

# Assignment Project Exam Help

## Today

Add WeChat powcoder

- Maximum Likelihood (cont'd)
- Classification

Assignment Project Exam Help

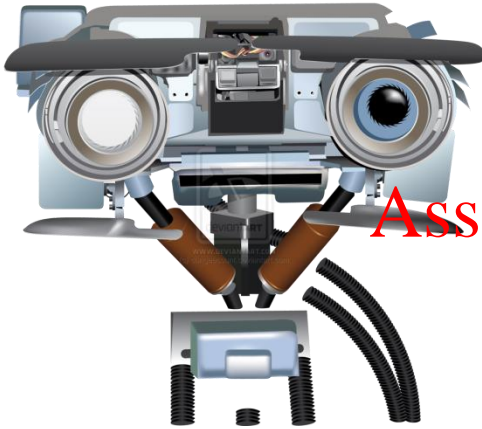
<https://powcoder.com>

Reminder: ps1 due at midnight

Add WeChat powcoder

Assignment Project Exam Help

Add WeChat powcoder



# Maximum Likelihood for Linear Regression

Assignment Project Exam Help

<https://powcoder.com>

---

Add WeChat powcoder

# Maximum likelihood way of estimating model parameters $\theta$

In general, assume data is generated by some distribution

$$U \sim p(U|\theta)$$

Observations ( $n$  i.i.d.)

$$D = \{u^{(1)}, u^{(2)}, \dots, u^{(m)}\}$$

Maximum likelihood estimate

$$\mathcal{L}(D) = \prod_{i=1}^m p(u^{(i)}|\theta)$$

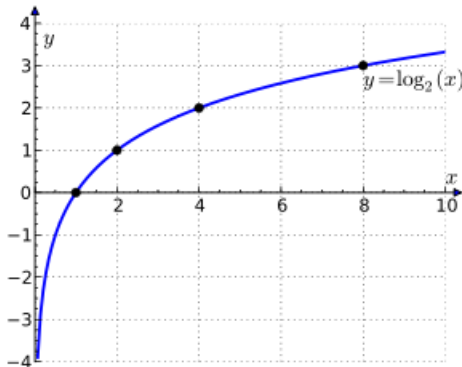
Likelihood

$$\theta_{ML} = \operatorname{argmax}_{\theta} \mathcal{L}(D)$$

Log likelihood

$$= \operatorname{argmax}_{\theta} \sum_{i=1}^m \log p(u^{(i)}|\theta)$$

Note:  $p$  replaces  $h$ !



$\log(f(x))$  is monotonic/increasing, same argmax as  $f(x)$

# Assignment Project Exam Help

## i.i.d. observations

Add WeChat powcoder

- independently identically distributed random variables

Assignment Project Exam Help

- If  $u^i$  are i.i.d. r.v.s, then

<https://powcoder.com>

$$p(u^1, u^2, \dots, u^m) = p(u^1)p(u^2) \dots p(u^m)$$

Add WeChat powcoder

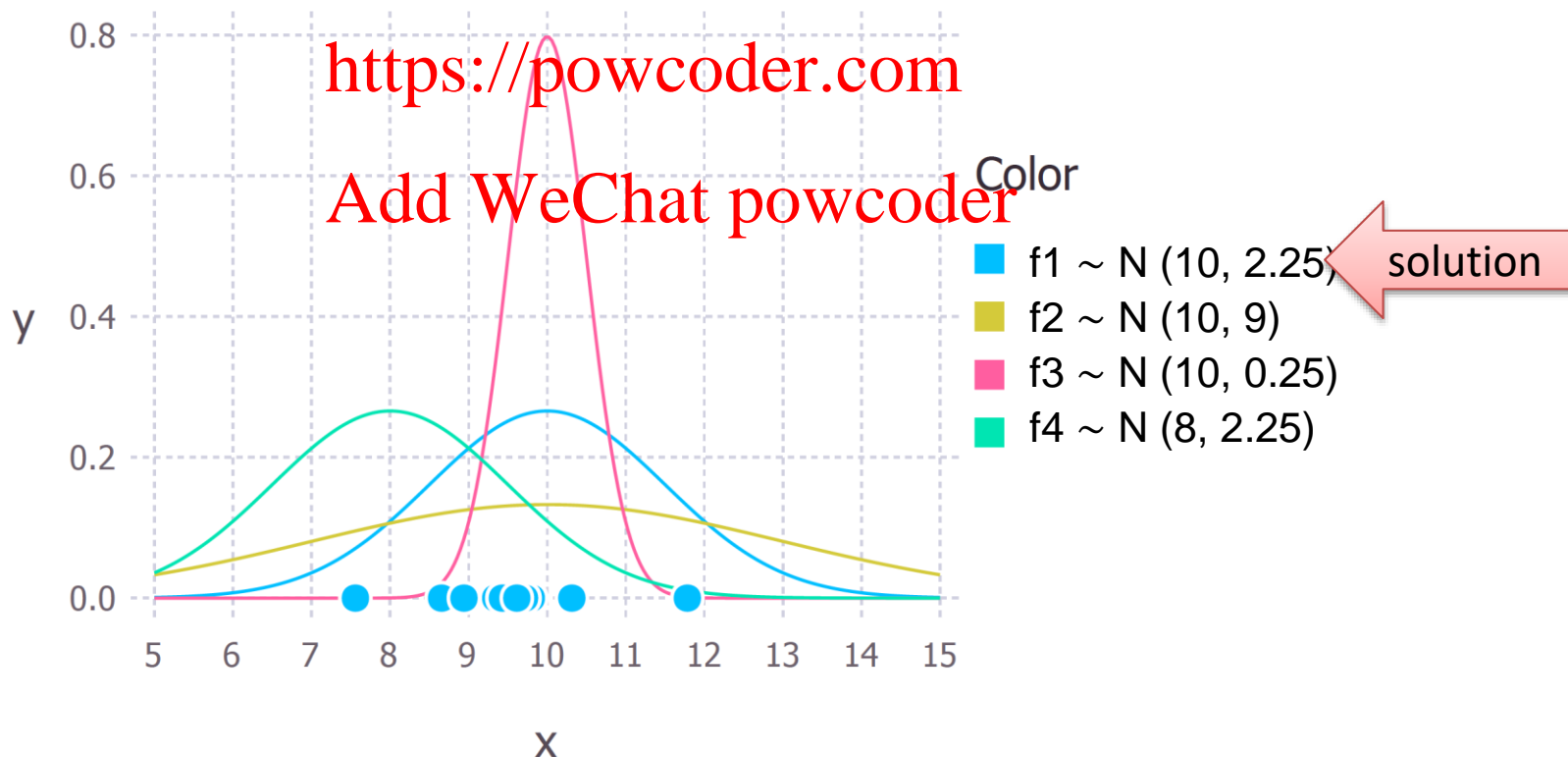
- A reasonable assumption about many datasets, but not always

