

# Assignment Project Exam Help

## Linear Regression

<https://powcoder.com>

MAST90083 Computational Statistics and Data Mining  
Dr Karim Seghouane  
School of Mathematics & Statistics  
The University of Melbourne

# Add WeChat powcoder



## Statistical Models

# Assignment Project Exam Help

- ▶ What is the simplest mathematical model that describes the relationship between two variables ?
- ▶ **Straight line**  
<https://powcoder.com>
- ▶ Statistical models are fitted for a variety of reasons:
- ▶ Explanation and prediction: Uncover causes by studying the relationship between an interested variable (the response) and a set of variables called the explanatory variables & use the model for prediction
- ▶ Add WeChat powcoder
- ▶ Examine and test scientific hypotheses

## Linear Models

# Assignment Project Exam Help

- ▶ Linear models have a long history in statistics, but even in today's computer era they are still important and widely used in supervised learning.
- ▶ They are simple and provide a picture of how the inputs affect the output
- ▶ For prediction purposes they can sometimes outperform fancier nonlinear models, particularly in small sample cases, low signal-to-noise ratio or sparse data
- ▶ We will study some of the key questions associated with the linear regression model

<https://powcoder.com>

Add WeChat powcoder





# Assignment Project Exam Help

$$y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + \epsilon_i \quad i = 1, \dots, N$$

$y_i = \mathbf{x}_i^T \boldsymbol{\beta} + \epsilon_i \quad i = 1, \dots, N$

- Taking  $\mathbf{X}$  as the  $N \times (p-1)$  matrix with each row an input vector and 1 in the first position and  $\mathbf{y}$  an  $N \times 1$  vector of responses:  $E(\mathbf{y}) = \mathbf{X}\boldsymbol{\beta}$  and  $\text{Cov}(\mathbf{y}) = \sigma^2 \mathbf{I}_N$ , so that

$$\mathbf{y} \sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_N)$$

## Note 1

# Assignment Project Exam Help

- $$\text{RSS}(\beta) = (\mathbf{y} - X\beta)^\top (\mathbf{y} - X\beta)$$
- Assuming that  $X$  has full column rank or  $X^\top X$  is positive definite gives a unique solution

Assuming that  $X$  has full column rank or  $X^T X$  is positive definite, the least squares solution is

# Add WeChat powcoder

- ▶ and the fitted values at the training inputs are ( **Note 2** )

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X} \left( \mathbf{X}^\top \mathbf{X} \right)^{-1} \mathbf{X}^\top \mathbf{y} = \mathbf{H} \mathbf{y}$$

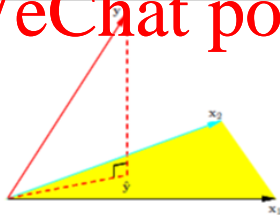


## Geometric Interpretation

The hat matrix  $H$  is square and satisfies:  $H^2 = H$  and  $H^T = H$  (Note 9)

- $H$  is the orthogonal projector onto  $V = Sp(X)$  (column space of  $X$  or the subspace of  $\mathcal{R}^N$  spanned by the column vectors of  $X$ )
- and  $\hat{\mathbf{y}}$  is the orthogonal projection of  $\mathbf{y}$  onto  $Sp(X)$
- The residual vector  $\mathbf{y} - \hat{\mathbf{y}}$  is orthogonal to this subspace

Add WeChat powcoder



## Statistical Properties

- Assuming the model  $\mathbf{y} \sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_N)$  gives

$$\hat{\boldsymbol{\beta}} \sim \mathcal{N}\left(\boldsymbol{\beta}, \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}\right)$$

or

$$\hat{\mathbf{y}} \sim \mathcal{N}(\mathbf{X}\hat{\boldsymbol{\beta}}, \sigma^2 \mathbf{H})$$

- where  $\sigma^2$  is estimated by

Add WeChat powcoder

$$\hat{\sigma}^2 = \frac{1}{N - p - 1} \sum_{i=1}^N (y_i - \hat{y}_i)^2 = \frac{1}{N - p - 1} (\mathbf{y} - \hat{\mathbf{y}})^\top (\mathbf{y} - \hat{\mathbf{y}})$$

- and  $(N - p - 1)\hat{\sigma}^2 \sim \sigma^2 \chi_{N-p-1}^2$
- Furthermore  $\hat{\sigma}^2$  and  $\hat{\boldsymbol{\beta}}$  are statistically independent. (Why ?)

