

Assignment Project Exam Help

Add WeChat powcoder

Assignment Project Exam Help

Linear Regression

Add WeChat powcoder

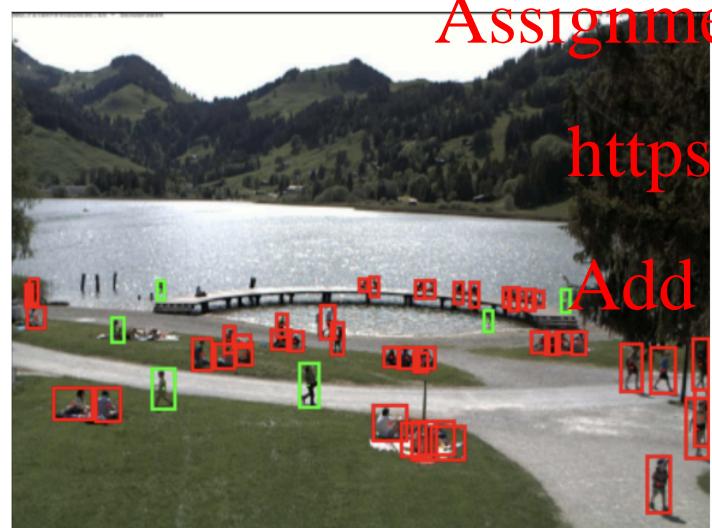
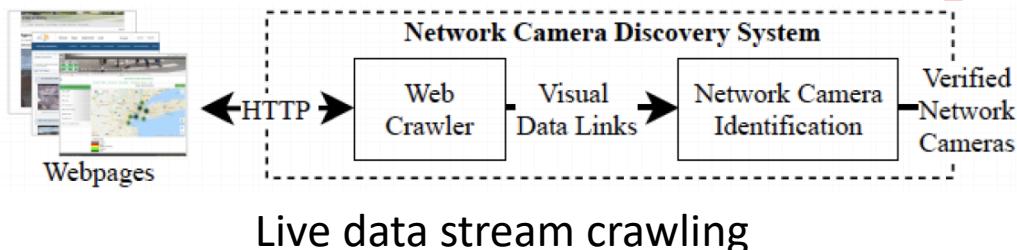
Liang Zheng

Australian National University

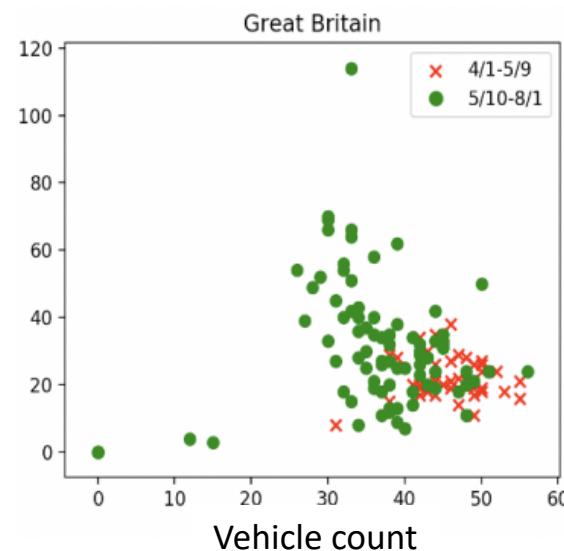
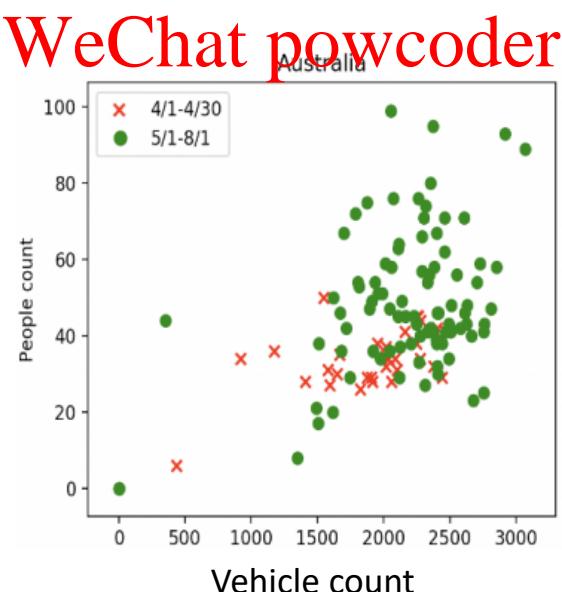
liang.zheng@anu.edu.au

Analyzing Worldwide Social Distancing through Large-Scale Computer Vision, Ghodgaonkar et al., Arxiv 2020

Add WeChat powcoder



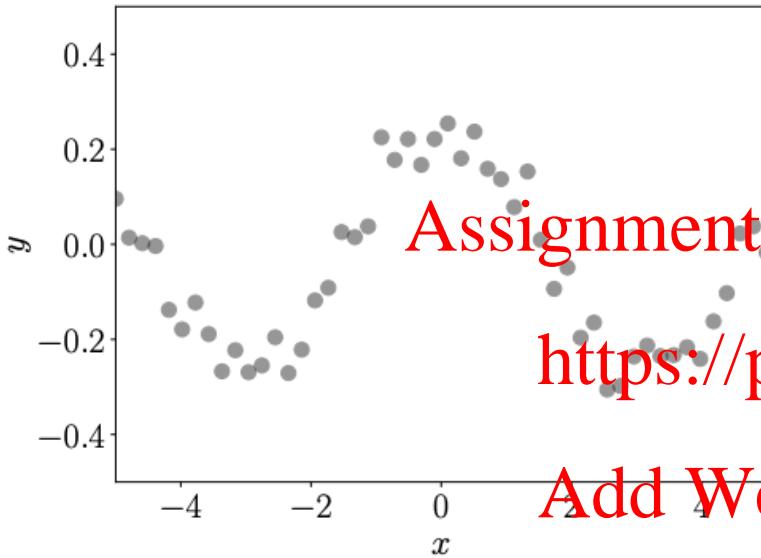
A park in Schwarzsee, Switzerland,
on May 22, 2020



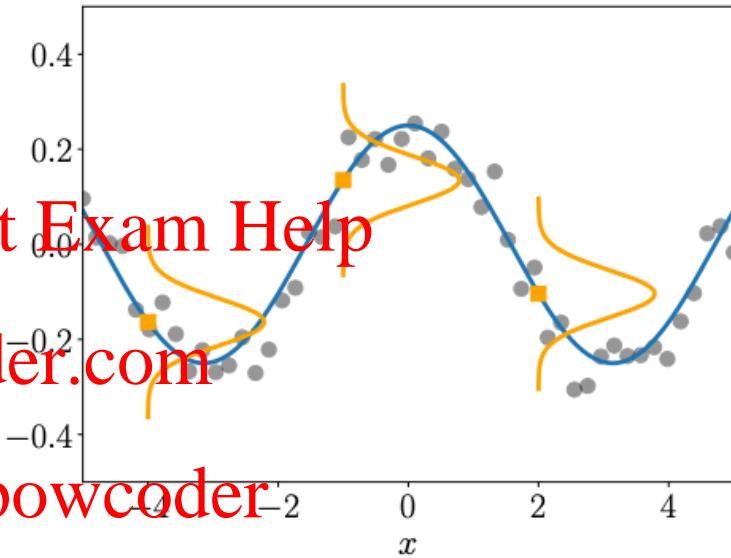
Regression problem

Add WeChat powcoder

- Regression is a fundamental problem in machine learning



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



Regression problem: observed noisy function values from which we wish to infer the underlying function that generated the data.

Regression solution: possible function that could have generated the data (blue) with indication of the measurement noise of the function value at the corresponding inputs (orange distributions)

