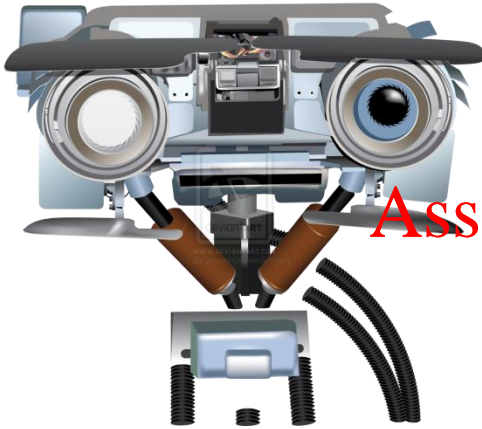


# Announcements

**Reminder:** ps4 self-grading form out, due Friday 10/30

## Assignment Project Exam Help

- pset 5 out Thursday 10/29 due 11/5 (1 week)  
<https://powcoder.com>
- Midterm grades will go up by Monday (don't discuss it yet)  
Add WeChat powcoder
- My Thursday office hours moved to 11am
- Lab this week – probabilistic models, ipython notebook examples



Assignment Project Exam Help

# Bayesian Methods

<https://powcoder.com>

---

Add WeChat powcoder

CS542 Machine Learning

# Bayesian Methods

- Before, we derived cost functions from maximum likelihood, then added regularization terms to these cost functions

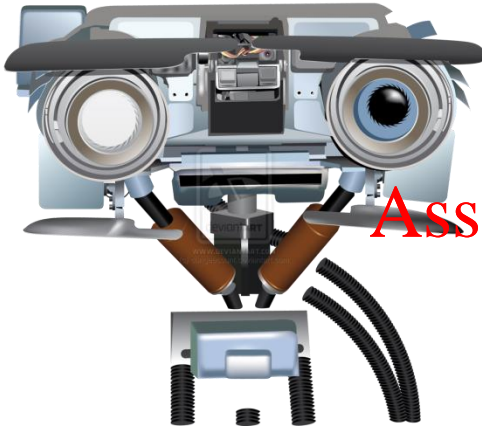
**Assignment Project Exam Help**

- Can we derive regularization directly from probabilistic principles?

**<https://powcoder.com>**

**Add WeChat powcoder**

- Yes! Use Bayesian methods



Assignment Project Exam Help

# Bayesian Methods

<https://powcoder.com>

---

Add WeChat powcoder

Motivation

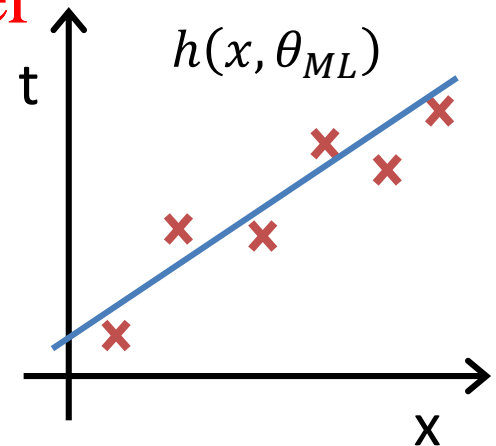
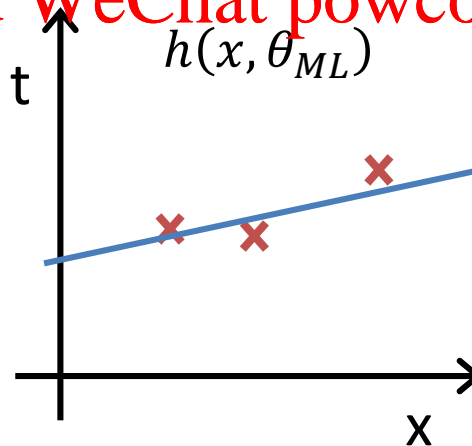
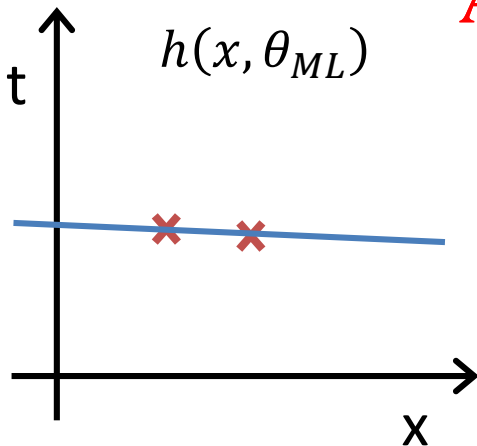
# Problem with Maximum Likelihood: Bias

- ML estimates are biased
- Especially a problem for small number of samples, or high input dimensionality
- Suppose we sample 2,3,6 points from the same dataset, use ML to fit regression parameters

Assignment Project Exam Help

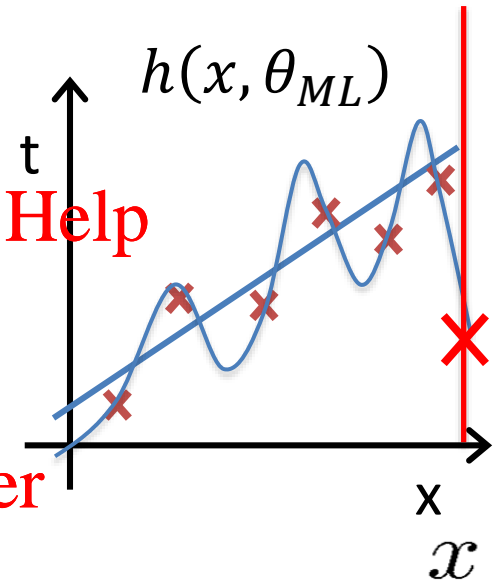
<https://powcoder.com>

Add WeChat powcoder



# Problem with Maximum Likelihood: Overfitting

- ML estimates cannot be used to choose complexity of model
  - E.g. suppose we want to estimate the number of basis functions
  - Choose  $K=1$ ?
  - Or  $K=15$ ?
- ML will always choose  $K$  that best fits training data (in this case,  $K=15$ )
- Solution: use a Bayesian method--define a prior distribution over the parameters (results in regularization)



# Bayesian vs. Frequentist

**Frequentist:** maximize data likelihood

$$p(D|model) = p(D|\theta)$$

Assignment Project Exam Help

**Bayesian:** treat  $\theta$  as random variable, maximize posterior

$$p(\theta|D) = \frac{p(D|\theta)p(\theta)}{p(D)}$$

$p(D|\theta)$  is the data likelihood,  $p(\theta)$  is the prior over the model parameters

# Bayesian Method

Treat  $\theta$  as random variable, maximize posterior

$$p(\theta|D) = \frac{p(D|\theta)p(\theta)}{p(D)}$$

Assignment Project Exam Help

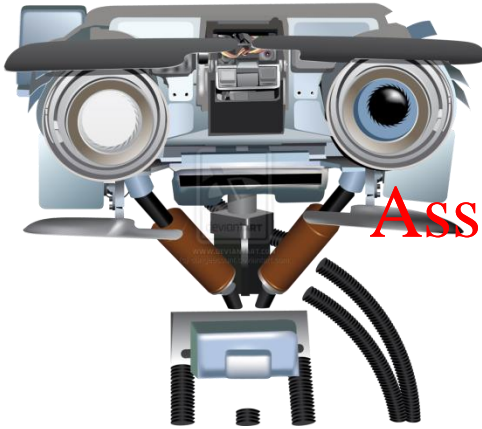
Likelihood  $p(D|\theta)$  is the same as before, as in Maximum Likelihood

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

**Prior**  $p(\theta)$  is a new distribution we model; specifies which parameters are more likely *a priori*, before seeing any data

$p(D)$  does not depend on  $\theta$ , constant when choosing  $\theta$  with the highest posterior probability





# Prior over Model Parameters

Assignment Project Exam Help

<https://powcoder.com>

---

Add WeChat powcoder  
Intuition

# Will he score?

Score!

Score!

Miss

Score!

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Your estimate of  
 $\theta = p(score)$ ?



# Will he score?

Score!

Score!

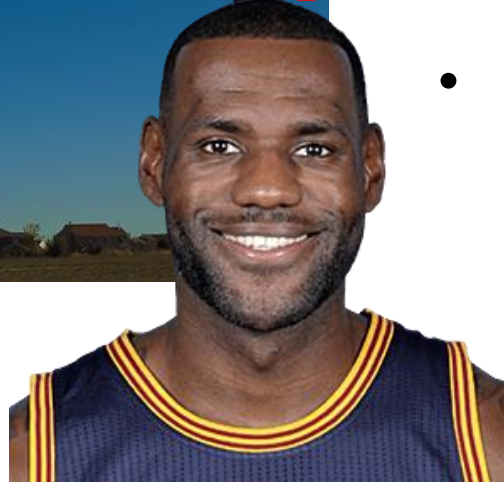
Miss

Score!

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



- Prior information:  
player= [LeBron James](#)
- Your estimate of  $\theta = p(\text{score})$ ?
- Prior  $p(\theta)$  reflects prior knowledge, e.g.,  $\theta \approx 1$

# Prior Distribution

Prior distributions  $p(\theta)$  are probability distributions of model parameters based on some a priori knowledge about the parameters.

**Assignment Project Exam Help**

Prior distributions are independent of the observed data.

**<https://powcoder.com>**

**Add WeChat powcoder**

# Coin Toss Example

---

What is the probability of heads ( $\theta$ )?