## Assessing the Accuracy of the coefficient estimates

- Approximate confidence set for the parameter vector  $\beta$

$$C_{\beta} = \{\beta | (\hat{\beta} - \beta)^{\top} X^{\top} X (\hat{\beta} - \beta) \leq \hat{\sigma}^2 \chi_{p+1}^{(1-\alpha)}\}$$

- Test the null hypothesis of  $H_0 : \beta_j = 0$  vs.  $H_1 : \beta_j \neq 0$  using

$$z_j = \frac{\hat{\beta}_j}{\hat{\sigma} \sqrt{v_j}}$$

- and  $v_j$  is the  $j$ th diagonal element of  $(X^{\top} X)^{-1}$ . Under  $H_0$ ,  $z_j$  is  $t_{N-p-1}$ .
- Testing for a group of variables,  $H_0$  : smaller model is correct

$$F = \frac{(RSS_0 - RSS_1)/(p_1 - p_0)}{RSS_1/(N - p_1 - 1)} \sim F_{p_1 - p_0, N - p_1 - 1}$$

Note 4



# Assignment Project Exam Help

- $y_0 = f(x_0) + \epsilon_0$  at input  $x_0$

$$E \left[ y_0 - \tilde{f}(x_0) \right]^2 = \sigma^2 + E \left[ x_0^\top \tilde{\beta} - f(x_0) \right]^2 = \sigma^2 + MSE \left[ \tilde{f}(x_0) \right]$$

## Note 5

# Assignment Project Exam Help

- $R^2$  Statistics

# Add WeChat powcoder

- ▶ measure the amount of variability ( $TSS = \sum (y_i - \bar{y})^2$ ) removed by the model

# Assignment Project Exam Help

- ▶ Correlation of the error terms
- ▶ Interactions or collinearity
- ▶ Categorical predictors and their interpretation (two or more categories).
- ▶ Non-linear effects of predictors
- ▶ Outliers and high-leverage points
- ▶ Multiple outputs





# Assignment Project Exam Help

- ▶ It is difficult to separate the individual effects of collinear variables on the response.

$$\text{Var}(\hat{\beta}) = \sigma^2 (X^\top X)^{-1}$$

- ▶ Collinearity has considerable effect on the precision of  $\hat{\beta}$  → large variances, wide confidence interval and low power of the tests
- ▶ It is important to identify and address potential collinearity problems

## Detection of collinearity

- ▶ Look at the correlation matrix of the variables to detect pair of highly correlated variables

- ▶ Collinearity between three or more variables compute *variance inflation factor*  $VIF$  for each variable

$$VIF(\hat{\beta}_j) = \frac{1}{1 - R_j^2}$$

- ▶ Geometrically  $1 - R_j^2$  measures how close  $x_j$  is to the subspace spanned by  $X_{-j}$

$$\text{Var}(\hat{\beta}_j) = \frac{\sigma^2}{1 - R_j^2} \leq \lambda_{\max} \lambda_{\min}^{-1} = \kappa(X)^2$$

- ▶ Examine the eigenvalues and eigenvectors **Note 7**

## Categorical predictors

- Also referred as categorical or discrete predictors or variables.
- Prediction task is called regression for quantitative output and classification for qualitative outputs

- Qualitative variables are represented by numerical codes

$$x_i = \begin{cases} 1 & \text{if the } i\text{th experiment is a success} \\ 0 & \text{if the } i\text{th experiment is a failure} \end{cases}$$

