

A NOTE ON BAYESIAN LINEAR REGRESSION [1–3]

Mohammad Shekaramiz and Todd K. Moon

Electrical and Computer Engineering Department and Information Dynamics Laboratory
Utah State University
{mohammad.shekaramiz@aggiemail.usu.edu, todd.moon@usu.edu}

Abstract– In this report, we briefly discuss Bayesian linear regression as well as the proof for the inference to perform prediction based on the training data using this technique.

1. MODEL DESCRIPTION

- **Training data:** input-output pairs $D = \{(\mathbf{x}_m, y_m) | m = 1, 2, \dots, M\}$
- Each input is a vector \mathbf{x}_m of dimension N
- Suppose that the training data D be a set of *i.i.d.* samples from some unknown distribution
- The standard probabilistic interpretation of linear regression states that

$$y_m = \boldsymbol{\theta}^T \mathbf{x}_m + \epsilon_m, \quad m = 1, 2, \dots, M, \quad (1)$$

where ϵ is the noise and $\epsilon_m \sim \mathcal{N}(0, \sigma^2)$

- For notational convenience, define

$$X := \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_m^T \end{pmatrix}, \mathbf{y} := \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{pmatrix}, \text{ and } \boldsymbol{\epsilon} := \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_m \end{pmatrix},$$

where $X \in \mathbf{R}^{M \times N}$, $\mathbf{y} \in \mathbf{R}^M$, and $\boldsymbol{\epsilon} \in \mathbf{R}^M$.

- In Bayesian linear regression, we assume that a “prior distribution” over parameters is given
- Prior over the weight vector $\boldsymbol{\theta}$:

$$\boldsymbol{\theta} \sim \mathcal{N}(\mathbf{0}, \tau^2 I_N) \quad (2)$$

2. POSTERIOR DISTRIBUTION OVER $\boldsymbol{\theta}$

- Apply Bayes’ theorem to obtain the posterior distribution on the weight set $\boldsymbol{\theta}$

$$p(\boldsymbol{\theta} | X, \mathbf{y}) \propto p(\mathbf{y} | X, \boldsymbol{\theta}) p(\boldsymbol{\theta}) \quad (3)$$

- Use equations (2) and (3) to find the posterior distribution over $\boldsymbol{\theta}$

$$\log p(\boldsymbol{\theta} | X, \mathbf{y}) \propto (\mathbf{y} - X\boldsymbol{\theta})^T (\sigma^2 I_M)^{-1} (\mathbf{y} - X\boldsymbol{\theta}) + \boldsymbol{\theta}^T (\tau^2 I_N)^{-1} \boldsymbol{\theta}$$

- Collect the terms dependent on $\boldsymbol{\theta}$

$$\log p(\boldsymbol{\theta} | X, \mathbf{y}) \propto \boldsymbol{\theta}^T \left(\frac{1}{\sigma^2} X^T X + \frac{1}{\tau^2} I_N \right) \boldsymbol{\theta} - 2\boldsymbol{\theta}^T \left(\frac{1}{\sigma^2} X^T \mathbf{y} \right)$$

- Therefore, the posterior distribution over $\boldsymbol{\theta}$ becomes

$$p(\boldsymbol{\theta} | X, \mathbf{y}) = \mathcal{N}(\boldsymbol{\mu}_\boldsymbol{\theta}, \Sigma_\boldsymbol{\theta}), \quad (4)$$

where

$$\boldsymbol{\mu}_\boldsymbol{\theta} = \left(\frac{1}{\sigma^2} X^T X + \frac{1}{\tau^2} I_N \right)^{-1} \frac{1}{\sigma^2} X^T \mathbf{y} \quad (5)$$

and

$$\Sigma_\boldsymbol{\theta} = \left(\frac{1}{\sigma^2} X^T X + \frac{1}{\tau^2} I_N \right)^{-1} \quad (6)$$

