction for r

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

SVD is a low rank approximation to any matrix

Assignment Project Exam Help

<https://eduassistpro.github.io>

- Single image as

matrix M

Add WeChat edu_assist_pro

- $M_{1200 \times 1200}$

Reconstruction for r

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

- Reconstruction:

$$\widetilde{M} = \sum_i^r \sigma_i u_i v_i^T$$

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://eduassistpro.github.io>

Add WeChat edu_assist_pro

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 $\mathbf{z}^T \mathbf{w}$.

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

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 $\mathbf{X}\mathbf{w}$.
- $\text{var}(X\mathbf{w}) = \mathbf{w}^T X^T X \mathbf{w} = \mathbf{w}^T S \mathbf{w}$ where $S = \text{var}(X)$.

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

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 $\mathbf{z}^T \mathbf{w}$.
- $\text{var}(X\mathbf{w}) = \mathbf{w}^T X^T X \mathbf{w} = \mathbf{w}^T S \mathbf{w}$ where $S = \text{var}(X)$.
- Seek

w^T <https://eduassistpro.github.io>

$$\underset{\mathbf{w}}{\text{argmax}} \quad \mathbf{w}^T S \mathbf{w} - \lambda(\mathbf{w}^T \mathbf{w})$$

Add WeChat edu_assist_pro

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 $\mathbf{z}^T \mathbf{w}$.
- $\text{var}(X\mathbf{w}) = \mathbf{w}^T X^T X \mathbf{w} = \mathbf{w}^T S \mathbf{w}$ where $S = \text{var}(X)$.
- Seek

w^T <https://eduassistpro.github.io>

$$\underset{\mathbf{w}}{\text{argmax}} \quad \mathbf{w}^T S \mathbf{w} - \lambda(\mathbf{w}^T \mathbf{w})$$

- Order eigenvalues $(\lambda_1, \dots, \lambda_p)$ λ_1 first

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 $z_1^n = \sum_{j=1}^p w_{ij} x_j^n$ written as $\mathbf{z}^n w$.
- $\text{var}(Xw) = w^T X^T X w = w^T S w$ where $S = \text{var}(X)$.
- Seek

w^T <https://eduassistpro.github.io>

$$\underset{w}{\operatorname{argmax}} \quad w^T S w - \lambda(w^T w)$$

- Order eigenvalues $(\lambda_1, \dots, \lambda_p)$ in first
- 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 $z_1^n = \sum_{j=1}^p w_{ij} x_j^n$ written as $\mathbf{z}^n w$.
- $\text{var}(Xw) = w^T X^T X w = w^T S w$ where $S = \text{var}(X)$.
- Seek

https://eduassistpro.github.io

$$\underset{w}{\operatorname{argmax}} \quad w^T S w - \lambda(w^T w)$$

- Order eigenvalues $(\lambda_1, \dots, \lambda_p)$ in first
- 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 λ_1 .

Assignment Project Exam Help

- Other eigenvectors w_k (from $S w_k = \lambda_k w_k$) represent combinations that create $z^n = w_{k,1}x^n + w_{k,2}x^n + \dots + w_{k,p}x^n$ decorrelated from z_1 .

cov(/

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

Assignment Project Exam Help

- Other eigenvectors w_k (from $S w_k = \lambda_k w_k$) represent combinations that create $z^n = w_{k,1}x^n + w_{k,2}x^n + \dots + w_{k,p}x^n$ decorrelated from z_1 .
 $\text{cov}($ /
• Large eigenvalues λ_k account for most of the variance in the data and

Add WeChat $\frac{(\sum_{k=1}^r \lambda_k)}{\lambda_k}$ edu_assist_pro

Assignment Project Exam Help

- Other eigenvectors w_k (from $Sw_k = \lambda_k w_k$) represent combinations that create $z^n = w_{k,1}x^n + w_{k,2}x^n + \dots + w_{k,p}x^n$ decorrelated from z_1 .
 $\text{cov}(\quad / \quad)$
 - Large eigenvalues indicate data and account for most variance.

Add WeChat_edu_assist_pr

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

Assignment Project Exam Help

To make P
standardi
them to unit variance.

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

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 (e

represent
more robust

<https://eduassistpro.github.io> and

Add WeChat edu_assist_pr

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

Assignment Project Exam Help

<https://eduassistpro.github.io>

Mean face

Add WeChat edu_assist_pro

Higher eigenfaces show only some random structure.

- In regression $y = Xw$

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

- 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):

<https://eduassistpro.github.io>

and t

$$X\hat{W} = X$$

Add WeChat edu_assist_pro

- 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):

<https://eduassistpro.github.io>

and t

$$X\hat{W} = X$$

- Choose a set of q orthonormal vectors {

$$z^n = \sum_{i=1}^q \alpha_i^n v_i \text{ minimise } \|x^n - z^n\|^2$$

Add WeChat edu_assist_pro

- 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):

<https://eduassistpro.github.io>

and t

$$X\hat{W} = X$$

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

Exercise

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

Assignment Project Exam Help

$$\left(x^{(n)} \right) = \frac{v_{1;1}}{v_{1;2}} \frac{v_{2;1}}{v_{2;2}} \left(\alpha^{(n)} \right)$$

(in v)

<https://eduassistpro.github.io>

$$\alpha_2^{(n)} = \frac{v_{1;2}}{v}$$

Add WeChat edu_assist_pro

Exercise

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

Assignment Project Exam Help

$$\left(x^{(n)} \right) = \frac{v_{1;1}}{v_{1;2}} \frac{v_{2;1}}{v_{2;2}} \left(\alpha^{(n)} \right)$$

(in v)

