

Foundations of Machine Learning AI2000 and AI5000

FoML-05

Maximum A Posteriori

Fully Bayesian treatment

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
- Maximum Likelihood Principle



Maximum A Posteriori



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



Maximum A Posteriori

- Given - Dataset of N independent observations D



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



Maximum A Posteriori

- Given - Dataset of N independent observations D
- ML estimate - w that maximizes the data likelihood

$$\mathbf{w}_{ML} = \underset{w}{\operatorname{argmax}} p(D|w)$$



Maximum A Posteriori

- Given - Dataset of N independent observations $D = \{x_1, x_2, \dots, x_N\}$
- MAP estimate - choose most probable w given data

Maximum A Posteriori

- Given - Dataset of N independent observations D
- MAP estimate - choose most probable w given data

$$\mathbf{w}_{MAP} = \underset{\mathbf{w}}{\text{argmax}} \quad P(\mathbf{w}/D)$$



MAP - Curve Fitting

- Given data D $D = \{(x_1, t_1), (x_2, t_2), \dots (x_N, t_N)\} = \{\mathbf{x}, \mathbf{t}\}$



MAP - Curve Fitting

- Given data D $D = \{(x_1, t_1), (x_2, t_2), \dots (x_N, t_N)\} = \{\mathbf{x}, \mathbf{t}\}$

- Model $p(t|x, \mathbf{w}, \beta) = \mathcal{N}(t|y(x, \mathbf{w}), \beta^{-1}) = \sqrt{\frac{\beta}{2\pi}} e^{-\frac{\beta}{2} (t - y(x, \mathbf{w}))^2}$

$p(\mathbf{w}|D)$

MAP - Curve Fitting

- Given data D $D = \{(x_1, t_1), (x_2, t_2), \dots (x_N, t_N)\} = \{\mathbf{x}, \mathbf{t}\}$

- Model $p(t|x, \mathbf{w}, \beta) = \mathcal{N}(t|y(x, \mathbf{w}), \beta^{-1})$

$$\mathbf{w}_{MAP} = \arg \max_{\mathbf{w}} p(\mathbf{w}|\mathbf{x}, \mathbf{t}, \beta)$$



MAP - Curve Fitting

$$p(\mathbf{w}|D) \propto \frac{p(D|\mathbf{w})p(\mathbf{w})}{\cancel{p(D|\mathbf{w})p(\mathbf{w})}}$$

- Given data $D = \{(x_1, t_1), (x_2, t_2), \dots, (x_N, t_N)\} = \{\mathbf{x}, \mathbf{t}\}$
- Model $p(t|x, \mathbf{w}, \beta) = \mathcal{N}(t|y(x, \mathbf{w}), \beta^{-1})$

$$\mathbf{w}_{MAP} = \arg \max_{\mathbf{w}} p(\mathbf{w}|\mathbf{x}, \mathbf{t}, \beta)$$

Given a prior $p(\mathbf{w}|\alpha)$ the posterior distribution becomes

$$p(\mathbf{w}|\mathbf{x}, \mathbf{t}, \beta, \alpha) = \frac{p(\mathbf{t}|\mathbf{w}, \mathbf{x}, \beta) p(\mathbf{w}|\alpha)}{p(\mathbf{t}|\mathbf{x}, \beta, \alpha)} \propto p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) p(\mathbf{w}|\alpha)$$



MAP - Curve Fitting

- MAP estimate - for convenience apply log

$$\mathbf{w}_{MAP} = \underset{\mathbf{w}}{\operatorname{argmax}} \left[\log p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) + \log p(\mathbf{w}|\alpha) - \underbrace{\log p(\mathbf{t}|\mathbf{x}, \beta, \alpha)}_{\text{independent of } \mathbf{w}} \right]$$



MAP - Curve Fitting

- Assuming Gaussian Prior and independence on parameters $\mathbf{w} \in \mathbb{R}^M$

$$\begin{aligned} p(\mathbf{w}|\alpha) &= \prod_{i=1}^M \mathcal{N}(\mathbf{w}_i | \mathbf{0}, \alpha^{-1}) \\ &= \prod_{i=1}^M \frac{1}{\sqrt{2\pi}} e^{-\frac{\alpha}{2} w_i^2} = \left(\frac{\alpha}{2\pi}\right)^{M/2} \prod_{i=1}^M e^{-\frac{\alpha}{2} w_i^2} \\ &= \left(\frac{\alpha}{2\pi}\right)^{M/2} e^{-\frac{\alpha}{2} \sum_{i=1}^M w_i^2} = \left(\frac{\alpha}{2\pi}\right)^{M/2} e^{-\frac{\alpha}{2} \mathbf{w}^T \mathbf{w}} \end{aligned}$$