Assignment Project Exam Help

Add WeChat powcoder

Assignment Project Exam Help

Linear Regression from a statistical point of view

<https://powcoder.com>

Add WeChat powcoder

Problem Formulation

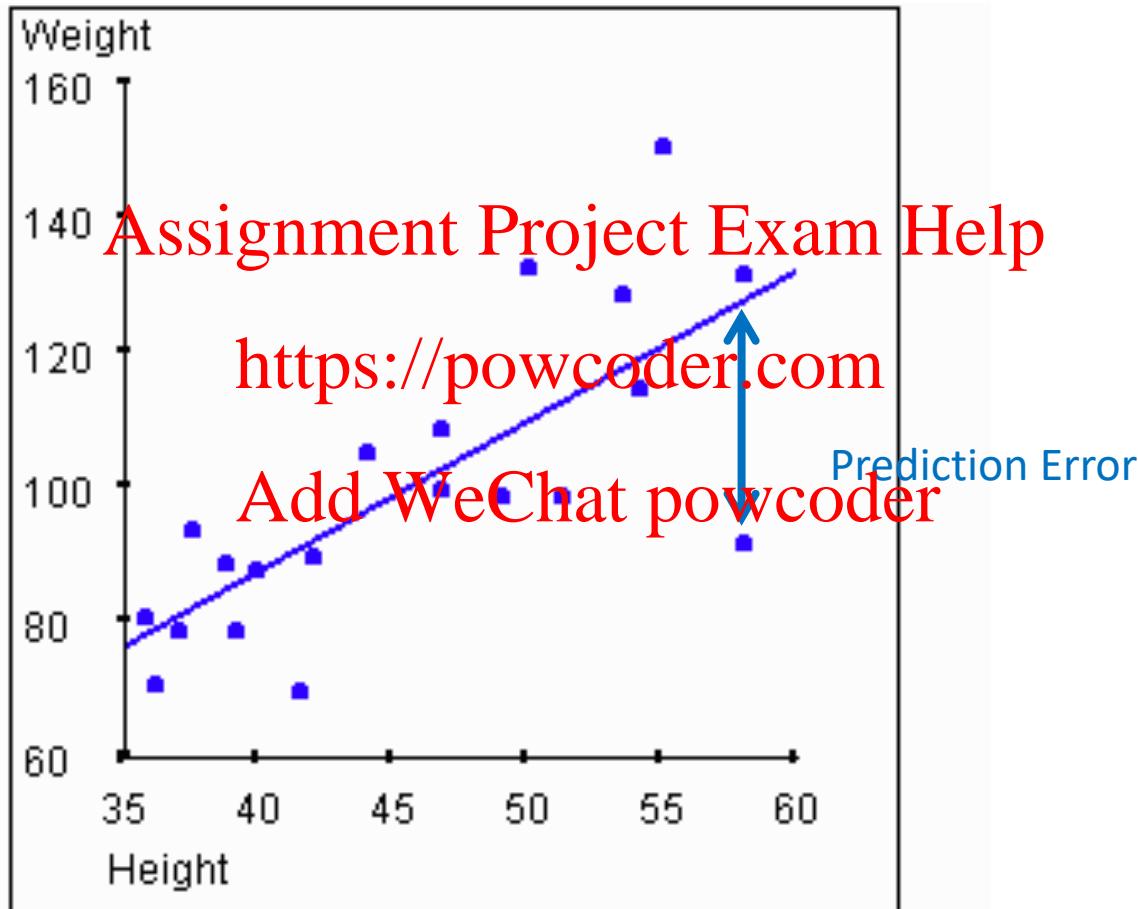
Add WeChat powcoder

- Regression belongs to supervised learning
- **Task.** Find function $f: \mathbb{R}^D \rightarrow \mathbb{R}$ such that $y \approx f(x; \theta)$
- **Experience.** Training set $\mathcal{D} = \{(x_1, y_1), \dots, (x_N, y_N)\}$, consisting of N inputs $x_N \in \mathbb{R}^D$ and corresponding observations/targets $y_n \in \mathbb{R}, n = 1, \dots, N$
- **Performance.** Estimated empirical risk on the test data

$R_{\text{emp}}(f, X_{\text{test}}, y_{\text{test}})$
Add WeChat powcoder

Working Example

Add WeChat powcoder



Loss Function for Training (review)

Add WeChat powcoder

- Under the i.i.d assumption, the empirical mean is a good estimate of the population mean.
- We can use the empirical mean of the loss on the training data
- Given a training set $\{(x_1, y_1), \dots, (x_N, y_N)\}$, we use the notation of an example matrix

$$\mathbf{X} := [x_1, \dots, x_N]^T \in \mathbb{R}^{N \times D}$$

and a label vector

$$\mathbf{y} = [y_1, \dots, y_N]^T \in \mathbb{R}^N$$

- The average loss is given by

$$\mathbf{R}_{emp}(f, \mathbf{X}, \mathbf{y}) = \frac{1}{N} \sum_{n=1}^N \ell(y_n, \hat{y}_n)$$

where $\hat{y}_n = f(\mathbf{x}_n, \boldsymbol{\theta})$. The above equation is called the **empirical risk**. The learning strategy is called **empirical risk minimization**.

Least Squares Linear Regression (revision)

- We use the squared loss function

$$\ell(y_n, \hat{y}_n) = (y_n - \hat{y}_n)^2$$

- The empirical risk becomes,

$$R_{emp}(f, X, y) = \frac{1}{N} \sum_{n=1}^N \ell(y_n, \hat{y}_n) = \frac{1}{N} \sum_{n=1}^N (y_n - f(x_n, \theta))^2$$

Assignment Project Exam Help

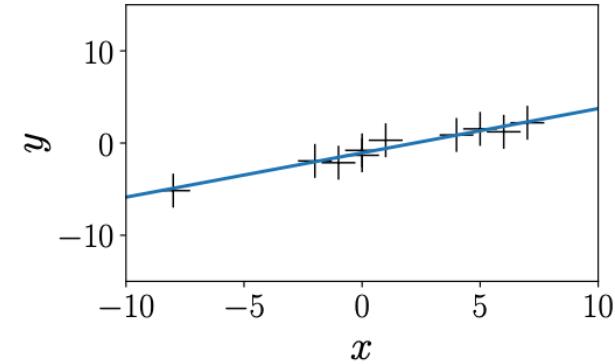
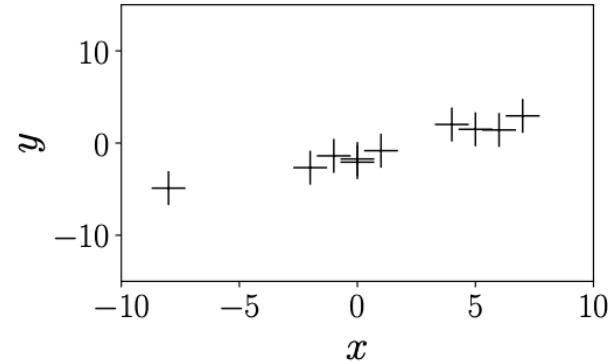
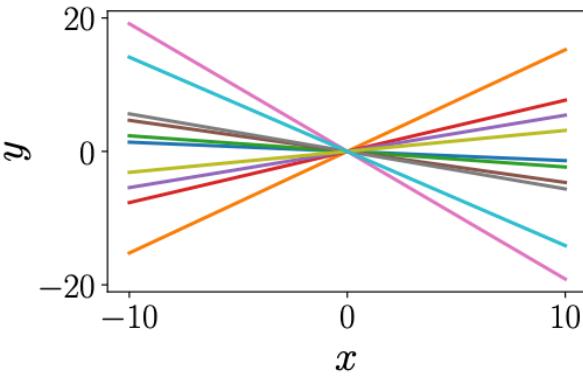
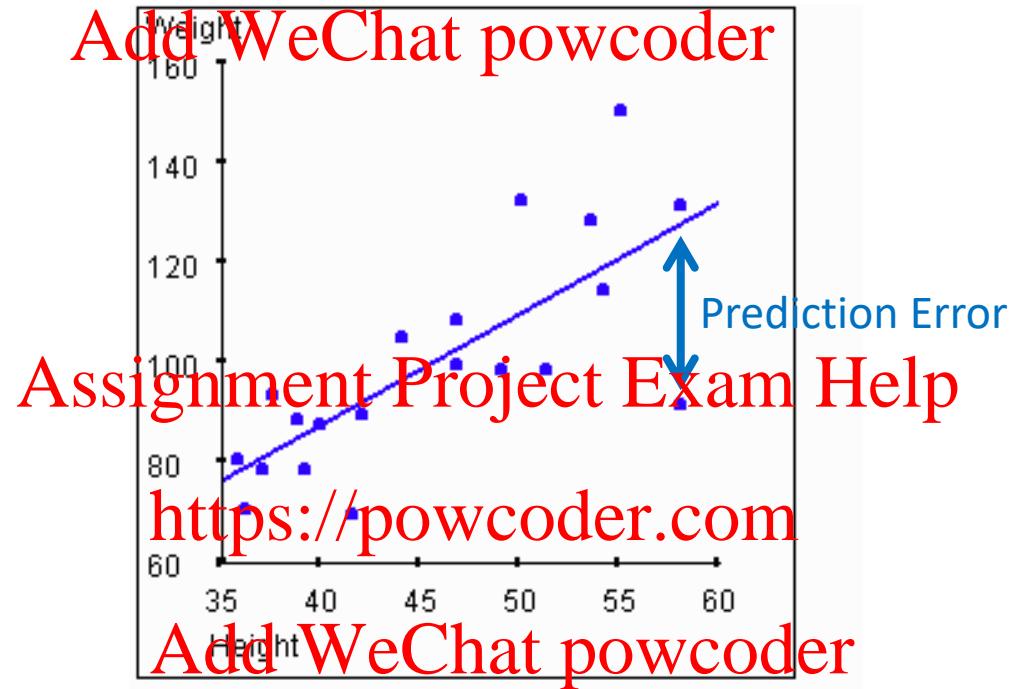
Using the linear predictor $f(x_n, \theta) = \theta^T x_n$, we obtain the empirical risk as

$$R_{emp}(f, X, y) = \frac{1}{N} \sum_{n=1}^N (y_n - f(x_n, \theta))^2$$

- This equation can be equivalently expressed in matrix form

$$\mathcal{L}(\theta) := R_{emp}(f, X, y) = \frac{1}{N} \|y - X\theta\|^2 = \frac{1}{N} (y - X\theta)^T (y - X\theta)$$

Working Example



A closed-form analytic solution

Add WeChat powcoder

- Loss function (mean square error)

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{N} \|y - X\boldsymbol{\theta}\|^2 = \frac{1}{N} (y - X\boldsymbol{\theta})^T (y - X\boldsymbol{\theta})$$

- We calculate the gradient of \mathcal{L} with respect to the parameters $\boldsymbol{\theta}$ as

$$\begin{aligned} \frac{d\mathcal{L}}{d\boldsymbol{\theta}} &= \frac{1}{N} \left(\frac{d}{d\boldsymbol{\theta}} (y - X\boldsymbol{\theta})^T (y - X\boldsymbol{\theta}) \right) \\ &= \frac{1}{N} \frac{d}{d\boldsymbol{\theta}} (y^T y - \boldsymbol{\theta}^T X^T y - y^T X \boldsymbol{\theta} + \boldsymbol{\theta}^T X^T X \boldsymbol{\theta}) \\ &= \frac{1}{N} \frac{d}{d\boldsymbol{\theta}} (y^T y - 2y^T X \boldsymbol{\theta} + \boldsymbol{\theta}^T X^T X \boldsymbol{\theta}) \\ &= \frac{1}{N} (-2y^T X + 2\boldsymbol{\theta}^T X^T X) \in \mathbb{R}^{1 \times D} \end{aligned}$$

- The minimum is attained when the gradient is zero.

$$\begin{aligned} \frac{d\mathcal{L}}{d\boldsymbol{\theta}} = \mathbf{0}^T &\Leftrightarrow \boldsymbol{\theta}^T X^T X = y^T X \\ &\Leftrightarrow \boldsymbol{\theta}^T = y^T X (X^T X)^{-1} \\ &\Leftrightarrow \boldsymbol{\theta} = (X^T X)^{-1} X^T y \end{aligned}$$

Assignment Project Exam Help

Add WeChat powcoder

.