<https://eduassistpro.github.io>

$$\alpha_2^{(n)} = \frac{v_{1;2}}{v}$$

- Extend to q-dimensional representation of x

Add WeChat edu_assist_pro

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

PCA: Low rank projection

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

$$\alpha_1^{(n)}$$

$$x_1^{(n)}$$

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

PCA: Low rank projection

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

$$\alpha_1^{(n)}$$

$$x_1^{(n)}$$

Assignment Project Exam Help

<https://eduassistpro.github.io>

- Since $(V_q V_q^\top)(V_q V_q^\top) = (V_q V_q^\top) P_q V_q V_q^\top$ n onto
q-dim subspace spanned by columns of V

Add WeChat edu_assist_pro

PCA: Low rank projection

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

$$\alpha_1^{(n)}$$

$$x_1^{(n)}$$

<https://eduassistpro.github.io>

- Since $(V_q V_q^\top)(V_q V_q^\top) = (V_q V_q^\top) P_q V_q V_q^\top$ n onto
q-dim subspace spanned by columns of V
- $Y = (X)_{(1)} \in V_q^\top X$ is the matrix of N

Add WeChat edu_assist_pro

PCA: Low rank projection

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

$$\alpha_1^{(n)}$$

$$x_1^{(n)}$$

Assignment Project Exam Help

<https://eduassistpro.github.io>

- Since $(V_q V_q^\top)(V_q V_q^\top) = (V_q V_q^\top) P_q V_q V_q^\top$ n onto q-dim subspace spanned by columns of V
- $Y = V_q^\top X$ is the matrix of N
- Among all rank q matrices A , quantity $\| A = V_q V_q^\top X \text{, q leading vectors (in order of singular values } \sigma_i \text{ in } \Sigma = \text{diag}(\sigma_1, \dots, \sigma_{\min(p, N)} \text{ in } V \text{ in } X = U \Sigma V^\top)$

PCA: Low rank projection

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

$$\alpha_1^{(n)}$$

$$x_1^{(n)}$$

Assignment Project Exam Help

<https://eduassistpro.github.io>

- Since $(V_q V_q^\top)(V_q V_q^\top) = (V_q V_q^\top) P_q V_q V_q^\top$ n onto q-dim subspace spanned by columns of V
- $Y = V_q^\top X$ is the matrix of N
- Among all rank q matrices A , quantity $\|A - V_q V_q^\top X\|$, q leading vectors (in order of singular values σ_i in $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_{\min(p, N)})$ in V in $X = U \Sigma V^\top$)
- $\|X - V_q V_q^\top X\|^2$ is the minimum residual.

Compare PCA pre-processed regression with regularisation

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

($\|w\|$)

$$w_{\text{ridge}} = \operatorname{argmin}_w \|y - Xw\|^2 + \lambda \|w\|^2$$

<https://eduassistpro.github.io>

$$= u_k y \frac{\|v_k\|}{2}$$

$$\sum_{k=1}^K$$

Add WeChat edu_assist_pro

Compare PCA pre-processed regression with regularisation

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

(μ)

$$w_{\text{ridge}} = \operatorname{argmin}_w \|y - Xw\|^2 + \lambda \|w\|^2$$

<https://eduassistpro.github.io>

$$= u_k y \frac{\|v_k\|_2}{\|v_k\|_2} v_k$$

$k=1$

- PCA drops all terms with small variance contribution.
regularisation performs a “soft” re-weighting.

Add WeChat `edu_assist_pro`

Compare PCA pre-processed regression with regularisation

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

(μ)

$$w_{\text{ridge}} = \operatorname{argmin}_w \|y - Xw\|^2 + \lambda \|w\|^2$$

<https://eduassistpro.github.io>

$$= u_k y \frac{\|v_k\|_2}{2}$$

$k=1$

- PCA drops all terms with small variance contribution. Lasso performs a “soft” re-weighting.
- **Lasso:** regularisation method that drops components

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

Compare PCA pre-processed regression with regularisation

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

($\|w\|_2$)

$$w_{\text{ridge}} = \operatorname{argmin}_w \|y - Xw\|^2 + \lambda \|w\|^2$$

<https://eduassistpro.github.io>

$$= u_k y \frac{\|v_k\|_2}{\|v_k\|_2} v_k$$

$k=1$

- PCA drops all terms with small variance contribution. Lasso performs a “soft” re-weighting.
- **Lasso:** regularisation method that drops components

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

Compare PCA pre-processed regression with regularisation

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

(μ)

$$w_{\text{ridge}} = \operatorname{argmin}_w \|y - Xw\|^2 + \lambda \|w\|^2$$

<https://eduassistpro.github.io>

$$= u_k y \frac{\|v_k\|_2}{2} v_k$$

$k=1$

- PCA drops all terms with small variance contribution. It performs a “soft” re-weighting.

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

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

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

- Using ideas from

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

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://eduassistpro.github.io>

Add WeChat edu_assist_pro

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://eduassistpro.github.io>

Add WeChat edu_assist_pro

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

- Using ideas from

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

- **Elastic net:** $\arg \min_w \|y - Xw\|^2 + \lambda_2 \|w\|^2 + \lambda_1 \|w\|_1.$

- **Spa**

$$A^T$$

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

ν_i sparse PCA components

<https://eduassistpro.github.io>

),

$$\lambda_1 \|\beta\|_1$$

Add WeChat edu_assist_pro

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

Assignment Project Exam Help

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

- **Sparsity:** $A^T \hat{\beta} = \lambda_1 \|\beta\|_1$,

$$A^T$$

<https://eduassistpro.github.io>

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

ν_i sparse PCA components

- Implement PCA via neural network with linear activation function, we arrive at an **autoencoder**, trained via backprop.