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



# Beta Prior for $\theta$

$$P(\theta) = \text{Beta}(\alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{(\alpha-1)} (1 - \theta)^{(\beta-1)}$$

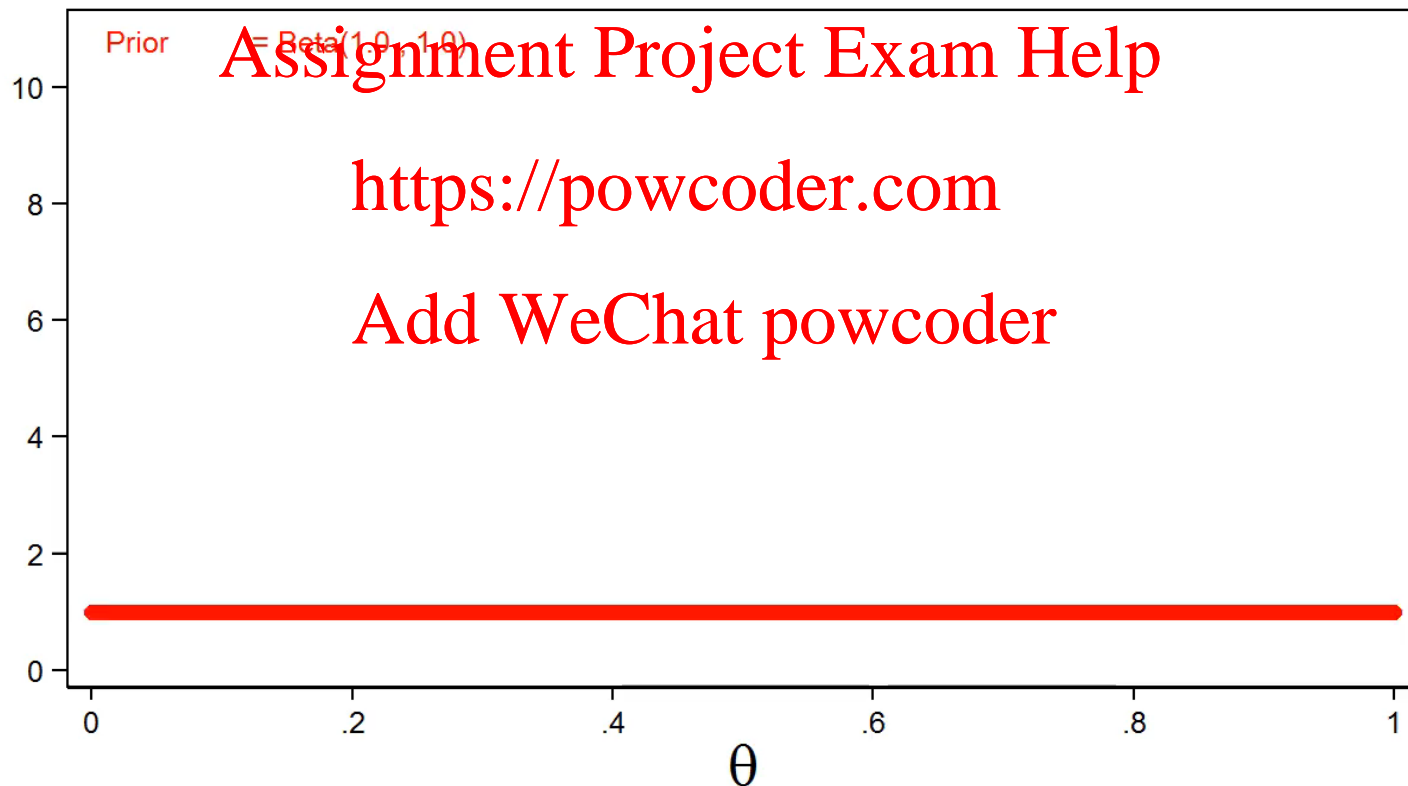
Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

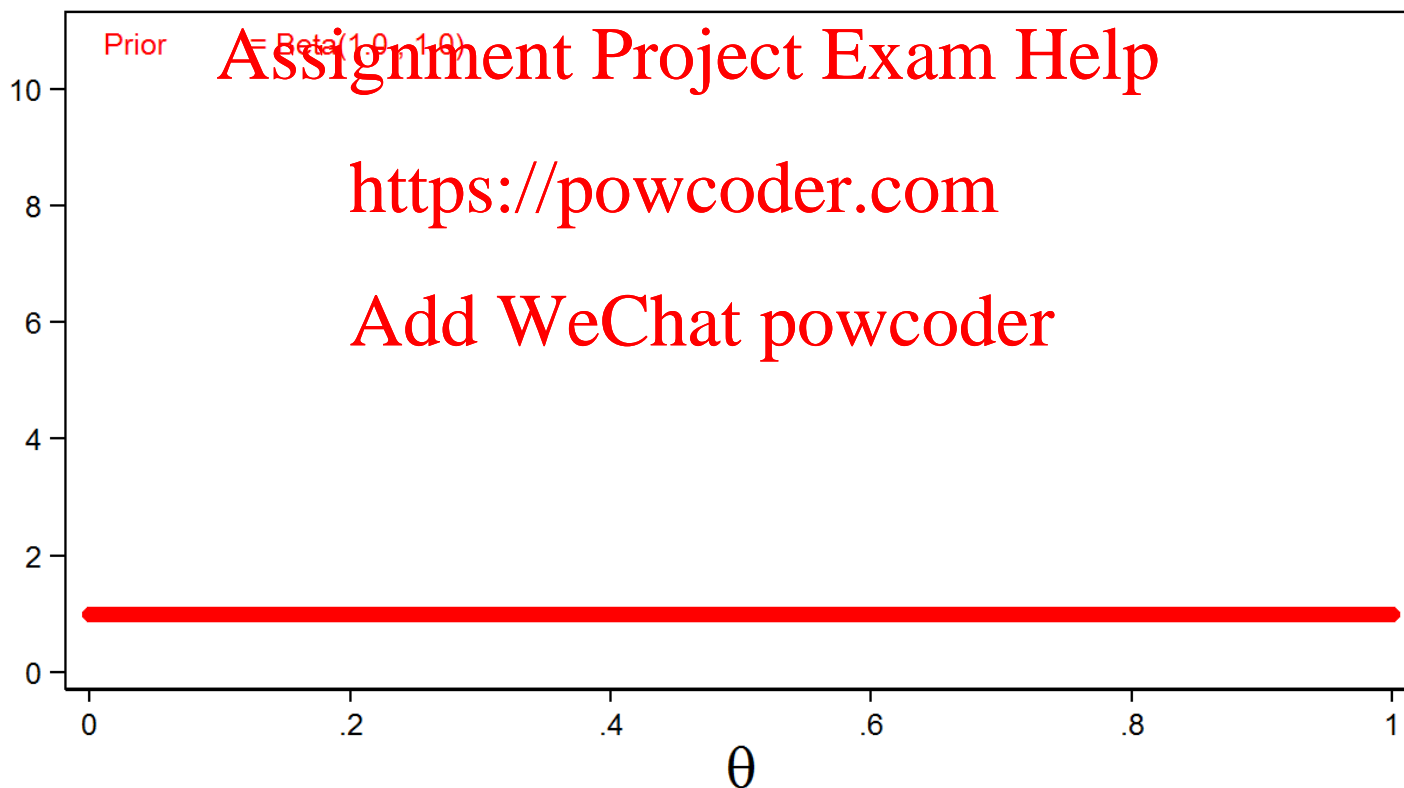
# Beta Prior for $\theta$

$$P(\theta) = \text{Beta}(\alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{(\alpha-1)}(1 - \theta)^{(\beta-1)}$$



# Uninformative Prior

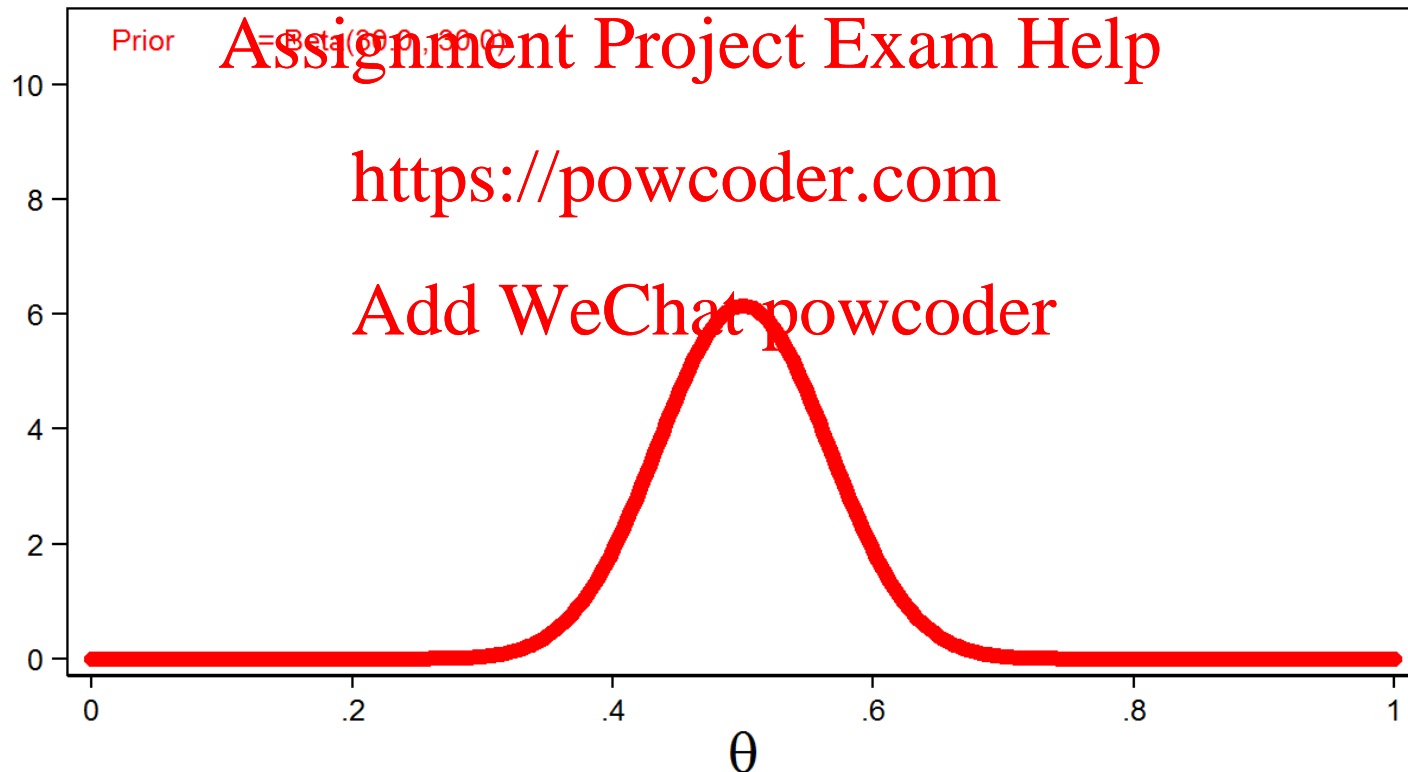
$$P(\theta) = \text{Beta}(\alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{(\alpha-1)}(1 - \theta)^{(\beta-1)}$$





