

Foundations of Machine Learning

AI2000 and AI5000

FoML-11

Bayesian Regression

Dr. Konda Reddy Mopuri

Department of AI, IIT Hyderabad

July-Nov 2025



భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad



So far in FoML

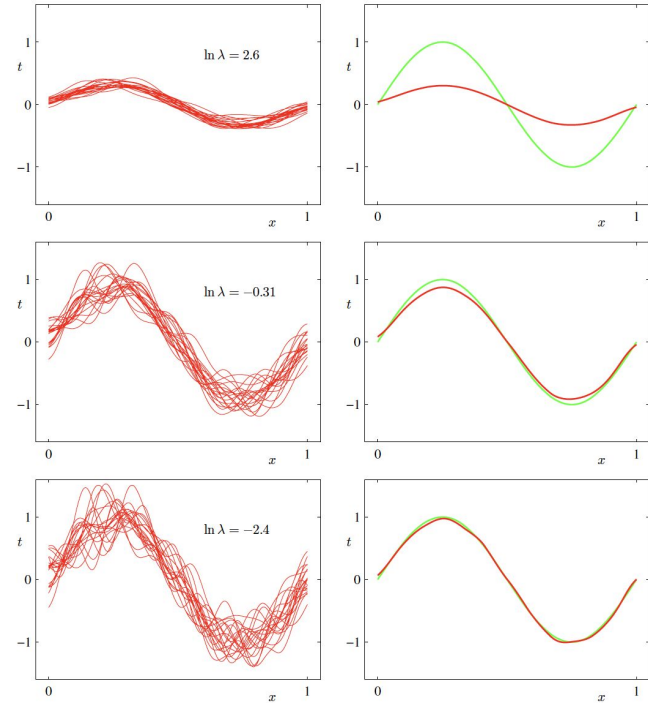
- What is ML and the learning paradigms
- Probability refresher
- MLE, MAP, and fully Bayesian treatment
- Linear Regression with basis functions - and regularization
- Model selection
- Bias-Variance Decomposition/Trade-off

Bayesian Regression



We have seen that

- Model averaging may be a good thing to do
 - Across different datasets



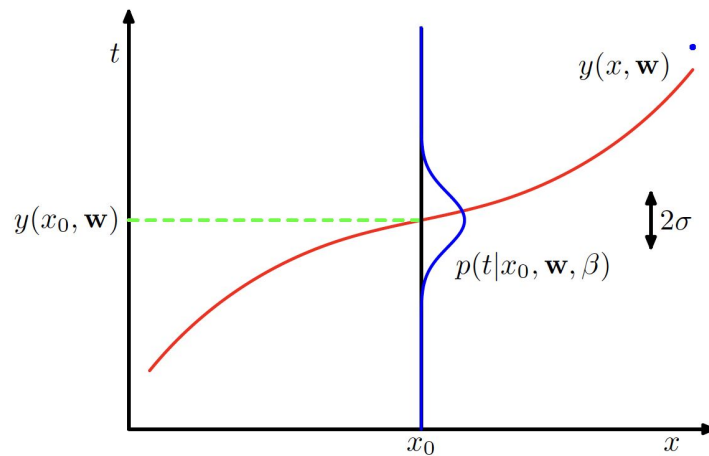
Bayesian Regression

- Instead of averaging over different datasets, we do it over different parameter sets



Bayesian Linear Regression

$$\text{Data } \mathbf{t} = (t_1, \dots, t_N)^T \quad \mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T$$



Bayesian Linear Regression

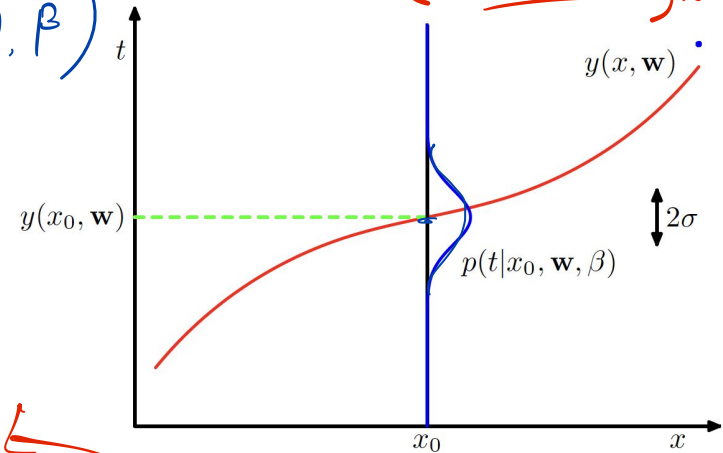
Data $\mathbf{t} = (t_1, \dots, t_N)^T$ $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T$

Likelihood: $p(t'|\mathbf{x}', \mathbf{w}, \beta) = \mathcal{N}(t' | \mathbf{w}^T \phi(x), \beta^{-1})$