- This results in the model

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i = \begin{cases} \beta_0 + \beta_1 + \epsilon_i & \text{if the } i\text{th exp. is a success} \\ \beta_0 + \epsilon_i & \text{if the } i\text{th exp. is a failure} \end{cases}$$

### Note 8

# Assignment Project Exam Help

- <https://powcoder.com>

In this case non-linearity is obtained by considering transformed versions of the predictors

20/61

# Assignment Project Exam Help

- The residual as an estimate of the error can be used to identify outliers by examination for extreme values

$$E[\mathbf{e}] = E[(I_p - H) \mathbf{y}] = (I_p - H) \delta$$

- 21/61

## High-leverage points

# Assignment Project Exam Help

- Observations with high leverage have an unusual value for  $x_i$
- Difficult to identify when there are multiple predictors
- To quantify an observation's leverage use the leverage statistic

<https://powcoder.com>

$$h_i = \frac{1}{N} + (N - 1)^{-1} (\mathbf{x}_i - \bar{\mathbf{x}}) \mathbf{S}^{-1} (\mathbf{x}_i - \bar{\mathbf{x}})$$

- $\mathbf{S}$  is the sample covariance matrix,  $\mathbf{x}_i$  the  $i$ th row of  $\mathbf{X}$  and  $\bar{\mathbf{x}}$  the average row
- The leverage statistic  $\frac{1}{n} \leq h_i \leq 1$  and the average is  $(p + 1)/n$ .
- If an observation has  $h_i$  greatly exceeds  $(p + 1)/n$ , then we may suspect that the corresponding point has high leverage.

Note 11

# Assignment Project Exam Help

- Multiple outputs  $y_1, \dots, y_K$  need to be predicted from  $x_1, \dots, x_p$  where a linear model is assumed for each output

$\mathbf{y}_k \neq \beta_{0k} + \sum_{j=1}^p \mathbf{x}_j \beta_{jk} + \epsilon_k = f_k(X) + \epsilon_k$

- In matrix notation  $\mathbf{Y} = \mathbf{XB} + \mathbf{E}$  where  $\mathbf{Y}$  is  $N \times K$ ,  $\mathbf{X}$  is  $N \times (p+1)$ ,  $\mathbf{B}$  is  $(p+1) \times K$  (matrix of parameters) and  $\mathbf{E}$  is  $N \times K$  matrix of errors.

$$RSS(\mathbf{B}) = \sum_{k=1}^K \sum_{i=1}^N (y_{ik} - f_k(x_i))^2 = \text{tr} \left[ (\mathbf{Y} - \mathbf{XB})^\top (\mathbf{Y} - \mathbf{XB}) \right]$$

# Assignment Project Exam Help

- <https://powcoder.com>
- ▶ In case of correlated errors  $\epsilon \sim \mathcal{N}(0, \Sigma)$  the multivariate criterion becomes

$$RSS(\mathbf{B}) = \sum_{i=1}^n (y_i - f(x_i))^{\top} \Sigma^{-1} (y_i - f(x_i))$$

## Note 11



# Why ?

# Assignment Project Exam Help

- ▶ The least squares estimates is in most cases not satisfied when a large number of potential explanatory variables are available
- ▶ Improving prediction accuracy: LSE often has low bias but large variance, sacrifice a little bit of bias to reduce the variance of the predicted values and improve overall prediction accuracy
- ▶ Interpretation: Do all the predictors help to explain  $y$ ? determine a smaller subset with strongest effects and sacrifices the small details



## Deciding on the Important Variables

# Assignment Project Exam Help

### Subset selection

- ▶ **All subsets** or **best subsets** regression (examine all potential combinations of variables)
- ▶ **Forward selection** - begin with intercept and iteratively add one variable.
- ▶ **Backward selection** - begin with the full model and iteratively remove one variable.
- ▶ **What is best for cases where  $p > n$ ?**