Assignment Project Exam Help

Linear Regression from a probability point of view

<https://powcoder.com>

Add WeChat powcoder

Assignment Project Exam Help

9.1 Problem formulation

Add WeChat powcoder

- A linear regression problem with likelihood function

$$p(y|x, \theta) = \mathcal{N}(y|x^T\theta, \sigma^2)$$

$$\Leftrightarrow y = x^T\theta + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2)$$

Assignment Project Exam Help

- The class of functions described by $y = x^T\theta + \varepsilon$ are straight lines that pass through the origin.
- The likelihood in $p(y|x, \theta) = \mathcal{N}(y|x^T\theta, \sigma^2)$ is the probability density function of y evaluated at $x^T\theta$.
- The only source of uncertainty originates from the observation noise ε (we assume ε is known)

Assignment Project Exam Help

9.2 Parameter estimation

Add WeChat powcoder

- Given a training set $\{(x_1, y_1), \dots, (x_N, y_N)\}$ where $x_n \in \mathbb{R}^D$ and corresponding observations/targets $y_n \in \mathbb{R}$
- y_i and y_j are conditionally independent given their inputs x_i and x_j
- we use the notation of set of training input

$$\mathcal{X} := \{x_1, \dots, x_N\}$$

and set of corresponding targets

$$y := \{y_1, \dots, y_N\}$$

Add WeChat powcoder

- The likelihood factorizes

$$\begin{aligned} p(y|\mathcal{X}, \boldsymbol{\theta}) &= p(y_1, \dots, y_N | x_1, \dots, x_N, \boldsymbol{\theta}) \\ \prod_{n=1}^N p(y_n | x_n, \boldsymbol{\theta}) &= \prod_{n=1}^N \mathcal{N}(y_n | x_n^T \boldsymbol{\theta}, \sigma^2) \end{aligned}$$

9.2.1 Maximum Likelihood Estimation

Add WeChat powcoder

- Intuitively, maximizing the likelihood means maximizing the predictive distribution of the training data given the model parameters. We obtain the maximum likelihood parameters as

$$\boldsymbol{\theta}_{\text{ML}} = \arg \max p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}) = \arg \max \prod_{n=1}^N p(y_n|x_n, \boldsymbol{\theta})$$

- Instead of maximizing the likelihood directly, we apply the log-transformation to the likelihood function and minimize the negative log-likelihood,

$$-\log p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}) \stackrel{\text{Add WeChat powcoder}}{=} -\log \prod_{n=1}^N p(y_n|x_n, \boldsymbol{\theta}) \stackrel{\text{Add WeChat powcoder}}{=} -\sum_{n=1}^N \log p(y_n|x_n, \boldsymbol{\theta})$$

- In the linear regression model $y = \mathbf{x}^T \boldsymbol{\theta} + \varepsilon$, $\varepsilon \sim \mathcal{N}(0, \sigma^2)$, the likelihood is Gaussian (due to the noise term ε), such that we arrive at

$$\log p(y_n|x_n, \boldsymbol{\theta}) = -\frac{1}{2\sigma^2} (y_n - \mathbf{x}^T \boldsymbol{\theta})^2 + \text{const}$$

Assignment Project Exam Help Maximum Likelihood Estimation

Add WeChat powcoder

- Negative log-likelihood:

$$-\log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) = -\log \prod_{n=1}^N p(y_n|x_n, \boldsymbol{\theta}) = -\sum_{n=1}^N \log p(y_n|x_n, \boldsymbol{\theta})$$

- We have

Assignment Project Exam Help

$$\log p(y_n|x_n, \boldsymbol{\theta}) = -\frac{1}{2\sigma^2} (y_n - \mathbf{x}_n^\top \boldsymbol{\theta})^2 + \text{const}$$

<https://powcoder.com>

- By substitution and discarding const terms, we have

$$\mathcal{L}(\boldsymbol{\theta}) := \frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \mathbf{x}_n^\top \boldsymbol{\theta})^2 = \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) = \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|^2$$

where \mathbf{X} is the example feature matrix

$$\mathbf{X} := [\mathbf{x}_1, \dots, \mathbf{x}_N]^\top \in \mathbb{R}^{N \times D}$$

and \mathbf{y} the label vector

$$\mathbf{y} = [y_1, \dots, y_N]^\top \in \mathbb{R}^N$$

Linear regression from Stats and Prob views

- From the Stats view:

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N \|y_i - \mathbf{x}_i^\top \boldsymbol{\theta}\|^2 = \frac{1}{N} (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})$$

- From the Prob view:

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) = \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|^2$$

- Both N and σ are known, so these two loss functions are equivalent and have the same closed-form solution

<https://powcoder.com>

Issues.

Add WeChat powcoder

- Need $\mathbf{X}^\top \mathbf{X}$ to be invertible

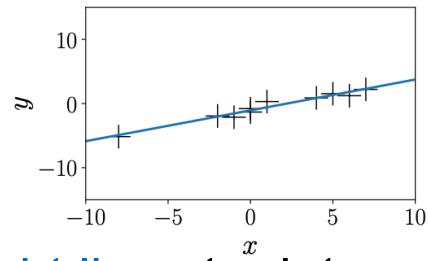
- Feature vectors $\mathbf{x}_1, \dots, \mathbf{x}_N$ must span \mathbb{R}^D
- Must have more data than features, $N \geq D$
- Use regularization if $\mathbf{X}^\top \mathbf{X}$ is not invertible

- What if $\mathbf{X}^\top \mathbf{X} \in \mathbb{R}^{d \times d}$ is a large matrix?

- Takes long time to invert
- Use stochastic gradient descent if $\mathbf{X}^\top \mathbf{X}$ is too large

Linear regression with features

Add WeChat powcoder



- So far, we have discussed linear regression which fits **straight lines** to data.
- However straight lines are not sufficiently expressive when fitting more complex data
- Fortunately, “linear regression” only refers to “linear in the parameters”
- we can perform an arbitrary nonlinear transformation $\phi(x)$ of the inputs x and then linearly combine the components of this transformation.

$$p(y|x, \theta) = \mathcal{N}(y|\phi^T(x)\theta, \sigma^2)$$

Add WeChat powcoder

$$\Leftrightarrow y = \phi^T(x)\theta + \epsilon = \sum_{k=0}^{n-1} \theta_k \phi_k(x) + \epsilon$$

where $\phi: \mathbb{R}^D \rightarrow \mathbb{R}^K$ is a (nonlinear) transformation of the inputs x and $\phi_k: \mathbb{R}^D \rightarrow \mathbb{R}$ is the k th component of the **feature vector** ϕ . Note that the model parameters θ still appear only linearly

Example - Polynomial Regression

Add WeChat powcoder

- We are concerned with a regression problem $y = \phi^T(x)\theta + \epsilon$, where $x \in \mathbb{R}$ and $\theta \in \mathbb{R}^K$. A transformation that is often used in this context is

$$\phi(x) = \begin{bmatrix} \phi_0(x) \\ \phi_1(x) \\ \vdots \\ \phi_{K-1}(x) \end{bmatrix} = \begin{bmatrix} 1 \\ x \\ x^2 \\ x^3 \\ \vdots \\ x^{K-1} \end{bmatrix} \in \mathbb{R}^K$$

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

- We create a K -dimensional feature from a 1-dimensional input
- With these features, we can model polynomials of degree $\leq K - 1$ within the framework of linear regression: A polynomial of degree $K - 1$ is

$$f(x) = \sum_{k=0}^{K-1} \theta_k x^k = \phi^T(x)\theta$$

where $\theta = [\theta_0, \dots, \theta_{K-1}]^T \in \mathbb{R}^K$ contains the (linear) parameters.

Assignment Project Exam Help

Check your understanding

Add WeChat powcoder

- Linear regression belongs to supervised learning.
- If we use a polynomial feature expression, linear regression is not longer *linear*.
- Linear regression is *linear* in that it uses a linear modeling of the input feature.
- If we have limited data but very high dimensional data, we will likely have overfitting, but it is still possible to use the equation $\theta = (X^T X)^{-1} X^T y$ to calculate the optimal parameters.
- You are running a restaurant and want to use linear regression to estimate your daily income. Your prior knowledge tells you that the income is cyclical with different seasons. Which feature is best for your model?
 - Polynomials of days
 - Cosine function of days
 - Exponential function of days
 - Days

Synthetic-to-Real Unsupervised Domain Adaptation for Scene Text Detection in the Wild. Wu et al, Arxiv 2020