# Assignment Project Exam Help

## ML: Another example

Add WeChat powcoder

- Observe a dataset of points  $D = \{x^i\}_{i=1:10}$
- Assume  $x$  is generated by Normal distribution,  $x \sim N(x|\mu, \sigma)$
- Find parameters  $\theta_{ML} = [\mu, \sigma]$  that maximize  $\prod_{i=1}^{10} N(x^i|\mu, \sigma)$



# ML for Linear Regression

Assignment Project Exam Help

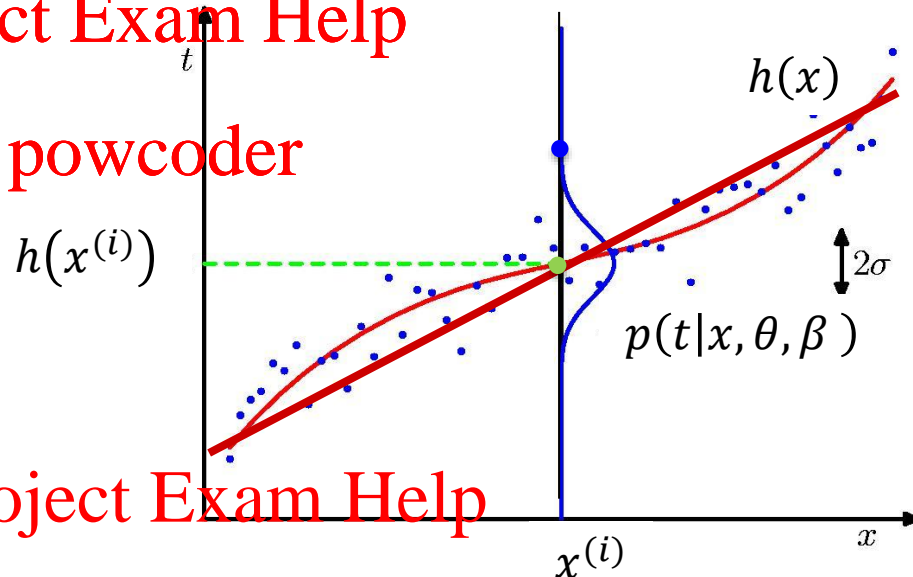
Add WeChat powcoder

Assume:

$$t = y + \epsilon = h(x) + \epsilon$$

$$\text{Noise } \epsilon \sim N(\epsilon|0, \beta^{-1}),$$

$$\text{where } \beta = \frac{1}{\sigma^2}$$



Assignment Project Exam Help

*we don't get to see  $y$  only  $t$*

Add WeChat powcoder

$$t_i \quad h(x^{(i)})$$

# ML for Linear Regression

Assignment Project Exam Help

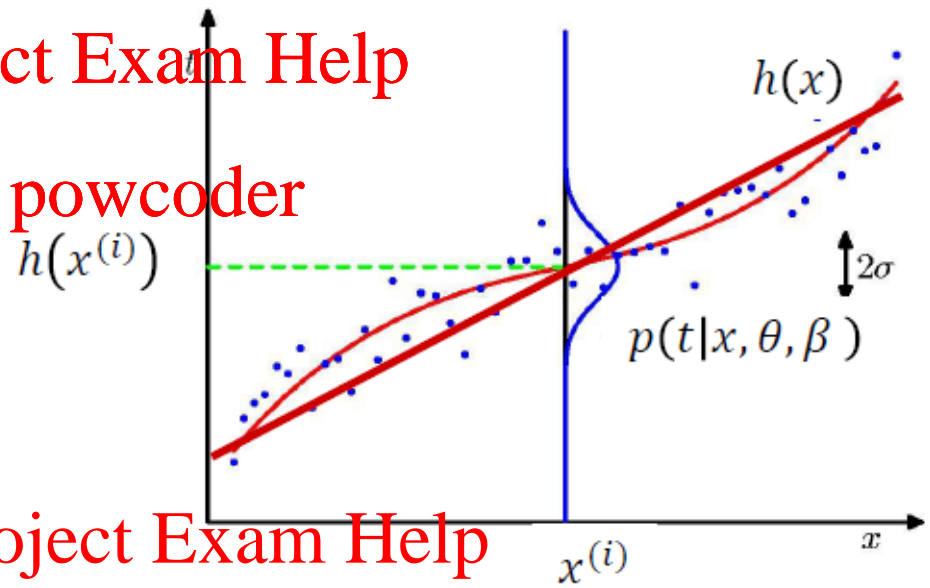
Assume:

$$t = y + \epsilon = h(x) + \epsilon$$

$$\text{Noise } \epsilon \sim N(\epsilon|0, \beta^{-1}),$$

$$\text{where } \beta = \frac{1}{\sigma^2}$$

Add WeChat powcoder



Assignment Project Exam Help

$$p(t|x, \theta, \beta) = N(t|h(x), \beta^{-1})$$

<https://powcoder.com>

Probability of one data point

Add WeChat powcoder

$$p(\mathbf{t}|\mathbf{x}, \theta, \beta) = \prod_{i=1}^m N(t^{(i)}|h(x^{(i)}), \beta^{-1})$$

Likelihood function

Max. likelihood solution

$$\theta_{ML} = \operatorname{argmax}_{\theta} p(\mathbf{t}|\mathbf{x}, \theta, \beta)$$

$$\beta_{ML} = \operatorname{argmax}_{\beta} p(\mathbf{t}|\mathbf{x}, \theta, \beta)$$

## Assignment Project Exam Help

Want to maximize

Add WeChat powcoder

$$p(\mathbf{t}|\mathbf{x}, \theta, \beta) = \prod_{i=1}^m N(t^{(i)} | h(x^{(i)}), \beta^{-1})$$

Easier to maximize  $\log()$

Assignment Project Exam Help

$$\ln p(\mathbf{t}|\mathbf{x}, \theta, \beta) =$$

<https://powcoder.com>

$$-\frac{\beta}{2} \sum_{i=1}^m (h(x^{(i)}) - t^{(i)})^2 + \frac{m}{2} \ln \beta - \frac{m}{2} \ln(2\pi)$$

Add WeChat powcoder



# Assignment Project Exam Help

Want to maximize w.r.t.  $\theta$

Add WeChat powcoder

$$\ln p(\mathbf{t}|\mathbf{x}, \theta, \beta) = -\frac{\beta}{2} \sum_{i=1}^m (h(x^{(i)}) - t^{(i)})^2 + \frac{m}{2} \ln \beta - \frac{m}{2} \ln(2\pi)$$

... but this is same as minimizing sum-of-squares cost<sup>1</sup>

$$\frac{1}{2m} \sum_{i=1}^m (h(x^{(i)}) - t^{(i)})^2$$

Add WeChat powcoder

... which is the same as our SSE cost from before!!

