



Regression  
1/9

D.E. Brown

Bayesian  
Regression

Problem Formulation

Regularization

Conjugate Approach

# Bayesian Least Squares Regression

Donald E. Brown

School of Data Science  
University of Virginia  
Charlottesville, VA 22904



# Bayesian Formulation

Regression  
3/9

D.E. Brown

Bayesian  
Regression

Problem Formulation  
Regularization  
Conjugate Approach

- Regression with response  $\mathbf{y}$ , data  $\mathbf{X}$ , and parameter  $\boldsymbol{\theta}$ :

$$\mathbf{y} = \mathbf{X}\boldsymbol{\theta} + \boldsymbol{\epsilon}$$

- The parameters are random variables, so  $\boldsymbol{\theta}$  is a random vector
- Given the prior is  $p(\boldsymbol{\theta})$  and data  $\{\mathbf{X}, \mathbf{y}\}$ , then the posterior is

$$p(\boldsymbol{\theta}|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{X}, \mathbf{y}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathbf{X}, \mathbf{y})}$$

- $\mathbf{X}$  is known (i.e., not random) &  $\mathbf{y}$  is a function of  $\mathbf{X}$  made random by  $\boldsymbol{\epsilon}$ . So will simply write  $p(\boldsymbol{\theta}|\mathbf{y})$
- Next estimate  $\boldsymbol{\theta}$



# Maximum A Posteriori (MAP) Estimate

Regression

4/9

D.E. Brown

Bayesian  
Regression

Problem Formulation

Regularization

Conjugate Approach

- Find the estimate that maximizes the posterior distribution for  $\theta$
- Conjugate prior approach:

$$p(\theta) \sim N(\theta_0, \Sigma_\theta)$$

$$f(\mathbf{y}|\theta, \sigma_\epsilon^2) \sim N(\mathbf{X}\theta, \sigma_\epsilon^2 \mathbf{I})$$

- The conditional posterior is Gaussian, and the marginal is a t distribution; for both, the mean is the max

$$E[\theta|\mathbf{y}] =$$

$$\theta_0 + (\Sigma_\theta^{-1} + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{y} - \mathbf{X}\theta_0)$$



# Bayesian Formulation as Regularization

Regression  
5/9

D.E. Brown

Bayesian  
Regression

Problem Formulation

Regularization

Conjugate Approach

- Prior acts to regularize the resulting estimate
- Suppose  $\Sigma_{\theta} = \sigma_{\theta}^2 \mathbf{I}$  and  $\theta_0 = \mathbf{0}$

$$\theta_{MAP} = (\lambda \mathbf{I} + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T y$$

where,  $\lambda = \sigma_{\epsilon}^2 / \sigma_{\theta}^2$

- This is ridge regression
- The ridge regression solution is

$$\theta_{Ridge} = (\lambda \mathbf{I} + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T y$$



# Regularization Parameter

Regression  
6/9

D.E. Brown

Bayesian  
Regression

Problem Formulation

Regularization

Conjugate Approach

- Choosing  $\lambda = \sigma_\epsilon^2 / \sigma_\theta^2$  determines the performance of the regularization
- Typically use grid search with cross-validation
- Bayesian formulation provides the basis for improved understanding through the posterior,  $p(\theta|\mathbf{y})$
- Instead of just the maximum, we may want all of  $p(\theta|\mathbf{y})$ , which means we also want the Bayesian denominator

$$f(\mathbf{y}) = \int f(\mathbf{y}|\theta)p(\theta)$$



# Conjugate Approach

Regression  
7/9

D.E. Brown

Bayesian  
Regression  
Problem Formulation  
Regularization  
Conjugate Approach

