# ANLY-601
## Advanced Pattern Recognition

Assignment Project Exam Help

*Spring 2018*

https://powcoder.com

Add WeChat powcoder

L11 – Map Estimates, Bayesian Inference, Hyperparameter Choice

GEORGETOWN
UNIVERSITY

# *Continuing with Bayesian Methods*

MAP Estimates, Bayesian Inference,
and Hyperparameter Choice

*GEORGETOWN
UNIVERSITY*

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

Consider the expected squared error of any estimator:

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

Assignment Project Exam Help

https://powcoder.com

$$= E\left[(\overline{\mu} - E[\overline{\mu}])^2\right] + (E[\overline{\mu}] - \mu)^2$$

Add WeChat powcoder

$$= \operatorname{var}(\overline{\mu}) + bias^2$$

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.

*GEORGETOWN UNIVERSITY*

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

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

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

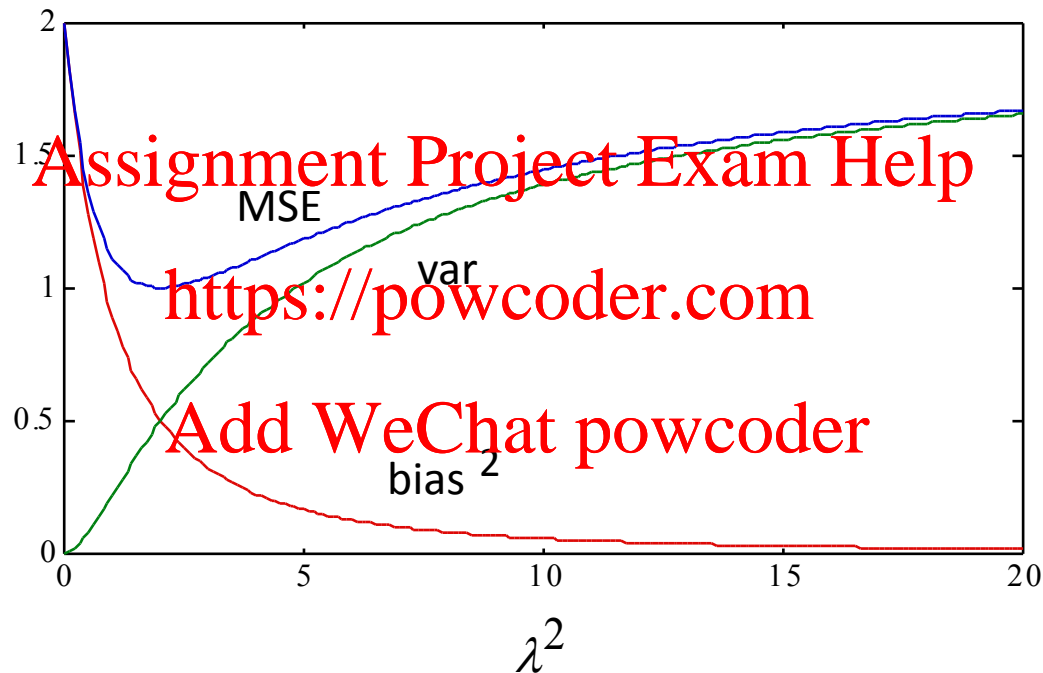The bias and variance are (show these!)

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

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

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 -- $\lambda_{opt}^2 = (\mu_0 - \mu)^2$

GEORGETOWN
UNIVERSITY

# MAP Estimates and Regularizers

The log of the posterior on the parameters is

$$\log\left(p(\Theta\,|\,D)\right) \;=\; \log\left(p(D\,|\,\Theta)\right) \;+\; \log\left(p(\Theta)\right) \;-\; \log\left(p(D)\right)$$

<span style="color:red">Assignment Project Exam Help</span>

log-likelihood – bare cost          log prior -- regularizer

<span style="color:red">https://powcoder.com</span>

<span style="color:red">Add WeChat powcoder</span>

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

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

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

One can also use this MAP value to estimate the target density

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

GEORGETOWN
UNIVERSITY

# *Example of Map Regression – Ridge*

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

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

The regularized cost function is thus

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

That for linear regression, *f(x,Θ)* 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).