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

<sup>1</sup>multiply by  $-\frac{1}{m\beta}$ , changing max to min, omit last two terms (don't depend on  $\theta$ )

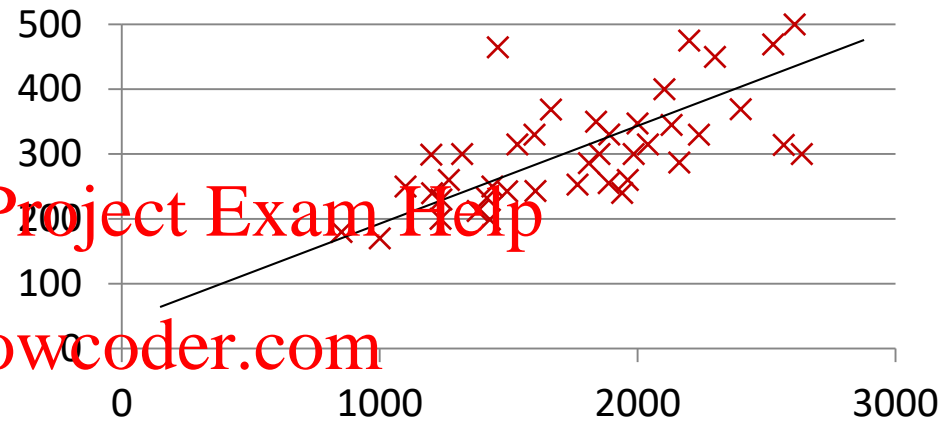
# Summary: Maximum Likelihood Solution for Linear Regression

Hypothesis:

$$h_{\theta}(x) = \theta^T x$$

$\theta$ : parameters

$D = (x^{(i)}, t^{(i)})$ : data



Likelihood:

$$p(\mathbf{t}|\mathbf{x}, \theta, \beta) = \prod_{i=1}^m N(t^{(i)} | h_{\theta}(x^{(i)}), \beta^{-1})$$

Goal: maximize likelihood, equivalent to

$$\operatorname{argmin}_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - t^{(i)})^2 \quad (\text{same as minimizing SSE})$$

# Assignment Project Exam Help

## Probabilistic Motivation for SSE

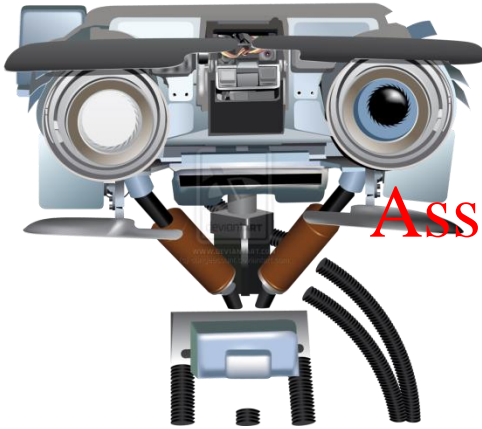
- Under the Gaussian noise assumption, maximizing the probability of the data points is the same as minimizing a sum of squares cost function

<https://powcoder.com>

- Also known as least squares method
- ML can be used for other hypotheses!
  - But linear regression has closed-form solution

Assignment Project Exam Help

Add WeChat powcoder



# Supervised Learning: Classification

Assignment Project Exam Help

<https://powcoder.com>

---

Add WeChat powcoder

# Assignment Project Exam Help

## Classification

Add WeChat powcoder

$$y \in \{0,1\}$$

0: “Negative Class” (e.g., benign tumor)

1: “Positive Class” (e.g., malignant tumor)

# Assignment Project Exam Help

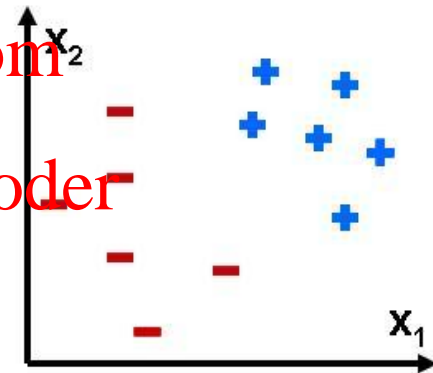
Tumor: Malignant / Benign?

Email: Spam / Not Spam?

Video: Viral / Not Viral?

<https://powcoder.com>

Add WeChat powcoder



# Assignment Project Exam Help

## Classification

Add WeChat powcoder

$$y \in \{0,1\}$$

0: "Negative Class" (e.g., benign tumor)

1: "Positive Class" (e.g., malignant tumor)

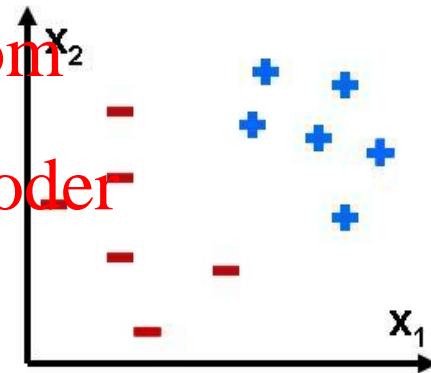
# Assignment Project Exam Help

Why not use least squares regression?

$$\operatorname{argmin}_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

<https://powcoder.com>

Add WeChat powcoder



# Assignment Project Exam Help

## Classification

Add WeChat powcoder

$$y \in \{0,1\}$$

0: “Negative Class” (e.g., benign tumor)

1: “Positive Class” (e.g., malignant tumor)

# Assignment Project Exam Help

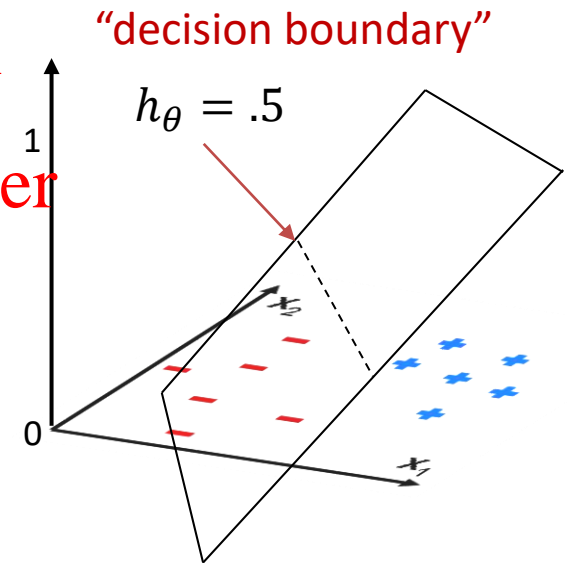
Why not use least squares regression?

