

Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro

L11 – Map Estimates, Bayesian  
Inference, Hyperparameter Choice

# *Continuing with Bayesian Methods*

Assignment Project Exam Help  
MAP Estimates, Bayesian Inference,  
and Hy <https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# *Why use a MAP Estimate, they're Biased?*

Consider the expected squared error of any estimator:

$$MSE = E[(\bar{\mu} - \mu)^2] = E[(\bar{\mu} - E[\bar{\mu}] + (E[\bar{\mu}] - \mu))^2]$$

$$= E[(\bar{\mu} - E[\bar{\mu}])^2] + (E[\bar{\mu}] - \mu)^2$$

$$= \text{var}(\bar{\mu}) + \text{bias}^2$$

Bias isn't the only consideration - variance is also important. There's usually a trade-off; increase the bias, and the variance drops, and vice-versa.

# Bias-Variance Trade-Off and MAP Estimates

Let's go back to our MAP estimate of the mean for Gaussian data:

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

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

The bias and variance are (show t

Add WeChat [edu\\_assist\\_pro](#)

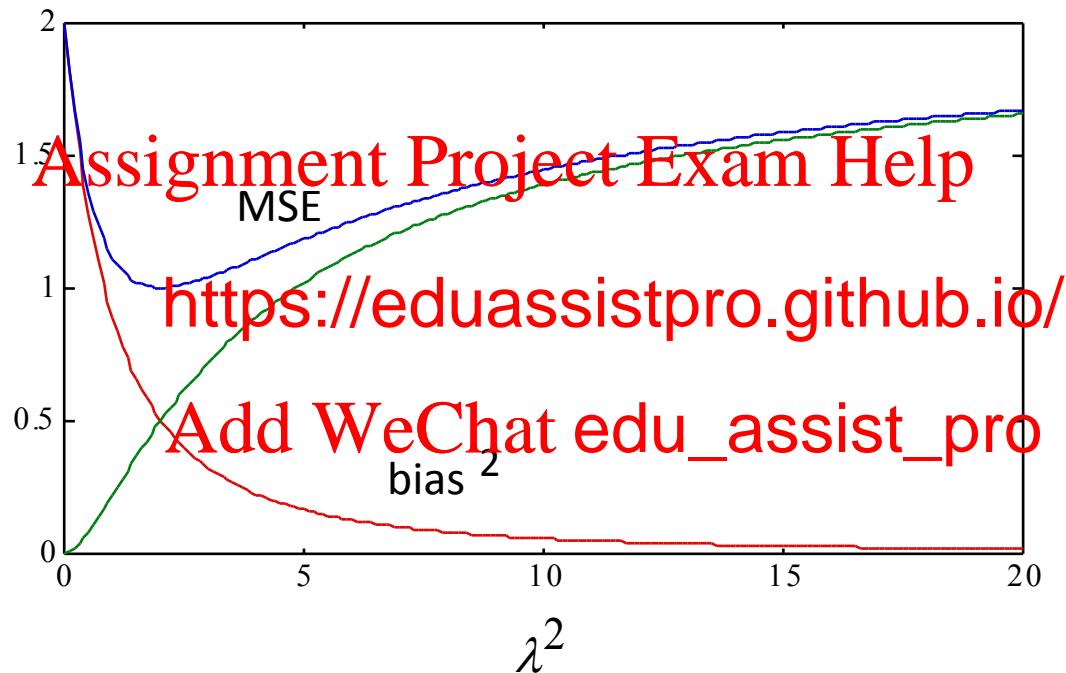
$$bias^2 = (E[\bar{\mu}] - \mu)^2 = \left( \frac{\sigma^2}{m\lambda^2 + \sigma^2} (\mu_0 - \mu) \right)^2$$

$$var(\bar{\mu}) = E[(\bar{\mu} - E[\bar{\mu}])^2] = \left( \frac{m\lambda^2}{m\lambda^2 + \sigma^2} \right)^2 \frac{\sigma^2}{m}$$

As  $m \rightarrow \infty$   
both go to zero

# *Bias-Variance Trade-Off and MAP Estimates*

The curves look like this



The curve of MSE has its minimum at a non-zero value of  $\lambda$ . Specifically --  $\lambda_{opt}^2 = (\mu_0 - \mu)$

# MAP Estimates and Regularizers

The log of the posterior on the parameters is

$$\log(p(\Theta | D)) = \log(p(D | \Theta)) + \log(p(\Theta)) - \log(p(D))$$

~~Assignment Project Exam Help~~

log-likelihood

log-posterior -- regularizer

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

We saw that maximizing the data log-likelihood is equivalent to minimizing some nice cost function -- e.g. the mean-squared-error.