*GEORGETOWN UNIVERSITY*

# *Bayesian Estimation*

Let's continue.  Suppose we have obtained the posterior on the parameters p(Θ| 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)

$$p(x|D) = \int p(x|\Theta)\, p(\Theta|D)\, d\Theta$$

 A Bayesian computes the mean of any function f(x) as

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

*GEORGETOWN UNIVERSITY*

# 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|Θ) about the maximum and substitute into the integral

$$p(x \mid D) = \int p(x \mid \Theta) p(\Theta \mid D) d\Theta = \int \left[ p(x \mid \hat{\Theta}) + \frac{dp(x \mid \Theta)}{d\Theta} \bigg|_{\hat{\Theta}} (\Theta - \hat{\Theta}) \right.$$

$$\left. + \frac{1}{2} \frac{d^2 p(x \mid \Theta)}{d\Theta^2} \bigg|_{\hat{\Theta}} (\Theta - \hat{\Theta})^2 + ... \right] p(\Theta \mid D) \ d\Theta$$

$$= p(x \mid \hat{\Theta}) + \frac{dp(x \mid \Theta)}{d\Theta} \bigg|_{\hat{\Theta}} E\big[(\Theta - \hat{\Theta}) \mid D\big] + \frac{1}{2} \frac{d^2 p(x \mid \Theta)}{d\Theta^2} \bigg|_{\hat{\Theta}} E\big[\big(\Theta - \hat{\Theta}\big)^2 \mid D\big] + ...$$

*GEORGETOWN UNIVERSITY*

# Bayesian and MAP Estimates

$$p(x \mid D) = \int p(x \mid \Theta) \, p(\Theta \mid D) \, d\Theta = \int \left[ p(x \mid \hat{\Theta}) + \left. \frac{dp(x \mid \Theta)}{d\Theta} \right|_{\hat{\Theta}} (\Theta - \hat{\Theta}) \right.$$

$$\left. + \frac{1}{2} \left. \frac{d^2 p(x \mid \Theta)}{d\Theta^2} \right|_{\hat{\Theta}} (\Theta - \hat{\Theta})^2 + ... \right] p(\Theta \mid D) \, d\Theta$$

$$= p(x \mid \hat{\Theta}) + \left. \frac{dp(x \mid \Theta)}{d\Theta} \right|_{\hat{\Theta}} E\left[ (\Theta - \hat{\Theta}) \mid D \right] + \frac{1}{2} \left. \frac{d^2 p(x \mid \Theta)}{d\Theta^2} \right|_{\hat{\Theta}} E\left[ (\Theta - \hat{\Theta})^2 \mid D \right] + ...$$

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

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

# *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 dataset with n points by $D^n$ = {x1, x2, …, xn}, and its likelihood by

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

Using the last expression, the posterior can be written

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

*GEORGETOWN UNIVERSITY*

# *Recursive Bayesian Estimation*

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

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

Starting with zero data, we take $\quad p(\Theta \mid D^0) = p(\Theta)$
and generate the sequence

$$p(\Theta \mid D^1),\ p(\Theta \mid D^2),\ ...$$

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

*GEORGETOWN UNIVERSITY*

# *How Does a Bayesian do Regression?*

Get a dataset $D \equiv \left\{ \left( x_i, t_i \right) \ i = 1, \ldots, m \right\}$

Choose a parameterized regression function f(x;Θ) to fit to the data.

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

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

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

Calculate the data likelihood and the posterior distribution of the parameters

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

*GEORGETOWN UNIVERSITY*

# *How Does a Bayesian Do Regression?*

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

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

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

- The conditional mean  $E[\, t \mid x, D \,]$  )  (called the <u>regressor</u>, and equal to

$$E[t \mid x, D] = \int t \; p(t \mid x, D) \, dt$$

$$= \int t \, p(t \mid x, \sigma^2, \Theta) \, dt \, d\sigma^2 \, d\Theta = \int f(x; \Theta) \, p(\Theta \mid D) \, d\Theta$$

  for our Gaussian model.)

- The most likely value(s) of  t      $\arg\max_{t} p(t \mid x, D)$

- The target variance var( t | x,D).