$$\operatorname{argmin}_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

<https://powcoder.com>

Add WeChat powcoder

- Indeed, this is possible!
  - Predict 1 if  $h_{\theta}(x) > .5$ , 0 otherwise
- However, outliers lead to problems...
- Instead, use **logistic regression**



# Least Squares vs. Logistic Regression for Classification

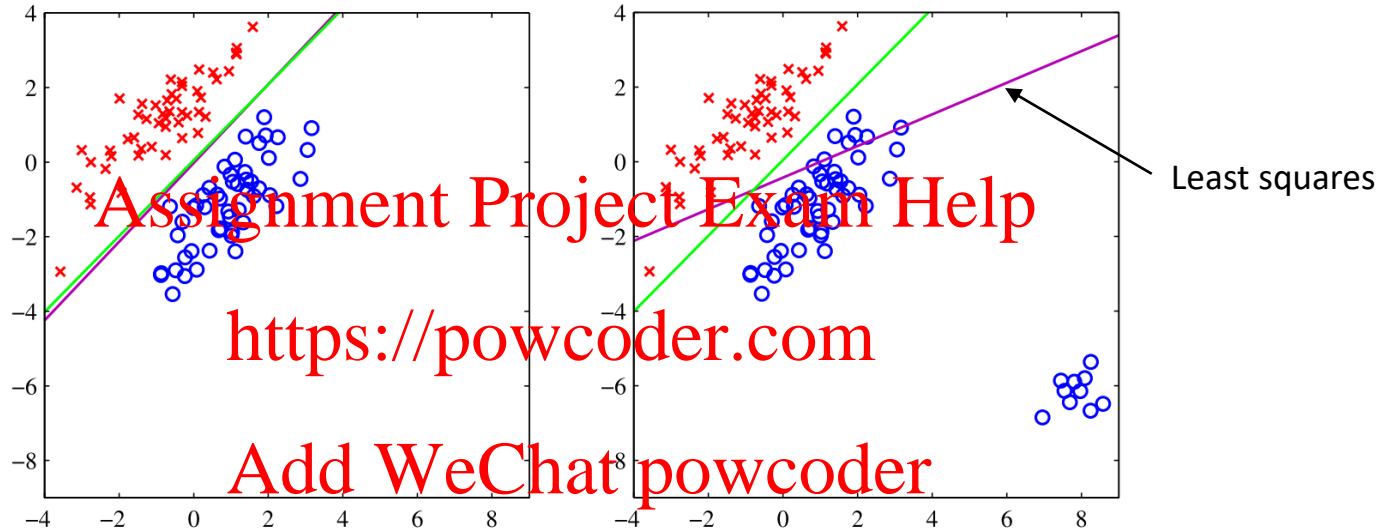


Figure 4.4 from Bishop. The left plot shows data from two classes, denoted by red crosses and blue circles, together with the decision boundary found by least squares (magenta curve) and also by the logistic regression model (green curve). The right-hand plot shows the corresponding results obtained when extra data points are added at the bottom left of the diagram, showing that **least squares is highly sensitive to outliers**, unlike logistic regression.

(see Bishop 4.1.3 for more details)



# Assignment Project Exam Help

## Logistic Regression

Add WeChat powcoder

$$0 \leq h_{\theta}(x) \leq 1$$

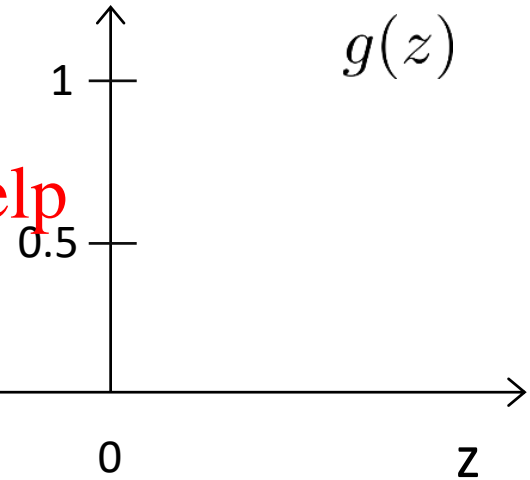
map to (0, 1) with “sigmoid” function

$$g(z) = \frac{1}{1 + e^{-z}}$$

$$h_{\theta}(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}}$$

https://powcoder.com

Add WeChat powcoder



$$h_{\theta}(x) = p(y = 1|x) \quad \text{“probability of class 1 given input”}$$

# Assignment Project Exam Help

# Logistic Regression

Add WeChat powcoder

Hypothesis:

$$h_{\theta}(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}}$$

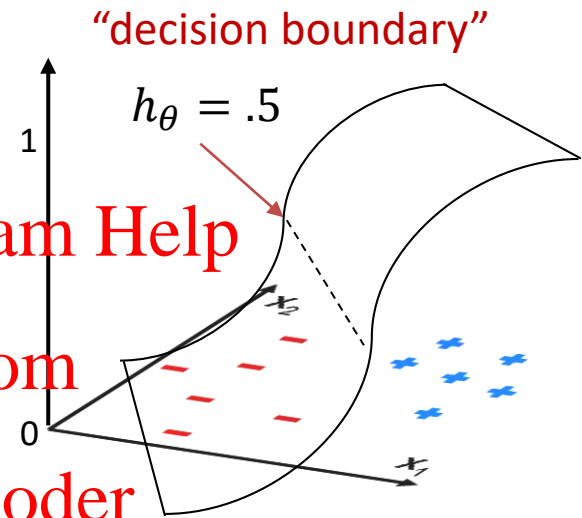
Assignment Project Exam Help

Predict “y = 1” if  $h_{\theta}(x) \geq 0.5$

Predict “y = 0” if  $h_{\theta}(x) < 0.5$

<https://powcoder.com>

Add WeChat powcoder



# Assignment Project Exam Help

# Logistic Regression Cost

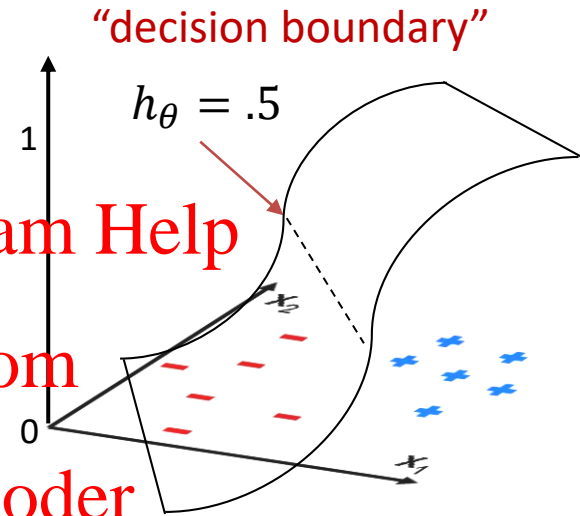
Add WeChat powcoder

Hypothesis:

$$h_{\theta}(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}}$$

$\theta$ : parameters

$D = (x^{(i)}, y^{(i)})$ : data



Cost Function: cross entropy

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_{\theta}(x^{(i)}), y^{(i)})$$
$$= -\frac{1}{m} \left[ \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\theta}(x^{(i)})) \right]$$