ICDAR2013



ICDAR2015



ICDAR2017 MLI



SynthText



VISD

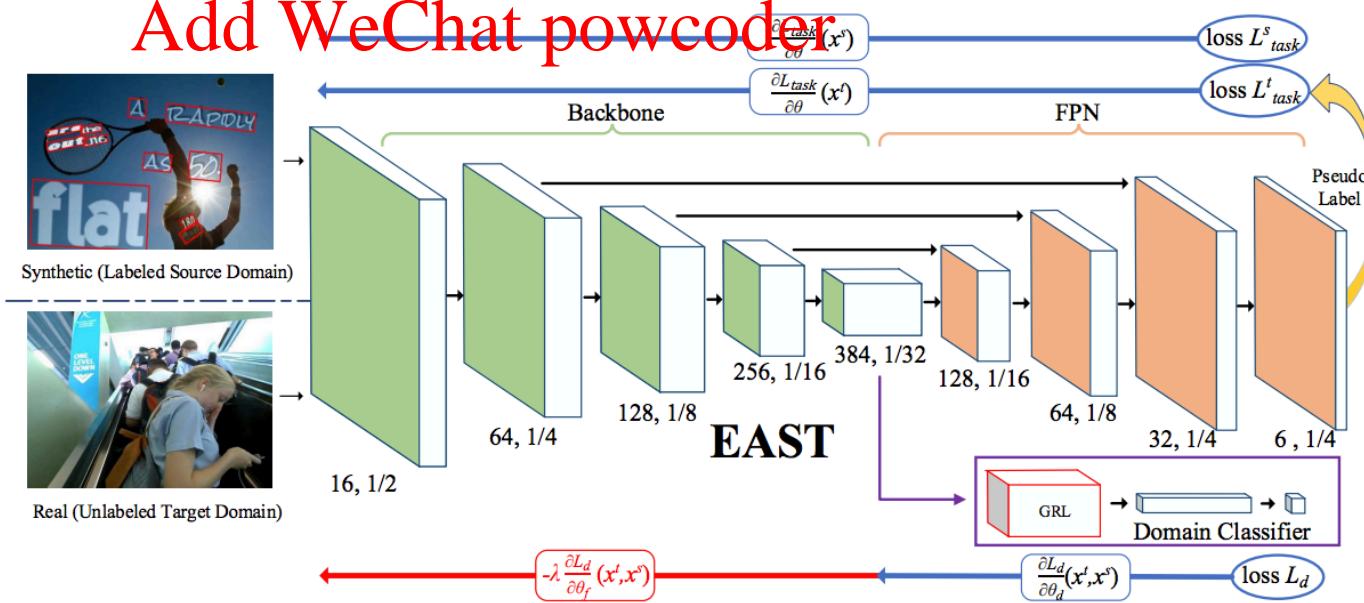


UnrealText

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

gradient reversal layer (GRL)

Add WeChat powcoder



Linear regression with features

- We consider training inputs $\mathbf{x}_n \in \mathbb{R}^D$ and targets $y_n \in \mathbb{R}, n = 1, \dots, N$ and define the **feature matrix (design matrix)** as

$$\Phi := \begin{bmatrix} \phi^T(\mathbf{x}_1) \\ \vdots \\ \phi^T(\mathbf{x}_N) \end{bmatrix} = \begin{bmatrix} \phi_0(\mathbf{x}_1) & \cdots & \phi_{K-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \cdots & \phi_{K-1}(\mathbf{x}_2) \\ \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \cdots & \phi_{K-1}(\mathbf{x}_N) \end{bmatrix} \in \mathbb{R}^{N \times K}$$

<https://powcoder.com>

where $\Phi_{ij} = \phi_j(\mathbf{x}_i)$ and $\phi_j: \mathbb{R}^D \rightarrow \mathbb{R}$

Add WeChat powcoder

- With this feature matrix Φ , the negative log-likelihood for the linear regression model can be written as

$$-\log p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}) = \frac{1}{2\sigma^2} (\mathbf{y} - \Phi\boldsymbol{\theta})^T (\mathbf{y} - \Phi\boldsymbol{\theta}) + \text{const.}$$

Linear regression with features

Add WeChat powcoder

- With this feature matrix Φ , the negative log-likelihood for the linear regression model can be written as

$$-\log p(y|\mathcal{X}, \theta) = \frac{1}{2\sigma^2} (\mathbf{y} - \Phi\theta)^T(\mathbf{y} - \Phi\theta) + \text{const} \quad (1)$$

- Recall we have

$$\mathcal{L}(\theta) := \frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \mathbf{x}_n^T \theta)^2 = \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\theta)^T(\mathbf{y} - \mathbf{X}\theta) = \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\theta\|^2 \quad (2)$$

<https://powcoder.com>

where \mathbf{X} is the matrix of the data points.

- Comparing (1) and (2), we immediately see we just need to replace \mathbf{X} with Φ .
- The solution to (2) was

$$\theta = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- we arrive immediately at the closed form solution to the Linear regression with features problem

$$\theta = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y}$$

Example - Feature Matrix for Second-order Polynomials

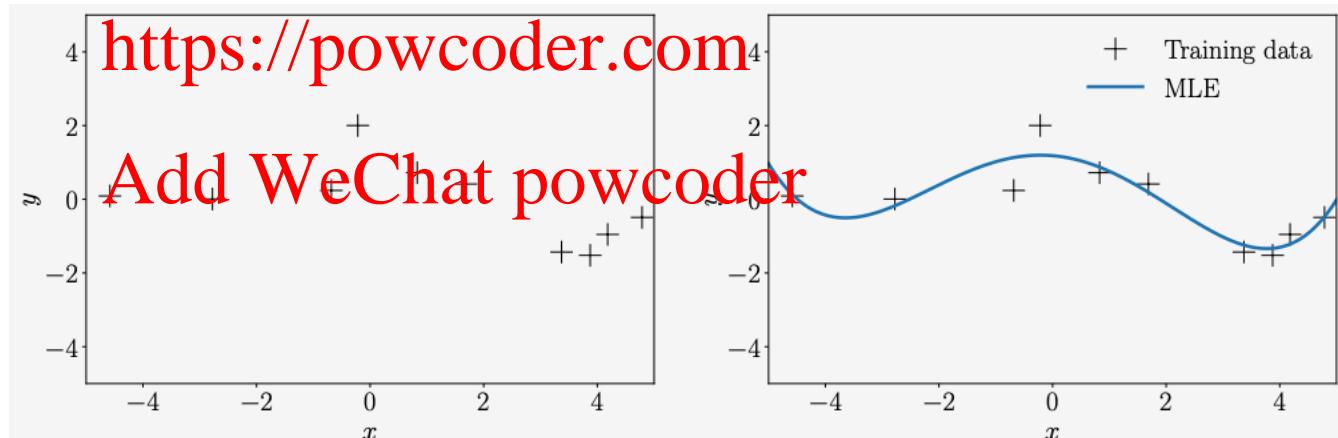
Add WeChat powcoder

- For a second-order polynomial and N training points $x_n \in \mathbb{R}, n = 1, \dots, N$, the feature matrix is

$$\Phi = \begin{bmatrix} 1 & x_1 & x_1^2 \\ 1 & x_2 & x_2^2 \\ \vdots & \vdots & \vdots \\ & & \end{bmatrix}_{N \times 3}$$

Add WeChat powcoder

- Another example:



- The dataset consists of $N = 20$ data pairs (x_n, y_n) , where $x_n \sim \mathcal{U}[-5, 5]$, and $y_n = -\sin(x_n/5) + \cos(x_n) + \varepsilon$, where $\varepsilon \sim \mathcal{N}(0, 0.2^2)$
- In this figure, we fit a polynomial of degree $K = 4$ using maximum likelihood estimation.

9.2.2 Overfitting in Linear Regression

Add WeChat powcoder

- Loss function (mean square error)

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{N} \|y - \Phi\boldsymbol{\theta}\|^2 = \frac{1}{N} \sum_{n=1}^N (y_n - \phi^T(x_n)\boldsymbol{\theta})^2$$

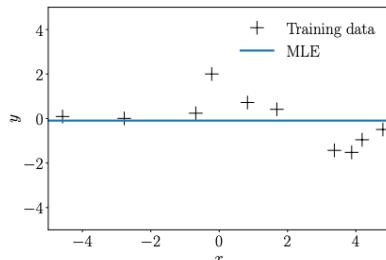
- Root mean square error

$$\mathcal{L}(\boldsymbol{\theta}) = \sqrt{\frac{1}{N} \|y - \Phi\boldsymbol{\theta}\|^2} = \sqrt{\frac{1}{N} \sum_{n=1}^N (y_n - \phi^T(x_n)\boldsymbol{\theta})^2}$$

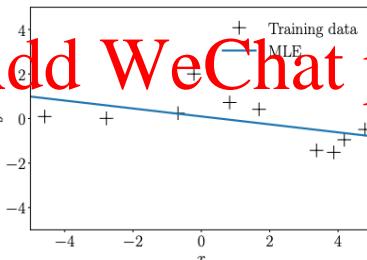
Add WeChat powcoder

- We have N data pairs and want to fit a polynomial model to the data
- We want to determine the polynomial degree M that yields a low training loss and generalizes well (low test error).

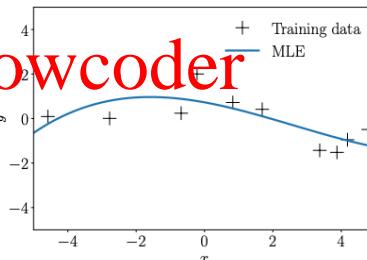
- To determine the best value of M (polynomial degree), we can perform a brute-force search and enumerate all (reasonable) values of M .



