# Bayesian Inference for Multiparameter Models

*Aleksandar Dogandžić*

*2017-03-23*

## Contents

READING: §11 in the textbook and [Gelman et al. 2014, §3.1 and 3.2 and pp. 108–109].

## Multiparameter Models

So far, we have mostly discussed Bayesian estimation for toy scenarios with single parameters. In most real applications, we have multiple parameters that need to be estimated.

For example, if $X[n]$ are independent, identically distributed (i.i.d.) $\mathcal{N}(\mu, \sigma^2)$ given $\mu$ and $\sigma^2$, *both* $\mu$ and $\sigma^2$ may be *unknown* parameters. Any signal-plus-noise model where we do not know a signal parameter ($\mu$ in the above example) and a noise parameter ($\sigma^2$ in the above example) is a multiparameter model.

☞    THE noise variance $\sigma^2$ is often considered a *nuisance parameter*. We are not interested in the value of $\sigma^2$, but it is not known and hence is a "nuisance".

Consider the case with two parameters $\theta_1$ and $\theta_2$ and assume that only $\theta_1$ is of interest.[1] An example would be the DC-signal-in-additive white Gaussian noise (AWGN) model, where $\theta_1 = \mu$ and $\theta_2 = \sigma^2$.

[1] $\theta_1$ and $\theta_2$ could be vectors, but we describe the scalar case for simplicity.

Our inference should be based on

$$f_{\Theta_1 | \boldsymbol{X}}(\theta_1 \mid \boldsymbol{x})$$

the marginal posterior probability density function (pdf) (or probability mass function (pmf)) of $\theta_1$, which accounts for the uncertainty due to the fact that $\theta_2$ is unknown. First, we start with the joint posterior pdf (pmf):

$$f_{\Theta_1, \Theta_2 | \boldsymbol{X}}(\theta_1, \theta_2 \mid \boldsymbol{x}) \propto f_{\boldsymbol{X} | \Theta_1, \Theta_2}(\boldsymbol{x} \mid \theta_1, \theta_2) f_{\Theta_1, \Theta_2}(\theta_1, \theta_2)$$

Put the quantities we know on the right and the quantities we do not know on the left.

and then, in the continuous (pdf) case, *integrate out* the nuisance parameter:

$$f_{\Theta_1 | \boldsymbol{X}}(\theta_1 \mid \boldsymbol{x}) = \int f_{\Theta_1, \Theta_2 | \boldsymbol{X}}(\theta_1, \theta_2 \mid \boldsymbol{x}) \, \mathrm{d}\theta_2$$

also discussed in §10.7 of the textbook

or, equivalently,

$$f_{\Theta_1 | \boldsymbol{X}}(\theta_1 \mid \boldsymbol{x}) = \int f_{\Theta_1 | \Theta_2, \boldsymbol{X}}(\theta_1 \mid \theta_2, \boldsymbol{x}) f_{\Theta_2 | \boldsymbol{X}}(\theta_2 \mid \boldsymbol{x}) \, \mathrm{d}\theta_2 .$$

☞ THE marginal posterior distribution of $\theta_1$ can be viewed as its conditional posterior distribution (conditioned on the nuisance parameter, in addition to the data) *averaged over* the marginal posterior pdf of the nuisance parameter. Hence, the uncertainty due to the unknown $\theta_2$ is taken into account.

✳ EXAMPLE: Predictive Distribution. Suppose that we wish to predict an observation $X_\star$ coming from the model

$$f_{\boldsymbol{X} | \Theta_1, \Theta_2}(\boldsymbol{x} \mid \theta_1, \theta_2).$$