Goal: minimize cost  $\min_{\theta} J(\theta)$

# Assignment Project Exam Help

## Cross Entropy Cost

Add WeChat powcoder

- Cross entropy compares distribution  $q$  to reference  $p$

$$H(p, q) = - \sum_x p(x) \log q(x)$$

<https://powcoder.com>

- Here  $q$  is predicted probability of  $y=1$  given  $x$ , reference distribution is  $p=y^{(i)}$ , which is either 1 or 0

$$-\frac{1}{m} \left[ \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\theta}(x^{(i)})) \right]$$

# Assignment Project Exam Help

## Gradient of Cross Entropy Cost

- Cross entropy cost

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_{\theta}(x^{(i)}), y^{(i)})$$
$$= -\frac{1}{m} \left[ \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\theta}(x^{(i)})) \right]$$

- its gradient w.r.t  $\theta$  is:

$$(h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \quad (\text{left as exercise})$$

- No direct closed-form solution

Assignment Project Exam Help

# Gradient descent for Logistic Regression

Add WeChat powcoder

## Cost

$$J(\theta) = -\frac{1}{m} \left[ \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\theta}(x^{(i)})) \right]$$

Assignment Project Exam Help

Want  $\min_{\theta} J(\theta)$ : <https://powcoder.com>

Repeat {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

Add WeChat powcoder

} (simultaneously update all  $\theta_j$ )

Assignment Project Exam Help

# Gradient descent for Logistic Regression

Add WeChat powcoder

## Cost

$$J(\theta) = -\frac{1}{m} \left[ \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\theta}(x^{(i)})) \right]$$

Assignment Project Exam Help

Want  $\min_{\theta} J(\theta)$ : <https://powcoder.com>

Repeat {

$$\theta_j := \theta_j - \alpha \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

Add WeChat powcoder

(simultaneously update all  $\theta_j$ )  
}

# Maximum Likelihood Derivation of Logistic Regression Cost

We can derive the Logistic Regression cost

$$\begin{aligned} J(\theta) &= \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_{\theta}(x^{(i)}), y^{(i)}) \\ &= -\frac{1}{m} \left[ \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\theta}(x^{(i)})) \right] \end{aligned}$$

using Maximum Likelihood, by writing down the likelihood function as

$$p(D|\theta) = \prod_{i=1}^m p(y = 1|x^{(i)}, \theta)^{y^{(i)}} (1 - p(y = 1|x^{(i)}, \theta))^{(1-y^{(i)})}$$

where

$$p(y = 1|x^{(i)}, \theta) = h_{\theta}(x^{(i)})$$

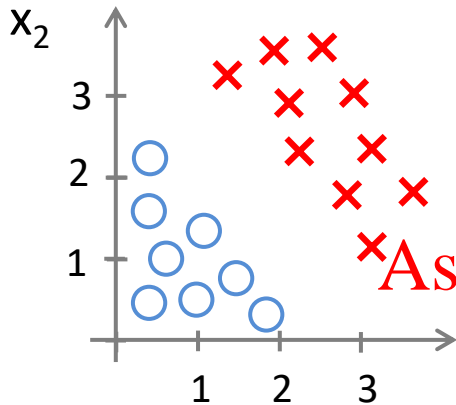
then taking the log.



# Assignment Project Exam Help

## Decision boundary

Add WeChat powcoder



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

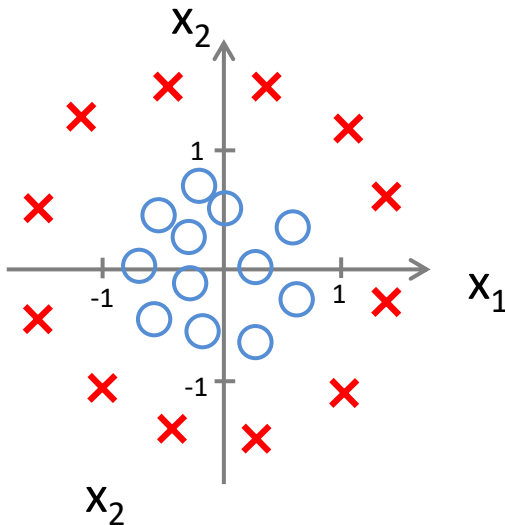
Predict “ $y = 1$ ” if  $-3 + x_1 + x_2 \geq 0$