## Best Subset

# Assignment Project Exam Help

- ▶ Retain a subset of predictor and eliminate the rest
- ▶ LSE is used to obtain the coefficients of the retained variables
- ▶ For each  $k \in \{0, 1, 2, \dots, p\}$  find the subset  $k$  that gives the smallest residual
- ▶ The choice of  $k$  is obtained using a criterion and involves a tradeoff between bias & variance
- ▶ Different criteria  $\leftarrow$  minimizes an estimate of the expected prediction error

Infeasible for large  $p$

**Note 12**

## Forward Selection

# Assignment Project Exam Help

- ▶ Sequential addition of predictors → forward stepwise selection
- ▶ Starts with the intercept and sequentially add the predictor that most improve the fit
- ▶ Add predictor producing the largest value of

$$F = \frac{\text{RSS}(\hat{\beta}_i) - \text{RSS}(\hat{\beta}_{i+1})}{\text{RSS}(\hat{\beta}_{i+1}) / (N - k - 2)}$$

- ▶ Use 90th or 95th percentile of  $F_{1, N-k-2}$  as  $F_e$

Note 14

# Assignment Project Exam Help

- Use  $F_d$  to choose the predictor to delete (smallest value)

## Add WeChat powcoder

## Add WeChat powcoder

## Alternative: Shrinkage Methods

# Assignment Project Exam Help

- ▶ Subset selection produces an interpretable model with possible lower prediction error than the full model
- ▶ The selection is discrete → often exhibits high variance
- ▶ Shrinkage methods are continuous and don't suffer as much from high variability
- ▶ We fit a model containing all  $p$  predictors using a technique that **constrains** or **regularizes** the coefficient estimates, or equivalently, that shrinks the coefficient estimates towards zero.
- ▶ Shrinking the coefficient estimates can significantly reduce their variance (not immediately obvious).

<https://powcoder.com>

Add WeChat powcoder

## Ridge Regression

- Recall that the least squares fitting procedure estimates  $\beta_0, \beta_1, \dots, \beta_p$  using the values that minimize

$$\text{RSS} = \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2$$

- In contrast, the ridge regression coefficient estimates  $\hat{\beta}^R$  are the values that minimize

$$\sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 = \text{RSS} + \lambda \sum_{j=1}^p \beta_j^2,$$

where  $\lambda \geq 0$  is a *tuning parameter*, to be determined separately.



## Ridge Regression

# Assignment Project Exam Help

$$\sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 = \text{RSS} + \lambda \sum_{j=1}^p \beta_j^2$$

<https://powcoder.com>

The shrinkage penalty.

The tuning parameter.

- It serves to control the relative impact of the penalty on the regression coefficient estimates.
- Selecting a good value for  $\lambda$  is critical.
- cross-validation is used for this.

- It is small when  $\beta_1, \beta_2, \dots, \beta_p$  are close to zero.

So it has the effect of *shrinking* the estimates of  $\beta_j$  towards zero.

Add WeChat powcoder

## Ridge Regression

# Assignment Project Exam Help

$$\sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2$$

<https://powcoder.com>

The shrinkage penalty.

Add WeChat powcoder

Accuracy ↔ Complexity

# Assignment Project Exam Help

- ▶ Ridge regression shrinks the regression coefficients by constraining their size
- ▶ This is the approach used in neural networks where it is known as weight decay
- ▶ The larger the value of  $\lambda$ , the greater the amount of shrinkage
- ▶  $\beta_0$  is left out of the penalty term

constraining their size  
This is the approach used in neural networks where it is known as weight decay  
The larger the value of  $\lambda$ , the greater the amount of shrinkage  
 $\theta_0$  is left out of the penalty term

## Ridge Regression

# Assignment Project Exam Help

- ▶ Because we have now the addition of the penalty term, the ridge regression coefficient estimates can change substantially when multiplying a given predictor by a constant!
- ▶ It is best to apply ridge regression after *standardizing the predictors*, using the formula

$$\hat{\lambda}_{ij} = \frac{x_{ij}}{\sqrt{\frac{1}{n} \sum_{i=1}^n (x_{ij} - \bar{x}_{ij})^2}}$$