If $X_\star$ and $\boldsymbol{x}$ are independent given $\theta_1$ and $\theta_2$, then, based on (7) from handout `bpred`, we obtain:

$$f_{X_\star | \boldsymbol{X}}(x_\star \mid \boldsymbol{x}) = \int \int f_{X_\star | \Theta_1, \Theta_2}(x_\star \mid \theta_1, \theta_2) f_{\Theta_1, \Theta_2 | \boldsymbol{X}}(\theta_1, \theta_2 \mid \boldsymbol{x}) \, \mathrm{d}\theta_1 \, \mathrm{d}\theta_2$$

$$= \int \int f_{X_\star | \Theta_1, \Theta_2}(x_\star \mid \theta_1, \theta_2) f_{\Theta_1 | \Theta_2, \boldsymbol{X}}(\theta_1 \mid \theta_2, \boldsymbol{x}) f_{\Theta_2 | \boldsymbol{X}}(\theta_2 \mid \boldsymbol{x}) \, \mathrm{d}\theta_1 \, \mathrm{d}\theta_2 .$$

The above integrals are often difficult to evaluate analytically and may be highly multidimensional if $\theta_2$ and $\theta_1$ are vectors. Typically, we need Monte Carlo methods to handle practical cases. If we just wish to find the mode of

$$f_{\Theta_1 | \boldsymbol{X}}(\theta_1 \mid \boldsymbol{x})$$

an expectation-maximization (EM) algorithm may suffice. The EM algorithm will be discussed in detail later.

Lack of analytical tractability is the reason why the Bayesian methodology had been considered impractical in the past. Sometimes,

Bayesians made analytically tractable but hard to justify constructs, which had made the Bayesian approach appear even more obscure.

The advent of computers and development of Monte Carlo methods seem to have changed the balance in favor of the Bayesian approach, which has become practical and more flexible than classical inference that still largely relies on analytical tractability or hard-to-justify asymptotic results.

In the following, we consider two classical cases where analytical Bayesian computations are possible.

## DC Level Estimation in AWGN with Unknown Variance

THE measurements $(X[n])_{n=0}^{N-1}$ are i.i.d. $\mathcal{N}(\mu, \sigma^2)$ given $\boldsymbol{\theta}$, where

$$\boldsymbol{\theta} = (\mu, \sigma^2)$$

are the unknown parameters.

Assume that $\mu$ and $\sigma^2$ are independent and use the standard noninformative Jeffreys' priors for each:

$$f_{\mu, \sigma^2}(\mu, \sigma^2) = f_\mu(\mu) f_{\sigma^2}(\sigma^2)$$
$$\propto 1 \frac{1}{\sigma^2} \mathbb{1}_{[0, +\infty)}(\sigma^2). \tag{1}$$

The likelihood function $f_{X|\Theta}(x \mid \boldsymbol{\theta})$ is

$$f_{X|\Theta}(x \mid \boldsymbol{\theta}) = (2\pi\sigma^2)^{-N/2}$$
$$\exp\left\{-\frac{1}{2\sigma^2} \sum_{n=0}^{N-1} [(x[n] - \bar{x}) + (\bar{x} - \mu)]^2\right\} \mathbb{1}_{[0, +\infty)}(\sigma^2)$$
$$= (2\pi\sigma^2)^{-N/2} \exp\left[-\frac{Ns^2(x) + N(\bar{x} - \mu)^2}{2\sigma^2}\right] \mathbb{1}_{[0, +\infty)}(\sigma^2) \tag{2a}$$

where

$$s^2(x) = s^2 = \frac{1}{N} \sum_{n=0}^{N-1} (x[n] - \bar{x})^2, \qquad \bar{x} = \frac{1}{N} \sum_{n=0}^{N-1} x[n] \tag{2b}$$

are the sufficient statistics for $\mu$ and $\sigma^2$. The joint posterior pdf of the parameters is the product of the prior pdf (1) and likelihood (2a):

$$f_{\Theta|X}(\boldsymbol{\theta} \mid x) = f(\mu, \sigma^2 \mid x)$$
$$\propto \frac{1}{\sigma^2} (\sigma^2)^{-N/2} \exp\left\{-\frac{1}{2\sigma^2} [Ns^2 + N(\bar{x} - \mu)^2]\right\}$$
$$\cdot \mathbb{1}_{[0, +\infty)}(\sigma^2). \tag{3}$$

Provided that $N \geq 2$, the posterior pdf (3) is proper.