*GEORGETOWN UNIVERSITY*

# *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 <u>hyperparameters chosen</u> $\mu_0, \ \lambda^2$

GEORGETOWN
UNIVERSITY

# *Hyperparameter Selection*

- We could calculate the likelihood function for particular values

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

  and choose the values of the hyperparameters that maximizes it.

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

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

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

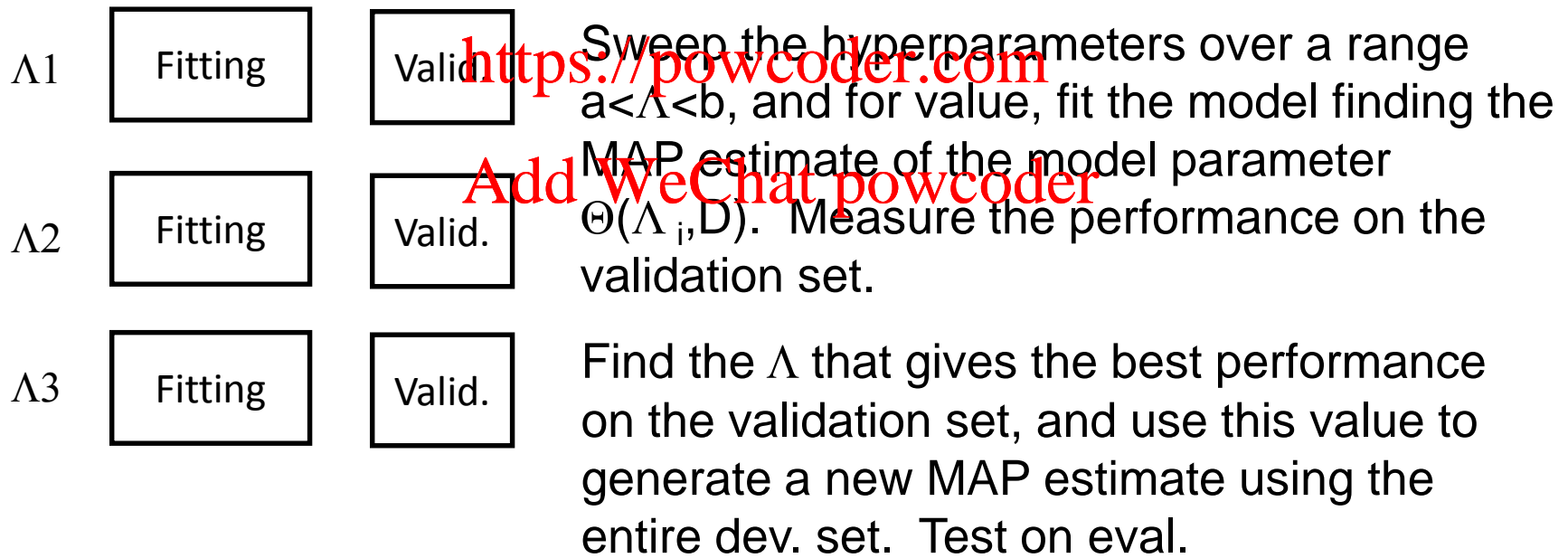- Use some sort of empirical technique.

*GEORGETOWN UNIVERSITY*

# *Empirical Hyperparameter Selection*

Using a 'validation' set and MAP estimates.
Divide data into two pieces, development and evaluation

| Development | | Eval. |
|---|---|---|

Further divide the development set into fitting and validation

$\Lambda 1$  | Fitting | Valid. |

Sweep the hyperparameters over a range $a<\Lambda<b$, and for value, fit the model finding the MAP estimate of the model parameter $\Theta(\Lambda_i, D)$. Measure the performance on the validation set.

$\Lambda 2$  | Fitting | Valid. |

$\Lambda 3$  | Fitting | Valid. |

Find the $\Lambda$ that gives the best performance on the validation set, and use this value to generate a new MAP estimate using the entire dev. set. Test on eval.

*GEORGETOWN UNIVERSITY*

Assignment Project Exam Help

https://powcoder.com

Add WeChat powcoder

GEORGETOWN UNIVERSITY