Assignment Project Exam Help

<https://powcoder.com>

## Non-linear decision boundaries

Add WeChat powcoder

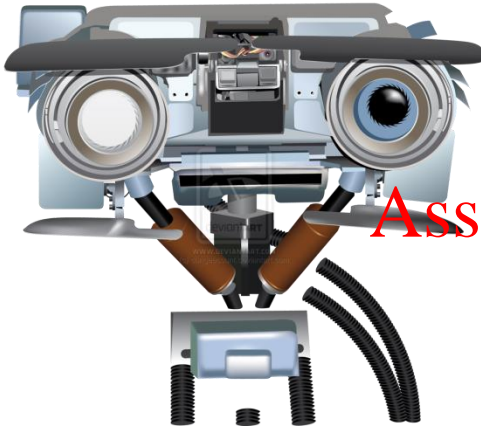


$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$$

Predict “ $y = 1$ ” if  $-1 + x_1^2 + x_2^2 \geq 0$

Assignment Project Exam Help

Add WeChat powcoder



Assignment Project Exam Help Supervised Learning II

<https://powcoder.com>

---

Add WeChat powcoder

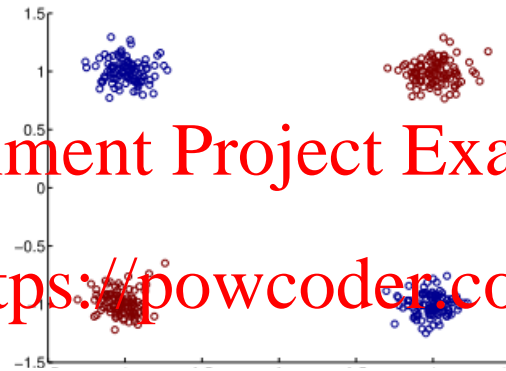
Non-linear features

# Assignment Project Exam Help

# What to do if data is nonlinear?

## Add WeChat powcoder

### Example of nonlinear classification

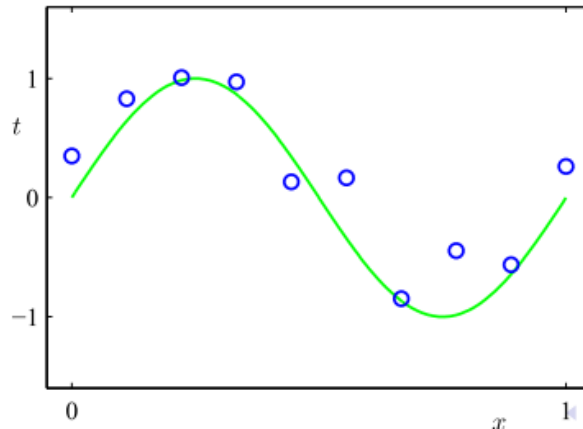


Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

### Example of nonlinear regression



# Assignment Project Exam Help

## Nonlinear basis functions

Add WeChat powcoder

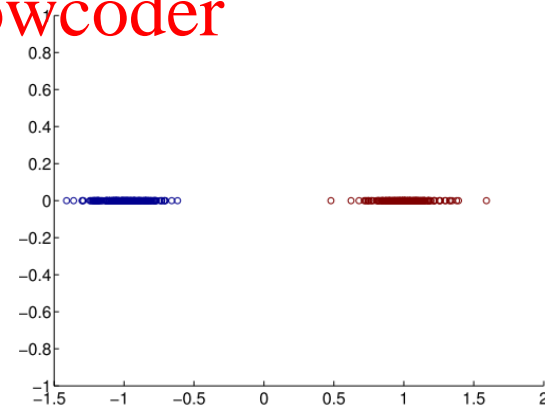
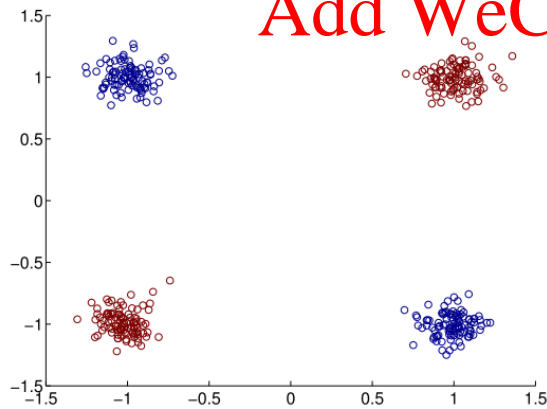
Transform the input/feature

$$\phi(x) : x \in \mathbb{R}^2 \rightarrow z = x_1 \cdot x_2$$

Assignment Project Exam Help

Transformed training data: <https://powcoder.com> linearly separable!

Add WeChat powcoder

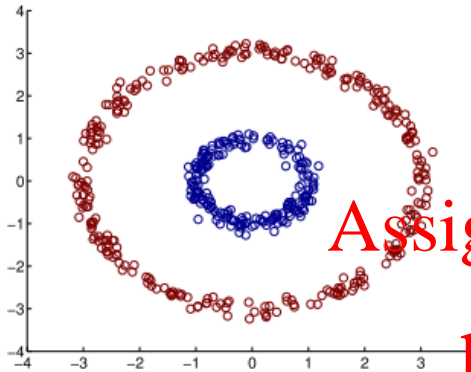


Assignment Project Exam Help

Another example

Add WeChat powcoder

How to transform the input/feature?



Assignment Project Exam Help

<https://powcoder.com>

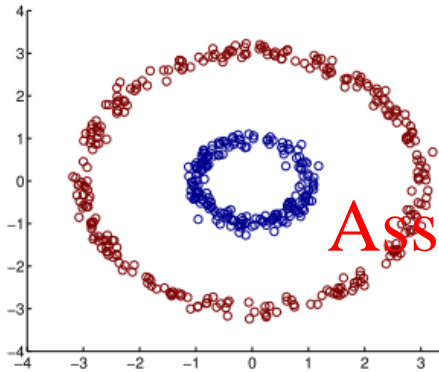
Add WeChat powcoder

# Assignment Project Exam Help

## Another example

Add WeChat powcoder

How to transform the input/feature?



Assignment Project Exam Help

$$\phi(x): x \in R^2 \rightarrow z = \begin{bmatrix} x_1^2 \\ x_1 \cdot x_2 \\ x_2^2 \end{bmatrix}$$

<https://powcoder.com>

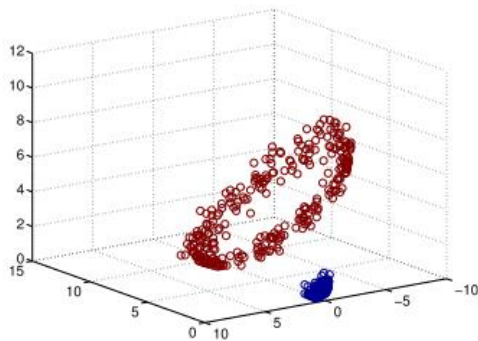
Transformed training data: linearly separable

Add WeChat powcoder

Intuition: suppose  $\theta = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}$

$$\text{Then } \theta^T z = x_1^2 + x_2^2$$

i.e., the sq. distance to the origin!



# Assignment Project Exam Help

## Non-linear basis functions

Add WeChat powcoder

- We can use a nonlinear mapping, or **basis function**

$$\phi(x) : x \in \mathbb{R}^N \rightarrow z \in \mathbb{R}^M$$

Assignment Project Exam Help

<https://powcoder.com>

- where  $M$  is the dimensionality of the new feature/input  $z$  (or  $\phi(x)$ )
  - Note that  $M$  could be either greater than  $D$  or less than, or the same
- Add WeChat powcoder