Conditional posterior pdf of $\mu$ given $\sigma^2$

THE conditional posterior pdf of $\mu$ given $\sigma^2$ is proportional to the joint posterior density with $\sigma^2$ held constant:

$$
\begin{aligned}
f(\mu \,|\, \sigma^2, \boldsymbol{x}) &\propto f(\mu, \sigma^2 \,|\, \boldsymbol{x}) \\
&\propto \frac{1}{\sigma^2}\, \mathbb{1}_{[0,+\infty)}(\sigma^2)(\sigma^2)^{-N/2} \exp\Big\{-\frac{1}{2\sigma^2}[Ns^2 + N(\bar{x}-\mu)^2]\Big\} \qquad \text{see (3)} \\
&\propto \exp\Big[-\frac{N}{2\sigma^2}(\bar{x}-\mu)^2\Big] \\
&\text{is the kernel of} \quad \mathcal{N}\big(\mu \,|\, \bar{x}, \sigma^2/N\big) \qquad\qquad\qquad (4) \qquad \text{keep only terms that have } \mu
\end{aligned}
$$

which agrees with (8) in handout `introBayes`, the case of estimating the DC level in AWGN with known variance.

We now write $f(\mu \,|\, \sigma^2, \boldsymbol{x})$ in (4) *with the normalizing constant*:

$$
f(\mu \,|\, \sigma^2, \boldsymbol{x}) = \frac{1}{\sqrt{2\pi\sigma^2/N}} \exp\Big[-\frac{(\mu-\bar{x})^2}{2\sigma^2/N}\Big]. \qquad (5) \qquad \text{see the table of distributions}
$$

Conditional posterior pdf of $\sigma^2$ given $\mu$

A reminder:

$$
\text{Inv-}\chi^2(\sigma^2 \,|\, v_0, \sigma_0^2) \propto \big(\sigma^2\big)^{-(v_0/2+1)} \exp\Big(-\frac{v_0\sigma_0^2}{2\sigma^2}\Big) \mathbb{1}_{[0,+\infty)}(\sigma^2) \qquad (6)
$$

and the full scaled inverted $\chi^2$ pdf is

$$
\begin{aligned}
\text{Inv-}\chi^2(\sigma^2 \,|\, v_0, \sigma_0^2) &= \frac{(v_0/2)^{v_0/2}}{\Gamma(v_0/2)}(\sigma_0^2)^{v_0/2}\big(\sigma^2\big)^{-(v_0/2+1)} \\
&\quad \cdot \exp\Big(-\frac{v_0\sigma_0^2}{2\sigma^2}\Big) \mathbb{1}_{[0,+\infty)}(\sigma^2) \qquad (7)
\end{aligned}
$$

see the table of distributions.

The conditional posterior pdf of $\sigma^2$ given $\mu$ is proportional to the joint posterior density with $\mu$ held constant:

$$
\begin{aligned}
f(\sigma^2 \,|\, \mu, \boldsymbol{x}) &\propto f(\mu, \sigma^2 \,|\, \boldsymbol{x}) \\
&\propto \frac{1}{\sigma^2}\, \mathbb{1}_{[0,+\infty)}(\sigma^2)(\sigma^2)^{-N/2} \exp\Big\{-\frac{1}{2\sigma^2}[Ns^2 + N(\bar{x}-\mu)^2]\Big\} \qquad \text{see (3)} \\
&\propto (\sigma^2)^{-(N/2+1)} \exp\Big\{-\frac{1}{2\sigma^2}[Ns^2 + N(\bar{x}-\mu)^2]\Big\} \mathbb{1}_{[0,+\infty)}(\sigma^2) \qquad \text{keep track of } \sigma^2 \\
&\text{is the kernel of} \quad \text{Inv-}\chi^2\big(\sigma^2 \,|\, N, s^2 + (\bar{x}-\mu)^2\big). \qquad (8)
\end{aligned}
$$