# Informative Prior

$$P(\theta) = \text{Beta}(\alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{(\alpha-1)}(1 - \theta)^{(\beta-1)}$$



# Coin Toss Experiment

- $n = 10$  coin tosses
- $y = 4$  number of heads

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



# Likelihood Function for the Data

$$P(y|\theta) = \textit{Binomial}(n, \theta) = \binom{n}{y} \theta^y (1 - \theta)^{(n-y)}$$

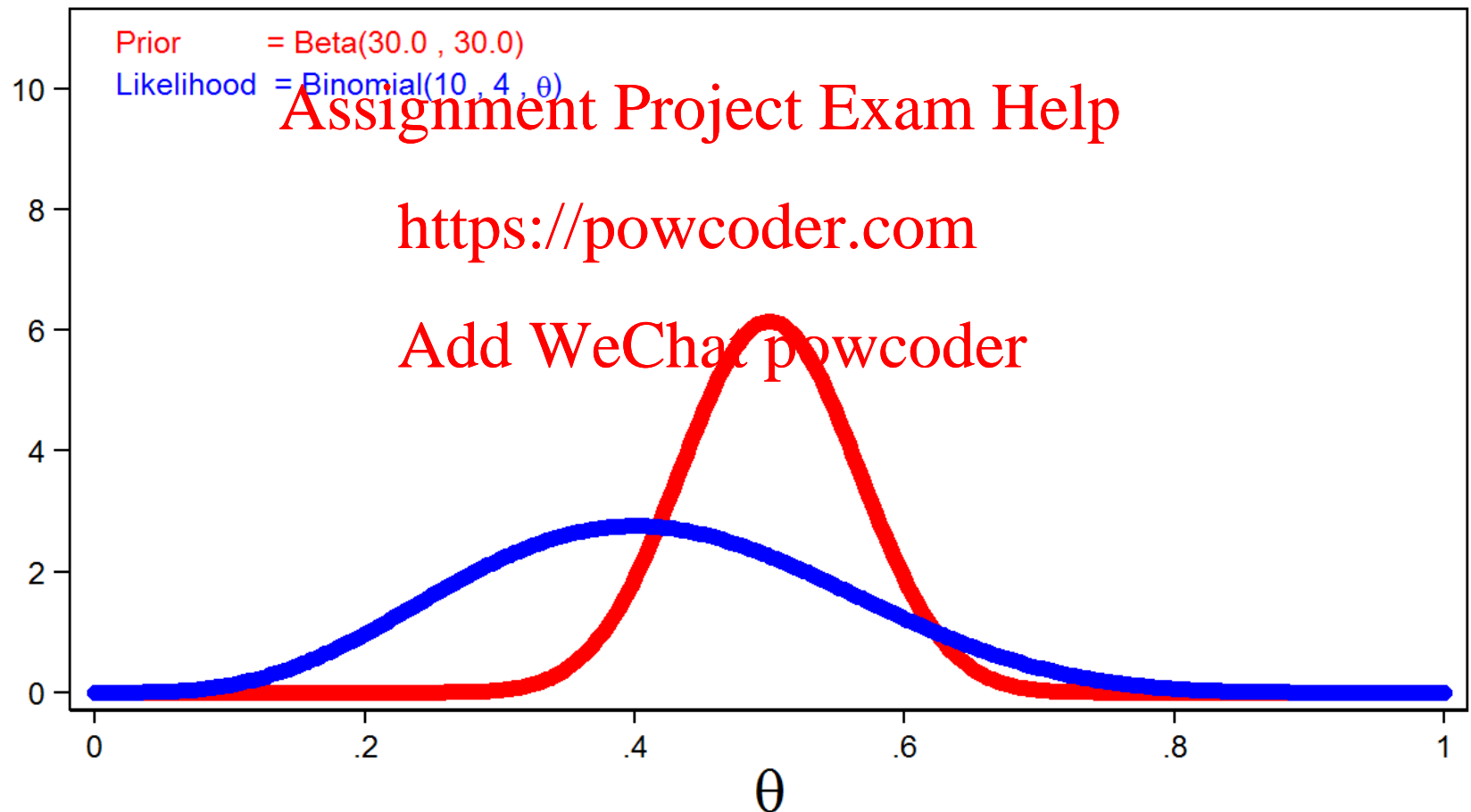
Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

# Prior and Likelihood

$$P(y|\theta) = \text{Binomial}(n, \theta) = \binom{n}{y} \theta^y (1 - \theta)^{(n-y)}$$



# Posterior Distribution

$$\text{Posterior} = \text{Prior} \times \text{Likelihood}$$

$$P(\theta|y) = \text{Assignment Project Exam Help}$$

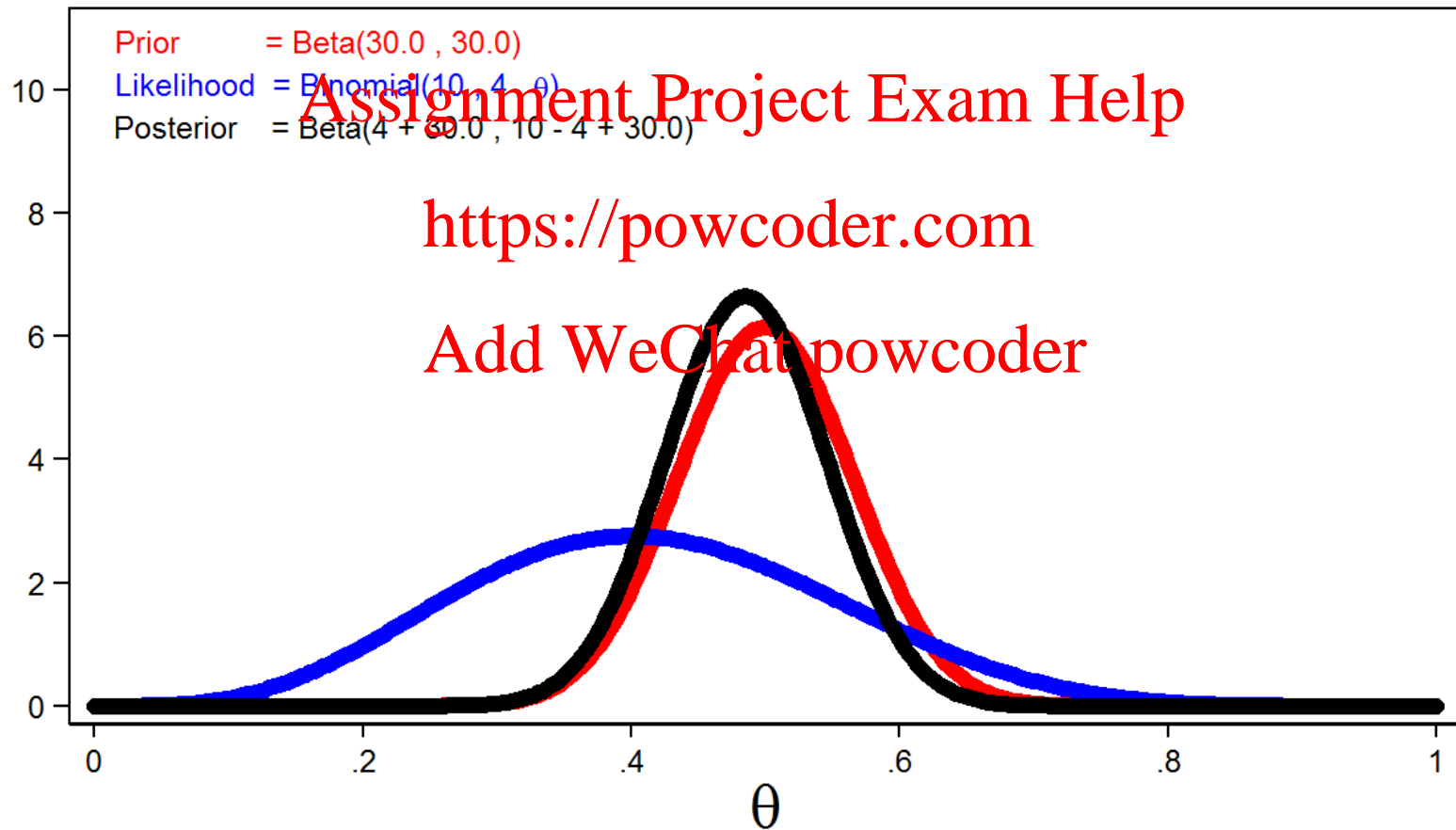
$$P(\theta|y) = \frac{\text{https://powcoder.com}}{\text{Beta}(\alpha, \beta)} \times \text{Binomial}(n, \theta)$$