- Define $A := \Sigma_\boldsymbol{\theta}^{-1}$ to be compatible with the notations used in [1]
- Posterior distribution over $\boldsymbol{\theta}$ becomes

$$p(\boldsymbol{\theta} | X, \mathbf{y}) = \mathcal{N}\left(\frac{1}{\sigma^2} A^{-1} X^T \mathbf{y}, A^{-1}\right) \quad (7)$$

3. PREDICTION USING BAYESIAN LINEAR REGRESSION

- Assume that there is the same noise model on the testing point, $\{\mathbf{x}_*, y_*\}$, as our training points

- Posterior predictive distribution over y_* :
Integrate out the weight vector θ

$$\begin{aligned} p(y_*|\mathbf{x}_*, X, \mathbf{y}) &= \int_{\theta} p(y_*|\mathbf{x}_*, \theta) p(\theta|X, \mathbf{y}) d\theta \\ &= \int_{\theta} \exp \left\{ - \left[(y_* - \mathbf{x}_*^T \theta)^T (\sigma^2 I)^{-1} (y_* - \mathbf{x}_*^T \theta) + \left(\theta - \frac{1}{\sigma^2} A^{-1} X^T \mathbf{y} \right)^T A \left(\theta - \frac{1}{\sigma^2} A^{-1} X^T \mathbf{y} \right) \right] \right\} d\theta \end{aligned} \quad (8)$$

- Matrix A is symmetric (covariance matrix), hence

$$\begin{aligned} p(y_*|\mathbf{x}_*, X, \mathbf{y}) &= \int_{\theta} \exp \left\{ - \left[\theta^T \left(A + \frac{1}{\sigma^2} \mathbf{x}_* \mathbf{x}_*^T \right) \theta - 2\theta^T \frac{1}{\sigma^2} (\mathbf{x}_* y_* + X^T \mathbf{y}) + \left(\frac{1}{\sigma^2} y_*^T y_* + \frac{1}{\sigma^4} \mathbf{y}^T X A^{-1} X^T \mathbf{y} \right) \right] \right\} d\theta \end{aligned}$$

- Results in

$$\begin{aligned} P(y_*|\mathbf{x}_*, X, \mathbf{y}) &= e^{-\left[\frac{1}{\sigma^2} (y_*^T y_* + \frac{1}{\sigma^2} \mathbf{y}^T X A^{-1} X^T \mathbf{y}) - \frac{1}{\sigma^4} (\mathbf{x}_* y_* + X^T \mathbf{y})^T \left(A + \frac{1}{\sigma^2} \mathbf{x}_* \mathbf{x}_*^T \right)^{-1} (-) \right]} \\ &\times \int_{\theta} e^{-\left[(\theta - \left(A + \frac{1}{\sigma^2} \mathbf{x}_* \mathbf{x}_*^T \right)^{-1} \frac{1}{\sigma^2} (\mathbf{x}_* y_* + X^T \mathbf{y}))^T \left(A + \frac{1}{\sigma^2} \mathbf{x}_* \mathbf{x}_*^T \right)^{-1} (-) \right]} d\theta \end{aligned}$$

- Collect terms that only depend on y_*

$$\begin{aligned} \log p(y_*|\mathbf{x}_*, X, \mathbf{y}) &\propto y_*^T \left(\frac{1}{\sigma^2} - \frac{1}{\sigma^4} \mathbf{x}_*^T \left(A + \frac{1}{\sigma^2} \mathbf{x}_* \mathbf{x}_*^T \right)^{-1} \mathbf{x}_* \right) y_* \\ &\quad - 2y_*^T \left(\frac{1}{\sigma^4} \mathbf{x}_*^T \left(A + \frac{1}{\sigma^2} \mathbf{x}_* \mathbf{x}_*^T \right)^{-1} X^T \mathbf{y} \right) \end{aligned}$$

