

# LECTURE 1: FOUNDATIONS IN DETERMINISTIC AND STOCHASTIC OPTIMIZATION

STEFAN VLASKI AND ALI H. SAYED

ABSTRACT. These notes cover foundational material in statistical learning theory as well as deterministic and stochastic optimization with a focus on single agent learners. It is expected that the reader will have had some exposure to concepts such as maximum likelihood and maximum a posteriori estimation, gradient descent as well as its stochastic variants, and as a result the exposition is kept brief. Nevertheless, we find it useful to collect this information here as it will form the foundation for the development of multi-agent learning algorithms in future lectures.

## 1. BAYESIAN INFERENCE

One of the most fundamental problems in statistics, signal processing, and machine learning, is the *inference* problem, where we wish to construct an estimate of some random quantity of interest  $\gamma \in \mathbb{R}^{M_\gamma}$ , given observations of a related random variable  $\mathbf{h} \in \mathbb{R}^{M_h}$ . The quantity of interest  $\gamma$  may take on continuous or discrete values, and is referred to as the “dependent variable,” “state of nature,” “class,” or “label” depending on the application. We will most commonly refer to it as the label. We will refer to the observed random variable  $\mathbf{h}$  generally as the feature, although in some applications it is known as “regressor” or “observation”.

If we are provided with the conditional distribution  $f_{\gamma|\mathbf{h}}(\gamma|\mathbf{h})$  along with a single realization of the feature  $\mathbf{h}$ , it is quite natural to estimate  $\gamma$  as the most likely outcome given  $\mathbf{h}$ . We can express this formally as:

$$(1) \quad \gamma^* \triangleq \arg \max_{\gamma \in \Gamma} f_{\gamma|\mathbf{h}}(\gamma|\mathbf{h})$$

The conditional distribution  $f_{\gamma|\mathbf{h}}(\gamma|\mathbf{h})$  is frequently referred to as the *posterior distribution* of  $\gamma$ , making  $\gamma^*$  the *maximum a posteriori (MAP) estimate* of  $\gamma$  given  $\mathbf{h}$ . We note that it is common in the literature on estimation theory to employ hat-notation to refer to estimates of a random quantity based on data. In this sense, we could have opted to denote the optimal solution to (1) by  $\hat{\gamma}$ , rather than  $\gamma^*$ . Nevertheless, as we transition from inference to learning and optimization, we will increasingly interpret estimates as solutions to generic optimization problems. In this context, it will be appropriate to employ the \*-notation to refer to an optimal solution. To preserve consistency throughout this text, we opt to use the \*-notation from the onset, with the understanding that within the Bayesian framework (1)

the optimal solution  $\gamma^*$  carries the interpretation of an estimate. Finally, we note that (1) is formulated for a particular realization  $h$  of the random variable  $\mathbf{h}$ . As we will see, most MAP estimators are derived for a particular realization of the feature vector, resulting in a deterministic estimate  $\gamma^*$  as a function of the realization  $h$ . We can also regard the MAP solution as a random variable, denoted in boldface notation  $\boldsymbol{\gamma}^*$ , when it is viewed as a function of the random observation  $\mathbf{h}$ :

$$(2) \quad \boldsymbol{\gamma}^* \triangleq \arg \max_{\gamma \in \Gamma} f_{\gamma|\mathbf{h}}(\gamma|\mathbf{h})$$

where the posterior distribution  $f_{\gamma|\mathbf{h}}(\gamma|\mathbf{h})$  is now a function of the random variable  $\mathbf{h}$ , rather than its realization  $h$ .

*Remark 1.1 ( **Bold font indicates random variables** ).* We note that we employ bold font to denote a random variables, for example  $\mathbf{h}$ , while regular font denotes its realization or a deterministic quantity, as in  $h$ . As a general rule, we use lowercase fonts to denote vectors, while uppercase letters denote matrices. In this way, a random matrix would be denoted by  $\mathbf{H}$ , while its realization would correspond to  $H$ .

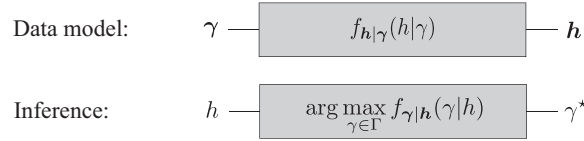


FIGURE 1. A general model underlying most inference problems. The top row illustrates the generative model for the observations  $\mathbf{h}$ , while the bottom row illustrates the Bayesian inference formulation for recovering the label variable.