(a) $M = 0$

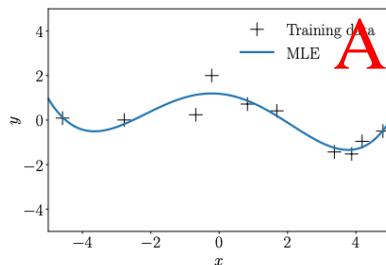


(b) $M = 1$

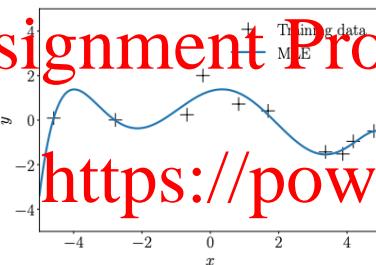


(c) $M = 3$

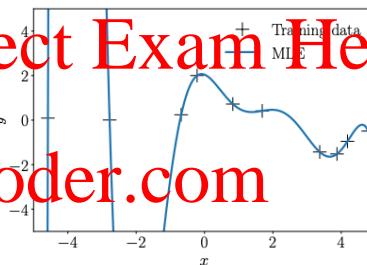
$N = 10$ Data points



(d) $M = 4$



(e) $M = 6$

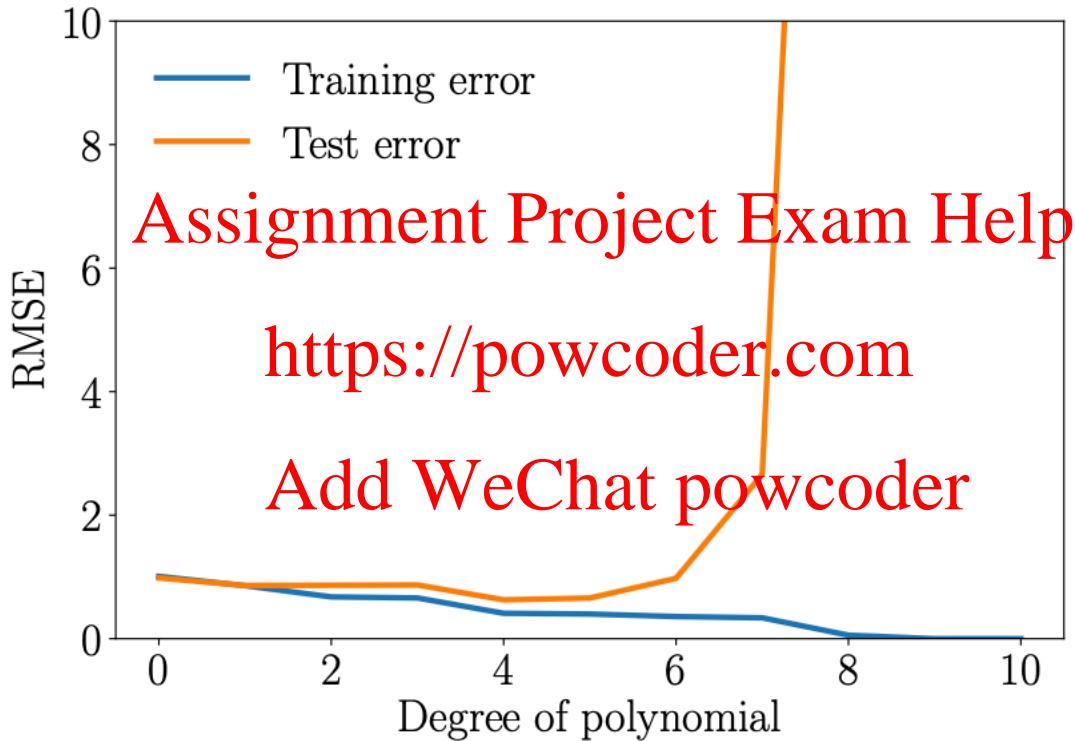


(f) $M = 9$

- When $M = 0, 1$, polynomials fit data poorly
- When $M = 3, 4$, the fits look plausible and smoothly interpolate the data
- When $M = 6, 9$, polynomials fit data better and better. In the extreme case of $M = N - 1 = 9$, the function will pass through every single data point. These high-degree polynomials oscillate wildly and are a poor representation of the underlying function that generated the data, such that we suffer from **overfitting**.

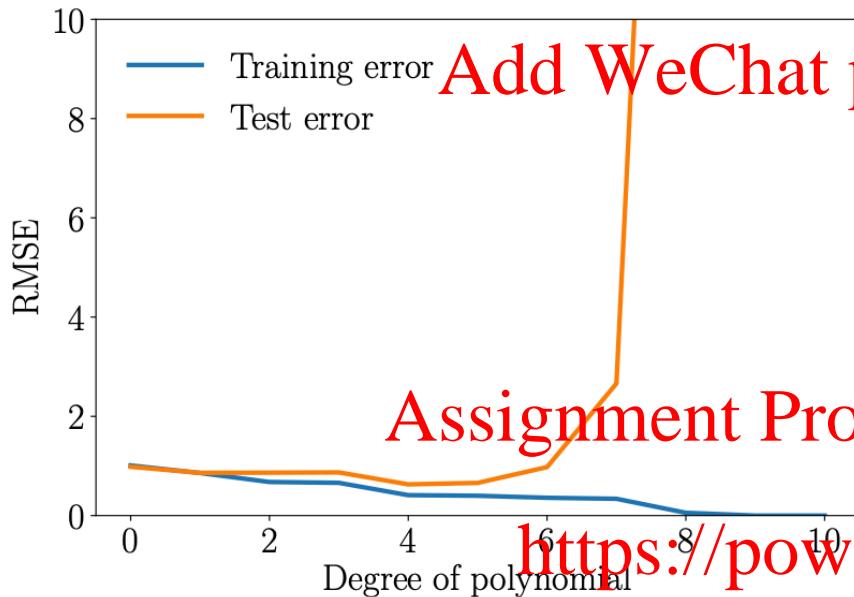
Add WeChat powcoder

- Evaluate whether a model generalizes well



A test set of 200 data points generated by the same procedure as the training set
For each choice of M , calculate the RMSE value.

- Evaluate whether a model generalizes well



A test set of 200 data points generated by the same procedure as the training set

For each choice of M , calculate the RMSE value.

<https://powcoder.com>

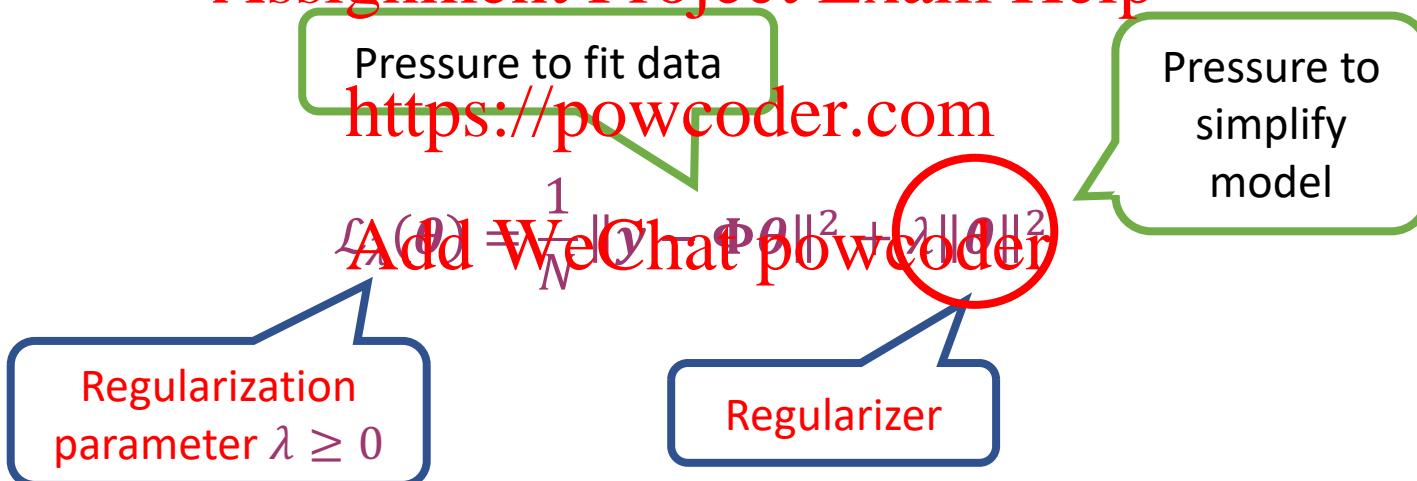
- Initially the test error decreases
- For fourth-order polynomials, the test error is relatively low and stays relatively constant up to degree 5
- From degree 6 onward the test error increases significantly, and high-order polynomials have very bad generalization properties.
- The training error never increases when the degree of the polynomial increases
- the best generalization (smallest test error) is obtained for $M = 4$.

9.2.4 Regularized Least Squares

Add WeChat powcoder

- Overfitting occurs because the model is too complex (θ has too many non-zero entries).
- We want to penalize the amplitude of parameters by regularization
- In regularized least squares, we consider the loss function

Assignment Project Exam Help



- First term: data-fit term; second term: regularizer
- We can use any p -norm $\|\cdot\|_p$
- smaller p leads to sparser solutions, i.e., many parameter values $\theta_d = 0$.

9.2.4 Regularized Least Squares

Add WeChat powcoder

- Loss function

$$\mathcal{L}_\lambda(\theta) = \frac{1}{N} \|y - \Phi\theta\|^2 + \lambda\|\theta\|^2$$