Maximizing the log-posterior is equivalent to minimizing a regularized cost function. **The effect of the regularizer is to reduce the parameter variance at the cost of adding parameter bias.**

# MAP Regression

One can use the MAP estimate of  $\Theta$ , and construct the regression function

$$E[t | x, D] = f(t | x, \hat{\Theta})$$

where  $\hat{\Theta}$  is the value that maximizes the posterior  $p(\Theta|D)$ .

One can also use the target density

$$p(t | x, \hat{\Theta}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp - \frac{1}{2\sigma^2} (t - f(t | x, \hat{\Theta}))^2$$

# Example of Map Regression – Ridge

Ridge regression uses a parameterized regressor  $f(x, \theta)$ , the familiar SSE cost function (Gaussian likelihood for the targets), and a Gaussian prior on the parameters, typically centered at zero

$$p(\theta) = \frac{1}{\sqrt{2\pi/\Lambda}} \exp\left(-\frac{1}{2\Lambda} |\theta|^2\right)$$

The regularized cost function is <https://eduassistpro.github.io/>

$$E(\Lambda, \theta) = \sum_{i=1}^M (t_i - f(x_i, \theta))^2 + \Lambda |\theta|^2$$

That for linear regression,  $f(x, \theta)$  is linear in  $\theta$  so  $E$  is quadratic in  $\theta$  and the cost function can be minimized in closed form (just like MLE estimation for linear regression).



# Bayesian Estimation

Let's continue. Suppose we have obtained the posterior on the parameters  $p(\Theta | D)$  and we wish to find the probability of a new data value  $x$ . A Bayesian says that you should calculate this from his version of the distribution  $p(x)$

Assignment Project Exam Help  
 $p(x|D) = \int p(x|\Theta) p(\Theta|D) d\Theta$

A Bayesian compute  $f(x)$  as  
<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

$$E[f | D] = \int f(x) p(x | D) dx = \iint f(x) p(x | \Theta) p(\Theta | D) d\Theta dx$$

# Bayesian and MAP Estimates

Relation to **MAP** Estimates: Suppose the posterior is sharply peaked up about its maximum value (the MAP estimate). Write a series expansion of  $p(x|\Theta)$  about the maximum and substitute into the integral

$$\begin{aligned}
 p(x|D) &= \int p(x|\Theta) p(\Theta|D) d\Theta = \int \left[ p(x|\hat{\Theta}) + \frac{dp(x|\Theta)}{d\Theta} \Big|_{\hat{\Theta}} (\Theta - \hat{\Theta}) \right. \\
 &\quad \left. + \frac{1}{2} \frac{d^2 p(x|\Theta)}{d\Theta^2} \Big|_{\hat{\Theta}} (\Theta - \hat{\Theta})^2 + \dots \right] p(\Theta|D) d\Theta \\
 &= p(x|\hat{\Theta}) + \frac{dp(x|\Theta)}{d\Theta} \Big|_{\hat{\Theta}} E[(\Theta - \hat{\Theta})|D] + \frac{1}{2} \frac{d^2 p(x|\Theta)}{d\Theta^2} \Big|_{\hat{\Theta}} E[(\Theta - \hat{\Theta})^2|D] + \dots
 \end{aligned}$$

