

QBUS6850

Lecture 2

Python Machine Learning

Assignment Project Exam Help

<https://powcoder.com> © Discipline of Business Analytics

Add WeChat powcoder

BUSINESS SCHOOL

QBUS6850 Team



THE UNIVERSITY OF
SYDNEY

□ Topics covered

- Machine learning model representation
- Cost/loss function
- Linear regression with single and multiple features
- Optimisation algorithm: gradient descent
- Model and feature selection techniques
- Learning curves
- Training, validation and test sets
- Cross validation

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



□ References

- Bishop (2006), Chapters 1.3 - 1.5; 3.2
- Friedman et al. (2001), Chapters 2.3.1, 3.1 - 3.2, 7.1 - 7.6, 7.10
- James et al., (2014), Chapters 2.1 - 2.2
- James et al., (2014) Chapter 5.1 (Cross Validation)

<https://powcoder.com>

Add WeChat powcoder

Learning Objectives

- ☐ Understand model representation and cost function
- ☐ Understand the matrix representation of linear regression with single and multiple features
- ☐ Understand how gradient descent algorithm works
- ☐ Understand overfitting and underfitting
- ☐ Understand bias and variance decomposition and be able to diagnose them
- ☐ Be able to interpret the learning curves
- ☐ Be able to do the polynomial order selection
- ☐ Understand the reason and process of Cross Validation

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



Assignment Project Exam Help

<https://powcoder.com>
ML Basic Concepts

and Workflow
Add WeChat powcoder

Terminology in ML

➤ Input/Feature Supervised learning:

- ❖ An object is usually characterized by a *feature* scalar or vector
- ❖ Denote by $\mathbf{x} = (x_1, x_2, \dots, x_d)^T$ where each component x_i is a specified feature for the object
- ❖ Each component x_i may be a quantitative value from \mathbb{R} (the set of all real numbers) or one of finite categorical values.

➤ Outcome/Target: <https://powcoder.com>

- ❖ An unknown system (to be learnt) which generates an *output/outcome/target* denoted by $\mathbf{t} = (t_1, t_2, \dots, t_m)^T$ for each object feature \mathbf{x}
- ❖ Each component t_j may be a quantitative value from \mathbb{R} (the set of all real numbers) or one of finite categorical values
- ❖ In most cases, we assume $m = 1$. We may assume $m = 1$ in this course. Thus t is a scalar
- ❖ As a measurement value, we always suppose there are some noises ϵ in t , i.e., the measurement is $t = y + \epsilon$ where y is the true target.

➤ Ask students for examples



Terminology in ML

➤ Training/Test Dataset:

- ❖ A pair of observed (\mathbf{x}, t) is called a training/test datum.
- ❖ In unsupervised learning case, there is no target observation t .
- ❖ All the available training data are collected together by a set of *training/test data*, denoted \mathcal{D} with or without target observation

$$\mathcal{D} = \{(\mathbf{x}_n, t_n)\}_{n=1}^N \text{ or } \{\mathbf{x}_n\}_{n=1}^N$$

➤ Learner or Model:

- ❖ Use this dataset \mathcal{D} to build a prediction model, or *learner*, which will enable us to predict the outcome for new unseen objects or characterize them if without outcomes.
- ❖ A good learner is one that accurately predicts such an outcome or make a right characterization.



Data Representation

- Machine learning algorithms are built upon data. There exist different types of data. Although the numeric data are widely seen in scientific world, the categorical data are more common in business world
- When we have a data set with N observations $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$, we will organise them into a matrix of size $N \times d$ such that each row corresponds to an observation (or a case). If we have the target/output for N observations $\{t_1, t_2, \dots, t_N\}$, we also organise them into a column (simulating a row corresponding to a case or an observation), denoted by as

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_N^T \end{pmatrix} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1d} \\ x_{21} & x_{22} & \dots & x_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \dots & x_{Nd} \end{pmatrix} \in \mathbb{R}^{N \times d}, \quad \mathbf{t} = \begin{pmatrix} t_1 \\ t_2 \\ \vdots \\ t_N \end{pmatrix} \in \mathbb{R}^N$$

Machine Learning Flow

- Learning is the process of estimating an unknown dependency between the input and output or structure of a system using a limited number of observations.

Assignment Project Exam Help

- A typical learning procedure consists of the following steps:

1. Statement of the Problem
2. Data Generation/Experiment Design
3. Data Collection and Pre-processing
4. Hypothesis Formulation and Objectives
5. Model Estimation and Assessment
6. Interpretation of the Model and Drawing Conclusions

<https://powcoder.com>

Add WeChat powcoder

- Many recent application studies tend to focus on the learning methods used (i.e., a neural network).

Example: Linear Regression

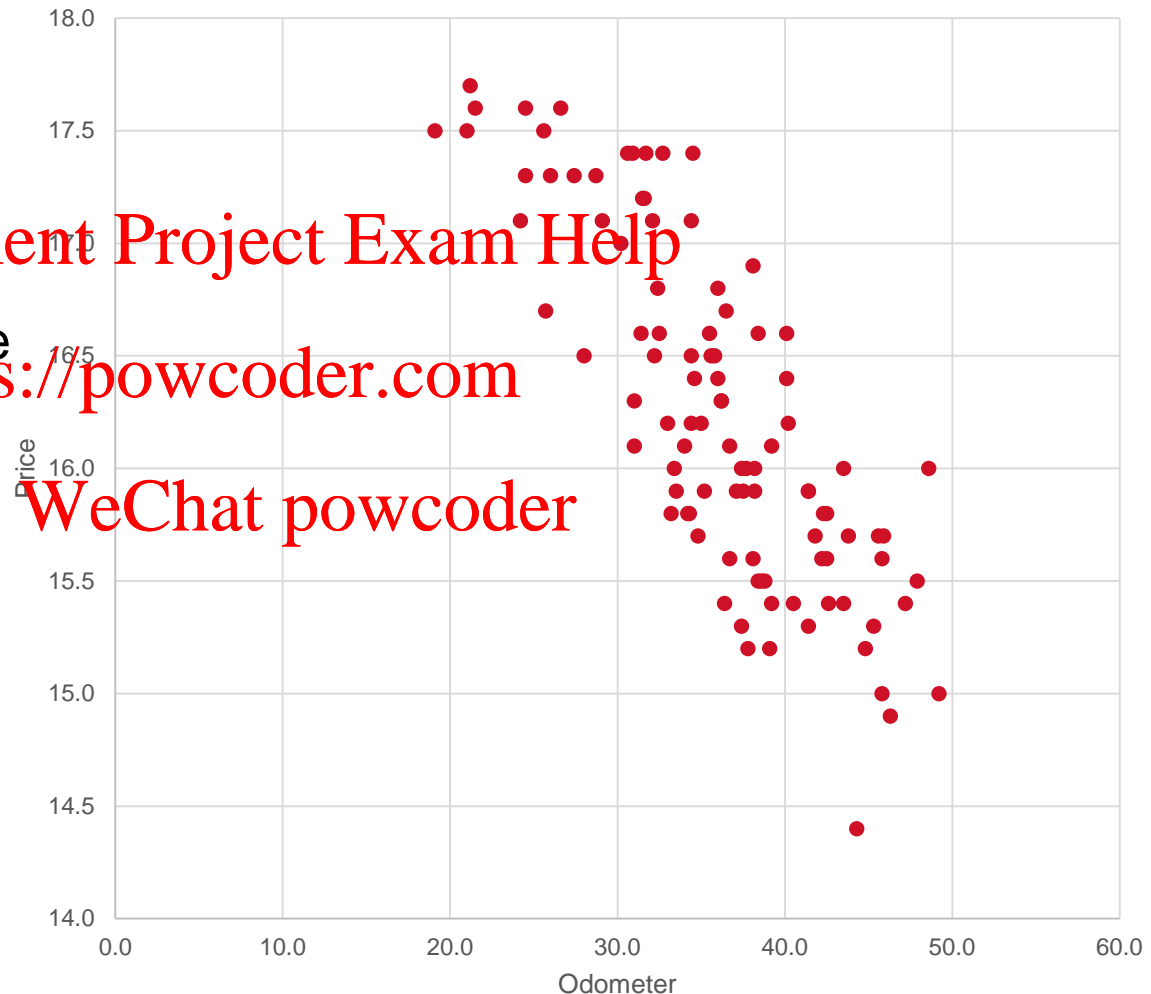
Supervised learning:

Step 1: Regression

'Problem': Predict car price

Steps 2&3: Completed

(Data collected and
labelled)



Example: Linear Regression

Step 4a: Linear Model Hypothesis

$$y = f(x; \beta) = \beta_0 + \beta_1 x$$

$f()$: simple (univariate) linear regression model;

$\beta = (\beta_0, \beta_1)^T$: model parameters.