- We calculate the gradient of \mathcal{L} with respect to the parameters θ as

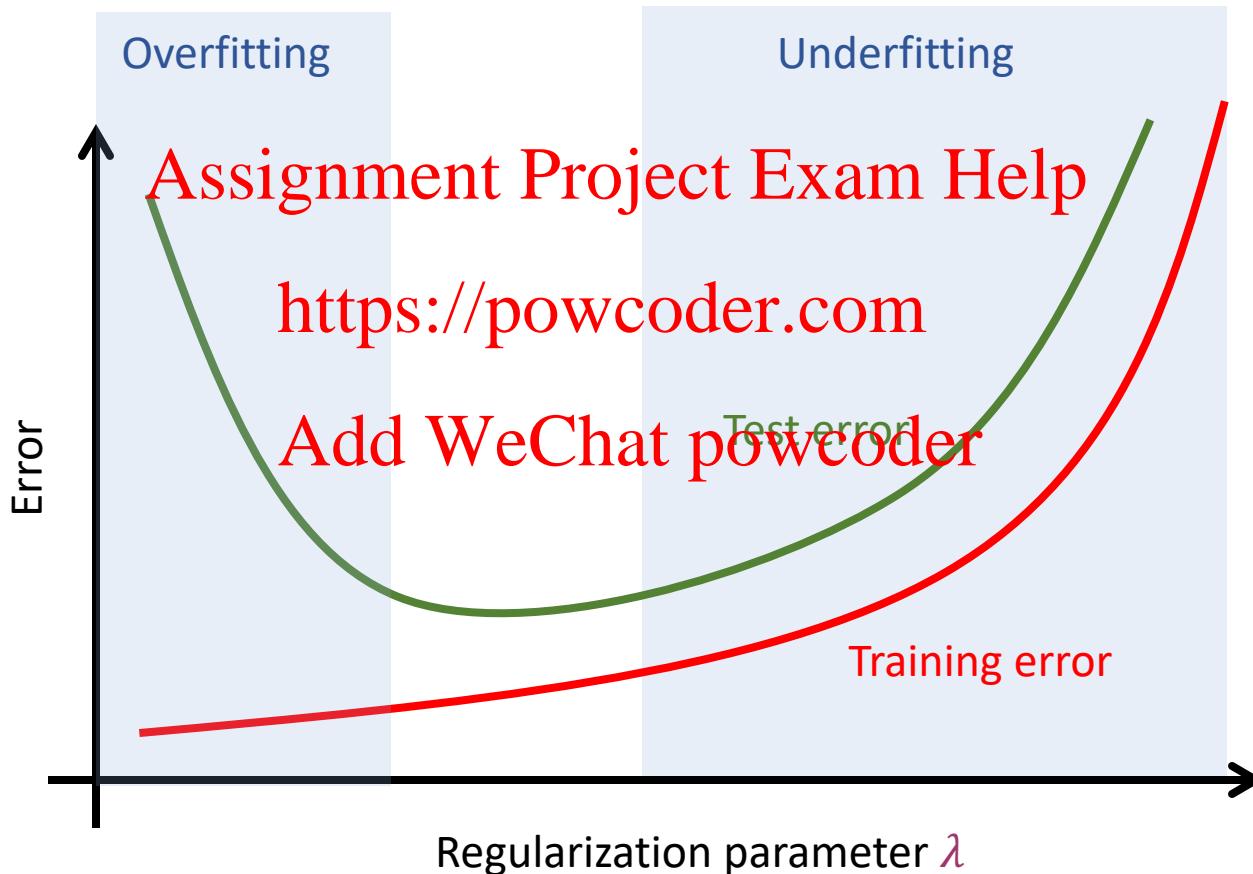
$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial \theta} &= \frac{1}{N} \left((y - \Phi\theta)^T (y - \Phi\theta) + \lambda\|\theta\|^2 \right) \\ &= \frac{1}{N} \left(-2y^T X + 2\theta^T X^T X \right) + \frac{1}{N} (\lambda\|\theta\|^2) \\ &= \frac{1}{N} (AddWeChatpowcoder \in \mathbb{R}^{1 \times D})\end{aligned}$$

- The minimum is attained when the gradient is zero.

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial \theta} = \mathbf{0}^T &\Leftrightarrow 2\lambda\theta^T + \frac{1}{N} 2\theta^T X^T X = \frac{1}{N} 2y^T X \\ &\Leftrightarrow \theta^T = y^T X (N\lambda I + X^T X)^{-1} \\ &\Leftrightarrow \theta = (N\lambda I + X^T X)^{-1} X^T y\end{aligned}$$

Effect of Regularization

Add WeChat powcoder



7. Continuous Optimization – analytic solution

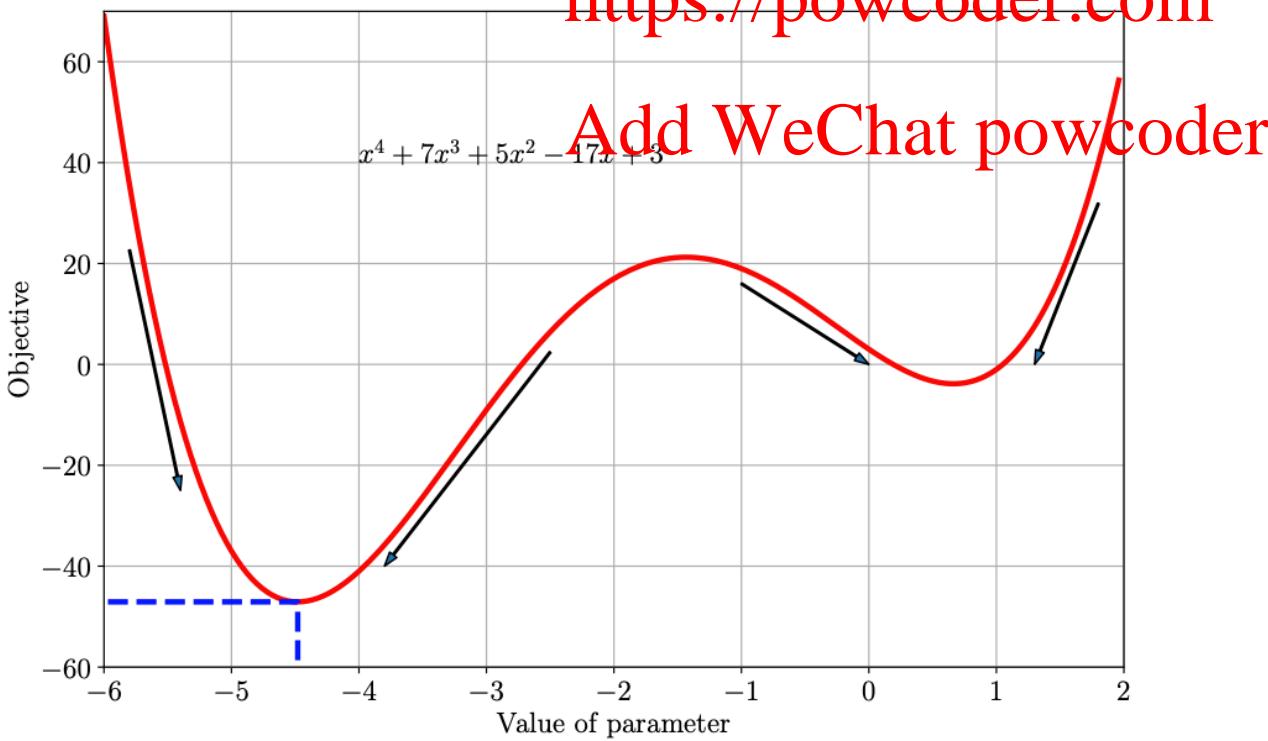
Add WeChat powcoder

- The function below has a global minimum **global minimum** around $x = -4.5$, where the function value is about **-47**
- There exists another **local minimum** around $x = 0.7$
- We can solve for all the stationary points of a function by calculating its **derivative** and setting it to **zero**

Assignment Project Exam Help

$$\ell(x) = x^4 + 7x^3 + 5x^2 - 17x + 3$$

<https://powcoder.com>



Assignment Project Exam Help

7. Continuous Optimization – analytic solution

Add WeChat powcoder

- Considering the following polynomial

$$\ell(x) = x^4 + 7x^3 + 5x^2 - 17x + 3$$

- We calculate the gradient

$$\frac{d\ell(x)}{dx} = 4x^3 + 21x^2 + 10x - 17$$

- A cubic equation, it has in general three solutions when set to zero
 - Two are minimums ($x \approx -4.5, 0.7$)
 - One is maximum ($x \approx -1.4$)
- To determine whether it is minimum or maximum, we calculate the second derivative

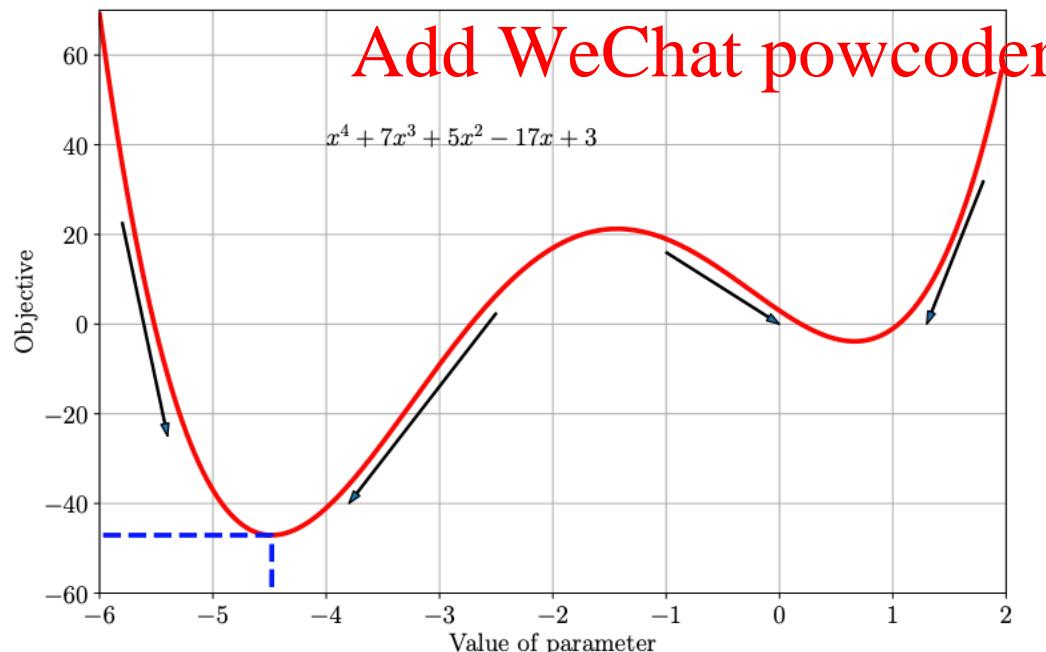
$$\frac{d^2\ell(x)}{dx^2} = 12x^2 + 42x + 10$$

- Substitute our visually estimated values of $x = -4.5, -1.4, 0.7$
- We observe that $x \approx -1.4$ is a maximum, because $\frac{d^2\ell(x)}{dx^2} < 0$
- $x \approx 4.5, 0.7$ are two minimums, because $\frac{d^2\ell(x)}{dx^2} > 0$

7. Continuous Optimization – Motivation for gradient descent

Add WeChat powcoder

- In general, we are unable to find analytic solutions
- We need to start at some value e.g., $x_0 = -10$, and follow the gradient
- The gradient points in the direction of steepest ascent of f .
- When we start from $x_0 < -10$, we know we should go right, but don't know how far we should go (step size)
- When starting from $x_0 > -1$, we might find the local minimum.