Add WeChat powcoder

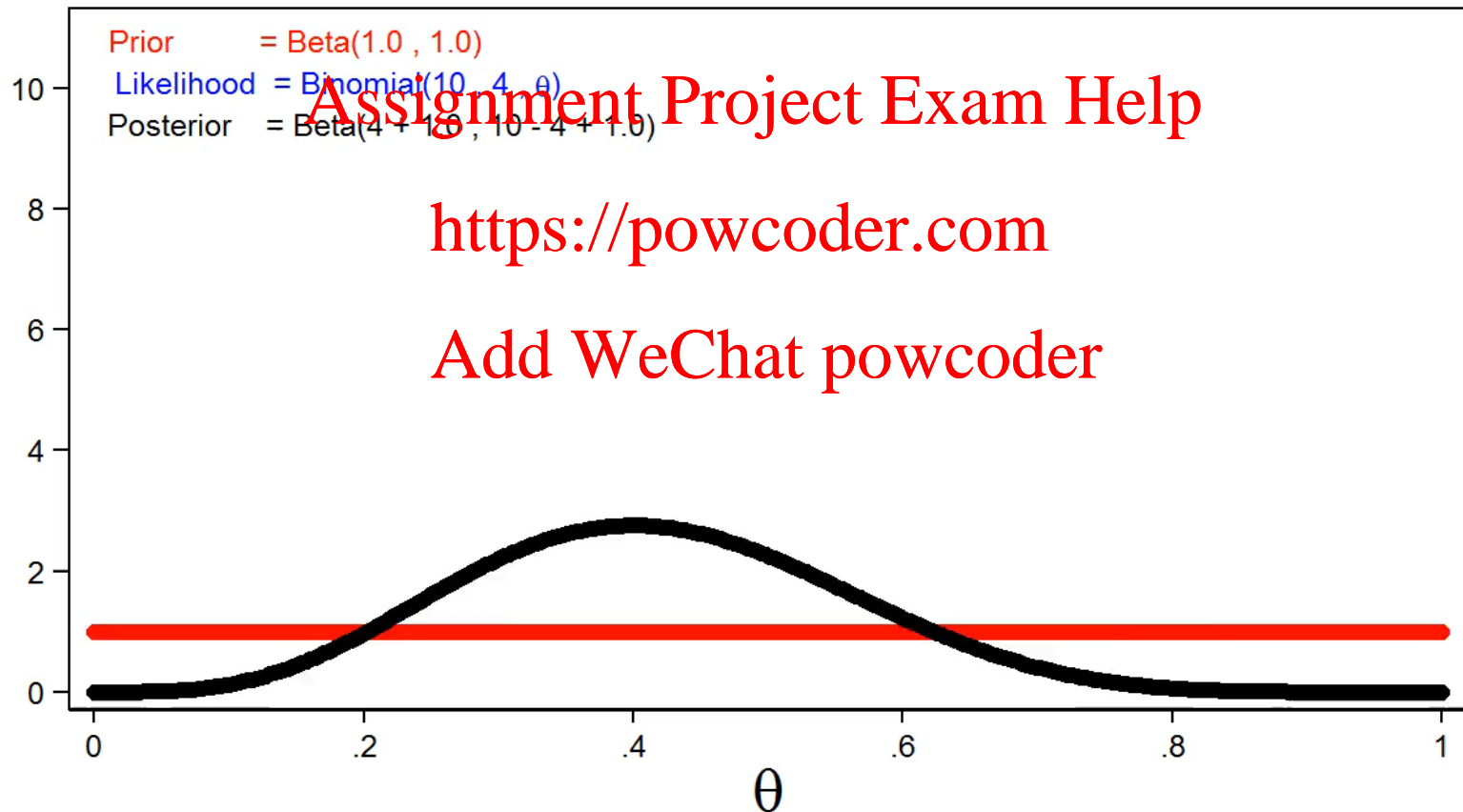
$$= \text{Beta}(y + \alpha, n - y + \beta)$$

This is why we chose the Beta distribution as our prior, posterior is also a Beta distribution: **conjugate prior**.

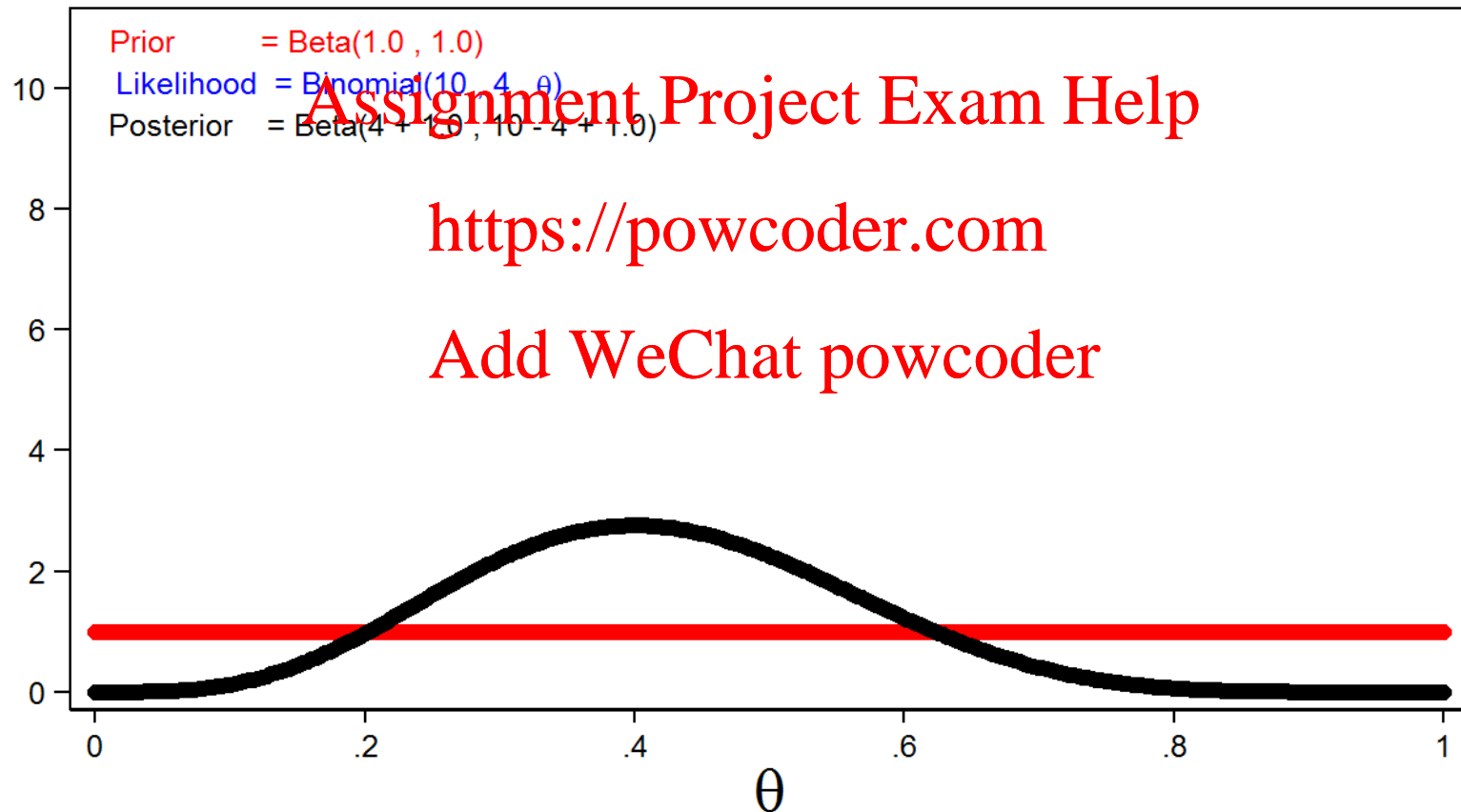
# Posterior Distribution



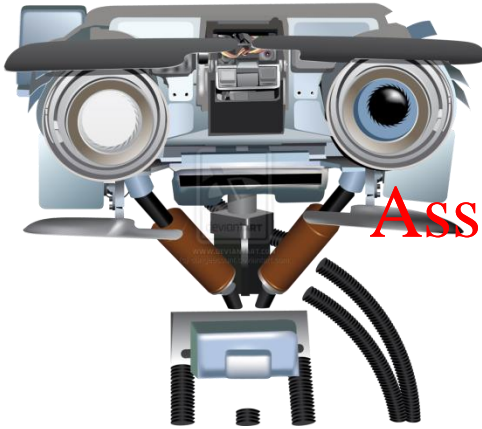
# Effect of Informative Prior



# Effect of Uninformative Prior







# Bayesian Linear Regression

Assignment Project Exam Help

<https://powcoder.com>

---

Add WeChat powcoder

# Bayesian Linear Regression

Let's now apply the Bayesian method to linear regression.

To do that, we must treat parameter  $\theta$  as a random variable, design a prior over it.

<https://powcoder.com>

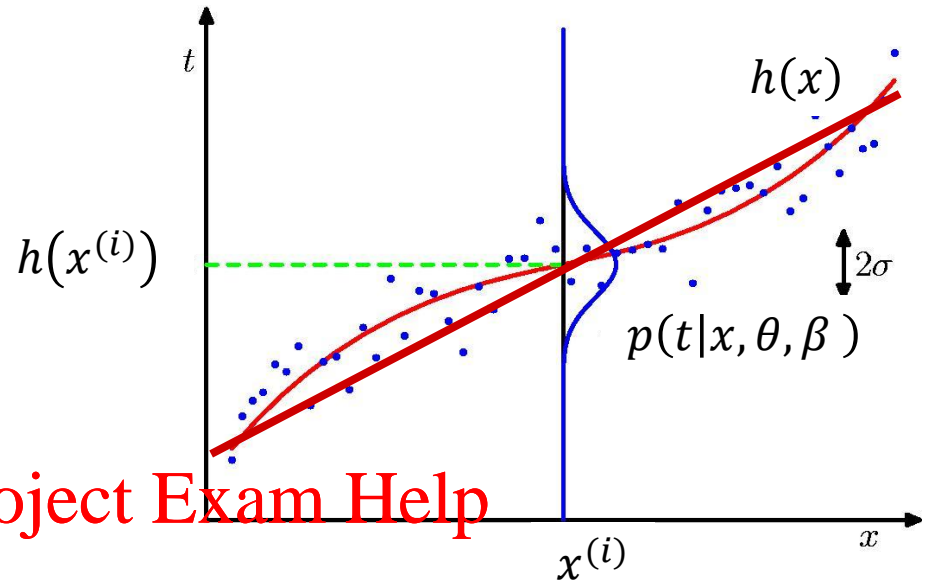
First, review maximum likelihood for linear regression.