# Bayesian and MAP Estimates

$$\begin{aligned}
 p(x|D) &= \int p(x|\Theta) p(\Theta|D) d\Theta = \int \left[ p(x|\hat{\Theta}) + \frac{1}{d\Theta} \left. \frac{dp(x|\Theta)}{d\Theta} \right|_{\hat{\Theta}} (\Theta - \hat{\Theta}) \right. \\
 &\quad \left. + \frac{1}{2} \frac{d^2 p(x|\Theta)}{d\Theta^2} \Big|_{\hat{\Theta}} (\Theta - \hat{\Theta})^2 + \dots \right] p(\Theta|D) d\Theta \\
 &= p(x|\hat{\Theta}) + \frac{1}{d\Theta} \left. \frac{dp(x|\Theta)}{d\Theta} \right|_{\hat{\Theta}} E[(\Theta - \hat{\Theta})|D] + \frac{1}{2} \frac{d^2 p(x|\Theta)}{d\Theta^2} \Big|_{\hat{\Theta}} E[(\Theta - \hat{\Theta})^2|D] + \dots
 \end{aligned}$$

Handwaving arg: As data increases, the posterior becomes more sharply peaked about the MAP value  $\hat{\Theta}$ , trailing terms will become small & integral is approximately

$$p(x|D) \approx p(x|\hat{\Theta})$$

# Recursive Bayesian Estimation

Back to Bayesian estimation of  $p(x|D)$

$$p(x|D) = \int p(x|\Theta) p(\Theta|D) d\Theta = \int p(x|\Theta) \frac{p(D|\Theta) p(\Theta)}{\int p(D|\Theta') p(\Theta') d\Theta'} d\Theta$$

Denote the data  $\{x_1, x_2, \dots, x_n\}$ ,  
and its likelihood

<https://github.com/eduassistpro>  
Add WeChat edu\_assist\_pro

$$p(D^n | \Theta) = \prod_{k=1}^n p(x_k | \Theta) = p(x_n | \Theta) p(D^{n-1} | \Theta)$$

Using the last expression, the posterior can be written

$$p(\Theta | D^n) = \frac{p(x_n | \Theta) p(\Theta | D^{n-1})}{\int p(x_n | \Theta') p(\Theta' | D^{n-1}) d\Theta'}$$

# *Recursive Bayesian Estimation*

We have written the posterior density for the n-sample data set as

$$p(\Theta | D^n) = \frac{p(x_n | \Theta) p(\Theta | D^{n-1})}{\int p(x_n | \Theta') p(\Theta' | D^{n-1}) d\Theta'}$$

Assignment Project Exam Help

Starting with zero and generate the  $p(\Theta | D^0) = p(\Theta)$   
<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](#)

and thus incrementally refine our estimate of the posterior density as more and more data becomes available.

# *How Does a Bayesian do Regression?*

Get a dataset  $D \equiv \{ (x_i, t_i) \mid i = 1, \dots, m \}$

Choose a parameterized regression function  $f(x; \Theta)$  to fit to the data.

Choose a model distribution function for the targets, e.g. a Gaussian

$$p(t \mid x, \sigma^2, \Theta) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[ -\frac{1}{2\sigma^2} (t - f(x; \Theta))^2 \right]$$

<https://eduassistpro.github.io/>  
Add WeChat edu\_assist\_pro

Choose a prior distribution on the parameters  $p(\Theta, \sigma^2)$ .

Calculate the data likelihood and the posterior distribution of the parameters

$$p(\Theta, \sigma^2 \mid D) = \frac{1}{p(D)} p(D \mid \Theta, \sigma^2) p(\Theta, \sigma^2)$$

# How Does a Bayesian Do Regression?

Calculate the target density as a function of  $x$  by integrating over the posterior distribution of the parameters

$$p(t | x, D) = \int p(t | x, \sigma^2, \Theta) p(\sigma^2, \Theta | D) d\sigma^2 d\Theta$$

From the distribution on  $t$ , we can calculate several quantities.

- The conditional mean  $E[t | x, D]$  and the conditional variance  $\text{var}(t | x, D)$  equal to

$$\begin{aligned} E[t | x, D] &= \int t p(t | x, D) dt \\ &= \int \int t p(t | x, \sigma^2, \Theta) p(\sigma^2, \Theta | D) d\sigma^2 d\Theta = \int f(x, \Theta) p(\Theta | D) d\Theta \end{aligned}$$

for our Gaussian model.)

- The most likely value(s) of  $t$   $\arg \max_t p(t | x, D)$
- The target variance  $\text{var}(t | x, D)$ .