Handwritten notes:
 $w_i \sim \mathcal{N}(0, \alpha^{-1})$
 $i = \{1, 2, \dots, m\}$



MAP - Curve Fitting

$$\mathbf{w}_{\text{MAP}} = \arg \min -\log \mathbf{p}(\mathbf{w}|\mathbf{x}, \mathbf{t}, \beta, \alpha) = \arg \min -\log \mathbf{p}(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) - \log \mathbf{p}(\mathbf{w}|\alpha)$$

$$= \arg \min_{\mathbf{w}} -\log \mathbf{p}(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) - \log \left(\left(\frac{\alpha}{2\pi} \right)^{M/2} e^{-\frac{\alpha}{2} \mathbf{w}^T \mathbf{w}} \right)$$

$$= \arg \min_{\mathbf{w}} -\log \mathbf{p}(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) - \frac{M}{2} \log \frac{\alpha}{2\pi} + \frac{\alpha}{2} \mathbf{w}^T \mathbf{w}$$

$\prod_{i=1}^N \sqrt{\frac{\beta}{2\pi}} e^{-\frac{\beta}{2} [t_i - y(x_i, \mathbf{w})]^2}$
Independent of \mathbf{w}

$$\mathbf{w}_{\text{MAP}} = \arg \min_{\mathbf{w}} \left\{ \frac{\beta}{2} \sum_{i=1}^N [t_i - y(x_i, \mathbf{w})]^2 + \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} \right\}$$



MAP - Curve Fitting

- Predictive distribution

$$p(t'|x', \beta) = \mathcal{N}(t' | \underbrace{y(w_{MAP}^* x_1)}_{\downarrow}, \beta)$$

w_{HL}^*

w_{MAP}^*



Bayesian Prediction



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



So far

- Our estimates for w have been point estimates
 - ML and MAP



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



So far

- Our estimates for w have been point estimates
 - ML and MAP
 - Regarded as frequentist because they discard 'uncertainty' about the w



Fully Bayesian

- An approach that relies on consistent application of sum and product rules of probability at all levels of modeling



Fully Bayesian

- Given a prior belief $p(\mathbf{w}|\alpha)$ over \mathbf{w} , and data D



Fully Bayesian

- Given a prior belief $p(\mathbf{w}|\alpha)$ over \mathbf{w} , and data D
- We are interested in the posterior

$$p(\mathbf{w}|D) = \frac{P(D|\mathbf{w}) P(\mathbf{w})}{P(D)}$$

Fully Bayesian

- The predictive distribution becomes

$$\underline{p(x'|D)} = \int p(x', w|D) \underline{dw} = \int p(x'|w, D) \underline{p(w|D)dw}$$



Fully Bayesian

- Curve fitting example



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



Fully Bayesian

- Curve fitting example
- Given training data (x, t)

Fully Bayesian

- Curve fitting example
- Given training data (x, t) and a test sample x

Fully Bayesian

- Curve fitting example
- Given training data (x, t) and a test sample x
- Goal - predict the value of t

Fully Bayesian

- Curve fitting example
- Given training data (\mathbf{x}, \mathbf{t}) and a test sample \mathbf{x}
- Goal - predict the value of \mathbf{t}

$$\mathcal{N}(\mathbf{t} | \underbrace{\mathbf{y}(\mathbf{w}_{ML}, \mathbf{x})}_{\mathbf{f}}, \mathbf{P})$$

We wish to evaluate the predictive distribution

$$\underline{p(\mathbf{t} | \mathbf{x}, \mathbf{x}, \mathbf{t})}$$



Fully Bayesian

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

$$\underline{p(t|x, \mathbf{x}, \mathbf{t})} = \int p(t|x, \mathbf{w}) p(\mathbf{w}|\mathbf{x}, \mathbf{t}) d\mathbf{w}.$$

$$= \int p(t, \mathbf{w} | \mathbf{x}, \mathbf{t}) d\mathbf{w}$$

$$= \int p(t | \mathbf{x}, \mathbf{t}, \mathbf{w}) p(\mathbf{w} | \mathbf{x}, \mathbf{t}) d\mathbf{w}$$

$$= \int p(t | \mathbf{x}, \mathbf{w}) p(\mathbf{w} | \mathbf{x}, \mathbf{t}) d\mathbf{w}$$



Fully Bayesian

- Advantages
 - Inclusion of the prior knowledge
 - Represents uncertainty in t' due to the target noise and uncertainty over w

Fully Bayesian

- Advantages
 - Inclusion of the prior knowledge
 - Represents uncertainty in t' due to the target noise and uncertainty over w
- Disadvantages
 - Posterior is hard to compute analytically
 - Prior is often a mathematical convenience

Rough work



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



Next Linear Models - Regression



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