# ML for Linear Regression

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

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

$$\text{where } \beta = \frac{1}{\sigma^2}, \quad h(x) = \theta^T x$$



Assignment Project Exam Help

Probability of one data point <https://powcoder.com>

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

Add WeChat powcoder

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

Maximum likelihood solution

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

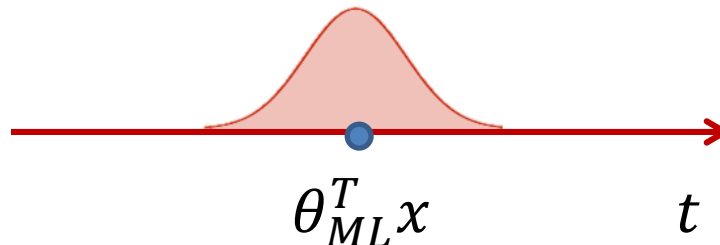
$$\beta_{ML} = \operatorname{argmax}_{\beta} p(\mathbf{t}|\mathbf{x}, \theta, \beta)$$

# What is $\beta$ useful for?

- Recall: we assumed observations  $t$  are Gaussian given  $h(x)$
- $\beta$  allows us to write down distribution over  $t$ , given new  $x$ , called predictive distribution

$$p(t|x, \theta_{ML}, \beta_{ML})$$

$$= N(t | \theta_{ML}^T x, \beta_{ML}^{-1})$$

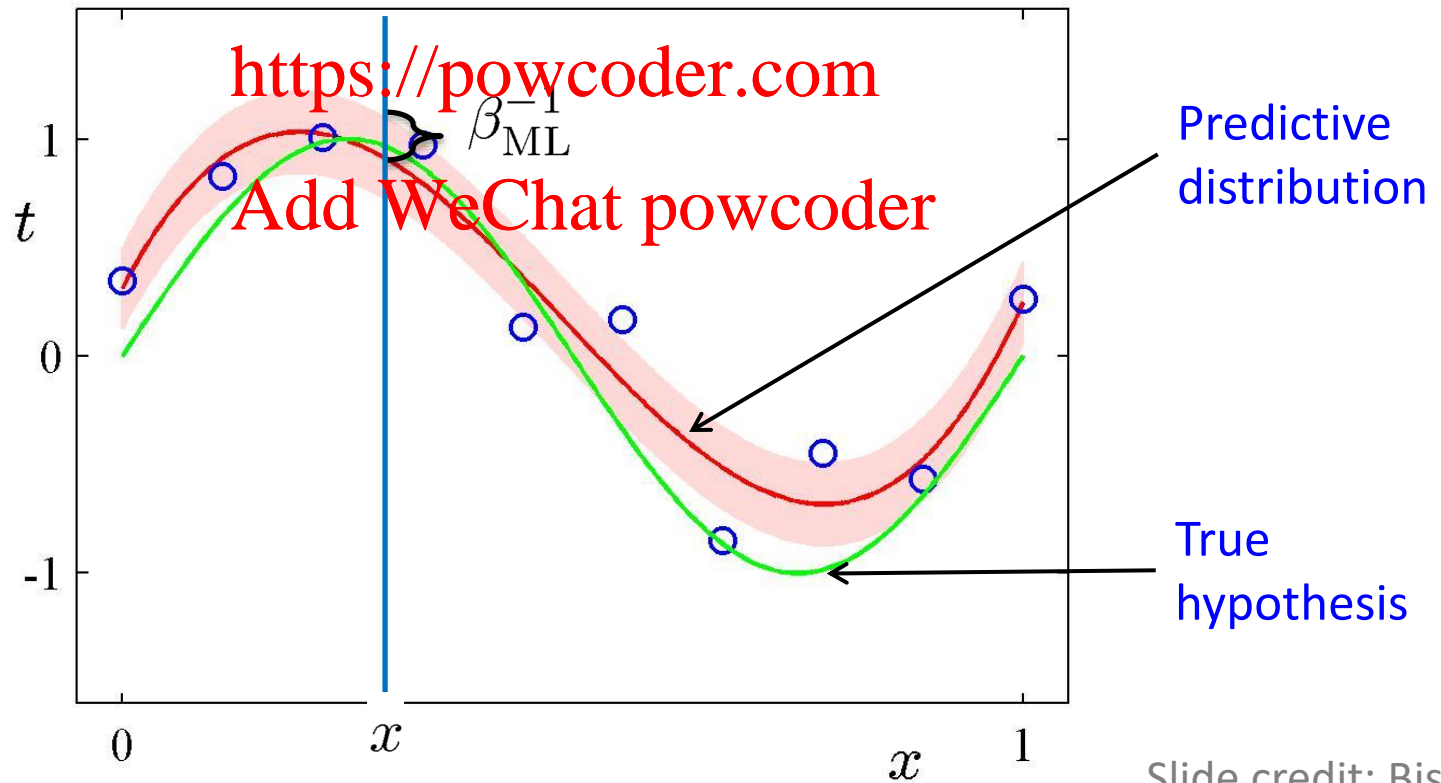


$\beta_{ML}^{-1}$  is the variance of this distribution

# Predictive Distribution

Given a new input point  $x$ , we can now compute a distribution over the output  $t$ :

$$p(t|x, \theta_{ML}, \beta_{ML}) = N(t|\theta_{ML}^T x, \beta_{ML}^{-1})$$



# Define a distribution over parameters

- Define **prior** distribution over  $\theta$  as

$$p(\boldsymbol{\theta}) = N(\boldsymbol{\theta} | \mathbf{m}_0, \mathbf{S}_0)$$

Assignment Project Exam Help

- Combining this with the **likelihood** function and using results for marginal and conditional Gaussian distributions<sup>1</sup>, gives the **posterior**

<https://powcoder.com>

Add WeChat powcoder

$$p(\boldsymbol{\theta} | \mathbf{t}) = N(\boldsymbol{\theta} | \mathbf{m}_N, \mathbf{S}_N)$$

- where

$$\mathbf{m}_N = \mathbf{S}_N (\mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \mathbf{X}^T \mathbf{t})$$

$$\mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta \mathbf{X}^T \mathbf{X}$$

<sup>1</sup>see Bishop 2.3.3

# A common choice for prior

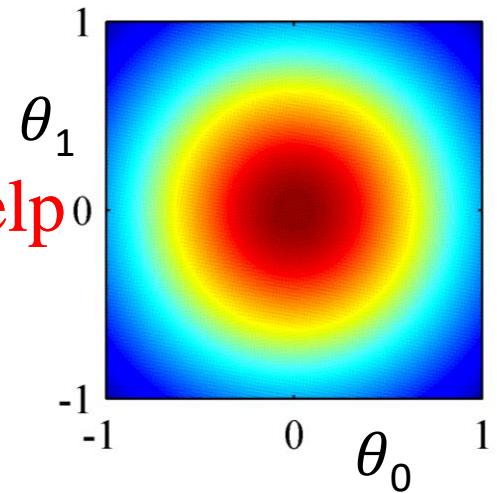
- A common choice for the prior is

$$p(\boldsymbol{\theta}) = N(\boldsymbol{\theta} | \mathbf{0}, \alpha^{-1} \mathbf{I})$$

- for which

$$\mathbf{m}_N = \beta \mathbf{S}_N \mathbf{X}^T \mathbf{t}$$

$$\mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \mathbf{X}^T \mathbf{X}$$



# Intuition: prefer $\theta$ to be simple

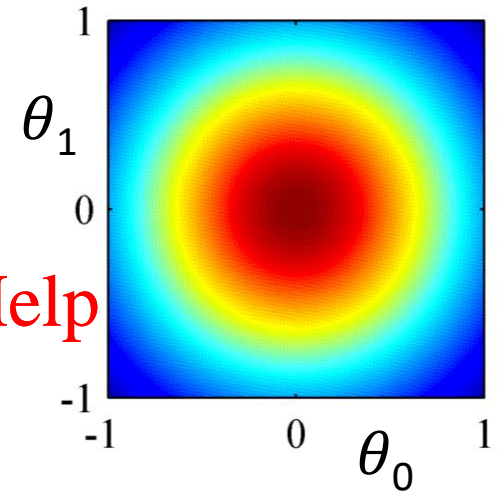
For a linear model for regression,  $\theta^T x$

What do we mean by  $\theta$  being simple?

Assignment Project Exam Help

$$p(\theta) = N(\theta | \mathbf{0}, \alpha^{-1} I)$$

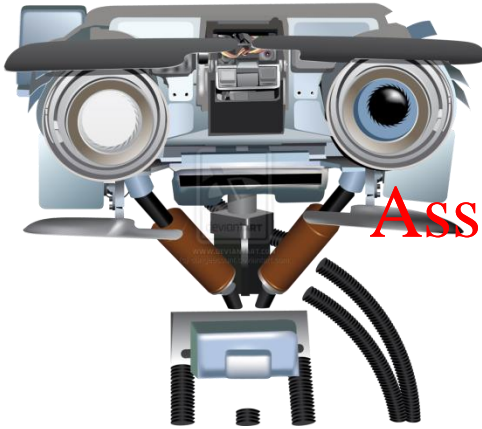
<https://powcoder.com>



Namely, put a prior ~~on  $\theta$~~ , which captures our belief that  $\theta$  is around zero, i.e., resulting in a simple model for prediction.