**1.1. From Inference to Learning.** Our previous discussion leading to (??) has uncovered how to perform MAP inference of a random quantity of interest  $\gamma$  given observations  $\{\mathbf{h}_n\}_{n=1}^N$ , models  $f_{\mathbf{h}|\gamma}(\cdot|\gamma)$ , and prior information  $f_{\gamma}(\gamma)$ . In many practical situations, we are not provided with prior knowledge about the model relating the label  $\gamma$  with the feature  $\mathbf{h}$ . Instead, we need to estimate  $f_{\mathbf{h}|\gamma}(\cdot|\gamma)$  from data before performing inference. We refer to the process of estimating statistical properties of data needed for inference, such as  $f_{\mathbf{h}|\gamma}(\cdot|\gamma)$ , as *learning*. As we will see in this section, the learning of models can be formalized using a Bayesian MAP framework analogous to (1). To this end, we will assume that the conditional likelihood  $f_{\mathbf{h}|\gamma}(\cdot|\gamma)$  is parameterized by a learnable parameter  $w \in \mathbb{R}^{M_w}$  and write instead  $f_{\mathbf{h}|\gamma,w}(\cdot|\gamma,w)$ . Under this parameterization, learning the conditional distribution of  $\gamma$  given  $\mathbf{h}$  is equivalent to learning the parameter  $w$  that parameterizes  $f_{\mathbf{h}|\gamma,w}(\cdot|\gamma,w)$ .

To formulate a procedure for learning  $w$  from pairs  $\{\mathbf{h}, \gamma\}$  we mirror the argument in Section ???. Suppose the model  $w$  is sampled once, yielding the

realization  $w^o$ . The *training data*  $\{\mathbf{h}_n, \gamma_n\}_{n=1}^N$  is sampled  $N$  times, where each pair  $\{\mathbf{h}_n, \gamma_n\}$  is sampled from  $f_{\mathbf{h}, \gamma | \mathbf{w}}(h, \gamma | w^o)$ . If we suppose that feature-label pairs  $\{\mathbf{h}_n, \gamma_n\}$  are identically and independently distributed after conditioning on the parameterization  $\mathbf{w}$ , we can factorize:

$$\begin{aligned} f_{\{\mathbf{h}_n, \gamma_n\}_{n=1}^N | \mathbf{w}} \left( \{\mathbf{h}_n, \gamma_n\}_{n=1}^N | \mathbf{w} \right) &\stackrel{(a)}{=} \prod_{n=1}^N f_{\mathbf{h}_n, \gamma_n | \mathbf{w}}(h_n, \gamma_n | \mathbf{w}) \\ (3) \qquad \qquad \qquad &\stackrel{(b)}{=} \prod_{n=1}^N f_{\mathbf{h}, \gamma | \mathbf{w}}(h_n, \gamma_n | \mathbf{w}) \end{aligned}$$

Step (a) holds by conditional independence and (b) holds by identical distribution of the pairs of random variables  $\{\mathbf{h}_n, \gamma_n\}$ . We can then define the MAP estimate of the weight vector  $\mathbf{w}$  as:

$$\begin{aligned} \mathbf{w}^* &\triangleq \arg \max_{\mathbf{w} \in \mathbb{R}^{M_w}} f_{\mathbf{w} | \{\mathbf{h}_n, \gamma_n\}_{n=1}^N} \left( \mathbf{w} | \{\mathbf{h}_n, \gamma_n\}_{n=1}^N \right) \\ &= \arg \max_{\mathbf{w} \in \mathbb{R}^{M_w}} \frac{f_{\{\mathbf{h}_n, \gamma_n\}_{n=1}^N | \mathbf{w}} \left( \{\mathbf{h}_n, \gamma_n\}_{n=1}^N | \mathbf{w} \right) \times f_{\mathbf{w}}(\mathbf{w})}{f_{\{\mathbf{h}_n, \gamma_n\}_{n=1}^N} \left( \{\mathbf{h}_n, \gamma_n\}_{n=1}^N \right)} \\ &= \arg \max_{\mathbf{w} \in \mathbb{R}^{M_w}} f_{\{\mathbf{h}_n, \gamma_n\}_{n=1}^N | \mathbf{w}} \left( \{\mathbf{h}_n, \gamma_n\}_{n=1}^N | \mathbf{w} \right) \times f_{\mathbf{w}}(\mathbf{w}) \\ (4) \qquad \qquad \qquad &= \arg \max_{\mathbf{w} \in \mathbb{R}^{M_w}} \left( \prod_{n=1}^N f_{\mathbf{h}, \gamma | \mathbf{w}}(h_n, \gamma_n | \mathbf{w}) \right) \times f_{\mathbf{w}}(\mathbf{w}) \end{aligned}$$