7.1 Optimization Using Gradient Descent

Add WeChat powcoder

- Given $f : \mathbb{R}^d \rightarrow \mathbb{R}$, we consider the problem of solving for the minimum of f

$$\min_{\mathbf{x}} f(\mathbf{x})$$

- Gradient descent is a first-order optimization algorithm.
- To find a local minimum of a function using gradient descent, one takes steps proportional to the negative of the gradient of the function at the current point.
- Gradient descent exploits the fact that $f(x_0)$ decreases fastest if one moves from x_0 in the direction of the negative gradient $-(\nabla f(x_0))^T$ of f at x_0 .
- Observation: if

$$\mathbf{x}_1 = \mathbf{x}_0 - \gamma (\nabla f(\mathbf{x}_0))^T$$

- For a small step size $\gamma \geq 0$, then $f(\mathbf{x}_1) \leq f(\mathbf{x}_0)$

Assignment Project Exam Help

7.1 Optimization Using Gradient Descent

Add WeChat powcoder

- A simple gradient descent algorithm:
- If we want to find a local optimum $f(\mathbf{x}_*)$ of a function $f : \mathbb{R}^n \rightarrow \mathbb{R}, \mathbf{x} \mapsto f(\mathbf{x})$, we start with an initial guess \mathbf{x}_0 (parameter we want to optimize), and then iterate according to

Assignment Project Exam Help

$$\mathbf{x}_{i+1} = \mathbf{x}_i - \gamma_i (\nabla f(\mathbf{x}_i))$$

- For suitable step-size γ_i , the sequence $f(\mathbf{x}_0) \geq f(\mathbf{x}_1) \geq \dots$ converges to a local minimum.

Add WeChat powcoder

Example – gradient descent

Add WeChat powcoder

- Consider a quadratic function in two dimensions

$$f \left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \right) = \frac{1}{2} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}^T \begin{bmatrix} 2 & 1 \\ 1 & 20 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} - \begin{bmatrix} 5 \\ 3 \end{bmatrix}^T \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

with gradient

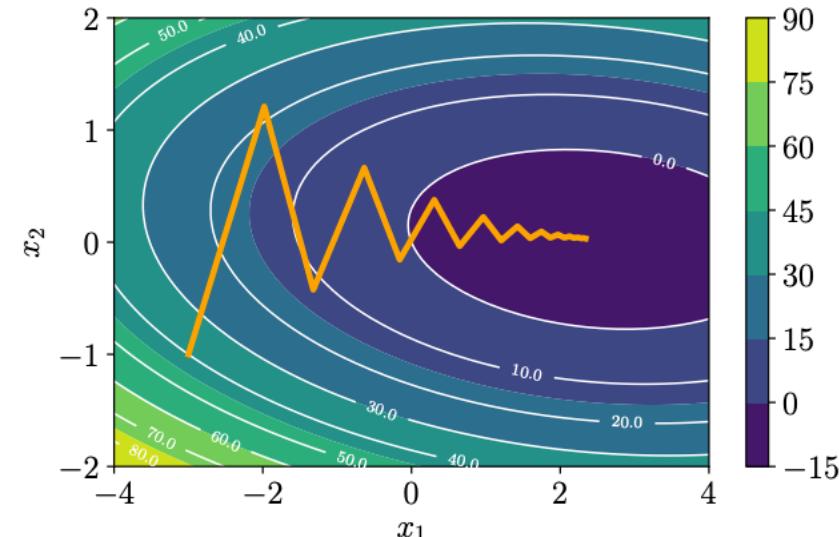
$$\nabla f \left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \right) = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \begin{bmatrix} 2 & 1 \\ 1 & 20 \end{bmatrix} - \begin{bmatrix} 5 \\ 3 \end{bmatrix}$$

<https://powcoder.com>

- Starting at the initial location $x_0 = [-3, 1]^T$, we apply gradient descent (iteratively) to obtain a sequence of estimates that converge to the minimum value

Add WeChat powcoder

- The gradient of x_0 points north and east
- leading to $x_1 = [-1.98, 1.21]^T$
- We can then have $x_2 = [-1.32, -0.42]^T$



Example - Solving a Linear Equation System

Add WeChat powcoder

- When we solve linear equations of the form $\mathbf{A}\mathbf{x} = \mathbf{b}$, in practice we aim to approximately find \mathbf{x}_* that minimizes the squares error

$$\mathcal{L}(\mathbf{x}) = \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2 = (\mathbf{A}\mathbf{x} - \mathbf{b})^T(\mathbf{A}\mathbf{x} - \mathbf{b})$$

- The gradient of $\mathcal{L}(\mathbf{x})$ with respect to \mathbf{x} is

Add WeChat powcoder

- We can use this gradient directly in a gradient descent algorithm

<https://powcoder.com>

Add WeChat powcoder

7.1.3 Stochastic Gradient Descent

Add WeChat powcoder

- Drawback of gradient descent
- Computing the gradient can be very time consuming, why?
- Given N data points $1, 2, \dots, N$, the loss function is a sum of losses \mathcal{L}_n incurred by each example n . That is,

Assignment Project Exam Help

$$\mathcal{L}_N(\boldsymbol{\theta}) = \frac{1}{N} \sum_{n=1}^N \mathcal{L}_n(\boldsymbol{\theta})$$

<https://powcoder.com>

- Example (regression)

$$\mathcal{L}_N(\boldsymbol{\theta}) = \frac{1}{N} \|\mathbf{y} - \Phi \boldsymbol{\theta}\|^2 = \frac{1}{N} \sum_{n=1}^N (y_n - \phi^T(\mathbf{x}_n) \boldsymbol{\theta})^2$$

where $\mathbf{x}_n \in \mathbb{R}^D$ are the training samples, y_n are the labels, and $\boldsymbol{\theta}$ are the parameters of the linear regression model we want to optimize.

- In standard gradient descent, optimization is performed using the **full training set**, by updating the vector of parameters according to

$$\boldsymbol{\theta}_{i+1} = \boldsymbol{\theta}_i - \gamma_i (\nabla \mathcal{L}(\boldsymbol{\theta}_i))^T = \boldsymbol{\theta}_i - \gamma_i \frac{1}{N} \sum_{n=1}^N \nabla \mathcal{L}_n(\boldsymbol{\theta})^T$$

7.1.3 Stochastic Gradient Descent

Add WeChat powcoder

$$\boldsymbol{\theta}_{i+1} = \boldsymbol{\theta}_i - \gamma_i (\nabla \mathcal{L}(\boldsymbol{\theta}_i))^T = \boldsymbol{\theta}_i - \gamma_i \frac{1}{N} \sum_{n=1}^N \nabla \mathcal{L}_n(\boldsymbol{\theta})^T$$

- Evaluating the sum gradient may require expensive evaluations of the gradients from all individual functions.
- Stochastic Gradient Descent
- Estimate the gradient by averaging over a smaller **minibatch** (subset of the training data).

Add WeChat powcoder

$$\nabla \mathcal{L}(\boldsymbol{\theta}) \approx \nabla \mathcal{L}_M(\boldsymbol{\theta}) = \frac{1}{M} \sum_{n \in \mathcal{B}_M} \nabla \mathcal{L}_n(\boldsymbol{\theta})^T$$

where $M \ll N$, and \mathcal{B}_M is a subset of the permuted indices of the samples in the minibatch. $|\mathcal{B}_M| = M$.

- For example, the whole dataset has $N = 4$ samples indexed with $n = 1, 2, 3, 4$. The minibatch contains $M = 2$ samples. \mathcal{B}_M could be $\{1, 3\}, \{2, 4\}, \dots$

Assignment Project Exam Help

Stochastic Gradient Descent

Add WeChat powcoder

1. Initialize θ randomly.
2. Select minibatch B_M of data from the training set at random.
$$\theta \leftarrow \theta - \gamma_i \nabla \mathcal{L}_M(\theta)$$
3. Repeat Step (2) until convergence.

Assignment Project Exam Help

<https://powcoder.com>

Gradient Descent

Add WeChat powcoder

1. Initialize θ randomly.
2. Update (use all the training set calculate the gradient)
$$\theta \leftarrow \theta - \gamma_i \nabla \mathcal{L}_N(\theta)$$
3. Repeat Step (2) until convergence.