Training set

N : number of training examples

X: “input” variable; **features**

t: “output” variable; “target” variable which is a noised version of model output y , i.e.,

$$t = y + \varepsilon$$

(x_i, t_i) i_{th} training example

Odometer (x)	Price (t)
37.4	16.0
44.8	15.2
45.8	15.0
30.9	17.4
34.7	17.4
34.0	16.1
45.9	15.7
...	...
41.4	15.3

$$(x_1, t_1) = (37.4, 16.0)$$

$$(x_3, t_3) = ??$$



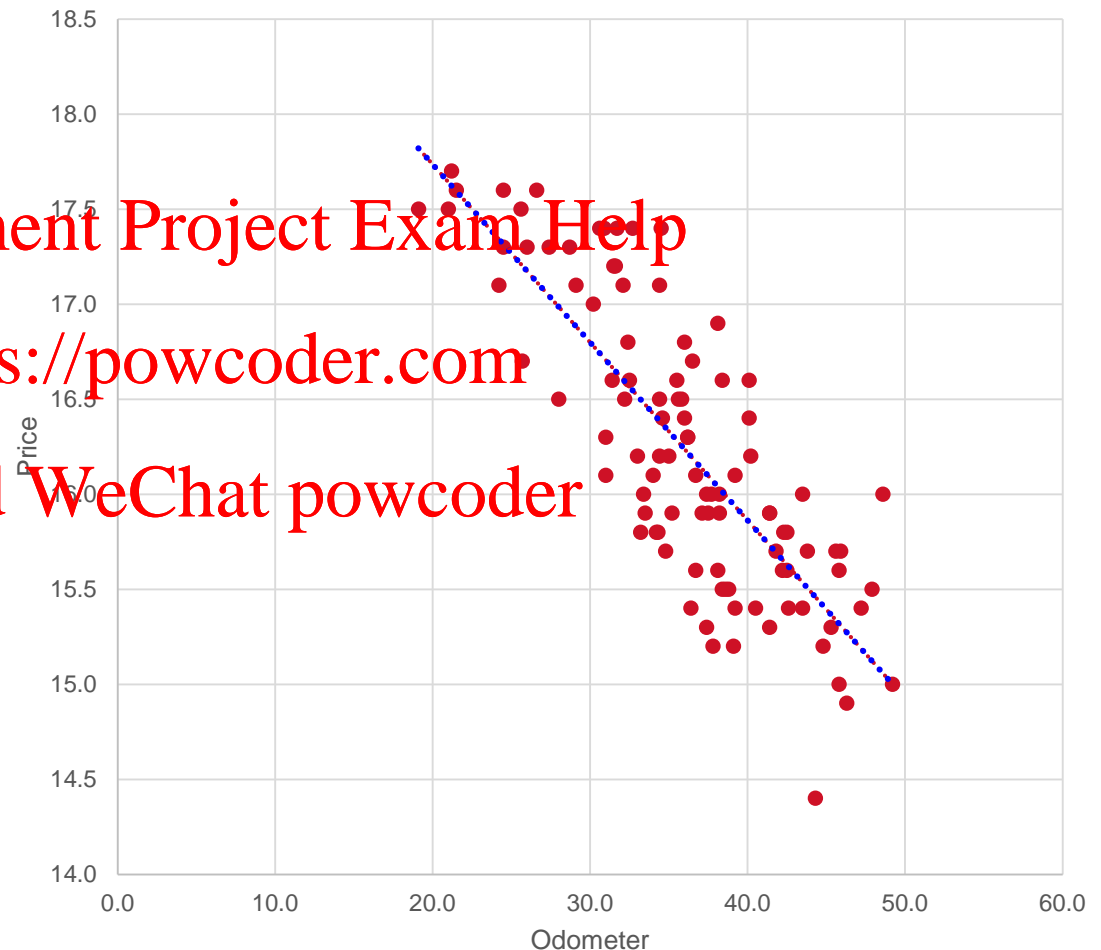
f() representation?

$$y = f(\mathbf{x}; \boldsymbol{\beta}) = \beta_0 + \beta_1 x$$

How to choose $\boldsymbol{\beta}$?

Or how can we estimate $\boldsymbol{\beta}$?

This is the task in Step 5



Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



Loss/cost function

Step 4b: Define an appropriate objective for the task in hand;

Purpose: to measure the error between the observed and the model. Loss function, also called a cost function, which is a single, overall measure of loss incurred in taking any of the available decisions or actions. (Bishop, C.M., 2006.)

Predicted y values

Observed t values

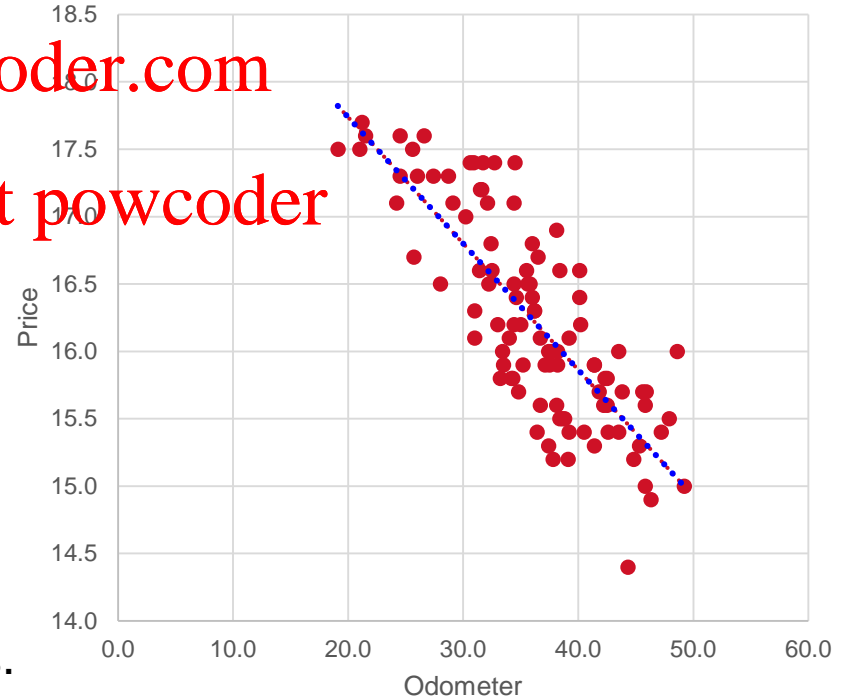
$$L(\beta_0, \beta_1) = \frac{1}{2N} \sum_{n=1}^N (f(x_n, \beta) - t_n)^2$$

For
computational
convenience

where $f(x; \beta) = \beta_0 + \beta_1 x$

$$\min_{\beta_0, \beta_1} L(\beta_0, \beta_1)$$

Choose parameters so that estimated linear regression line is close to our training examples.
This is also call **argmin**



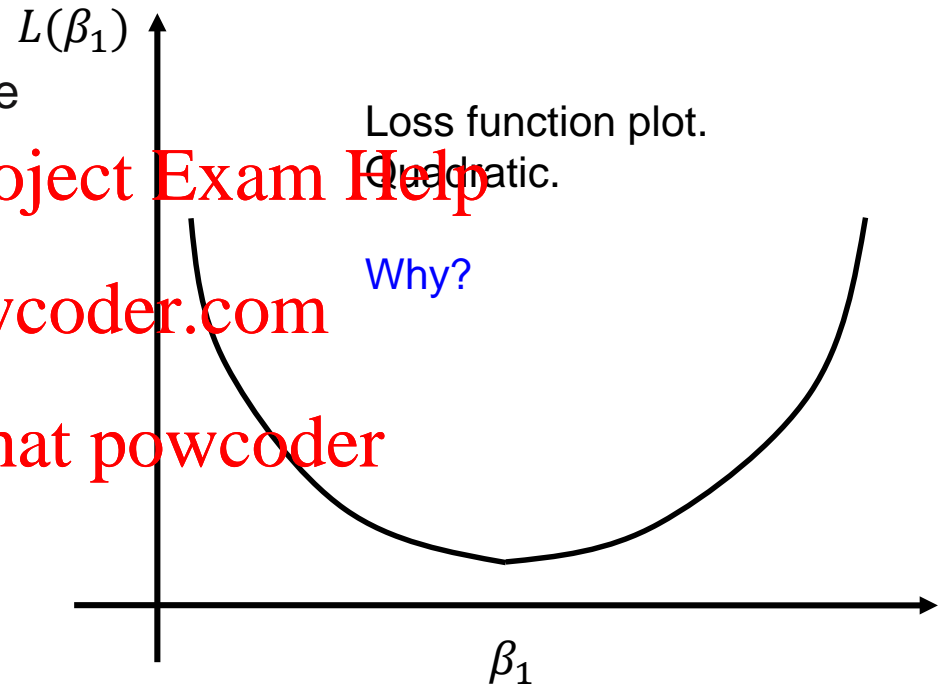
Optimisation Algorithm

Step 5: Find the model parameters;

For demo, we assume $\beta_0 = 0$, hence the loss function becomes

$$L(\beta_1) = \frac{1}{2N} \sum_{n=1}^N (f(x_n, \beta) - t_n)^2$$

Sometime this term is added for computational convenience, but will not change the estimation results of parameters

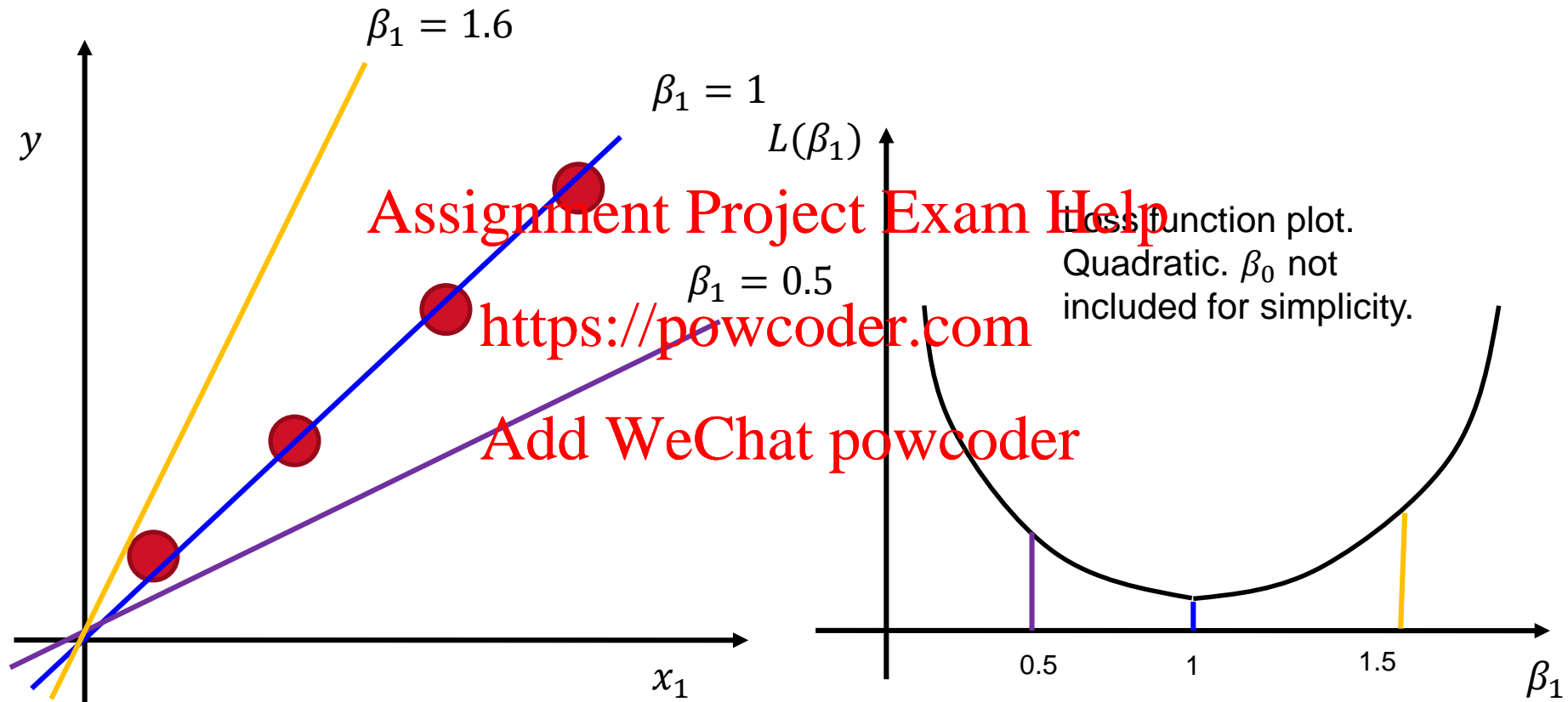


Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

β_0 not included for simplicity



If β_0 is included, how will the loss function $L(\beta_0, \beta_1)$ plot look like?



Gradient descent

- Have some random starting point for β_1 ;
- Keep updating β_1 to decrease the loss function $L(\beta_1)$ value;
- Repeat until achieving minimum (convergence).
- $\alpha > 0$ is called the learning rate: in empirical study, we can try many α values, and select the one generates least $L(\beta_1)$
- Gradient descent can converge to a **local minimum**

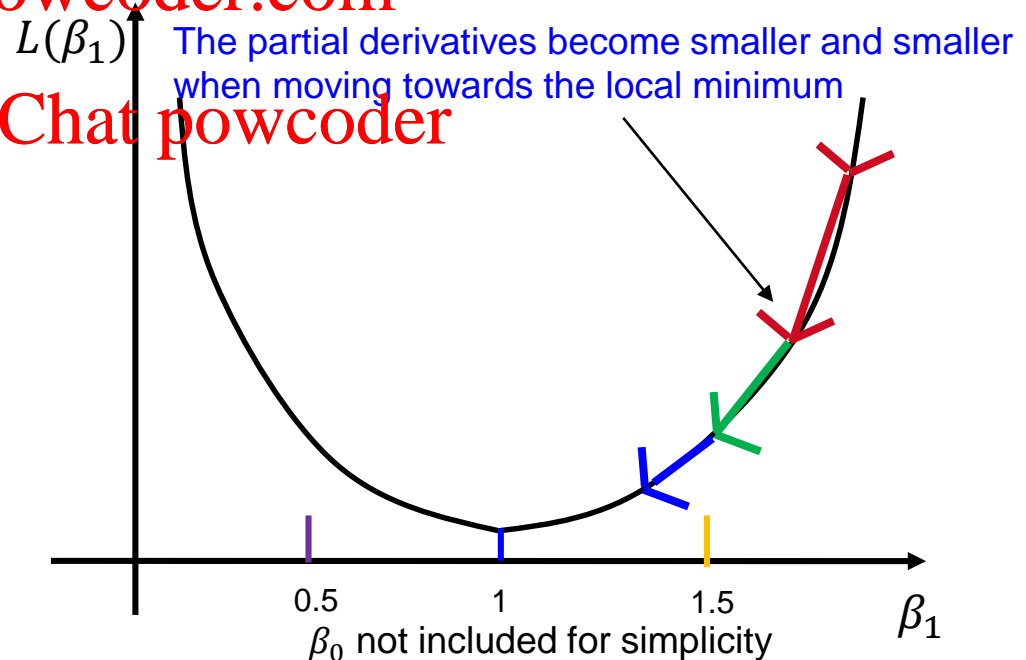
Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

$$\beta_1 := \beta_1 - \alpha \frac{\partial L(\beta_1)}{\partial \beta_1}$$

Assignment notation: keep updating β_1 based on calculations to the right hand side of this notation



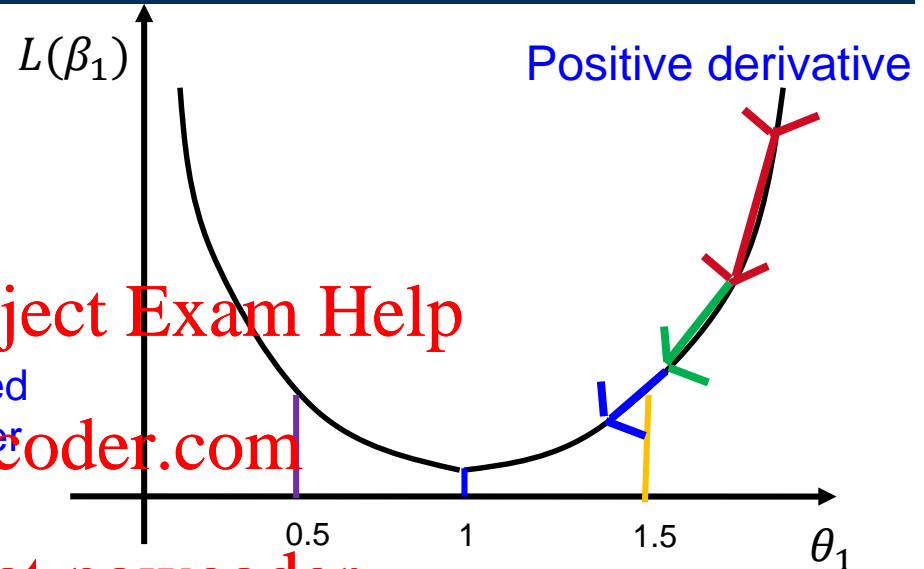


If starting point of β_1 is to the **right** of the local minimum:

$$\frac{\partial(L(\beta_1))}{\partial\beta_1} > 0$$

$$\beta_1 := \beta_1 - \alpha \frac{\partial(L(\beta_1))}{\partial\beta_1}$$

β_1 is updated
to be smaller
and smaller

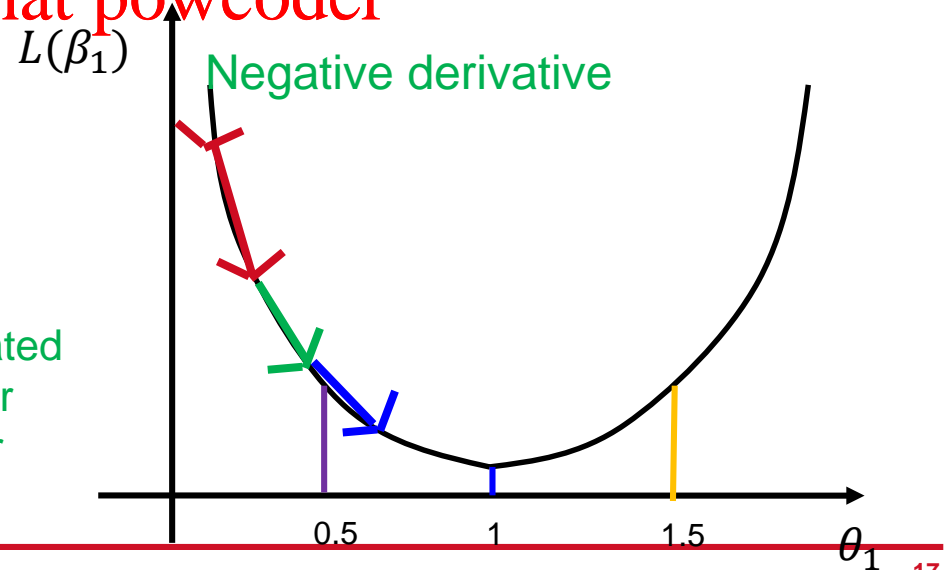


If starting point of β_1 is to the **left** of the local minimum:

$$\frac{\partial(L(\beta_1))}{\partial\beta_1} < 0$$

$$\beta_1 := \beta_1 - \alpha \frac{\partial(L(\beta_1))}{\partial\beta_1}$$

β_1 is updated
to be larger
and larger





Gradient descent of linear regression

- Have some random starting points for β_0 and β_1 ;
- Keep updating β_0 and β_1 (**simultaneously**) to decrease the loss function $L(\beta_0, \beta_1)$ value;
- Repeat until achieving minimum (convergence).

Assignment Project Exam Help

$$L(\beta_0, \beta_1) = \frac{1}{2N} \sum_{n=1}^N (f(x_n, \beta) - t_n)^2 = \frac{1}{2N} \sum_{n=1}^N (\beta_0 + \beta_1 x_{n1} - t_n)^2$$

Add WeChat powcoder

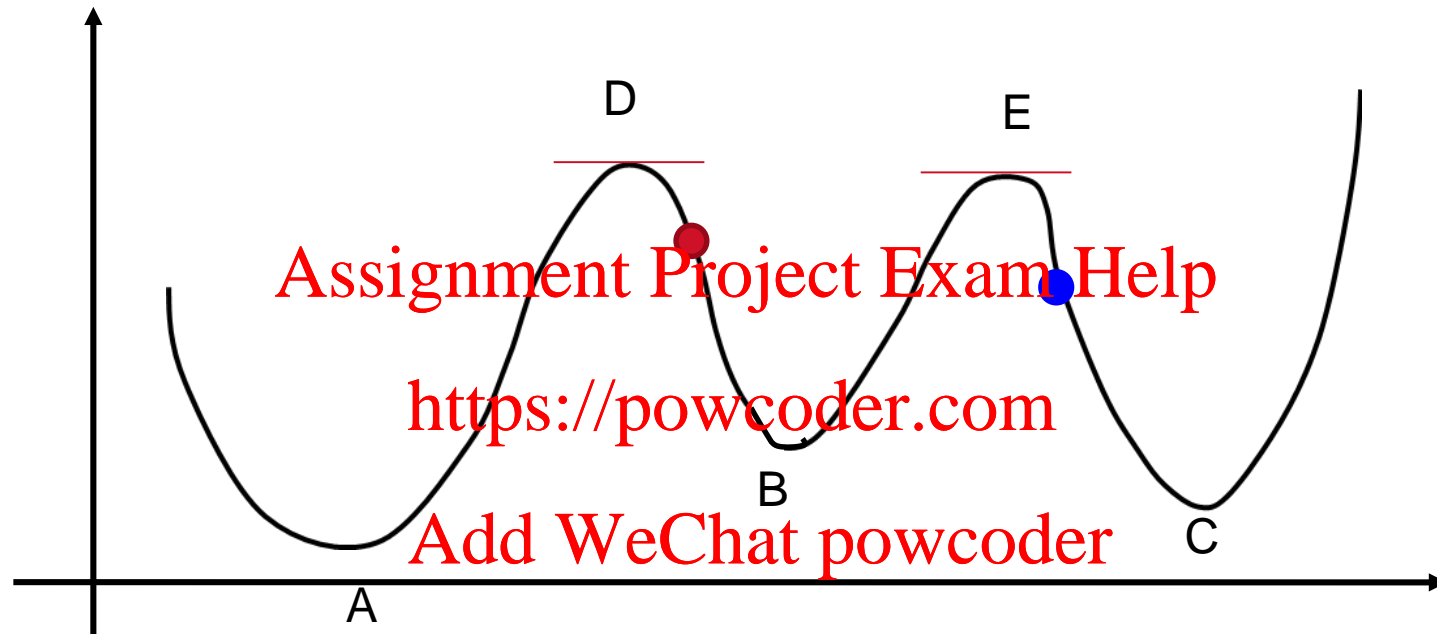
$$\beta_0 := \beta_0 - \alpha \frac{\partial L(\beta_0, \beta_1)}{\partial \beta_0} = \beta_0 - \alpha \frac{1}{N} \sum_{n=1}^N (\beta_0 + \beta_1 x_{n1} - t_n)$$

$$\beta_1 := \beta_1 - \alpha \frac{\partial L(\beta_0, \beta_1)}{\partial \beta_1} = \beta_1 - \alpha \frac{1}{N} \sum_{n=1}^N (\beta_0 + \beta_1 x_{n1} - t_n) x_{n1}$$

Update
simultaneously



Why local minimum?



- If the starting point is the red dot, then gradient descent can only converge to local minimum B
- If the starting point is the blue dot, then gradient descent can only converge to local minimum C
- The derivatives at D and E are 0
- The derivatives at A, B or C are also 0



Assignment Project Exam Help

<https://powcoder.com>
Linear regression with
multiple features
Add WeChat powcoder



Multiple features

Number (x_1)	Nearest (x_2)	Office (x_3)	Enrolment (x_4)	Income (x_5)	Distance (x_6)	Margin (y)
3203	4.2	54.9	8.0	40	4.3	55.5
2810	2.8	49.6	17.5	38	23.0	33.8
2890	2.4	25.4	20.0	38	4.2	49.0
3422	3.3	43.4	15.5	41	19.4	31.9
2687	0.9	57.8	15.5	46	11.0	57.4
3759	2.9	63.5	19.0	36	17.3	49.0
2341	2.3	58	23.0	31	11.8	46.0
3021	1.7	57.2	8.5	45	8.8	50.2

N: number of training examples

d: number of features

x: “input” variable; d features $\mathbf{x} = (x_1, x_2, \dots, x_d)^T \in \mathbb{R}^d$

y: “output” variable; “target” variable. We consider a single output

For dataset, we use notation

$x_{nj} \rightarrow n_{\text{th}}$ training example of j_{th} feature $\rightarrow x_{32} = 2.4$

For this example, the linear model is

$$f(\mathbf{x}, \boldsymbol{\beta}) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 + \beta_5 x_5 + \beta_6 x_6$$

d=6

Matrix Representation

In general with d features

$$f(\mathbf{x}, \boldsymbol{\beta}) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 + \cdots + \beta_d x_d$$

Define a special feature $x_0 = 1$ always taking value 1

Assignment Project Exam Help

So, new feature variable of a vector of d dimension

<https://powcoder.com>

$$\mathbf{x} = (x_0, x_1, x_2, \dots, x_d)^T \in \mathbb{R}^{d+1}$$

Add WeChat powcoder

Think about why this T

Collect all parameters into a vector as

$$\boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_d \end{bmatrix} \Longrightarrow f(\mathbf{x}, \boldsymbol{\beta}) = \mathbf{x}^T \boldsymbol{\beta}$$



Multiple features loss function

$$f(\mathbf{x}, \boldsymbol{\beta}) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 + \cdots + \beta_d x_d$$

$$f(\mathbf{x}, \boldsymbol{\beta}) = \mathbf{x}^T \boldsymbol{\beta}$$

Consider an input feature data $\mathbf{x}_n = (x_{n0}, x_{n1}, x_{n2}, \dots, x_{nd})$ and its corresponding output (or target) y_n . The squared model error is

$$e_n^2 = (t_n - f(\mathbf{x}_n, \boldsymbol{\beta}))^2 = (t_n - \mathbf{x}_n^T \boldsymbol{\beta})^2$$

The overall “mean” error is

$$L(\boldsymbol{\beta}) = \frac{1}{2N} \sum_{n=1}^N e_n^2 = \frac{1}{2N} \sum_{n=1}^N (t_n - \mathbf{x}_n^T \boldsymbol{\beta})^2$$

Note $\frac{1}{2}$ here is for mathematical convenience

Multiple features loss function

Another way to write the loss function

Denote

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_N^T \end{pmatrix} = \begin{pmatrix} x_{10} & x_{11} & x_{12} & \dots & x_{1d} \\ x_{20} & x_{21} & x_{22} & \dots & x_{2d} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_{N0} & x_{N1} & x_{N2} & \dots & x_{Nd} \end{pmatrix} \in \mathbb{R}^{N \times (d+1)}, \mathbf{t} = \begin{pmatrix} t_1 \\ t_2 \\ \vdots \\ t_N \end{pmatrix} \in \mathbb{R}^N$$

Assignment Project Exam Help
<https://powcoder.com>
Add WeChat powcoder

It is easy to prove that

$$L(\boldsymbol{\beta}) = \frac{1}{2N} (\mathbf{t} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{t} - \mathbf{X}\boldsymbol{\beta})$$

Question: What is the meaning of, for example, the second column vector of data matrix \mathbf{X} , called the design matrix? What is the 10th row of \mathbf{X} ?



Closed-Form Solution: Normal equation

Normal equation is an analytical solution:

- Compared with the gradient descent: no need to choose learning rate α and do not have to run a loop
- Can be slow when d is large

Assignment Project Exam Help

<https://powcoder.com>

$$(d+1) \times 1 \xrightarrow{\quad} \beta = (X^T X)^{-1} X^T t \xleftarrow{\quad} N \times 1$$

Add WeChat powcoder

$(d+1) \times N$

$$\begin{pmatrix} \mathbf{X}^T \mathbf{X} \end{pmatrix} \quad \text{Non-invertable?}$$



$\mathbf{X}^T \mathbf{X}$ non-invertable?

Reason: Multiconlinearity problem or redundant features.

Rank and determinant of $\mathbf{X}^T \mathbf{X} = ?$

Assignment Project Exam Help

Solution: Drop one or more highly correlated features from the model or collect more data

<https://powcoder.com>

Reason: The number of features is too large, e.g. ($N \ll d$).

Add WeChat powcoder

Rank and determinant of $\mathbf{X}^T \mathbf{X} = ?$

Solution: Drop some features or collect more data;
Add "regularization" term into the model

For real matrices \mathbf{X} , $\text{rank}(\mathbf{X}^T \mathbf{X}) = \text{rank}(\mathbf{X} \mathbf{X}^T) = \text{rank}(\mathbf{X}) = \text{rank}(\mathbf{X}^T)$



Gradient descent of linear regression with multiple features

- Have some random starting points for all β_i ;
- Keep updating all β_i (simultaneously) to decrease the loss function $L(\boldsymbol{\beta})$ value;
- Repeat until achieving minimum (convergence).

Assignment Project Exam Help

$$L(\boldsymbol{\beta}) = \frac{1}{2N} \sum_{n=1}^N (t_n - f(\mathbf{x}_n, \boldsymbol{\beta}))^2 = \frac{1}{2N} \sum_{n=1}^N (\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_d x_d - t_n)^2$$

<https://powcoder.com>

$$\beta_0 := \beta_0 - \alpha \frac{\partial L(\boldsymbol{\beta})}{\partial \beta_0} = \beta_0 - \alpha \frac{1}{N} \sum_{n=1}^N (\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_d x_d - t_n)$$

$$\beta_1 := \beta_1 - \alpha \frac{\partial L(\boldsymbol{\beta})}{\partial \beta_1} = \beta_1 - \alpha \frac{1}{N} \sum_{n=1}^N (\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_d x_d - t_n) x_{n1}$$

...

$$\beta_d := \beta_d - \alpha \frac{\partial L(\boldsymbol{\beta})}{\partial \beta_d} = \beta_d - \alpha \frac{1}{N} \sum_{n=1}^N (\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_d x_d - t_n) x_{nd}$$

Update
simultaneously



Gradient Descent in Matrix Form

- We can write the Gradient Descent for linear regression with multiple features in a matrix form
- The matrix form looks much more simple. First define

$$\mathbf{f}(\mathbf{X}, \boldsymbol{\beta}) := \begin{bmatrix} f(\mathbf{x}_1, \boldsymbol{\beta}) \\ f(\mathbf{x}_2, \boldsymbol{\beta}) \\ \vdots \\ f(\mathbf{x}_N, \boldsymbol{\beta}) \end{bmatrix}; \quad \mathbf{t} = \begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_N \end{bmatrix}; \quad \text{and} \quad \frac{\partial L(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} := \begin{bmatrix} \frac{\partial L(\boldsymbol{\beta})}{\partial \beta_0} \\ \frac{\partial L(\boldsymbol{\beta})}{\partial \beta_1} \\ \vdots \\ \frac{\partial L(\boldsymbol{\beta})}{\partial \beta_d} \end{bmatrix}$$

Then it can be proved that

$$\frac{\partial L(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = \frac{1}{N} \mathbf{X}^T (\mathbf{f}(\mathbf{X}, \boldsymbol{\beta}) - \mathbf{t})$$

- Hence gradient descent is

$$\boldsymbol{\beta} := \boldsymbol{\beta} - \alpha \frac{\partial L(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = \boldsymbol{\beta} - \frac{\alpha}{N} \mathbf{X}^T (\mathbf{f}(\mathbf{X}, \boldsymbol{\beta}) - \mathbf{t})$$



Feature Scaling

Target: transform features to be on a similar scale.

Results: faster convergence of the optimisation algorithm

Number (x_1)	Nearest (x_2)	Office (x_3)	Enrolment (x_4)	Income (x_5)	Distance (x_6)	Margin (y)
3203	4.2	54.9	8.0	40	4.3	55.5
2810	2.8	49.6	17.5	38	23.0	33.8
2890	2.4	25.4	20.0	38	4.2	49.0
3422	3.3	43.4	15.5	41	19.4	31.9
2687	0.9	67.8	15.5	46	11.0	57.4
3759	2.9	63.5	19.0	36	17.3	49.0
2341	2.3	68	23.0	31	11.8	46.0
3021	1.7	57.2	8.5	45	8.8	50.2

Number (x_1): 1613 to 4214

Nearest (x_2): 0.1 to 4.2

...