This Bayesian way of thinking is to regard  $\theta$  as a random variable, and we will use the observed data  $D$  to update our prior belief on  $\theta$





# Bayesian Linear Regression Example

Assignment Project Exam Help

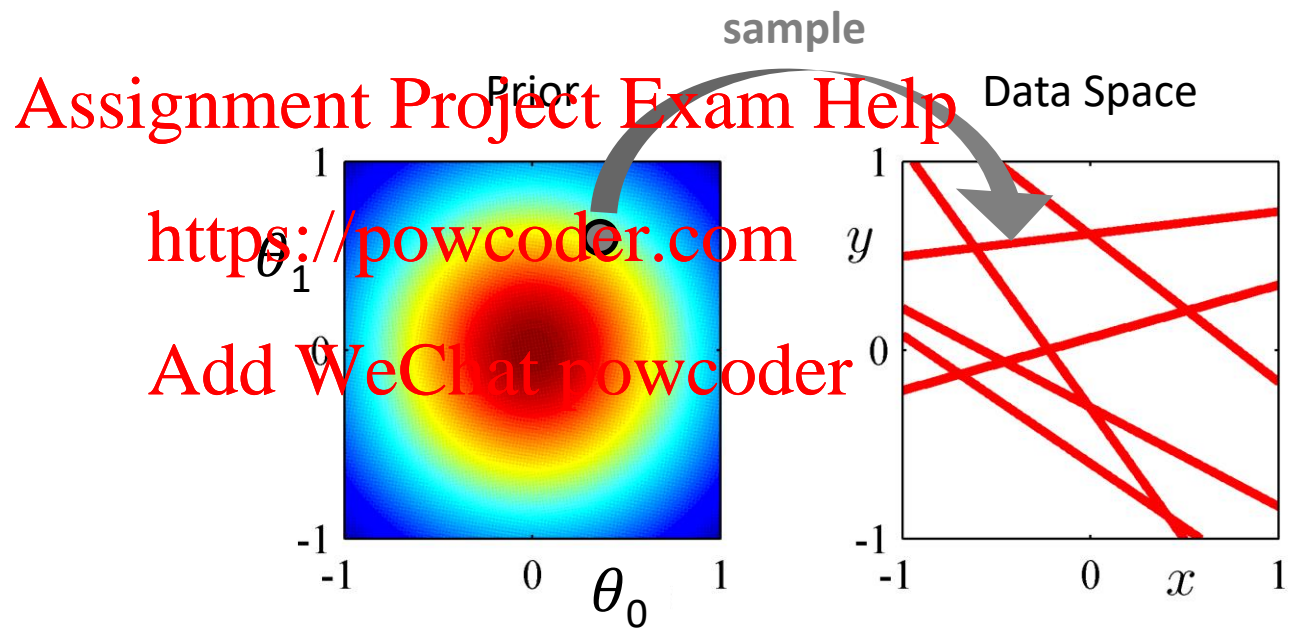
<https://powcoder.com>

---

Add WeChat powcoder

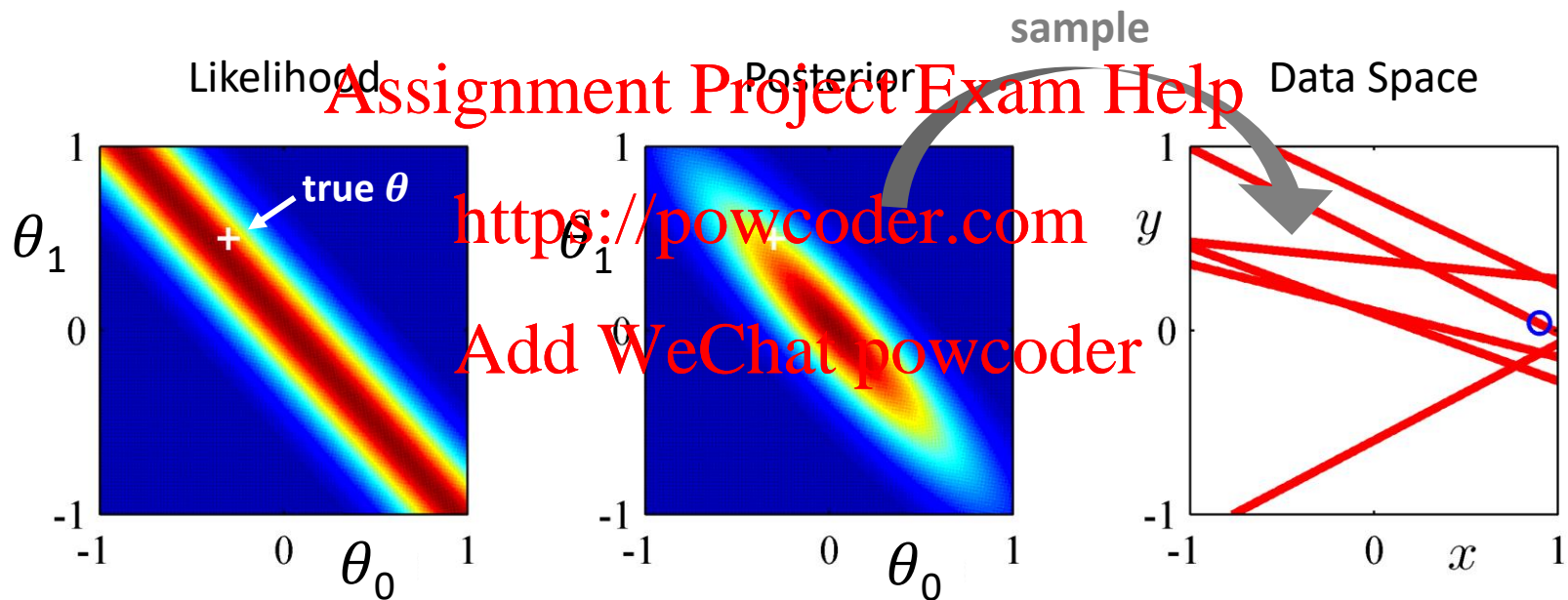
# Bayesian Linear Regression Example

0 data points observed



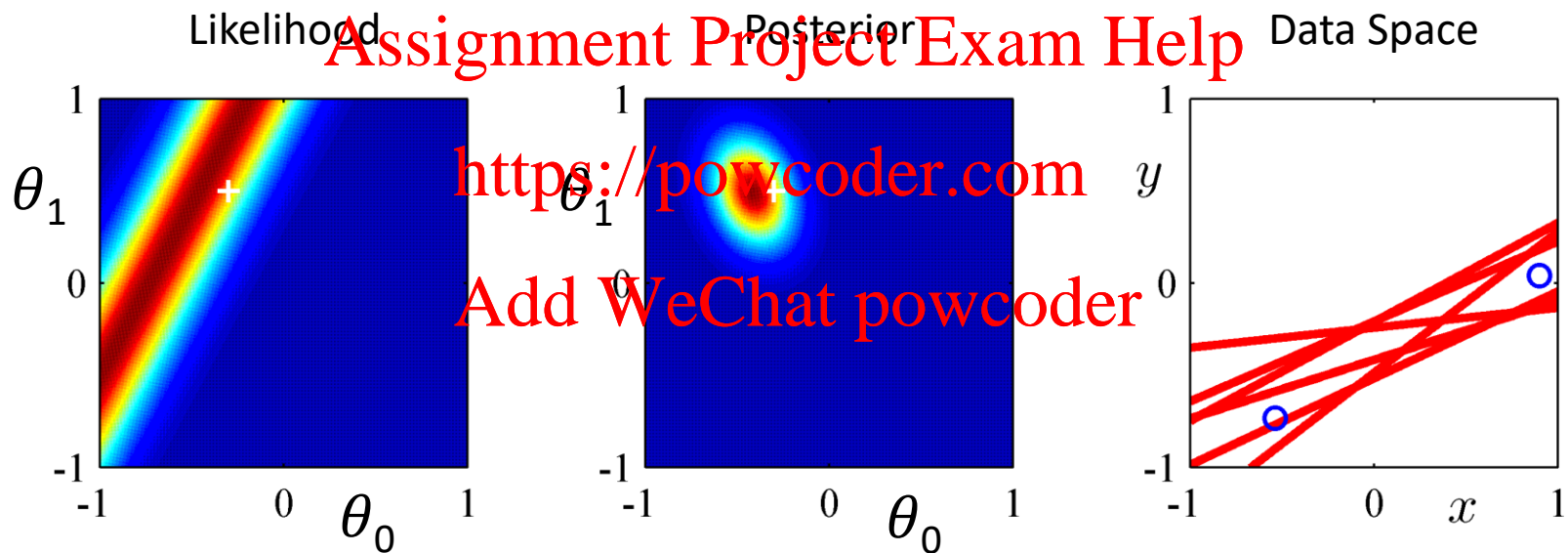
# Bayesian Linear Regression Example

1 data point observed



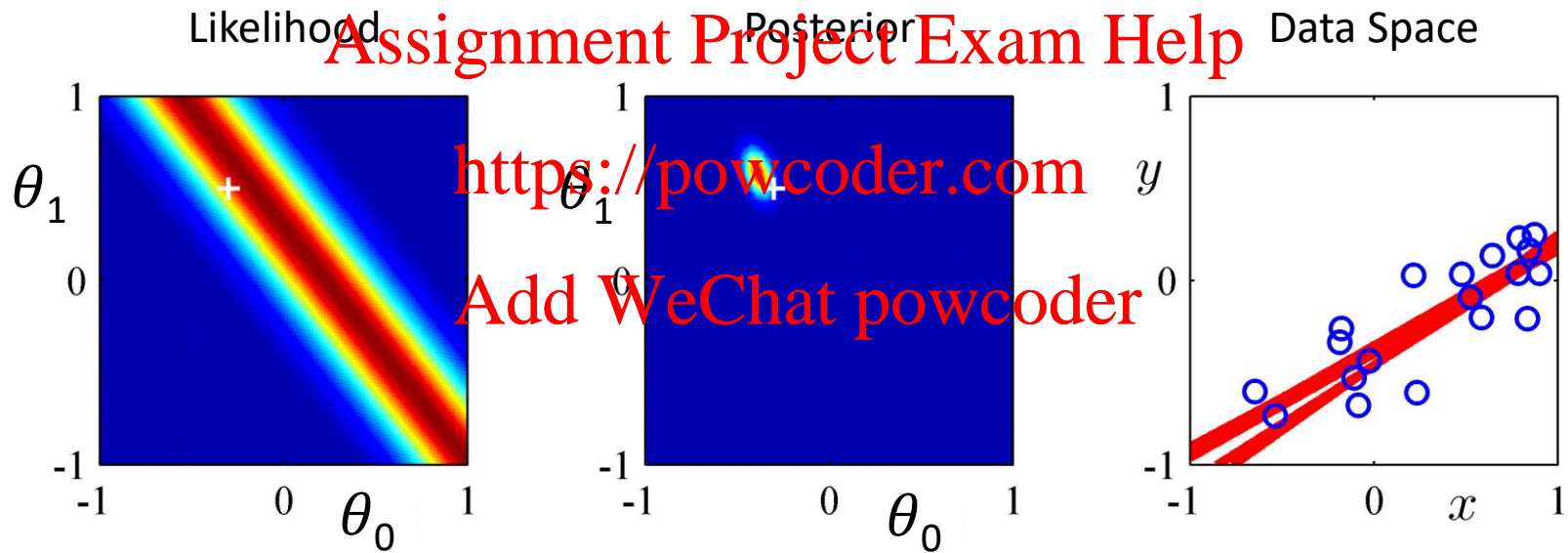
# Bayesian Linear Regression Example

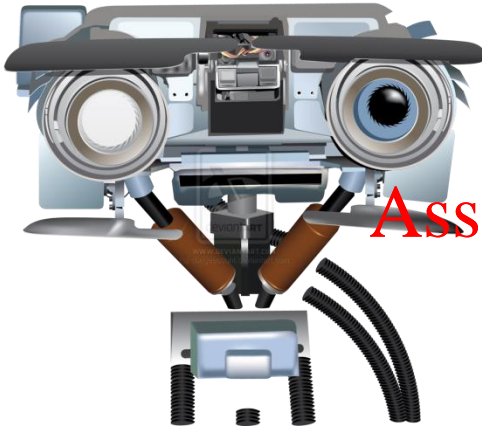
2 data points observed



# Bayesian Linear Regression Example

20 data points observed





# Bayesian Linear Regression

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Prediction

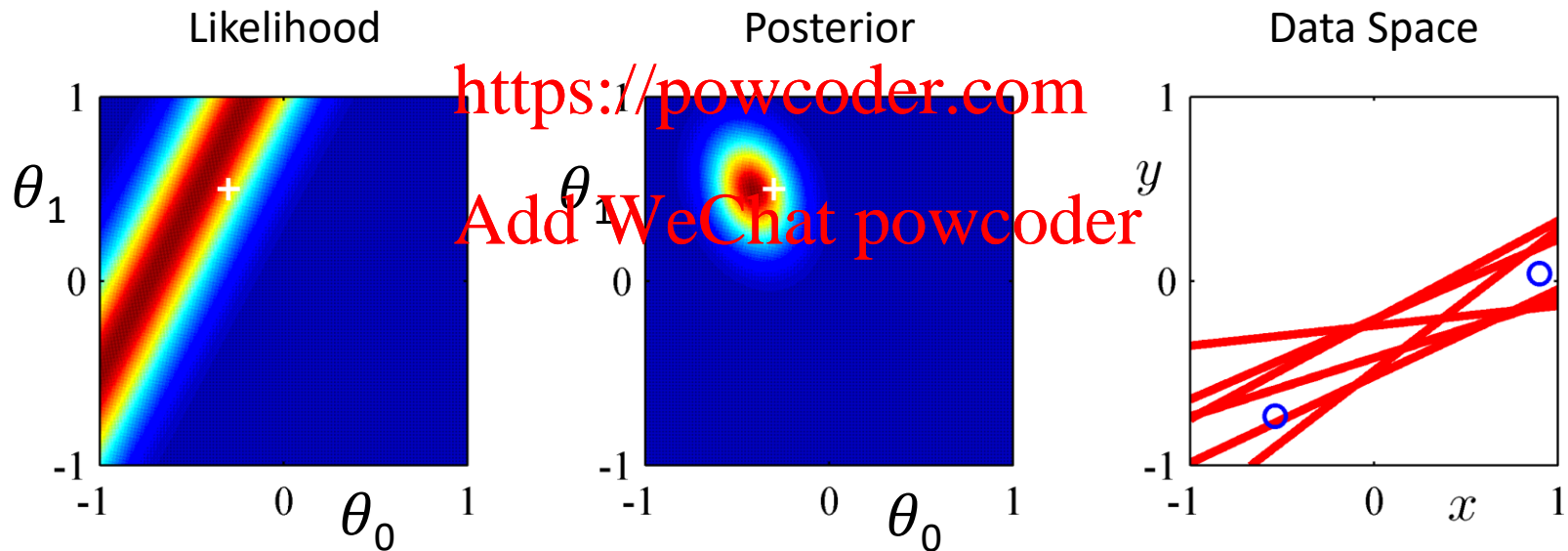
# Prediction

- Now that we have a Bayesian model, how do we use it to make predictions for new data points?

Assignment Project Exam Help

<https://powcoder.com>

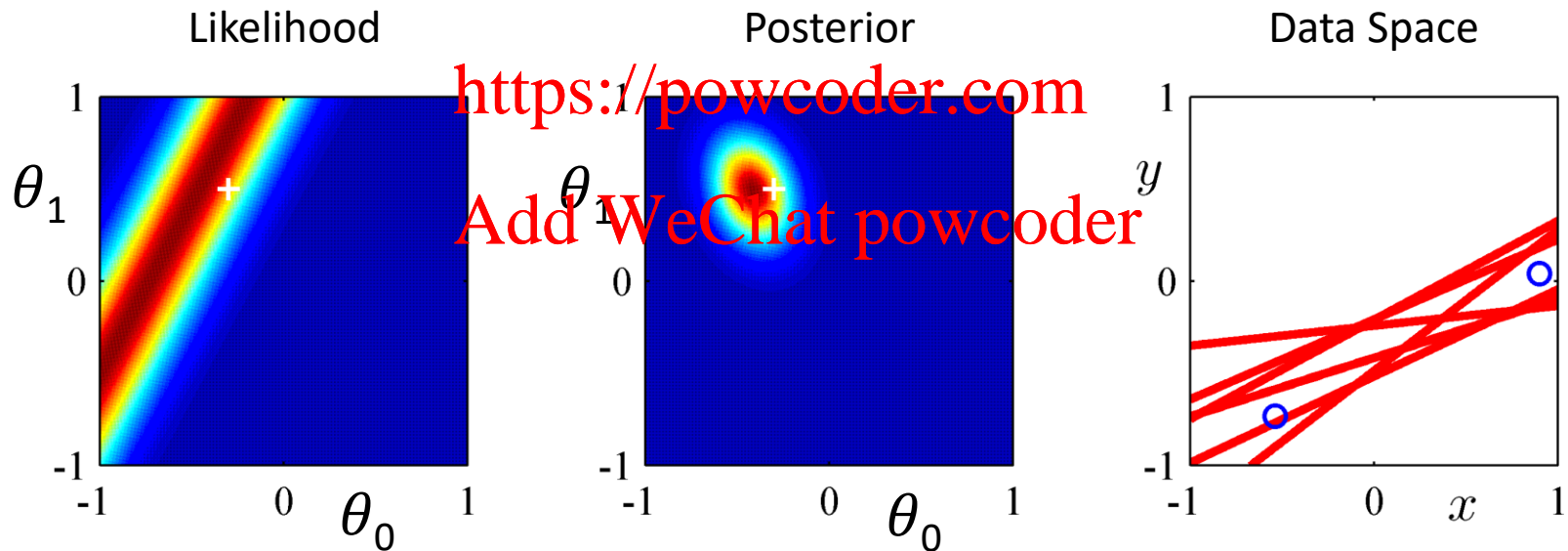
Add WeChat powcoder



# Prediction

- One way is to maximize the **posterior** to get an estimate of  $\theta_*$
- Then, plug  $\theta_*$  into the predictive distribution
- This is known as the **maximum a posteriori** estimate

Assignment Project Exam Help





# Maximum A Posteriori (MAP)

Output the parameter that maximizes its posterior distribution given the data

$$\theta_{MAP} = \underset{\theta}{\operatorname{argmax}} p(\theta | \mathbf{t})$$

Recall: for our prior  $p(\theta) = N(\theta | \mathbf{0}, \alpha^{-1} \mathbf{I})$ ,

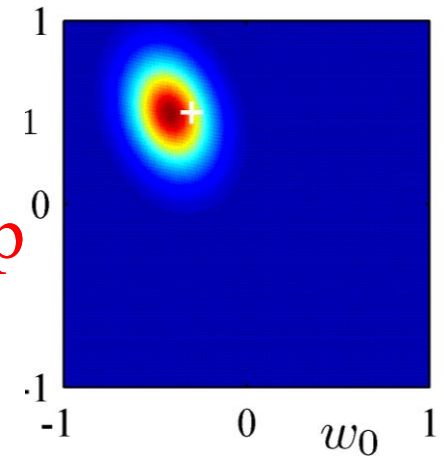
the posterior is  $p(\theta | \mathbf{t}) = N(\theta | \mathbf{m}_N, \mathbf{S}_N^{-1})$ ,

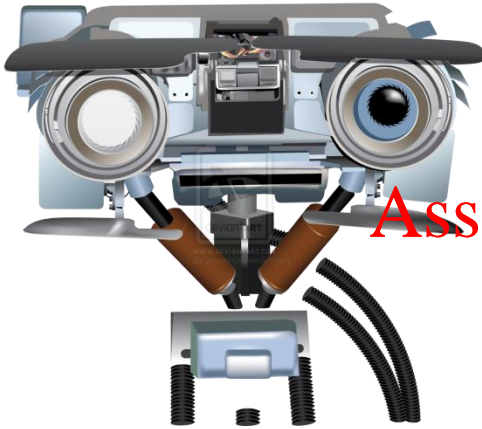
$$\text{where } \mathbf{m}_N = \beta \mathbf{S}_N \mathbf{X}^T \mathbf{t}, \quad \mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \mathbf{X}^T \mathbf{X}.$$

$$\text{Therefore, } \theta_{MAP} = \underset{\theta}{\operatorname{argmax}} p(\theta | \mathbf{t}) = \left( \mathbf{X}^T \mathbf{X} + \frac{\alpha}{\beta} \mathbf{I} \right)^{-1} \mathbf{X}^T \mathbf{y}$$

Same as solution to regularized regression with  $\|\theta\|^2$  term.