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = \prod_{i=1}^N \mathcal{N}(t_i | y(x_i, \mathbf{w}), \beta^{-1}) = \mathcal{N}(\mathbf{t} | \Phi \mathbf{w}, \beta^{-1} \mathbf{I})$$

$$\Phi = \begin{bmatrix} \phi_1^T & \dots & \phi_N^T \end{bmatrix}$$

$$\Phi = \begin{bmatrix} \text{---} \\ \text{---} \\ \vdots \\ \text{---} \end{bmatrix}_{N \times m}$$



Posterior distribution: $p(\mathbf{w}|\mathbf{t}, \mathbf{X}) = \frac{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta)p(\mathbf{w})}{p(\mathbf{t}|\mathbf{X}, \beta)}$

Bayesian Linear Regression

Conjugate Prior: $p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | m_0, \mathbf{S}_0)$

↳ prior; before observing any data

$$p(\mathbf{w} | \mathbf{t}, \mathbf{X}) = \frac{p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta) p(\mathbf{w})}{p(\mathbf{t} | \mathbf{X}, \beta)} = \mathcal{N}(\underline{\mathbf{w}} | \underline{m}_N, \underline{\mathbf{S}}_N)$$

↳ posterior; after observing
N data samples $\mathbf{X}_{N \times m}$

$$\mathcal{N}(\mathbf{0}, \mathbf{I})$$

↓

m_0

This is what
we worked with
during MAP
discussion
(FomL-05)

Bayesian Linear Regression

Conjugate Prior: $p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0)$

Makes it possible for
a comfortable posterior
(Gaussian in this
case)

$$p(\mathbf{w} | \mathbf{t}, \mathbf{X}) = \frac{p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta) p(\mathbf{w})}{p(\mathbf{t} | \mathbf{X}, \beta)} = \mathcal{N}(\mathbf{w} | \mathbf{m}_N, \mathbf{S}_N)$$

$$\checkmark \mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta \Phi^T \Phi$$

$$\checkmark \mathbf{m}_N = \mathbf{S}_N (\mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \Phi^T \mathbf{t})$$

If we want a point estimate,
 $\Rightarrow \underline{w}_{MAP} = \underline{m}_N$



Bayesian Linear Regression

- Simple prior: $p(\mathbf{w}|\alpha) = \mathcal{N}(\underline{\mu} / \underline{\sigma}, S_0)$

$$\mathbf{m}_0 = \mathbf{0}$$

$$\mathbf{S}_0 = \alpha^{-1} \mathbf{I}$$



Bayesian Linear Regression

- Simple prior: $p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I})$

$$\mathbf{m}_0 = \mathbf{0}$$

$$\mathbf{S}_0 = \alpha^{-1}\mathbf{I}$$

- Posterior

$$p(\mathbf{w}|\mathbf{t}, \mathbf{X}, \alpha, \beta) = \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N)$$

$$\checkmark \quad \mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta \Phi^T \Phi = (\alpha \mathbf{I} + \beta \Phi^T \Phi)$$

$$\checkmark \quad \mathbf{m}_N = \mathbf{S}_N (\mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \Phi^T \mathbf{t}) = (\alpha \mathbf{I} + \beta \Phi^T \Phi) (\mathbf{0} + \beta \Phi^T \mathbf{t})$$
$$= \beta \mathbf{S}_N \cdot \Phi^T \mathbf{t}$$



Bayesian Linear Regression

- Special prior: Infinitely **broad** prior (no restriction) on \mathbf{w}

$$p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}) \quad \underline{\alpha \rightarrow 0}$$

- Posterior

$$\underline{\mathbf{S}_N^{-1}} = \mathbf{S}_0^{-1} + \beta \Phi^T \Phi = \alpha \rightarrow 0 \quad \left[\alpha \mathbf{I} + \beta \Phi^T \Phi \right] = \beta \Phi^T \Phi$$

$$\underline{\mathbf{m}_N} = \underline{\mathbf{S}_N} (\mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \Phi^T \mathbf{t}) = (\beta \Phi^T \Phi)^{-1} (\beta \Phi^T \mathbf{t})$$

$$= (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

$$\mathbf{S}_0 = \frac{1}{\alpha} \mathbf{I}$$

$$\mathbf{m}_0 = \mathbf{0}$$



Bayesian Linear Regression

- Special prior: Infinitely **narrow** prior on \mathbf{w}

$$p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}) \quad \alpha \rightarrow \inf$$

- Posterior

$$\mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta \Phi^T \Phi = \lim_{\alpha \rightarrow \infty} [\alpha \mathbf{I} + \beta \Phi^T \Phi] \Rightarrow \mathbf{S}_N = \mathbf{0}$$

$$\mathbf{m}_N = \mathbf{S}_N (\mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \Phi^T \mathbf{t}) = \lim_{\alpha \rightarrow \infty} [\alpha \mathbf{I} + \beta \Phi^T \Phi]^{-1} \beta \Phi^T \mathbf{t}$$

$$= \lim_{\alpha \rightarrow \infty} \frac{\beta}{\alpha} \Phi^T \mathbf{t} = \mathbf{0}$$

$$\mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0)$$

$$\mathbf{S}_0 = [\alpha^{-1} \mathbf{I}]$$

$$\alpha \rightarrow 0$$



Next Decision Theory

