

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

L10 – Parameter Estimation

# Parameter Estimation

*Maximum Likelihood and Bayes Estimates*

The following lectures expand on our earlier discussion of parameter estimates, introducing some formal grounding. (A good supplemental source for this is chapter 2 of *Neural Networks for Pattern Recognition* by Pate ishop.)

<https://eduassistpro.github.io/>

We'll discuss parameter estimation in more detail, including mixture models for den

# Parametric Density Models

A model of specific functional form.

A small number of parameters that are estimated from data.

Assignment Project Exam Help

e.g. Normal distri

$$p(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} (x - \mu)^2\right)$$

Add WeChat edu\_assist\_pro

Data -  $D = \{x_1, x_2, x_3, \dots, x_m\}$

Parameter estimates

$$\hat{\mu} = \frac{1}{m} \sum_{i=1}^m x_i, \quad \overline{\quad}_{m-1}$$

but where did these forms for the estimators come from?

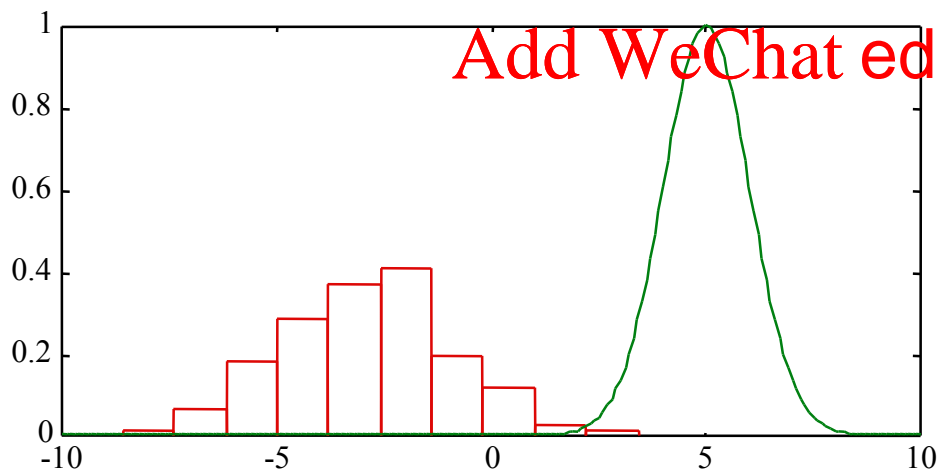
# Maximum Likelihood Estimation

Question -- what's the probability that the dataset D occurs, given the form of the model density?

We assume each of the  $x_i$  are sampled independently from the underlying (normal in this example) distribution, then

$$p(D | \mu, \sigma^2) = p(\{x_1, \dots, x_m\} | \mu, \sigma^2) = \prod_{i=1}^m p(x_i | \mu, \sigma^2)$$

Assignment Project Exam Help  
<https://eduassistpro.github.io/>



ta (histogram) and model look like this, the data will have **low probability** given the model (data at the tails of the model distribution).

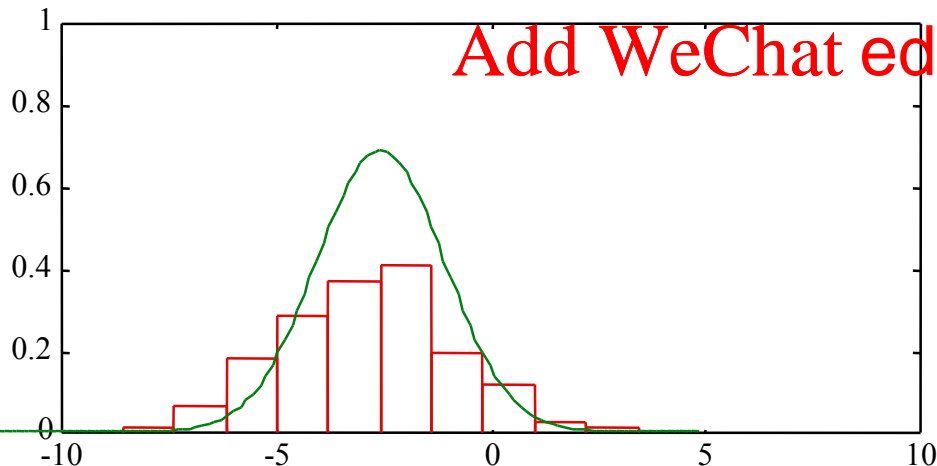
# Maximum Likelihood Estimation

Question -- what's the probability that the dataset D occurs, given the form of the model density?