- Therefore, the posterior of y_* is Gaussian with the following terms

$$\Sigma_{y_*} = \left(\frac{1}{\sigma^2} - \frac{1}{\sigma^4} \mathbf{x}_*^T \left(A + \frac{1}{\sigma^2} \mathbf{x}_* \mathbf{x}_*^T \right)^{-1} \mathbf{x}_* \right)^{-1} \quad (9)$$

and

$$\mu_{y_*} = \Sigma_{y_*} \frac{1}{\sigma^4} \mathbf{x}_*^T \left(A + \frac{1}{\sigma^2} \mathbf{x}_* \mathbf{x}_*^T \right)^{-1} X^T \mathbf{y}. \quad (10)$$

- Remark: Matrix Inversion Lemma for equation

$$\begin{aligned} (A + BCD)^{-1} &= \\ A^{-1} - A^{-1} B (C^{-1} + DA^{-1} B)^{-1} DA^{-1} \end{aligned} \quad (11)$$

- Simplifying covariance matrix Σ_{y_*} in (9) using (11)

$$\begin{aligned} \Sigma_{y_*} &= \left((\sigma^2 + \mathbf{x}_*^T A^{-1} \mathbf{x}_*)^{-1} \right)^{-1} \\ \Sigma_{y_*} &= \sigma^2 + \mathbf{x}_*^T A^{-1} \mathbf{x}_* \end{aligned} \quad (12)$$

- Simplifying the mean μ_{y_*} in (10) using (11):

$$\begin{aligned} \mu_{y_*} &= \Sigma_{y_*} \frac{1}{\sigma^4} \mathbf{x}_*^T \left(A + \frac{1}{\sigma^2} \mathbf{x}_* \mathbf{x}_*^T \right)^{-1} X^T \mathbf{y} \\ &= \frac{1}{\sigma^4} (\sigma^2 + \mathbf{x}_*^T A^{-1} \mathbf{x}_*) \mathbf{x}_*^T \left(A + \frac{1}{\sigma^2} \mathbf{x}_* \mathbf{x}_*^T \right)^{-1} X^T \mathbf{y} \\ &= \frac{1}{\sigma^4} (\sigma^2 + \mathbf{x}_*^T A^{-1} \mathbf{x}_*) \mathbf{x}_*^T \left(A^{-1} - \frac{A^{-1} \mathbf{x}_* \mathbf{x}_*^T A^{-1}}{\sigma^2 + \mathbf{x}_*^T A^{-1} \mathbf{x}_*} \right) X^T \mathbf{y} \\ &= \frac{1}{\sigma^4} (\sigma^2 + \mathbf{x}_*^T A^{-1} \mathbf{x}_*) \mathbf{x}_*^T A^{-1} - \mathbf{x}_*^T A^{-1} \mathbf{x}_* \mathbf{x}_*^T A^{-1} X^T \mathbf{y} \end{aligned}$$

- Therefore,

$$\mu_{y_*} = \frac{1}{\sigma^2} \mathbf{x}_*^T A^{-1} X^T \mathbf{y} \quad (13)$$

4. SUMMARY ON THE PREDICTION

$$p(y_*|\mathbf{x}_*, X, \mathbf{y}) = \mathcal{N}(\mu_{y_*}, \Sigma_{y_*}),$$

where

$$\mu_{y_*} = \frac{1}{\sigma^2} \mathbf{x}_*^T A^{-1} X^T \mathbf{y}$$

and

$$\Sigma_{y_*} = \sigma^2 + \mathbf{x}_*^T A^{-1} \mathbf{x}_*.$$

5. REFERENCES

- [1] C. Rasmussen and C. Williams, *Gaussian Processes in Machine Learning*. MIT Press, 2006.
- [2] C. B. Do, "Gaussian processes," Dec. 2007. [Online]. Available: <http://see.stanford.edu/materials/aimlcs229/cs229-gp.pdf/>
- [3] M. Shekaramiz, *Sparse Signal Recovery Based on Compressive Sensing and Exploration Using Multiple Mobile Sensors*. PhD Dissertation, Utah State University, Digitalcommons, 2018.