# Assignment Project Exam Help

## Example with regression

### Add WeChat powcoder

Polynomial basis functions

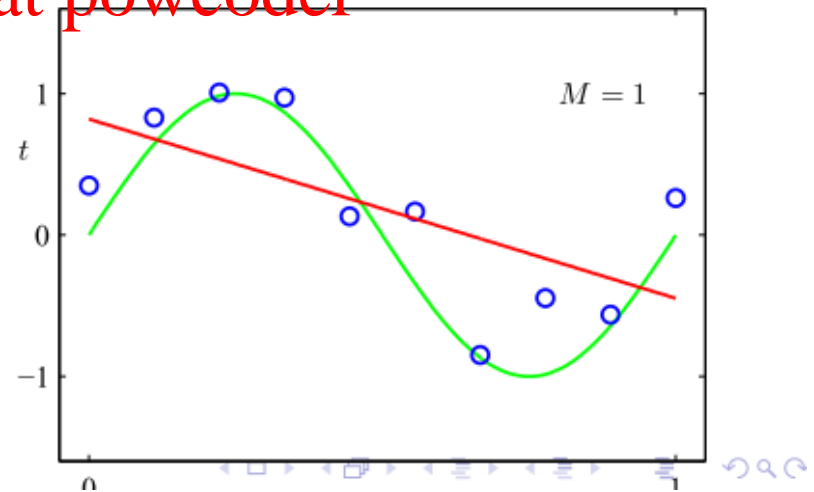
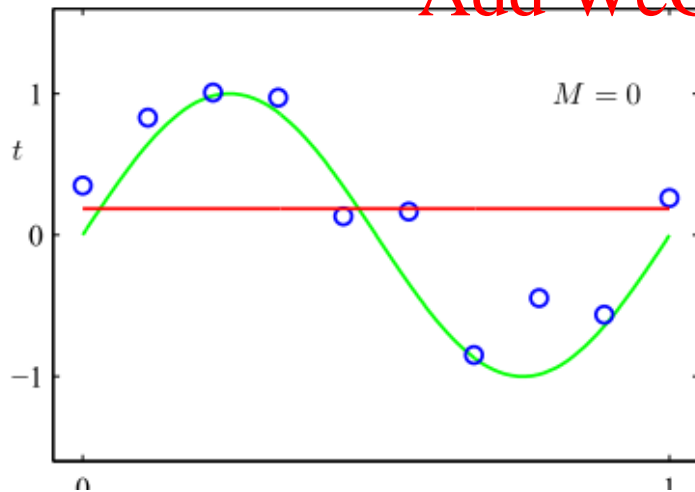
$$\phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^M \end{bmatrix}$$

Assignment Project Exam Help

<https://powcoder.com>

Fitting samples from a sine function: *underrfitting* as  $f(x)$  is too simple

Add WeChat powcoder





# Add more polynomial basis functions

Assignment Project Exam Help  
Add WeChat powcoder

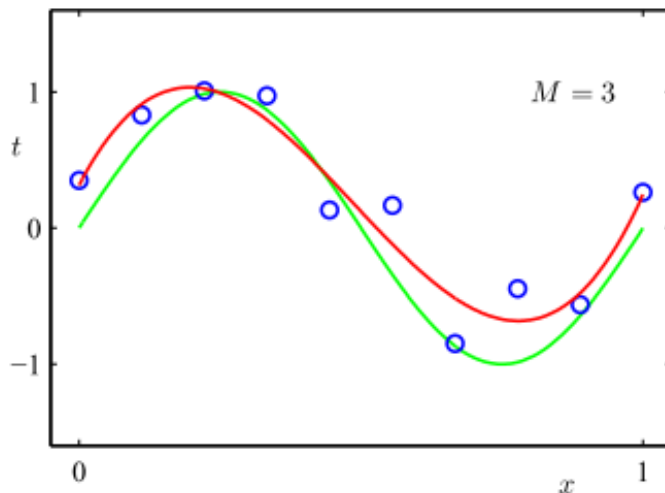
## Polynomial basis functions

$$\phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^M \end{bmatrix}$$

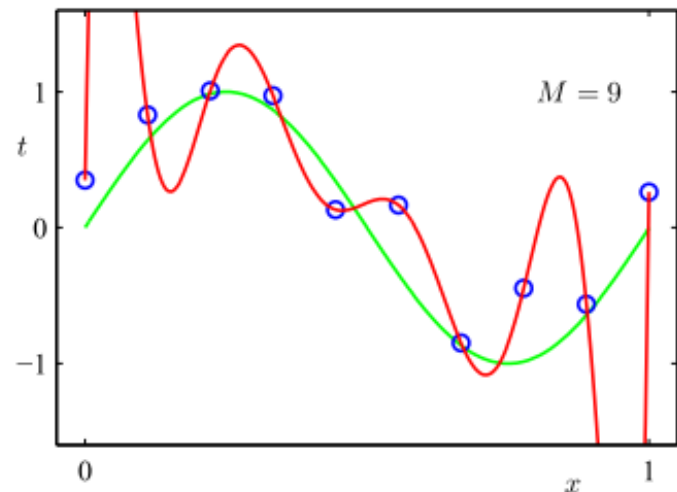
Being too adaptive leads to better results on the training data, but not so great on data that has not been seen!

<https://powcoder.com>

**M=3** *good fit*

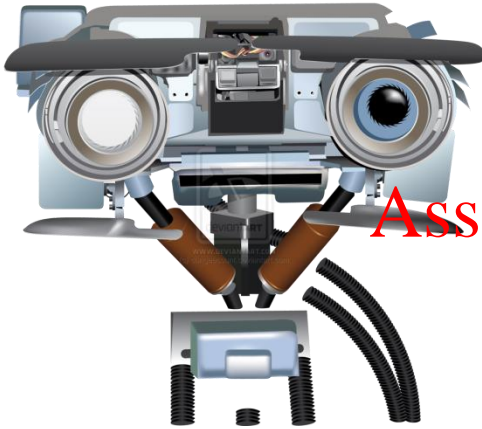


**M=9** *overfitting*



Assignment Project Exam Help

Add WeChat powcoder



Assignment Project Exam Help Supervised Learning II

<https://powcoder.com>

---

Add WeChat powcoder

Overfitting

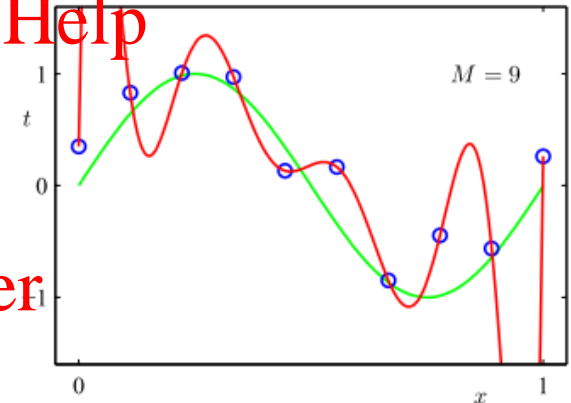
# Assignment Project Exam Help

## Overfitting

Add WeChat powcoder

Parameters for higher-order polynomials are very large

	M = 0	M = 1	M = 3	M = 9	M = 9: <i>overfitting</i>
$\theta_0$	0.19	0.82	0.31	0.35	
$\theta_1$		-1.27	7.99	232.37	
$\theta_2$			-25.43	-5321.83	
$\theta_3$			17.37	48568.31	
$\theta_4$				-231639.30	
$\theta_5$				640042.26	
$\theta_6$				-1061800.52	
$\theta_7$				1042400.18	
$\theta_8$				-557682.99	
$\theta_9$				125201.43	