Therefore,

$$
\begin{aligned}
f(\sigma^2 \,|\, \mu, \boldsymbol{x}) &= \frac{(\frac{N}{2})^{N/2}}{\Gamma(N/2)}[s^2 + (\bar{x}-\mu)^2]^{N/2} \\
&\quad \cdot \big(\sigma^2\big)^{-(N/2+1)} \exp\Big(-\frac{Ns^2 + N(\bar{x}-\mu)^2}{2\sigma^2}\Big) \mathbb{1}_{[0,+\infty)}(\sigma^2)
\end{aligned}
$$

see (7).

Marginal posterior pdf of $\sigma^2$

.

THE marginal posterior pdf of $\sigma^2$ is a scaled inverted $\chi^2$:

$$f(\sigma^2 \mid \boldsymbol{x}) = \text{Inv-}\chi^2\left(N - 1, \frac{1}{N-1}\sum_{n=0}^{N-1}(x[n] - \bar{x})^2\right) \qquad (9)$$

[Gelman et al. 2014, §3.2] performs the integration whereas here we present an algebraic derivation without explicit integral computation based on [Gelman et al. 2014, pp. 108-109]

i.e., $\{\sigma^2 \mid \boldsymbol{x} = \boldsymbol{x}\}$ can be simulated as follows:

$$\{\sigma^2 \mid \boldsymbol{x} = \boldsymbol{x}\} \sim \frac{\sum_{n=0}^{N-1}(x[n] - \bar{x})^2}{Z}$$

where $Z$ is a $\chi^2_{N-1}$ random variable, see (2b) and (12) from handout `introBayes`.

To derive this result, we apply an algebraic approach for integrating $\mu$ out. The key to this approach is the fact that the conditional posterior pdf $f(\mu \mid \sigma^2, \boldsymbol{x})$ is known exactly (*including the normalizing constant*), e.g., it belongs to a family of pdfs (pmfs) that appear in our table of distributions.

We derive (9):

$$
f(\sigma^2 \mid \boldsymbol{x}) = \frac{\overbrace{f(\mu, \sigma^2 \mid \boldsymbol{x})}^{\text{see (3)}}}{\underbrace{f(\mu \mid \sigma^2, \boldsymbol{x})}_{\mathcal{N}(\mu \mid \bar{x}, \sigma^2/N), \text{ see (5)}}}
$$

not a function of $\mu$

keep track of both $\mu$ and $\sigma^2$

$$
\propto \frac{(\sigma^2)^{-N/2-1} \exp\{-\frac{1}{2\sigma^2}[Ns^2 + N(\bar{x}-\mu)^2]\}\, \mathbb{1}_{[0,+\infty)}(\sigma^2)}{(\sigma^2)^{-1/2} \exp[-\frac{(\mu-\bar{x})^2}{2\sigma^2/N}]}
$$

$$
\propto \frac{(\sigma^2)^{-N/2-1} \exp[-Ns^2/(2\sigma^2)]\, \mathbb{1}_{[0,+\infty)}(\sigma^2)}{(\sigma^2)^{-1/2}}
$$

plug in $\mu = \bar{x}$

$$
\propto (\sigma^2)^{-[(N-1)/2-1]} \exp\left(-\frac{Ns^2}{2\sigma^2}\right) \mathbb{1}_{[0,+\infty)}(\sigma^2)
$$

$$
\propto (\sigma^2)^{-[(N-1)/2-1]} \exp\left\{-\frac{\sum_{n=0}^{N-1}(x[n]-\bar{x})^2}{2\sigma^2}\right\} \mathbb{1}_{[0,+\infty)}(\sigma^2)
$$

rearrange to look like (6)

is the kernel of $\mathrm{Inv}\text{-}\chi^2\left(\sigma^2 \,\middle|\, N-1, \dfrac{1}{N-1}\sum_{n=0}^{N-1}(x[n]-\bar{x})^2\right)$

∗ CHECK: $\mu$ must cancel out—important for verifying the correctness of our computations.