- Again, the conditional prior and likelihood are Gaussian  
 $p(\boldsymbol{\theta}) \sim N(\boldsymbol{\theta}_0, \boldsymbol{\Sigma}_\theta)$  &  $f(\mathbf{y}|\boldsymbol{\theta}) \sim N(\mathbf{X}\boldsymbol{\theta}, \sigma_\epsilon^2 \mathbf{I})$
- So,  $\mathbf{y}$ ,  $p(\mathbf{y}|\sigma_\epsilon^2) \sim N(\mathbf{X}\boldsymbol{\theta}_0, \sigma_\epsilon^2 \mathbf{I} + \mathbf{X}\boldsymbol{\Sigma}_\theta \mathbf{X}^T)$
- The conditional posterior is Gaussian,  
 $p(\boldsymbol{\theta}|\mathbf{y}, \sigma_\epsilon^2) \sim N(\boldsymbol{\mu}_{\theta|\mathbf{y}}, \boldsymbol{\Sigma}_{\theta|\mathbf{y}})$  where

$$\boldsymbol{\mu}_{\theta|\mathbf{y}, \sigma_\epsilon^2} = \boldsymbol{\theta}_0 + \frac{1}{\sigma_\epsilon^2} \left( \boldsymbol{\Sigma}_\theta^{-1} + \frac{1}{\sigma_\epsilon^2} \mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}_0)$$

$$\boldsymbol{\Sigma}_{\theta|\mathbf{y}, \sigma_\epsilon^2} = \left( \boldsymbol{\Sigma}_\theta^{-1} + \frac{1}{\sigma_\epsilon^2} \mathbf{X}^T \mathbf{X} \right)^{-1}$$

- This gives the distribution for the estimate,  $\hat{\boldsymbol{\theta}}$



# Prediction

Regression  
8/9

D.E. Brown

Bayesian  
Regression

Problem Formulation

Regularization

Conjugate Approach

- Goal is to predict  $y$  given new  $\mathbf{x}$
- From the conditional posterior,  $p(\boldsymbol{\theta}|\mathbf{y}, \sigma_\epsilon^2)$ , obtain

$$p(y|\mathbf{x}, \mathbf{y}, \sigma_\epsilon^2) = \int p(y|\mathbf{x}, \theta, \sigma_\epsilon^2) p(\boldsymbol{\theta}|\mathbf{y}) d\boldsymbol{\theta}$$

- From the regression model

$$p(y|\mathbf{x}, \theta, \sigma_\epsilon^2) \sim N(\mathbf{x}\boldsymbol{\theta}, \sigma_\epsilon^2)$$



# Distribution of the Prediction

Regression

9/9

D.E. Brown

Bayesian  
Regression

Problem Formulation

Regularization

Conjugate Approach

- For  $\Sigma_{\theta} = \sigma_{\theta}^2 \mathbf{I}$  the posterior parameters are

$$\mu_{\theta|y, \sigma_{\epsilon}^2} = \theta_0 + \frac{1}{\sigma_{\epsilon}^2} \left( \frac{1}{\sigma_{\theta}^2} \mathbf{I} + \frac{1}{\sigma_{\epsilon}^2} \mathbf{x}^T \mathbf{x} \right)^{-1} \mathbf{x}^T (\mathbf{y} - \mathbf{x} \theta_0)$$

$$\sigma_{\theta|y, \sigma_{\epsilon}^2}^2 = \left( \frac{1}{\sigma_{\theta}^2} \mathbf{I} + \frac{1}{\sigma_{\epsilon}^2} \mathbf{x}^T \mathbf{x} \right)^{-1}$$

- So  $p(y|\mathbf{x}, \mathbf{y}) \sim N(\mu_y, \sigma_y^2)$

$$\mu_y = \mathbf{x} \mu_{\theta|y, \sigma_{\epsilon}^2}$$

$$\sigma_y^2 = \sigma_{\epsilon}^2 + \sigma_{\epsilon}^2 \sigma_{\theta|y, \sigma_{\epsilon}^2}^2 \mathbf{x} \left( \sigma_{\epsilon}^2 \mathbf{I} + \sigma_{\theta|y, \sigma_{\epsilon}^2}^2 \mathbf{x}^T \mathbf{x} \right)^{-1} \mathbf{x}$$