## Ridge Regression

# Assignment Project Exam Help

$$\hat{\beta}^{ridge} = (X^T X + \lambda I)^{-1} X^T y$$

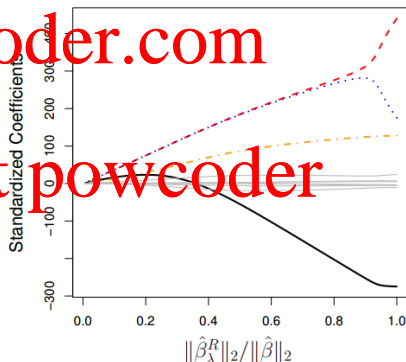
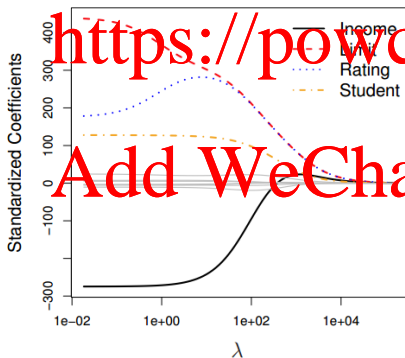
- ▶ The ridge solution is a linear function of  $y$
- ▶ Avoid singularity when  $X^T X$  is not full by adding a positive constant to the diagonal of  $X^T X$ .
- ▶ For orthogonal predictor  $\hat{\beta}^{ridge} = \gamma \hat{\beta}$  where  $0 \leq \gamma \leq 1$
- ▶ The effective degrees of freedom of the ridge regression fit is

$$df(\lambda) = \text{tr} \left( X [X^T X + \lambda I]^{-1} X^T \right)$$

Note 17

## Ridge regression - credit data example

*balance* ~  
*age, cards, education, income, limit, rating, gender, student,*  
*status, ethnicity*

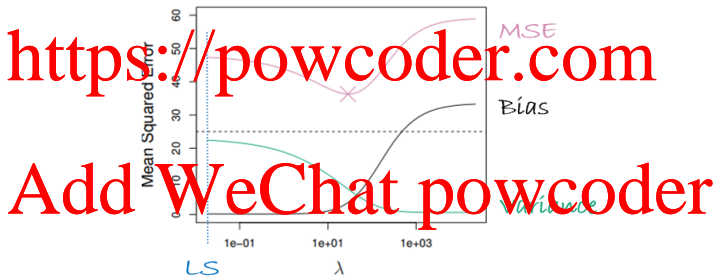


## Ridge regression vs. LS

- Shrinkage methods can reduce the variance of the estimate  $\hat{f}(x)$

Assignment Project Exam Help

$n \uparrow \rightarrow \# \text{ of (influential) variables} \downarrow$   
 $\rightarrow \text{Variance} \downarrow \quad \xi \quad \text{Bias} \uparrow$



- Ridge regression ( $\xi$  Lasso) improves on LS!
- The MSE is reduced
- The variance is much smaller at the expense of a small increase in bias. ↻ ↺ ↻

## Lasso

# Assignment Project Exam Help

- ▶ **Ridge regression** disadvantage: includes all  $p$  predictors (some of them with minor influence)
- ▶ **Lasso**, in contrast, select subset.
- ▶ The lasso coefficients,  $\hat{\beta}_{\lambda}^L$ , minimize the quantity

$$\sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p |\beta_j| = \text{RSS} + \lambda \sum_{j=1}^p |\beta_j|$$



## The Variable Selection Property of the Lasso

Why is it that the lasso, unlike ridge regression, results in coefficient estimates that are exactly equal to zero?

One can show that the lasso and ridge regression coefficient estimates solve the problems

$$\underset{\beta}{\text{minimize}} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \quad \text{subject to} \quad \sum_{j=1}^p |\beta_j| \leq s$$

and

$$\underset{\beta}{\text{minimize}} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \quad \text{subject to} \quad \sum_{j=1}^p \beta_j^2 \leq s,$$

respectively.

## Some Remarks on Lasso

# Assignment Project Exam Help

- ▶ Making  $s$  sufficiently small will cause some of the coefficients to be exactly zero  $\rightarrow$  continuous subset selection
- ▶ If  $s = \sum_{i=1}^p \|\hat{\beta}_j^{ls}\|$ , then the lasso estimates are the  $\hat{\beta}_j^{ls}$ 's.
- ▶  $s$  should be adaptively chosen to minimize an estimate of expected prediction error.

<https://powcoder.com>

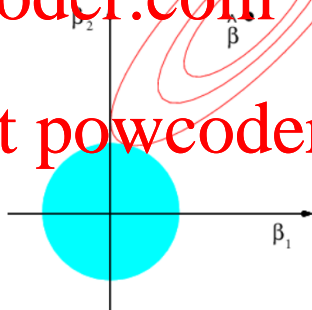
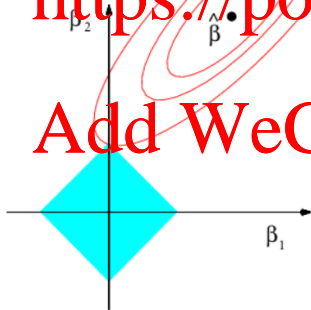
Add WeChat powcoder

## The Variable Selection Property of the Lasso

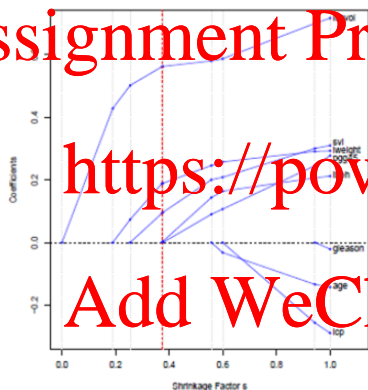
Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



## Profile of Lasso Coefficients



Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

- Profiles of lasso coefficients as the tuning parameter  $t$  is varied. The coefficients are plotted versus  $t = s / \sum_{i=1}^p \|\hat{\beta}_j^{ls}\|$

# Assignment Project Exam Help

- But, why not transform the predictors (to a lower dimension) and then fit the least squares model using the transformed predictors?

## Add WeChat powcoder

# Assignment Project Exam Help

- <https://powcoder.com>

- ▶  $\mathbf{v}_\ell^T \mathbf{S} \mathbf{x} = 0$  ensure  $\mathbf{z}_m = \mathbf{X} \mathbf{v}_m$  is uncorrelated with all previous linear combinations  $\mathbf{z}_\ell = \mathbf{X} \mathbf{v}_\ell$ ,  $\ell = 1, \dots, m-1$
- ▶  $\mathbf{y}$  is regressed on  $\mathbf{z}_1, \dots, \mathbf{z}_M$  for  $M \leq p$

## Principal Components Regression

- ▶ Since  $\mathbf{z}_m$  are orthogonal, this regression is just a sum of univariate regressions

$$\hat{\mathbf{y}}^{pcr} = \bar{\mathbf{y}} + \sum_{m=1}^M \hat{\theta}_m \mathbf{z}_m$$

$$\hat{\theta}_m = \langle \mathbf{z}_m, \mathbf{y} \rangle / \langle \mathbf{z}_m, \mathbf{z}_m \rangle$$

$$\hat{\boldsymbol{\beta}}^{pcr} = \sum_{m=1}^M \hat{\theta}_m \mathbf{v}_m$$

- ▶ if  $M = p$ ,  $\hat{\mathbf{y}}^{pcr} = \hat{\mathbf{y}}^{LS}$  since the columns of  $Z = UD$  span the column space of  $X$
- ▶ PCR discards the  $p - M$  smallest eigenvalue components.