We assume each of the  $x_i$  are sampled independently from the underlying (normal in this example) distribution, then

$$p(D | \mu, \sigma^2) = p(\{x_1, \dots, x_m\} | \mu, \sigma^2) = \prod_{i=1}^m p(x_i | \mu, \sigma^2)$$

**Assignment Project Exam Help**  
<https://eduassistpro.github.io/>



**Add WeChat edu\_assist\_pro**

adjust the model mean so that the data has higher probability under the model. But this model still has low tails where there's plenty of data, because the peak is too sharp.

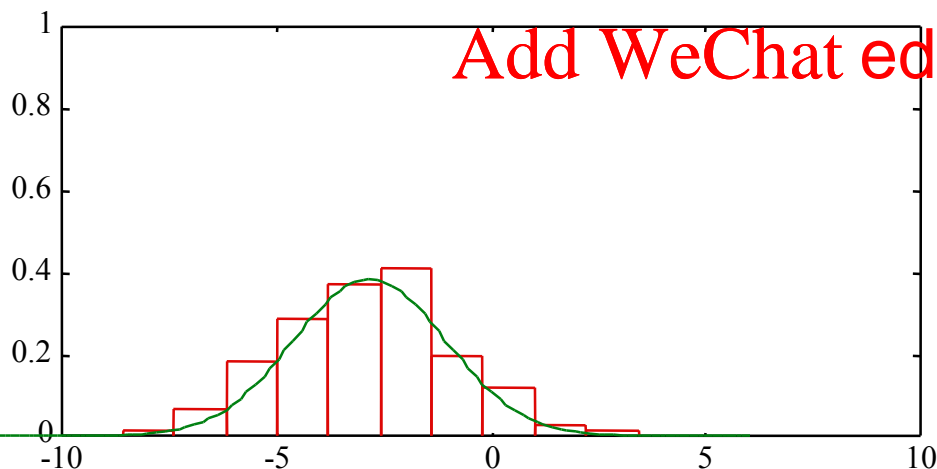
# Maximum Likelihood Estimation

Question -- what's the probability that the dataset  $D$  occurs, given the form of the model density?

We assume each of the  $x_i$  are sampled independently from the underlying (normal in this example) distribution, then

$$p(D | \mu, \sigma^2) = p(\{x_1, \dots, x_m\} | \mu, \sigma^2) = \prod_{i=1}^m p(x_i | \mu, \sigma^2)$$

**Assignment Project Exam Help**  
<https://eduassistpro.github.io/>



**Add WeChat edu\_assist\_pro**

Gaussian model that chooses the mean and variance so as to maximize the data likelihood under the model.

# Maximum Likelihood Estimation

So, we adjust the model parameters to maximize the data likelihood. Since the log is monotonic in its arguments, and we often deal with model distributions from the exponential family, it's convenient to maximize the log-likelihood.

Assignment Project Exam Help

$$L = \ln p(D | \mu, \sigma^2) = \sum_{i=1}^m \ln p(x_i | \mu, \sigma^2) = \sum_{i=1}^m \left[ -\frac{1}{2} \ln(2\pi \sigma^2) - \frac{1}{2\sigma^2} (x_i - \mu)^2 \right]$$

<https://eduassistpro.github.io/>

$$\left. \frac{\partial L}{\partial \mu} \right| = 0 = \frac{1}{\sigma^2} \sum_{i=1}^m (x_i - \mu) \Rightarrow \hat{\mu} = \frac{1}{m} \sum_{i=1}^m x_i$$

Add WeChat edu\_assist\_pro

$$\left. \frac{\partial L}{\partial \sigma^2} \right| = 0 \Rightarrow \hat{\sigma}^2 = \frac{1}{m} \sum_{i=1}^m (x_i - \hat{\mu})^2 \quad \text{Note } \sigma^2 \text{ is biased.}$$

# *Data Distributions and Cost Functions*

Regression - Minimizing mean square error between the data and a regression curve is equivalent to maximizing the data likelihood under the assumption that the fitting error is Gaussian.