☞ WE have computed

$$
f(\sigma^2 \mid \boldsymbol{x}) = \int_{-\infty}^{+\infty} f(\mu, \sigma^2 \mid \boldsymbol{x})\, \mathrm{d}\mu
$$

algebraically *without* performing the actual integration.

Marginal posterior pdf of $\mu$

APPLY the algebraic approach to integrate $\sigma^2$ out:

$$f(\mu \mid \boldsymbol{x}) = \frac{\overbrace{f(\mu, \sigma^2 \mid \boldsymbol{x})}^{\text{see (3)}}}{\underbrace{f(\sigma^2 \mid \mu, \boldsymbol{x})}}$$

[Gelman et al. 2014, §3.2] performs the integration whereas here we present an algebraic derivation without explicit integral computation based on [Gelman et al. 2014, pp. 108-109].

$$\propto \frac{(\sigma^2)^{-N/2-1} \exp\{-[Ns^2 + N(\bar{x}-\mu)^2]/(2\sigma^2)\}\, \mathbb{1}_{[0,+\infty)}(\sigma^2)}{[s^2 + (\bar{x}-\mu)^2]^{N/2}(\sigma^2)^{-N/2-1} \exp\{-[Ns^2 + N(\bar{x}-\mu)^2]/(2\sigma^2)\}\, \mathbb{1}_{[0,+\infty)}(\sigma^2)}$$

$\text{Inv-}\chi^2(\sigma^2 \mid N, s^2+(\bar{x}-\mu)^2)$, see (9)

not a function of $\sigma^2$

$$\propto [s^2 + (\bar{x}-\mu)^2]^{-N/2}$$

$$\propto \left[1 + \frac{1}{N-1}\frac{(\mu-\bar{x})^2}{s^2/(N-1)}\right]^{-(N-1+1)/2}$$

kernel of the Student-$t$ distribution keep track of both $\mu$ and $\sigma^2$

☞ CHECK: $\sigma^2$ must cancel out.

From the table of distributions, we find:

$$f(\mu \mid \boldsymbol{x}) = t_{N-1}\left(\mu \,\middle|\, \bar{x}, \frac{s^2}{N-1}\right)$$

which is a Student-$t$ pdf with $N-1$ degrees of freedom and parameters

$$\bar{x} \quad \text{mean, see (2b)} \quad \text{and} \quad \frac{s^2}{N-1} = \frac{\sum_{n=0}^{N-1}(x[n]-\bar{x})^2}{N(N-1)} \quad \text{scale.}$$

☞ WE have computed

$$f(\mu \mid \boldsymbol{x}) = \int_{-\infty}^{+\infty} f_{\mu,\sigma^2\mid \boldsymbol{X}}(\mu, \sigma^2 \mid \boldsymbol{x})\, d\sigma^2$$

algebraically, *without* performing the actual integration.

✳ A credible set for the DC-level-in-AWGN example. The typical 95 % highest posterior density (HPD) credible interval for $\mu$ based on

$$f_{\mu\mid \boldsymbol{X}}(\mu \mid \boldsymbol{x}) = t_{N-1}\left(\mu \,\middle|\, \bar{x}, \frac{s^2}{N-1}\right)$$

coincides with the common 95 % confidence interval based on the classical $t$ test, taught in STAT 101. However, the Bayesian interpretation is different.

☞ EXAMPLE 11.4 in the textbook is equivalent to the conjugate prior case for $(\mu, \sigma^2)$ in [Gelman et al. 2014, §3.3], which is not practically relevant.