Following the same argument that led to (??), we arrive at:

$$(5) \quad \mathbf{w}^* = \arg \min_{\mathbf{w} \in \mathbb{R}^M} \left\{ -\frac{1}{N} \sum_{n=1}^N \log f_{\mathbf{h}, \gamma | \mathbf{w}}(h_n, \gamma_n | \mathbf{w}) - \frac{1}{N} \log f_{\mathbf{w}}(\mathbf{w}) \right\}$$

**Remark 1.2 (Distinction between the true model and the MAP estimate).** Observe that we make a deliberate distinction between the *true* model  $w^o$ , which parameterizes the distribution  $f_{\mathbf{h}, \gamma | \mathbf{w}}(h, \gamma | w^o)$  from which the samples  $\{\mathbf{h}_n, \gamma_n\}_{n=1}^N$  are generated, and the MAP *estimate*  $\mathbf{w}^*$ , which maximizes the posterior distribution of  $\mathbf{w}$  after observing the samples  $\{\mathbf{h}_n, \gamma_n\}_{n=1}^N$ . In general, there will be a difference between the MAP model  $\mathbf{w}^*$  and the true model  $w^o$ . The expectation is that as the size of the data set  $N$  grows, the MAP model  $\mathbf{w}^*$  will become increasingly accurate and approach  $w^o$  as  $N \rightarrow \infty$ . We will verify that this is the case further below.  $\square$

Note that in order to be able to *learn* an optimal parameterization  $\mathbf{w}^*$  according to (5), we need to be able to collect a batch of feature-label pairs  $\{\mathbf{h}_n, \gamma_n\}_{n=1}^N$ . In the context of machine learning this step is frequently referred to as *training*, and the labeled data  $\{\mathbf{h}_n, \gamma_n\}_{n=1}^N$  is referred to as the *training data*. Once  $\mathbf{w}^*$  is determined, we can use the learned parameterization  $\mathbf{w}^*$  on an unlabeled feature vector  $h^{\text{test}}$ , to compute an approximate

MAP estimate for its label:

$$(6) \quad \gamma^{\text{test}\star} \triangleq \arg \max_{\gamma \in \Gamma} f_{\gamma|\mathbf{h},\mathbf{w}}(\gamma|\mathbf{h}^{\text{test}}, \mathbf{w}^\star)$$

We emphasize that (6) is only an *approximate* MAP estimate for the true label  $\gamma^{\text{test}}$ , since it is generated using an estimate of the posterior distribution  $f_{\gamma|\mathbf{h},\mathbf{w}}(\gamma|\mathbf{h}^{\text{test}}, \mathbf{w}^\star)$  using the learned parameterization  $\mathbf{w}^\star$ . In general, and in particular for finite sample sizes  $N$  of the training data, the learned parameterization  $\mathbf{w}^\star$  will be different from the true parameterization  $\mathbf{w}^o$  that actually generated the data. The difference between predictions made using the true model  $\mathbf{w}^o$  and the learned model  $\mathbf{w}^\star$  is known as a *generalization error*.

**Remark 1.3 (Compact notation for feature-label pairs).** As is evident from (5), during learning, we are provided with feature-label pairs  $\{\mathbf{h}, \gamma\}$  or their realizations  $\{h_n, \gamma_n\}_{n=1}^N$ . To simplify the notation, we will collect features  $\mathbf{h} \in \mathbb{R}^{M_h}$  and labels  $\gamma \in \mathbb{R}^{M_\gamma}$  into a single augmented data vector  $\mathbf{x} \in \mathbb{R}^{M_h+M_\gamma}$ , such that:

$$(7) \quad \mathbf{x} \triangleq \text{col}\{\mathbf{h}, \gamma\} = \begin{pmatrix} \mathbf{h} \\ \gamma \end{pmatrix}$$

Similarly, we will collect realizations into  $x_n \triangleq \text{col}\{h_n, \gamma_n\}$ . In this manner, we can write more compactly:

$$(8) \quad f_{\mathbf{h},\gamma|\mathbf{w}}(h, \gamma|\mathbf{w}) = f_{\mathbf{x}|\mathbf{w}}(x|\mathbf{w})$$

The MAP learning problem (5) then becomes:

$$(9) \quad \mathbf{w}^\star = \arg \min_{\mathbf{w} \in \mathbb{R}^M} \left\{ -\frac{1}{N} \sum_{n=1}^N \log f_{\mathbf{x}|\mathbf{w}}(x_n|\mathbf{w}) - \frac{1}{N} \log f_{\mathbf{w}}(\mathbf{w}) \right\}$$