Assignment Project Exam Help

The data is the sequence of  $(x, y)$  coordinates. The data  $y$  values are assumed Gaussian  $g(x)$ . That is

<https://eduassistpro.github.io/>

$$y = g(x) + \epsilon$$

Add WeChat edu\_assist\_pro

The data likelihood is

$$p(\{y_i\} \mid \{x_i\}; g(x), \sigma^2) = \prod_{i=1}^m \frac{1}{\sqrt{(2\pi\sigma^2)}} \exp\left[-\frac{1}{2\sigma^2} (y_i - g(x_i))^2\right]$$



# *Data Distributions and Cost Functions*

## *Regression*

Maximizing the data log-likelihood  $L$  with respect to  $g(x)$

Assignment Project Exam Help

$$L = \log(p(\{y\} \mid \{x\}; g(x), \sigma^2)) = \sum_{i=1}^m \left[ -\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (y_i - g(x_i))^2 \right]$$

<https://eduassistpro.github.io/>

is equivalent to minimizing the squared fitting error with respect to  $g(x)$ .

$$E = \sum_{i=1}^m \frac{1}{2\sigma^2} (y_i - g(x_i))^2$$

# *Data Distributions and Cost Functions*

## *Classification*

For a (two-class) classification problem, it's natural to write the data likelihood as a product of Bernouli distributions (since the target values are  $y = 0$  or  $1$  for each example)

$$L = p(\{y_i\} | \{x_i\}; \alpha(x)) = \prod_{i=1}^n \alpha(x_i)^{y_i} (1 - \alpha(x_i))^{(1-y_i)}$$

where  $\alpha(x)$  is the prob  
(rather than 0).

<https://eduassistpro.github.io/>  
ctor  $x$ , the class label is  $1$

Add WeChat edu\_assist\_pro

Maximizing this data likelihood is equivalent to minimizing its  $-\log$ , the cross entropy error

$$E = \sum_{i=1}^m y_i \log(\alpha(x_i)) + (1 - y_i) \log(1 - \alpha(x_i))$$

# Bayesian Estimation and Parameter Posterior Distributions

Maximum likelihood estimation --- there exists an actual value of the parameters  $\Theta_0$ , that we *estimate* by maximizing the probability of the data conditioned on the parameters

$$\Theta_0 = \arg \max_{\Theta} p(D | \Theta)$$

An arguably more naive approach is to assume that the parameters are the most probable values of the parameters, conditioned on the available data. Thus, parameters are regarded as random with their own distribution -- the posterior distribution

$$p(\Theta | D) = \frac{p(D | \Theta) P(\Theta)}{p(D)}$$

where  $P(\Theta)$  is the *prior* on  $\Theta$

# Maximum A Posterior Estimation

Maximizing the log of the posterior, with respect to the parameters, gives the maximum a posterior (or MAP) estimate

$$\hat{\Theta} = \arg \max_{\Theta} [ \log(p(D | \Theta)) + \log(P(\Theta)) ]$$

Assignment Project Exam Help

The prior distribution  $P(\Theta)$  is independent of  $\Theta$  (flat prior) at if the prior is independent of  $\Theta$  (flat prior) then the maximum likelihood estimates are the same.

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

Convenient to choose the prior distribution  $P(\Theta)$  so that the consequent posterior  $p(\Theta | D)$  has the same functional form as  $P(\Theta)$ . (The proper form depends, of course, on the form of the likelihood function  $p(D | \Theta)$ .)

Called conjugate priors.

# *Example: Posterior Distribution of the Gaussian Model's Mean*

Suppose the data is Gaussian, and the variance is known, but we just want to estimate the mean. The conjugate prior for this is a Gaussian. The posterior on the mean is

$$\begin{aligned} p(\mu | D, \sigma^2) &= \frac{p(D | \mu, \sigma^2) p(\mu)}{p(D)} \\ &= \frac{1}{p(D) \left(\sqrt{2\pi\sigma^2}\right)^m \sqrt{2\pi\lambda^2}} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^m (x_i - \mu)^2 + \frac{1}{2\lambda^2} (\mu - \mu_0)^2\right\} \end{aligned}$$

where,  $\lambda$  and  $\mu_0$  are the variance and mean of the prior distribution on  $\mu$ .