## Computing MAP Estimates

MAXIMUM *a posteriori* (MAP) estimation corresponds to the minimization problem:

$$\widehat{\boldsymbol{\theta}}_{\text{MAP}} = \arg\min_{\boldsymbol{\theta}} V(\boldsymbol{\theta})$$

where

$$V(\boldsymbol{\theta}) = -\ln f_{\boldsymbol{\Theta}|\boldsymbol{X}}(\boldsymbol{\theta} \mid \boldsymbol{x}).$$

## Newton-Raphson iteration

ITERATE

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \left(H^{(t)}\right)^{-1} \boldsymbol{g}^{(t)}$$

where $t$ denotes the iteration index and

$$\boldsymbol{g}^{(t)} = \left. \frac{\partial V(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}^{(t)}} \qquad \text{gradient}$$

$$H_i = \left. \frac{\partial^2 V(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}\, \partial \boldsymbol{\theta}^T} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}^{(t)}}.$$

A damped Newton-Raphson variation:

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \mu^{(t)}\left(H^{(t)}\right)^{-1}\boldsymbol{g}^{(t)} \tag{10}$$

see handout `ml`.

✳  COMMENTS:

- Newton-Raphson iteration is not guaranteed to converge but

- its convergence is fast in the neighborhood of the MAP estimate.

Upon convergence (i.e., as $t \nearrow \infty$) and if we reach the global optimum, we have

$$\boldsymbol{g}^{(+\infty)} = \left. \frac{\partial V(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}^{(+\infty)}=\widehat{\boldsymbol{\theta}}_{\text{MAP}}} = \boldsymbol{0}. \tag{11}$$

The EM algorithm for marginal MAP estimation is discussed in handout `emBayes`.

## MAP Estimation for Multiple Parameters

CONSIDER again the case of two parameter vectors, denoted by $\boldsymbol{\theta}$ and $\boldsymbol{u}$; then, the set of all unknown parameters is $(\boldsymbol{\theta}, \boldsymbol{u})$. The joint posterior pdf for $\boldsymbol{\theta}$ and $\boldsymbol{u}$ is

$$f_{\boldsymbol{\Theta},U|\boldsymbol{X}}(\boldsymbol{\theta}, \boldsymbol{u} \mid \boldsymbol{x}) = f_{U|\boldsymbol{\Theta},\boldsymbol{X}}(\boldsymbol{u} \mid \boldsymbol{\theta}, \boldsymbol{x}) f_{\boldsymbol{\Theta}|\boldsymbol{X}}(\boldsymbol{\theta} \mid \boldsymbol{x}).$$

We wish to estimate both $\boldsymbol{\theta}$ and $\boldsymbol{u}$.

Approach I

MAXIMIZE the marginal posterior pdfs (pmfs)

$$\widehat{\boldsymbol{\theta}} = \arg\max_{\boldsymbol{\theta}} f_{\boldsymbol{\Theta}|\boldsymbol{X}}(\boldsymbol{\theta} \mid \boldsymbol{x})$$

and

$$\widehat{\boldsymbol{u}} = \arg\max_{\boldsymbol{u}} f_{\boldsymbol{U}|\boldsymbol{X}}(\boldsymbol{u} \mid \boldsymbol{x})$$

which take into account the uncertainties about the other parameter. We can perform the two optimizations separately.

But, what if we cannot easily obtain these two marginal posterior pdfs (pmfs)? Suppose now that we can obtain $f_{\boldsymbol{\Theta}|\boldsymbol{X}}(\boldsymbol{\theta} \mid \boldsymbol{x})$, but not $f_{\boldsymbol{U}|\boldsymbol{X}}(\boldsymbol{u} \mid \boldsymbol{x})$.

Approach II

1. FIRST, find the marginal MAP estimate of $\boldsymbol{\theta}$:

$$\widehat{\boldsymbol{\theta}} = \arg\max_{\boldsymbol{\theta}} f_{\boldsymbol{\Theta}|\boldsymbol{X}}(\boldsymbol{\theta} \mid \boldsymbol{x})$$

   which, as desired, takes into account the uncertainty about $\boldsymbol{u}$ by integrating $\boldsymbol{u}$ out from the joint posterior pdf.