Mean Normalization

$$x_j^{(i)} = \frac{x_j^{(i)} - \bar{x}_j}{s}$$

Goal: have all the features to be

approximately 0 mean and 1 variance



Polynomial Regression

Question:

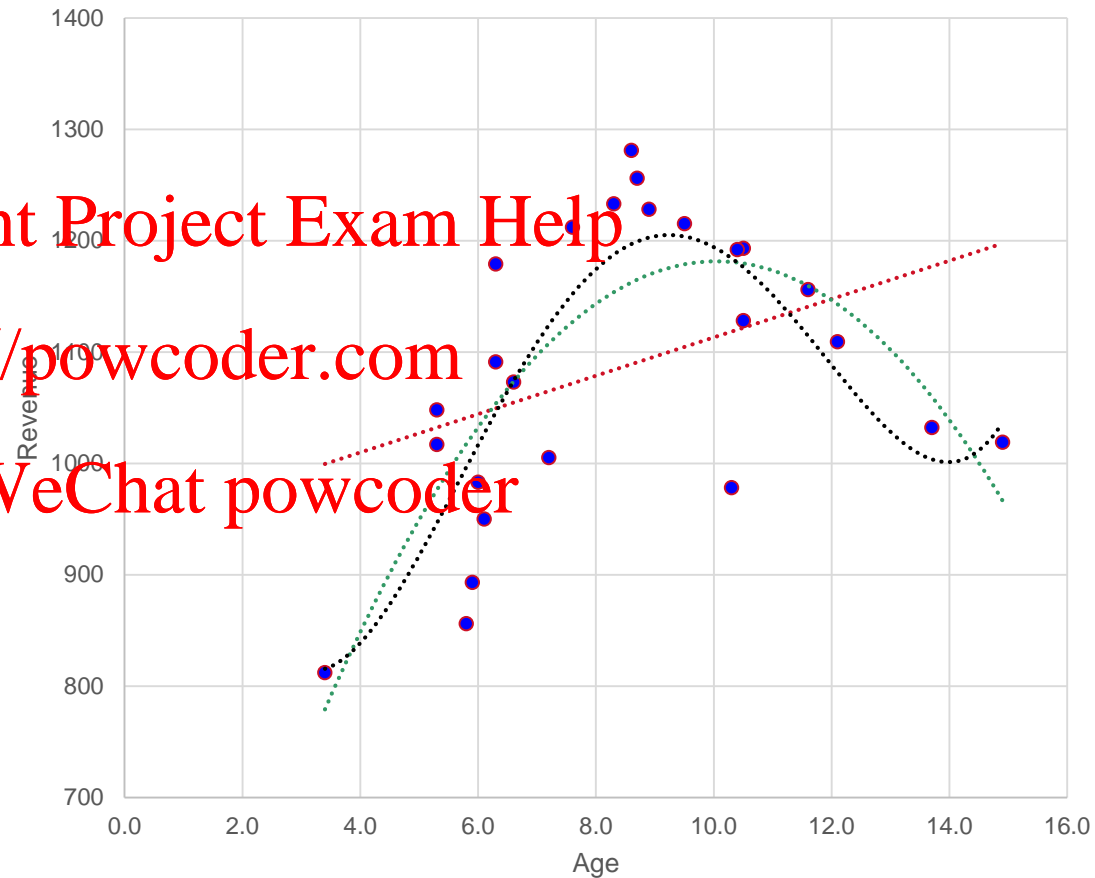
Which model is the best model?

$$f(\mathbf{x}, \boldsymbol{\beta}) = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2$$

We added a quadratic feature as $x_2 = x_1^2$ constructed from the first feature; Similarly

$$f(\mathbf{x}, \boldsymbol{\beta}) = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2 + \beta_3 x_1^3$$

$$f(\mathbf{x}, \boldsymbol{\beta}) = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2 + \beta_3 x_1^3 + \beta_4 x_1^4$$



Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



fit a 4th order polynomial regression model

from sklearn.preprocessing import PolynomialFeatures

poly = PolynomialFeatures(4)

poly_4 = poly.fit_transform(np.reshape(x_data, (1200, 1)))

poly_4 = np.asmatrix(poly_4)

lr_obj_4 = LinearRegression()

lr_obj_4.fit(poly_4, y_data)

print(lr_obj_4.intercept_) # This is the intercept \beta_0 in our notation

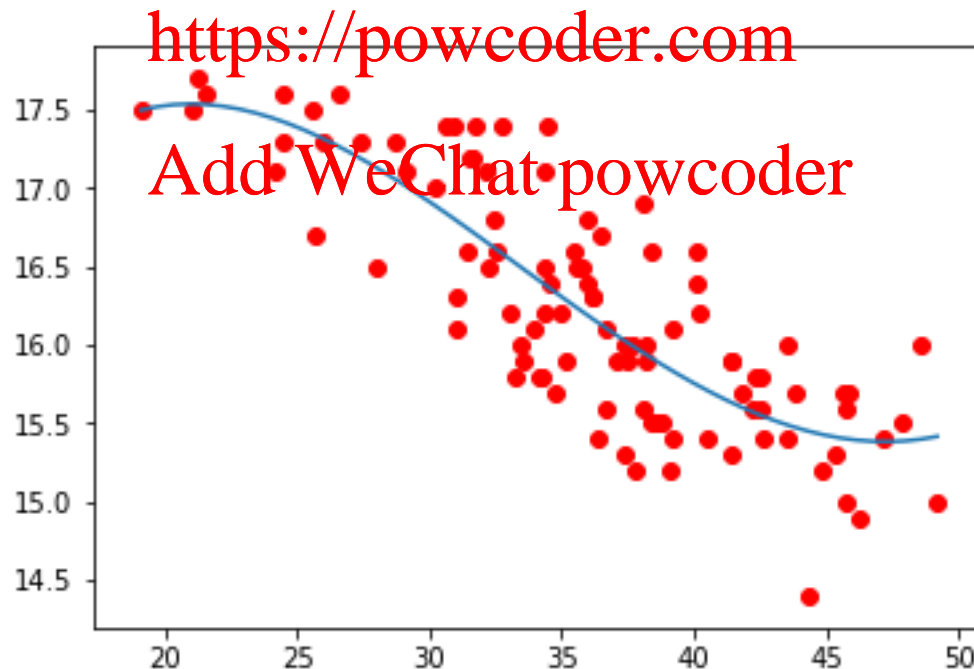
print(lr_obj_4.coef_)

```
...: print(lr_obj_4.intercept_) # This is the intercept \beta_0 in our notation
...: print(lr_obj_4.coef_)      # They are \beta_1, \beta_2, \beta_3, \beta_4 in our
notation
[ 9.70355382]
[[ 0.00000000e+00  9.14531510e-01 -3.44838143e-02  4.46054761e-04
 -1.52446757e-06]]
```



```
# plot the fitted 4th order polynomial regression line  
x_temp = np.reshape(np.linspace(np.min(x_data), np.max(x_data), 50), (50,1))  
poly_temp0_4 = poly.fit_transform(np.reshape(x_temp, (50,1)))  
y_temp = lr_obj.predict(poly_temp0_4)
```

```
plt.plot(x_temp, y_temp)  
plt.scatter(odometer, car_price, label = "Observed Points", color = "red")
```



Is this a better
model?



Assignment Project Exam Help

Model selection

<https://powcoder.com>

Add WeChat powcoder



Model Selection and Assessment

Step 5

Model Selection:

estimate the performance of different models in order to choose the (approximate) best one

Model Assessment:

after chosen the “best” model, estimate its prediction error (generalization error) on new data. (Friedman et al., 2001).

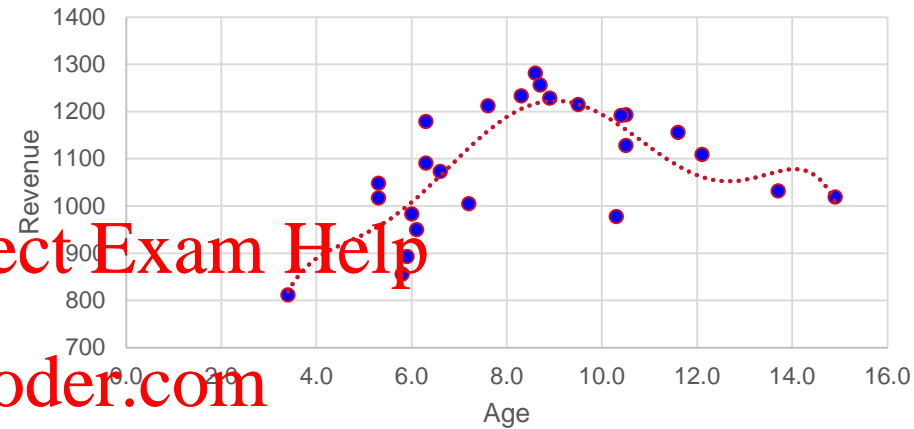
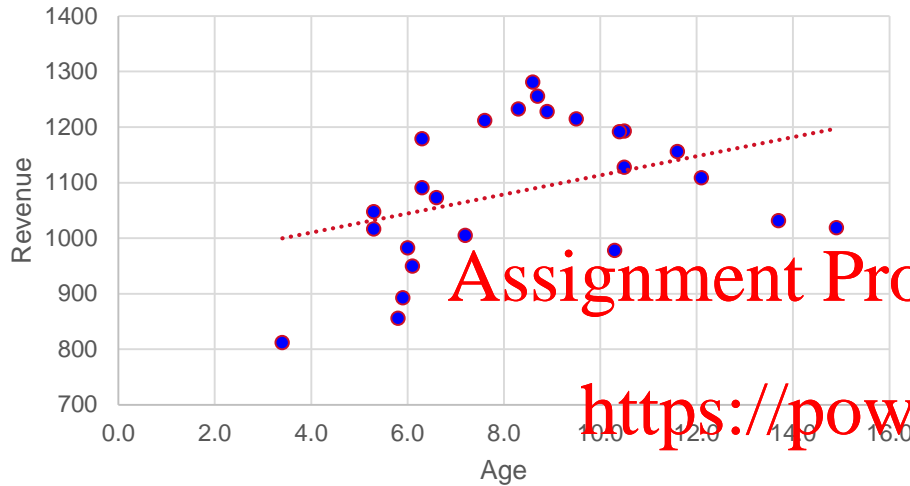
<https://powcoder.com>
Add WeChat powcoder

In general, we shall divide the given dataset into **three** parts:

- **Training dataset** (60%): used to estimate a model (or models)
 - **Validation dataset** (20%): used to select an appropriate model
 - **Test dataset** (20%): used to assess the performance of the selected model. This set of data should be hidden from training and validating process.
- Some academics use 50%, 25%, 20% split

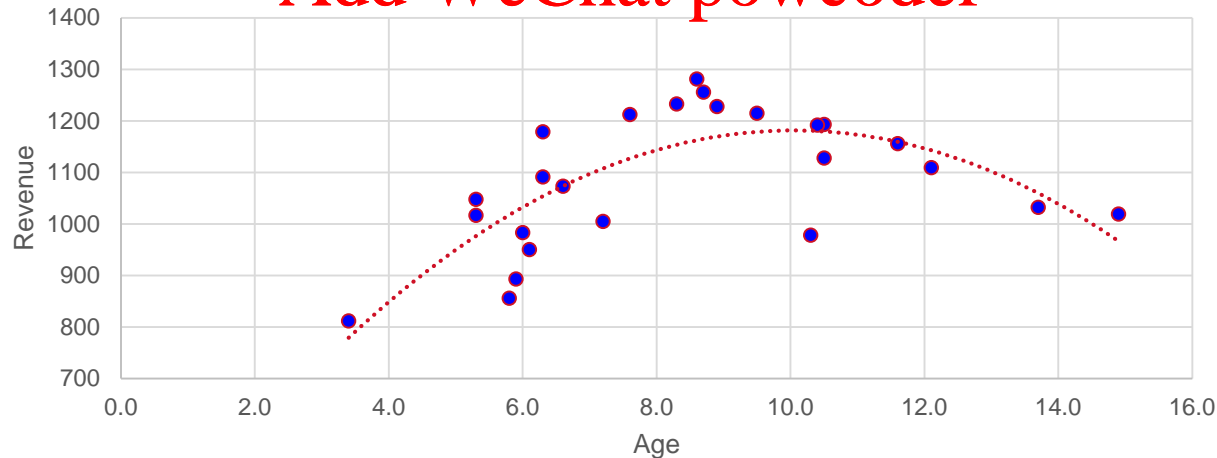


Underfitting & Overfitting



Underfitting: High bias, low variance

Overfitting: Low bias, high variance



Best model



Underfitting and Overfitting

Why is overfitting bad?

- Low training error, high generalization error
- Poor predictive performance
- Overreacts to minor fluctuations in training data

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

How to address overfitting:

☐ Drop some features

- Model selection algorithm
- Manually select features to keep

☐ Regularization

- Keep all features, reduce the magnitude/values of parameters



Training, validation and test sets

$$L(\boldsymbol{\beta}) = \frac{1}{2N} \sum_{n=1}^N (t_n - f(\mathbf{x}_n, \boldsymbol{\beta}))^2$$

This loss function issued for training, validation and test sets respectively

Assignment Project Exam Help

Odometer (x)	Price (t)
37.4	16.0
44.8	15.2
45.8	15.0
30.9	17.4
31.7	17.4
34.0	16.1
45.9	15.7
19.1	17.5
40.1	16.6
40.2	16.2

<https://powcoder.com>

Add WeChat powcoder

Training set: 60%

Validation set: 20%

Test set: 20%

Cost function

$L_{train}(\boldsymbol{\beta})$

$L_v(\boldsymbol{\beta})$

$L_{test}(\boldsymbol{\beta})$



Polynomial Order Selection

Training set	Validation set	Test
Estimate the parameters	Select the best model	Estimate the generalization error

Which one to use?

Assignment Project Exam Help

$$f(\mathbf{x}; \boldsymbol{\beta}) = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2$$

$$f(\mathbf{x}; \boldsymbol{\beta}) = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2 + \beta_3 x_1^3$$

Cannot use training set to select the best model.

Not a fair competition

$$f(\mathbf{x}; \boldsymbol{\beta}) = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2 + \beta_3 x_1^3 + \beta_4 x_1^4$$

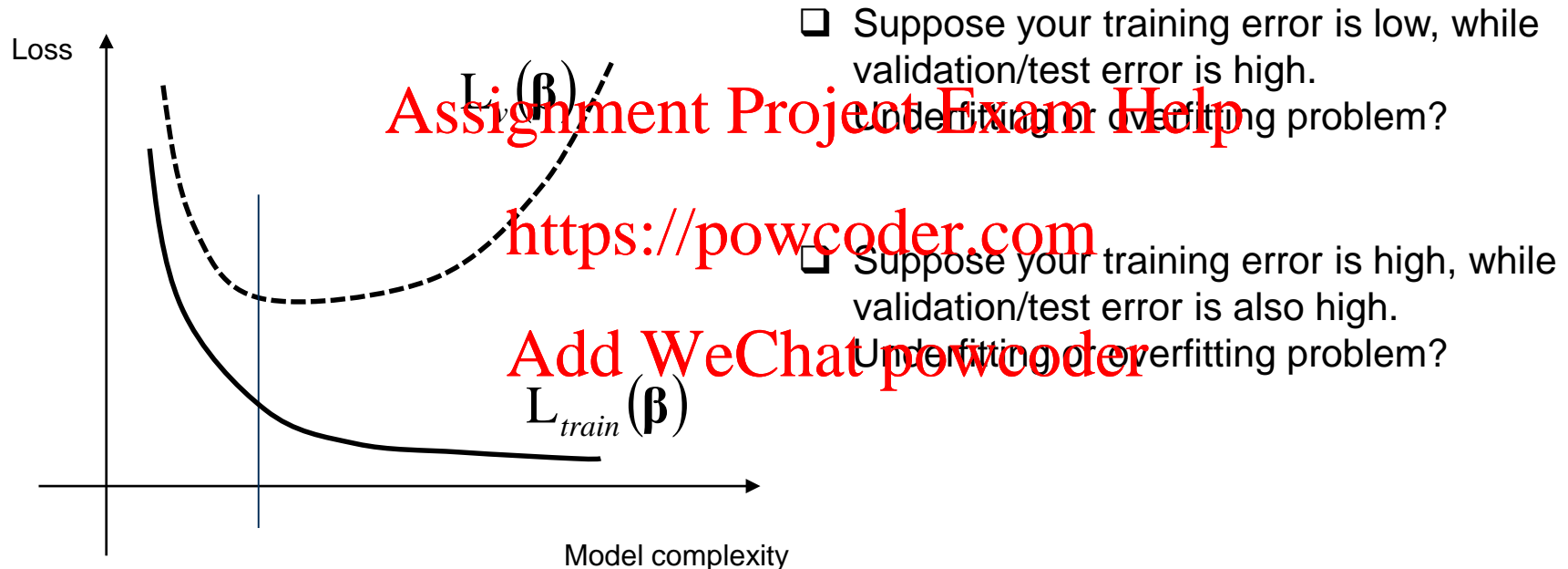
Since the model minimize this training data set, not necessarily minimize the new date sets.

- ☐ Optimize the parameters $\boldsymbol{\theta}$ employing the training set for each polynomial degree
- ☐ Find the polynomial degree d with the smallest error using the validation set
- ☐ Estimate the generalization error using the test set.



Diagnosing Learning Curve

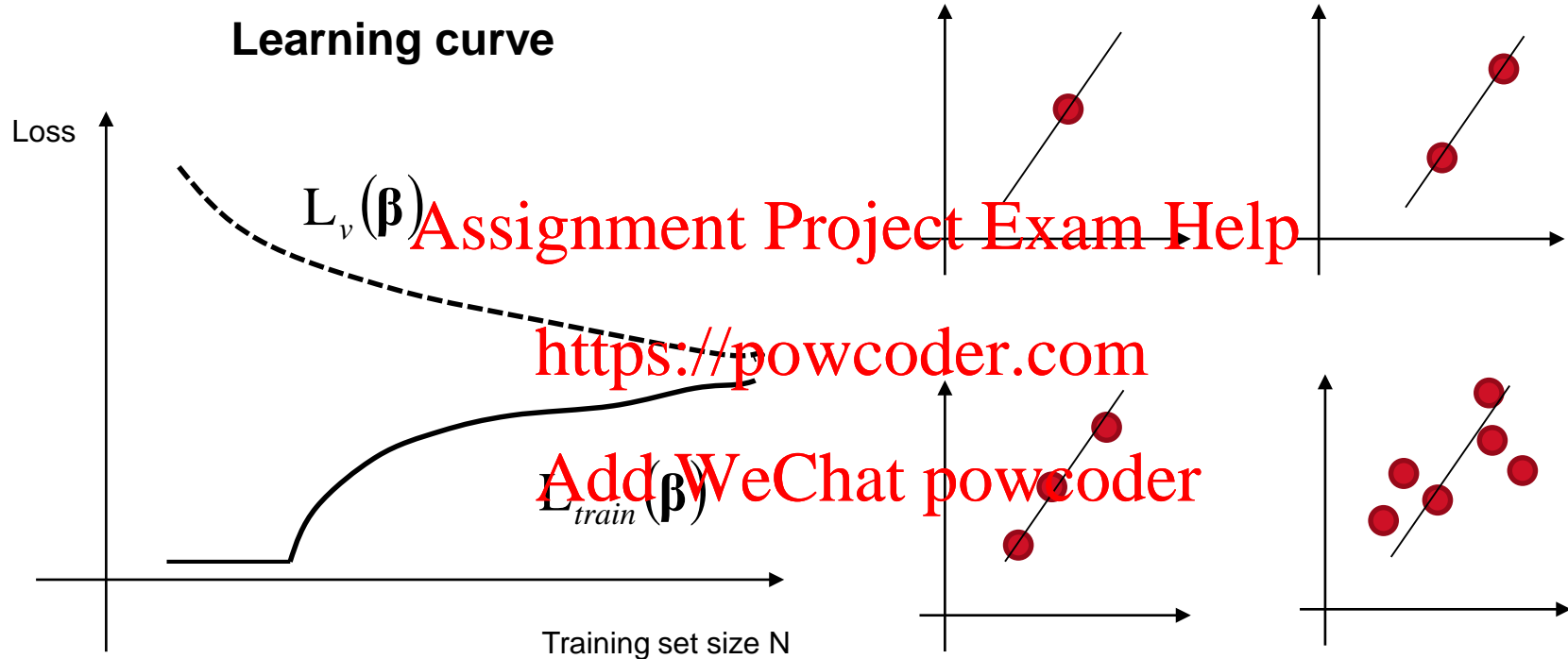
Learning Curve



Underfitting: training error is high, validation error is **slightly** > training error;
Overfitting: training error is low, validation error is **significantly** > training error.



Diagnosing Learning Curve



Why is this?

The impact of training set size
on loss function



Assignment Project Exam Help

Cross Validation

<https://powcoder.com>

Add WeChat powcoder

K-Fold Cross-Validation

- ❑ If we had enough data, we would set aside a validation set and use it to assess the performance of our prediction model

Assignment Project Exam Help

- ❑ However, data are often scarce, this is usually not possible

<https://powcoder.com>

- ❑ Particularly when we do not have sufficient labelled data

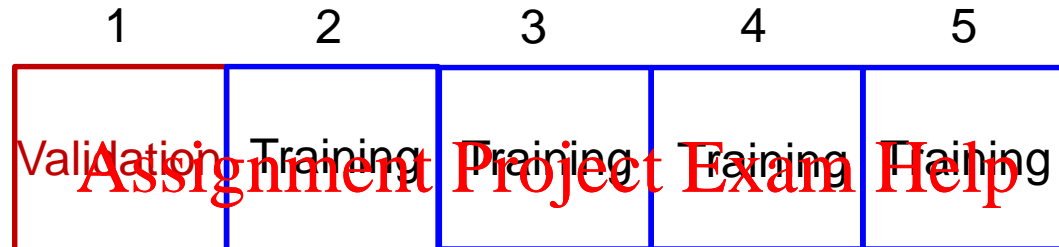
Add WeChat powcoder

- ❑ K-fold cross-validation (CV) uses part of the available data to fit the model, and a different part to test it, then iterate/repeat this process



5-fold Cross-Validation

K=5



<https://powcoder.com>

The original sample is **randomly** partitioned into K equal size subsamples; For each iteration, of the k subsamples, a single subsample is retained as the validation data for testing the model, and the remaining K-1 subsamples are used as training data.

Specifically, **at iteration 1:**

Data set **1** is chosen as validation set, and **2, 3, 4, 5** are chosen as training sets. Estimate the parameters β_1 using training sets and calculate the validation error $L_v(\beta_1)$



At iteration 2:

Data set **2** is chosen as validation set, and **1, 3, 4, 5** (K-1 sets) are chosen as training sets. Estimate the parameters β_2 using training sets and calculate the validation error $L_v(\beta_2)$

Assignment Project Exam Help

1	2	3	4	5
Validation	Training	Training	Training	Training

<https://powcoder.com>

Add WeChat powcoder

Repeat until iteration $K=5$. Estimate the parameters β_5 using training sets and calculate the validation error $L_v(\beta_5)$

Output mean validation error $(L_v(\beta_1) + L_v(\beta_2) + \dots + L_v(\beta_5))/5$ on validation sets and select the model that generates the least error.



Cross-Validation potential issues:

❑ Computational cost

- You must train each model K times.
- The K training sets (and hence the trained models) are highly correlated (see, e.g., Bengio Y & Grandvalet Y (2004))

<https://powcoder.com>

No Unbiased Estimator of the Variance of K -Fold Cross-Validation

Yoshua Bengio

*Dept. IRO, Université de Montréal
C.P. 6128, Montreal, Qc, H3C 3J7, Canada*

BENGIOY@IRO.UMONTREAL.CA

Yves Grandvalet

*Heudiasyc, UMR CNRS 6599
Université de Technologie de Compiègne, France*

YVES.GRANDVALET@UTC.FR

Scikit-learn Workflow

See Lecture02_Example01.py and Lecture02_Example02.py

- Python scikit-learn package provides facilities for most popular machine learning algorithms
- The best way to learn how to use scikit-learn functionalities to learn from examples and read user guide
- Workflow:
 - ❖ A typical machine learning task starts with data preparation: For example, loading data from database/files (e.g. using pandas); data cleaning; feature extraction, feature scaling and dimensionality reduction etc; some of these can be done with scikit-learn, some rely on other packages
 - ❖ Following data preparation there will be a step to define a machine learning model, for example, linear regression etc.
 - ❖ Scikit-learn introduces the concept of pipeline that chains all the steps in a linear sequence and automates the cross-validation
 - ❖ Read examples here <https://machinelearningmastery.com/automate-machine-learning-workflows-pipelines-python-scikit-learn/>