Assignment Project Exam Help

7.1.1 Step-size

Add WeChat powcoder

- Choosing a good step-size (learning rate) is important in gradient descent
- If the step-size is too small, gradient descent can be slow
- If the step-size is chosen too large, gradient descent can overshoot, fail to converge, or even diverge

Assignment Project Exam Help

<https://powcoder.com>

- There are several heuristics to adapt the step size
- When the function value increases after a gradient step, the step-size was too large. Undo the step and decrease the step-size
- When the function value decreases the step could have been larger. Try to increase the step-size.
- Typically, we choose a learning rate that starts big and ends small, e.g. $\gamma_i = 1/(i + 1)$

Assignment Project Exam Help

7.1.1 Step-size

Add WeChat powcoder

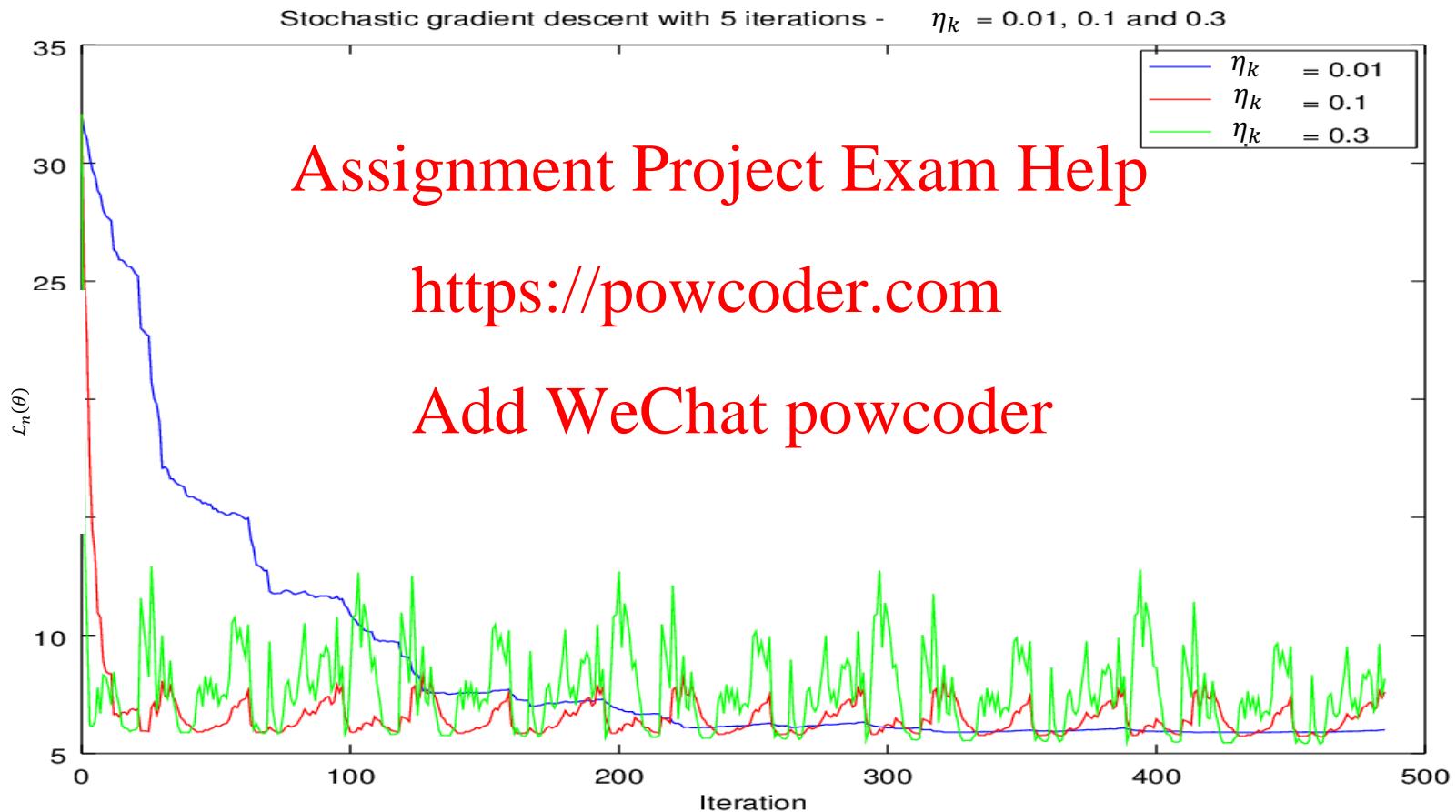
- There has been quite a few papers studying the step size/learning rate.
- [1] Accurate, Large Minibatch SGD: Training ImageNet in 1 Hour
- [2] Large Batch Training of Convolutional Networks
- [3] A Closer Look at Deep Learning Heuristics: Learning Rate Restarts, Warmup and Distillation
- [4] Deep Residual Learning for Image Recognition

Add WeChat powcoder

- A recent practice in deep learning is to start training with a small learning rate, switch to a large learning rate, and then decrease it gradually.

Impact of step size

Add WeChat powcoder



7.1.2 Gradient Descent With Momentum

Add WeChat powcoder

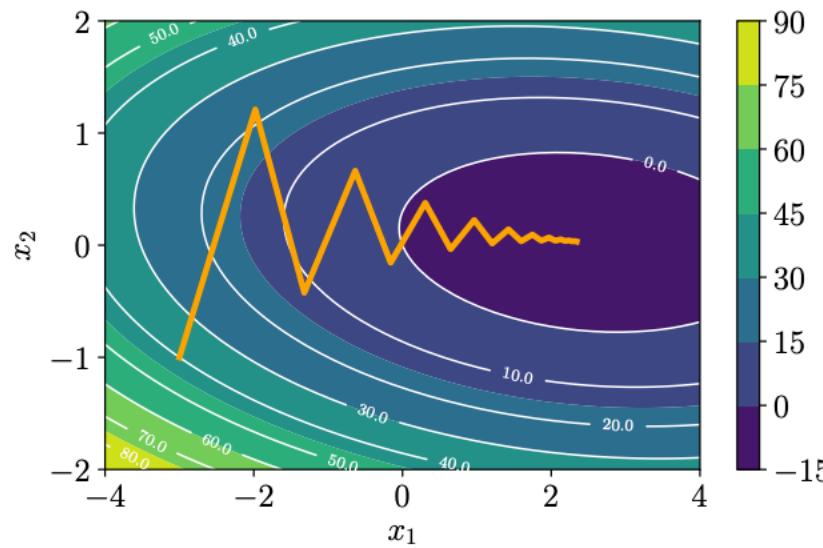
- The convergence of gradient descent may be very slow if the curvature of the optimization surface is such that there are regions that are poorly scaled
- The proposed method to improve convergence is to give gradient descent some memory

Add WeChat powcoder

- Gradient descent with momentum is a method that introduces an additional term to remember what happened in the previous iteration.

<https://powcoder.com>

- This memory dampens oscillations and smooths out the gradient updates
- The idea is to have a gradient update with memory to implement a moving average



7.1.2 Gradient Descent With Momentum

Add WeChat powcoder

- The momentum-based method remembers the update Δx_i at each iteration i and determines the next update as a linear combination of the current and previous gradients

$$x_{i+1} = x_i - \gamma_i (\nabla f(x_i))^T + \alpha \Delta x_i$$

Assignment Project Exam Help
 $\Delta x_i = x_i - x_{i-1} - \gamma_{i-1} (\nabla f(x_{i-1}))^T$

- Where $\alpha \in [-1,1]$. <https://powcoder.com>
- Sometimes we only know the gradient approximately.
- In such cases, the momentum term is useful since it averages out different noisy estimates of the gradient.

Linear regression with gradient descent

Add WeChat powcoder

- Loss function

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{N} \|y - X\boldsymbol{\theta}\|^2 = \frac{1}{N} (y - X\boldsymbol{\theta})^T (y - X\boldsymbol{\theta})$$

- Gradient

Assignment Project Exam Help

$$\nabla \mathcal{L}_N = \frac{1}{N} (-2y^T X + 2\boldsymbol{\theta}^T X^T X)$$

- Gradient descent

<https://powcoder.com>

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \gamma_i \left[\frac{2}{N} \boldsymbol{\theta}^T X^T X - \frac{2}{N} y^T X \right]$$

- Stochastic gradient descent

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \gamma_i \left[\frac{2}{M} \boldsymbol{\theta}^T X^T X - \frac{2}{M} y^T X \right]$$

- Why not use the exact solution? $\boldsymbol{\theta} = (X^T X)^{-1} X^T$

$X^T X \in \mathbb{R}^{D \times D}$ could be a large matrix

- Takes long time to invert

Assignment Project Exam Help

Regularized Linear regression with gradient descent

Add WeChat powcoder

- Loss function

$$\mathcal{L}_\lambda(\boldsymbol{\theta}) = \frac{1}{N} \|y - \mathbf{X}\boldsymbol{\theta}\|^2 + \lambda \|\boldsymbol{\theta}\|^2$$

- Gradient

$$\nabla \mathcal{L}_N = \frac{1}{N} (-2y^T \mathbf{X} + 2\boldsymbol{\theta}^T \mathbf{X}^T \mathbf{X}) + 2\lambda \boldsymbol{\theta}$$

- Gradient descent

<https://powcoder.com>

$$\boldsymbol{\theta} \leftarrow (1 - 2\gamma\lambda)\boldsymbol{\theta} - \gamma_i \left[\frac{2}{N} \boldsymbol{\theta}^T \mathbf{X}^T \mathbf{X} - \frac{2}{N} y^T \mathbf{X} \right]$$

- Stochastic gradient descent

$$\boldsymbol{\theta} \leftarrow (1 - 2\gamma_i\lambda)\boldsymbol{\theta} - \gamma_i \left[\frac{2}{M} \boldsymbol{\theta}^T \mathbf{X}^T \mathbf{X} - \frac{2}{M} y^T \mathbf{X} \right]$$

Check your understanding

Add WeChat powcoder

- When you increase the step size, the training time will be shortened.
- The step size is a hyperparameter.
- “Using a small step size all the time” vs “First use a big step size, then use a small step size”: both methods give you similar optimization results, when starting from the same random seed.
- We can use gradient descent to solve k-means.
<https://powcoder.com>
- SGD gives you a more precise optimization result than standard GD.
- In both GD and SGD, the loss function will decrease with each iteration.
- In linear regression, the training loss $\frac{1}{N} \|\mathbf{y} - \Phi\boldsymbol{\theta}\|^2$ is usually greater when using regularization than not using regularization.
- In $\mathcal{L}_\lambda(\boldsymbol{\theta}) = \frac{1}{N} \|\mathbf{y} - \Phi\boldsymbol{\theta}\|^2 + \lambda \|\boldsymbol{\theta}\|^2$, we can instead put λ in front of $\frac{1}{N} \|\mathbf{y} - \Phi\boldsymbol{\theta}\|^2$, which will give us similar optimization results.