## Principal Components Regression

- Let  $Z_1, Z_2, \dots, Z_M$  represent  $M < p$  *linear combinations* of our original  $p$  predictors. That is,

$$Z_m = \sum_{j=1}^p \phi_{mj} X_j \quad (1)$$

for some constants  $\phi_{m1}, \dots, \phi_{mp}$ .

- We can then fit the linear regression model,

$$y_i = \theta_0 + \sum_{m=1}^M \theta_m z_{im} + \epsilon_i \quad i = 1, \dots, n, \quad (2)$$

using ordinary least squares.

- Note that in model (2), the regression coefficients are given by  $\theta_0, \theta_1, \dots, \theta_M$ . If the constants  $\phi_{m1}, \dots, \phi_{mp}$  are chosen wisely, then such dimension reduction approaches can often outperform OLS regression.



## Dimension Reduction Methods

- Notice that from definition (1);

$$\sum_{m=1}^M \theta_m z_{im} = \sum_{m=1}^M \theta_m \sum_{j=1}^p \phi_{mj} x_{ij} = \sum_{j=1}^p \sum_{m=1}^M \theta_m \phi_{mj} x_{ij} = \sum_{j=1}^p \beta_j x_{ij},$$

where

$$\beta_j = \sum_{m=1}^M \theta_m \phi_{mj}. \quad (3)$$

- Hence model (2) can be thought of as a special case of the original linear regression model.
- Dimension reduction serves to constrain the estimated  $\beta_j$  coefficients, since now they must take the form (3).
- Can win in the bias-variance tradeoff.

# Assignment Project Exam Help

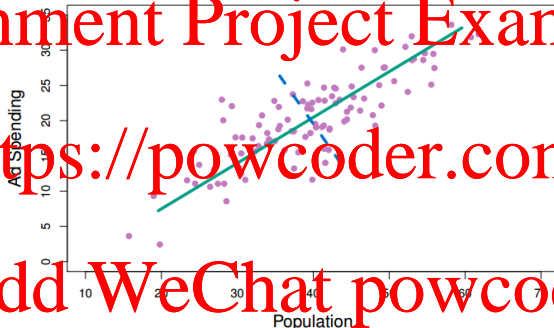
- And so on...
- Add WeChat powcoder**
- Many times we can explain most of the variation with only
- few principal components

## Principal Components

Assignment Project Exam Help

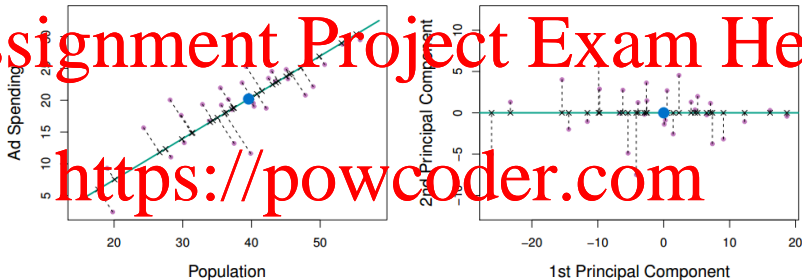
<https://powcoder.com>

Add WeChat powcoder



The population size (**pop**) and ad spending (**ad**) for 100 different cities are shown as purple circles. The green solid line indicates the first principal component, and the blue dashed line indicates the second principal component.

## Principal Components



**Add WeChat powcoder**

*A subset of the advertising data. Left: The first principal component, chosen to minimize the sum of the squared perpendicular distances to each point, is shown in green. These distances are represented using the black dashed line segments.*

*Right: The left-hand panel has been rotated so that the first principal component lies on the x-axis.*

# Assignment Project Exam Help

- We call it **Principal Components Regression (PCR)**
- Note, these directions are identified in an *unsupervised* way

# Assignment Project Exam Help

- PLS identifies these new features using the response  $Y$  (supervised way).