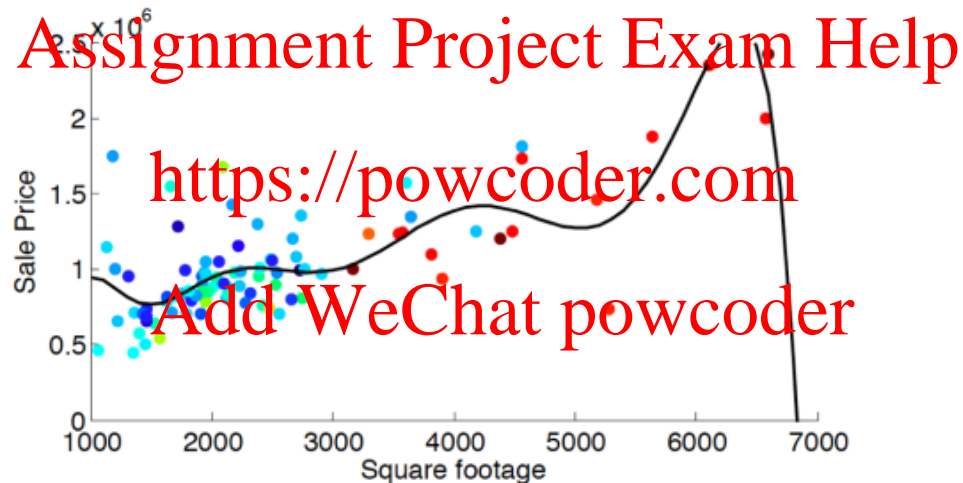


Assignment Project Exam Help

# Overfitting disaster

Add WeChat powcoder

Fitting the housing price data with  $M = 3$



Note that the price would go to zero (or negative) if you buy bigger houses!  
This is called poor generalization/overfitting.

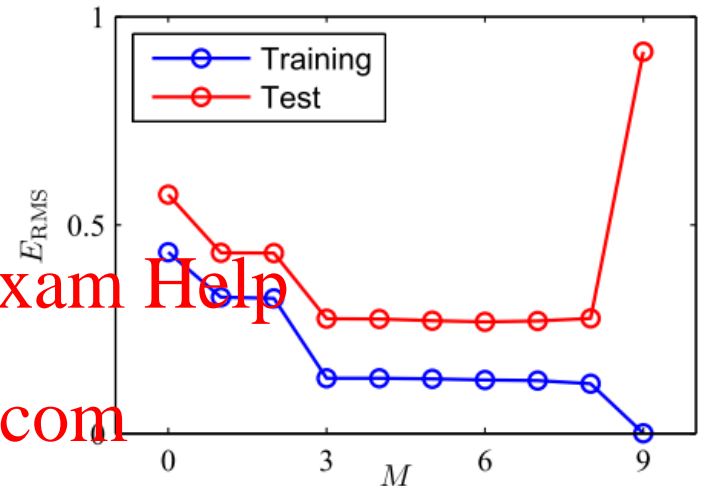
# Assignment Project Exam Help

## Detecting overfitting

Add WeChat powcoder

Plot model complexity versus  
objective function on test/train data

As model becomes more complex,  
performance on training keeps  
improving while on test data it increases



<https://powcoder.com>

**Horizontal axis:** measure of model complexity  
In this example, we use the maximum order of the polynomial basis  
functions.

**Vertical axis:** For regression, it would be SSE or mean SE (MSE)  
For classification, the vertical axis would be classification error rate or  
cross-entropy error function

Assignment Project Exam Help

# Overcoming overfitting

Add WeChat powcoder

- Basic ideas

- Use more training data.

Assignment Project Exam Help

- Regularization methods

- Cross-validation

<https://powcoder.com>

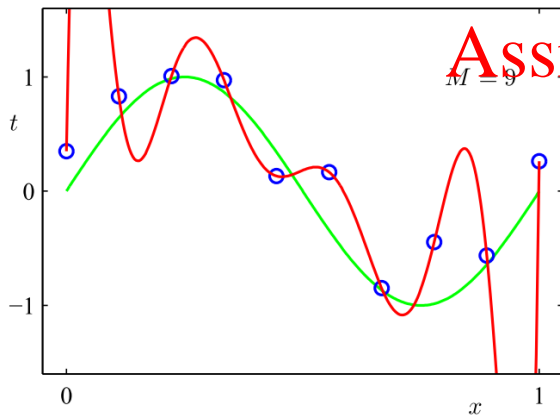
Add WeChat powcoder

# Assignment Project Exam Help

## Solution: use more data

Add WeChat powcoder

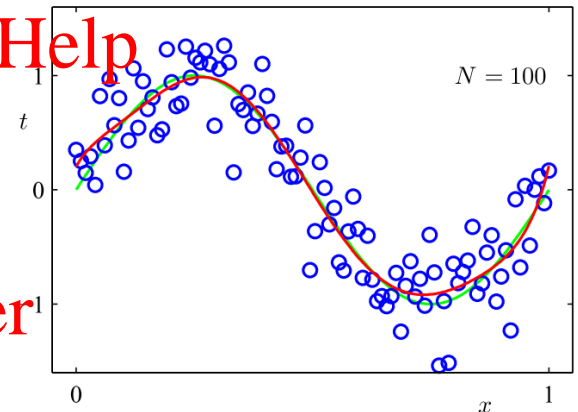
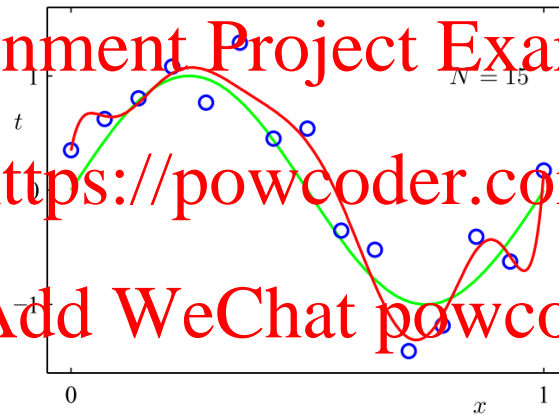
$M=9$ , increase  $N$



Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



What if we do not have a lot of data?

Assignment Project Exam Help

# Overcoming overfitting

Add WeChat powcoder

- Basic ideas

- Use more training data

Assignment Project Exam Help

- Regularization methods

<https://powcoder.com>

- Cross-validation

Add WeChat powcoder



Assignment Project Exam Help

Next Class

Add WeChat powcoder

## **Supervised Learning 3: Regularization**

more logistic regression, regularization; bias-variance

Assignment Project Exam Help

<https://powcoder.com>

**Reading:** Bishop 3.1, 3.2

Add WeChat powcoder

**Discussion/Lab this week:** Intro to Numpy

*PSet 2 out on Thursday*