2. Then, find the conditional MAP estimate of $\boldsymbol{u}$ by maximizing $f_{\boldsymbol{U}|\boldsymbol{\Theta},\boldsymbol{X}}(\boldsymbol{u} \mid \boldsymbol{\theta}, \boldsymbol{x})$:

$$\widehat{\boldsymbol{u}} = \arg\max_{\boldsymbol{u}} f_{\boldsymbol{U}|\boldsymbol{\Theta},\boldsymbol{X}}(\boldsymbol{u} \mid \widehat{\boldsymbol{\theta}}, \boldsymbol{x}).$$

Finally, what if we cannot easily obtain either of the two marginal posterior pdfs (pmfs)?

Approach III

FIND the joint MAP estimate $(\widehat{\boldsymbol{\theta}}, \widehat{\boldsymbol{u}})$,

$$(\widehat{\boldsymbol{\theta}}, \widehat{\boldsymbol{u}}) = \arg\max_{\boldsymbol{\theta}, \boldsymbol{u}} f_{\boldsymbol{\Theta},\boldsymbol{U}|\boldsymbol{X}}(\boldsymbol{\theta}, \boldsymbol{u} \mid \boldsymbol{x}).$$

This estimation is sometimes performed as follows: iterate between

1. finding the conditional MAP estimate of $\boldsymbol{\theta}$ by maximizing $f_{\boldsymbol{\Theta}|\boldsymbol{U},\boldsymbol{X}}(\boldsymbol{\theta} \mid \boldsymbol{u}^{(t)}, \boldsymbol{x})$:

$$\boldsymbol{\theta}^{(t+1)} = \arg\max_{\boldsymbol{\theta}} f_{\boldsymbol{\Theta}|\boldsymbol{U},\boldsymbol{X}}(\boldsymbol{\theta} \mid \boldsymbol{u}^{(t)}, \boldsymbol{x})$$

   and

2.  finding the conditional MAP estimate of $u$ by maximizing $f_{U|\Theta,X}(u \mid \theta^{(t+1)}, x)$:

$$u^{(t+1)} = \arg\max_{u} f_{U|\Theta,X}(u \mid \theta^{(t+1)}, x)$$

known as the *iterated conditional modes (ICM) algorithm*. This is simply an application of the *stepwise-ascent approach to optimization*. A general ICM algorithm is not restricted to two components ($\theta$ and $u$ in this example) and can employ many more; the components can be scalars or vectors.

✳   COMMUNICATIONS example. In communications, $\theta$ and $u$ can be (blocks of) symbols. Then, approach III corresponds to the Viterbi algorithm and Approach I to the Bahl-Cocke-Jelinek-Raviv (BCJR) algorithm [Bahl et al. 1974].

## Acronyms

*AWGN*  additive white Gaussian noise. 1, 4, 7

*BCJR*  Bahl-Cocke-Jelinek-Raviv. 10

*EM*  expectation-maximization. 2, 8

*HPD*  highest posterior density. 7

*i.i.d.*  independent, identically distributed. 1, 3

*ICM*  iterated conditional modes. 10

*MAP*  maximum *a posteriori*. 8–10

*pdf*  probability density function. 2–5, 7–9

*pmf*  probability mass function. 2, 5, 9

## References

Bahl, L., J. Cocke, F. Jelinek, and J. Raviv (1974). "Optimal decoding of linear codes for minimizing symbol error rate". In: *IEEE Trans. Inf. Theory* 20.2, pp. 284–287 (cit. on p. 10).

Gelman, Andrew, John B. Carlin, Hal S. Stern, David B. Dunson, Aki Vehtari, and Donald B. Rubin (2014). *Bayesian Data Analysis*. 3rd ed. Boca Raton, FL: Taylor & Francis (cit. on pp. 1, 3, 5, 7).