## Partial Least Squares (PLS)

# Assignment Project Exam Help

- ▶ Uses  $\mathbf{y}$  to construct linear combinations of the inputs.
- ▶ The inputs are weighted by the strength of their univariate effect on  $\mathbf{y}$
- ▶ Regress  $\mathbf{y}$  on  $\mathbf{z}_m \rightarrow \hat{\mathbf{y}}_m$  and orthogonalize with respect to  $\mathbf{z}_m$
- ▶ Continue the process until  $M < p$  directions are obtained
- ▶ PLS seeks directions that have high variance and have high correlation with the response

$$\max_{\|\alpha\|=1} \text{Corr}^2(\mathbf{y}, X\alpha) \text{ Var}(X\alpha)$$

$$\mathbf{v}_\ell^\top S\alpha = 0, \quad \ell = 1, \dots, m-1$$

Assignment Project Exam Help

56/61



## Illustrating the connection

The connection between these methods can be seen through the optimisation criterion they use to define projection directions

- ▶ PCR extracts components that explain the variance of the predictor space

$$\max_{\|\alpha\|=1} \text{Var}(X\alpha)$$

$$\mathbf{v}_\ell^\top \mathbf{S} \alpha = 0, \quad \ell = 1, \dots, m-1$$

- ▶ PLS extracts components that have a high covariance with

$$\max_{\|\alpha\|=1} \text{Corr}^2(\mathbf{y}, X\alpha) \text{Var}(X\alpha)$$

$$\mathbf{v}_\ell^\top \mathbf{S} \alpha = 0, \quad \ell = 1, \dots, m-1$$

- ▶ Both methods are similar in their aim to extract  $m$  components from the predictor space  $X$

# Assignment Project Exam Help

- ▶ at expressing the solution in lower dimensional subspace  $\beta = Vz$  where  $V$  is an  $p \times m$  matrix of orthonormal columns
- ▶ Using this basis for the subspace, an alternative approximate minimization problem is considered

# Add WeChat powcoder

- 58/61

## Illustrating the connection

Considering

- ▶ The singular value decomposition  $X = UDV^T$  where  $U$  is  $N \times p$ ,  $D = \text{diag}(d_1, \dots, d_p)$  is  $p \times p$  and  $V$  is  $p \times p$
- ▶ The columns of  $U$  and  $V$  are orthogonal such that  $U^T U = I_p$  and  $V^T V = I_p$
- ▶ The least squares solution takes the form

$$\hat{\beta} = \sum_{i=1}^p \frac{\mathbf{u}_i^T \mathbf{y}}{d_i} \mathbf{v}_i = \sum_{i=1}^p \beta_i$$

- ▶ The other estimator are shrinkage estimators and can be expressed as

$$\hat{\beta} = \sum_i^p f(d_i) \beta_i$$

## Multiple Outcome Shrinkage

- ▶ When the output are not correlated  $\rightarrow$  apply a univariate technique individually to each outcome or work with each column output individually
- ▶ Other approaches exploit correlations in the different responses  $\rightarrow$  canonical correlation analysis
- ▶ CCA find a sequence of linear combinations  $\mathbf{X}v_m$  and  $\mathbf{Y}u_m$  such that the correlations are maximized

Add WeChat  $\text{Corr}^2(\mathbf{Y}u_m, \mathbf{X}v_m)$  powcoder

- ▶ Reduced rank regression

$$\hat{\mathbf{B}}^{rr}(m) = \arg \min_{\text{rank}(\mathbf{B})=m} \sum_{i=1}^N (\mathbf{y}_i - \mathbf{B}\mathbf{x}_i)^\top \Sigma^{-1} (\mathbf{y}_i - \mathbf{B}\mathbf{x}_i)$$

Note 19

# Assignment Project Exam Help

- ▶ Summaries on LMS.
- ▶ <https://powcoder.com> Chapter 3, 5 & 14.5 from 'The elements of statistical learning' book.
- ▶ Chapters 3, 6 & 10.2 from 'An introduction to statistical learning' book.

# Add WeChat powcoder