# Hyperparameters and Model Selection

Our prior on model parameters is itself a parameterized distribution. Recall for our Gaussian density model, we put a prior on the distribution of the mean

$$p(\mu \mid \mu_0, \lambda^2) = \frac{1}{\sqrt{2\pi\lambda^2}} \exp\left[-\frac{1}{2\lambda^2} (\mu - \mu_0)^2\right]$$

But how were the hyperparameters chosen  $\mu_0, \lambda^2$



# Hyperparameter Selection

- We could calculate the likelihood function for particular values

$$p(D | \mu_0, \lambda^2) = \int p(D | \mu) p(\mu | \mu_0, \lambda^2) d\mu$$

and choose the values of the hyperparameters that maximizes it.

- We could set up a hyperprior on the hyperparameters and choose maximum hyperparameters by maximizing

$$p(\mu_0, \lambda^2 | D) \propto p(D | \mu_0, \lambda^2) p(\mu_0, \lambda^2)$$

(but the hyperprior is going to have its own parameters ...).

- Use some sort of empirical technique.

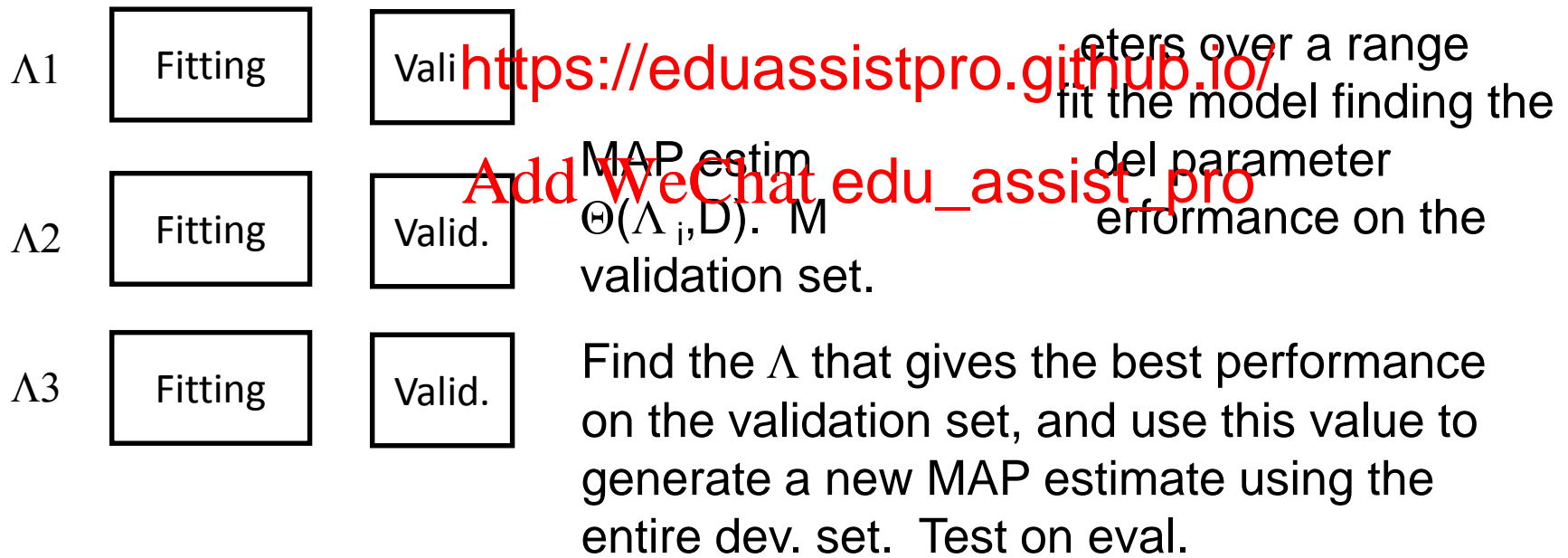
# Empirical Hyperparameter Selection

Using a 'validation' set and MAP estimates.

Divide data into two pieces, development and evaluation



Further divide the development set into fitting and validation



Assignment Project Exam Help

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

Add WeChat edu\_assist\_pro