Note, this is the **mode** of the distribution





# Bayesian Regression

Assignment Project Exam Help

<https://powcoder.com>

---

Add WeChat powcoder

Connection to Regularized  
Linear Regression

# Maximizing posterior leads to regularized cost function

Joint likelihood of both training data and parameter

$$\log p(\mathcal{D}, \theta) = \sum_n \log p(y_n | \mathbf{x}_n, \theta) + \log p(\theta)$$

Assignment Project Exam Help

$$= -\frac{\sum_n (\theta^T \mathbf{x}_n - y_n)^2}{2\beta} - \sum_d \frac{1}{2\alpha^{-2}} \theta_d^2 + \text{const}$$

where  $\beta^{-2}$  is the noise variance and  $\alpha^{-2}$  is the prior covariance parameter

**Maximum a posterior (MAP) estimate:** we seek to maximize

$$\theta_{MAP} = \arg \max_{\theta} \log p(\theta | \mathcal{D}) \propto \log p(\mathcal{D}, \theta)$$

that is, the most likely  $\theta$  *conditioning* on observed training data  $\mathcal{D}$ .

# Maximizing posterior leads to regularized cost function

Can re-write the optimization in the same form as the **regularized linear regression** cost:

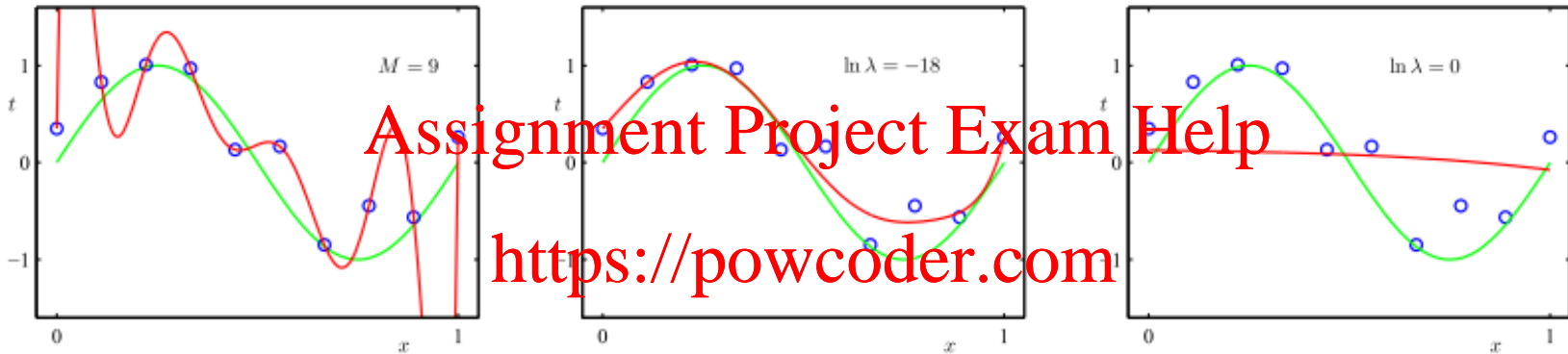
$$L(\boldsymbol{\theta}) = \sum (\boldsymbol{\theta}^T \mathbf{x}_n - y_n)^2 + \lambda \|\boldsymbol{\theta}\|^2$$

where  $\lambda = \beta^{-2} / \alpha^{-2}$  corresponds to the regularization hyperparameter.

- Intuitively, as  $\lambda \rightarrow +\infty$ , then  $\beta^{-2} \gg \alpha^{-2}$ . That is, the variance of noise is far greater than what our prior model can allow for  $\theta$ . In this case, our prior would be more accurate than what data can tell us, so we are getting a simple model, where  $\theta_{MAP} \rightarrow 0$ .
- If  $\lambda \rightarrow 0$ , then  $\beta^{-2} \ll \alpha^{-2}$ , and we trust our data more, so the MAP solution approaches the maximum likelihood solution, i.e.  $\theta_{MAP} \rightarrow \theta_{ML}$ .

# Effect of lambda

Overfitting is reduced from complex model to simpler one with the help of increasing regularizers

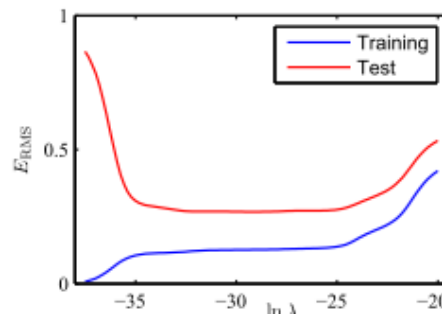


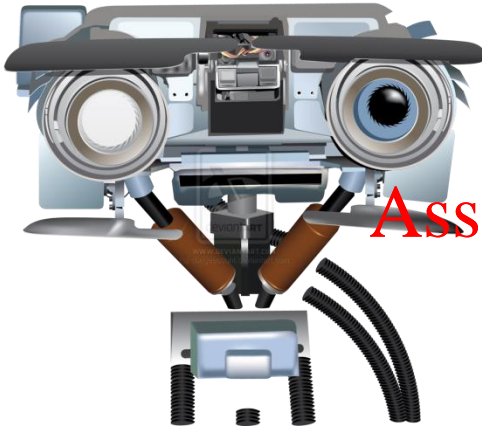
Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

$\lambda$  vs. residual error shows the difference of the model performance on training and testing dataset





# Bayesian Predictive Distribution

Assignment Project Exam Help

<https://powcoder.com>

---

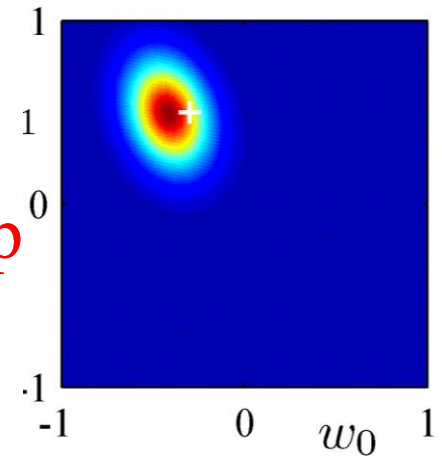
Add WeChat powcoder

# Maximum A Posteriori (MAP)

- Output the parameter that maximizes its posterior distribution given the data

$$\theta_{MAP} = \underset{\theta}{\operatorname{argmax}} p(\theta | t)$$

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



- Note, this is the ~~mode~~ of the distribution
- However, sometimes we may want to hedge our bets and **average (integrate)** over all possible parameters, e.g. if the posterior is multi-modal

# Bayesian Predictive Distribution

- Predict  $t$  for new values of  $x$  by integrating over  $\theta$  :

$$p(t|x, \mathbf{t}, \alpha, \beta) = \int p(t|\theta, \beta) p(\theta|\mathbf{t}, x, \beta) d\theta$$

$$= N(t|m_N^T x, \sigma_N^2(x))$$

- where

$$\sigma_N^2(x) = \frac{1}{\beta} + x^T S_N x$$

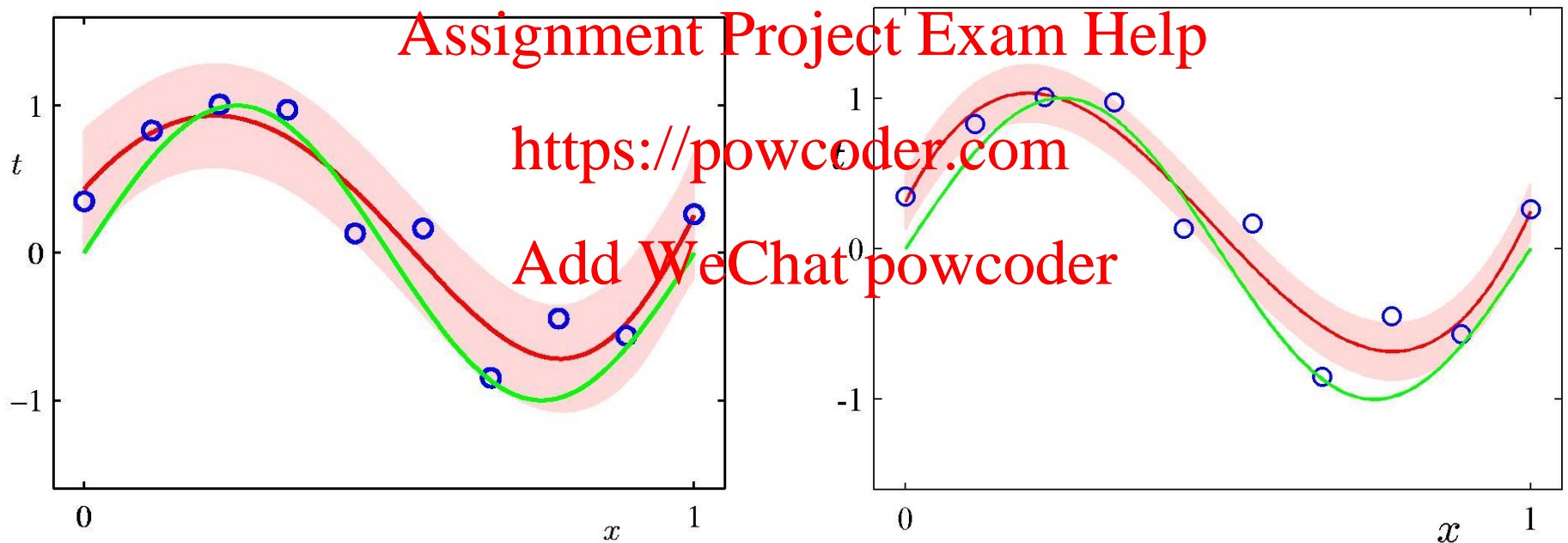


# What does it look like?

Compare to **Maximum Likelihood**:

$$p(t|x, \mathbf{x}, \mathbf{t}) = N(t|m_N^T x, \sigma_N^2)$$

$$p(t|x, \theta_{ML}, \beta_{ML}) = N(t|\theta_{ML}^T x, \beta_{ML}^{-1})$$



# Next Class

## Support Vector Machines I

maximum margin methods; support vector machines; primal vs dual SVM formulation; Hinge loss vs. cross-entropy loss

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

**Reading:** Bishop Ch 7.1.1-7.1.2