# Example: Posterior Distribution of the Gaussian Model's Mean

After some algebraic manipulation, we can rewrite the posterior dist. as:

$$p(\mu | D, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma_\mu^2}} \exp\left(-\frac{1}{2\sigma_\mu^2}(\mu - \bar{\mu})^2\right)$$

<https://eduassistpro.github.io/>

with the posterior mean  $\bar{\mu}$  (show this!)

Add WeChat edu\_assist\_pro

$$\bar{\mu} = \frac{m\lambda^2}{m\lambda^2 + \sigma^2} \frac{1}{m} \sum_{i=1}^m x_i + \frac{\mu_0}{m\lambda^2 + \sigma^2}$$

$$\sigma_\mu^2 = \frac{\sigma^2 \lambda^2}{m\lambda^2 + \sigma^2}$$

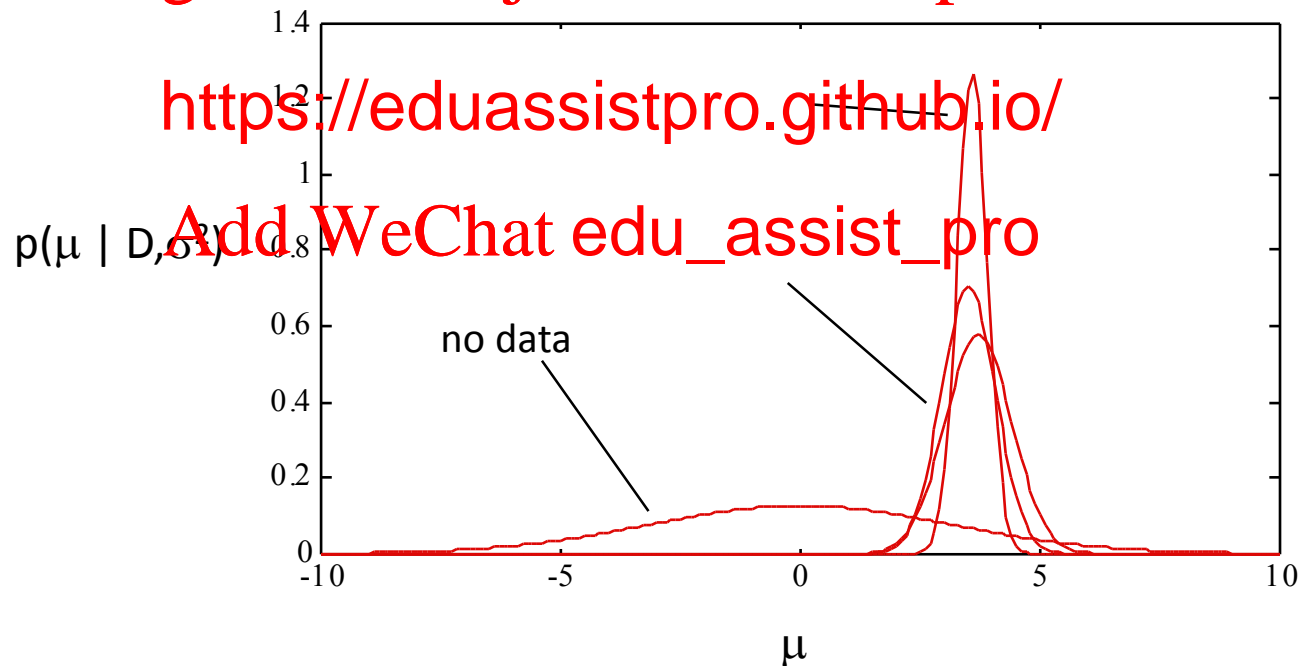
# *Example: Posterior Distribution of the Gaussian Model's Mean*

$$\bar{\mu} = \frac{m\lambda^2}{m\lambda^2 + \sigma^2} \frac{1}{m} \sum_{i=1}^m x_i + \frac{\sigma^2}{m\lambda^2 + \sigma^2} \mu_0$$
$$\sigma_{\mu}^2 = \frac{\sigma^2 \lambda^2}{m\lambda^2 + \sigma^2}$$

Note that for  $m \gg \sigma^2 / \lambda^2$  the posterior mean approaches the sample mean (the ML estimate), and the posterior variance becomes small.

# Example: Posterior Distribution of the Gaussian Model's Mean

Without data,  $m=0$ , the posterior is just the original prior on  $\mu$ . As we add samples, the posterior remains Gaussian (that's the point of a conjugate prior) but its mean and variance change in response to the data.





Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro