

1 Inference, Learning and Optimization

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 γ , given observations of a related random variable \mathbf{h} . The quantity of interest γ may take on continuous or discrete values, 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 “predictor”.

1.1 BAYSIAN INFERENCE

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

$$\gamma^* \triangleq \arg \max_{\gamma \in \Gamma} f_{\gamma|\mathbf{h}}(\gamma|h) \quad (1.1)$$

We note that it is more common in the literature on estimation 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 (1.1) as $\hat{\gamma}$, rather than γ^* . 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 more appropriate to employ *-notation to denote an optimal solution. To preserve consistency throughout this textbook, we opt to use *-notation from the onset, with the understanding that within the Bayesian framework (1.1) the optimal solution γ^* carries the interpretation of an estimate.

Example 1.1 (Rainy day) Suppose one steps outside in the morning and observes that the grass is wet. The question then arises whether it had rained overnight. The fact that the grass is wet (the observation \mathbf{h}) naturally contains information about the unobserved weather overnight (the state of nature γ).

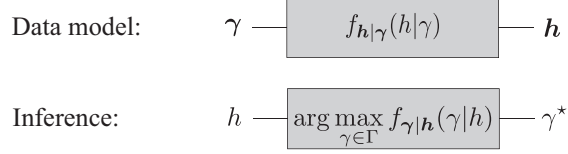


Figure 1.1 A general model underlying most inference problems.

This is formalized through the conditional likelihood:

$$\Pr(\mathbf{h} = \text{wet grass} | \gamma = \text{rainy night}) = 0.9 \quad (1.2)$$

$$\Pr(\mathbf{h} = \text{dry grass} | \gamma = \text{rainy night}) = 0.1 \quad (1.3)$$

$$\Pr(\mathbf{h} = \text{wet grass} | \gamma = \text{dry night}) = 0.05 \quad (1.4)$$

$$\Pr(\mathbf{h} = \text{dry grass} | \gamma = \text{dry night}) = 0.95 \quad (1.5)$$

The conditional distribution (or likelihood) above captures the fact that the state of the grass observed in the morning is highly correlated with the weather overnight. Suppose we live in an area of even distribution between rainy and dry nights, implying a non-informative prior:

$$\Pr(\gamma = \text{rainy night}) = \Pr(\gamma = \text{dry night}) = 0.5 \quad (1.6)$$

For the marginal probabilities we hence find:

$$\begin{aligned} \Pr(\mathbf{h} = \text{wet}) &= \Pr(\mathbf{h} = \text{wet grass} | \gamma = \text{rainy night}) \cdot \Pr(\gamma = \text{rainy night}) \\ &\quad + \Pr(\mathbf{h} = \text{wet grass} | \gamma = \text{dry night}) \cdot \Pr(\gamma = \text{dry night}) \\ &= 0.9 \cdot 0.5 + 0.05 \cdot 0.5 = 0.45 + 0.025 = 0.475 \end{aligned} \quad (1.7)$$

$$\begin{aligned} \Pr(\mathbf{h} = \text{dry}) &= \Pr(\mathbf{h} = \text{dry grass} | \gamma = \text{rainy night}) \cdot \Pr(\gamma = \text{rainy night}) \\ &\quad + \Pr(\mathbf{h} = \text{dry grass} | \gamma = \text{dry night}) \cdot \Pr(\gamma = \text{dry night}) \\ &= 0.1 \cdot 0.5 + 0.95 \cdot 0.5 = 0.05 + 0.475 = 0.525 \end{aligned} \quad (1.8)$$

Bayes' rule then allows us to find the desired conditional probabilities via:

$$\Pr(\gamma = \gamma | \mathbf{h} = h) = \frac{\Pr(\mathbf{h} = h | \gamma = \gamma) \cdot \Pr(\gamma = \gamma)}{\Pr(\mathbf{h} = h)} \quad (1.9)$$

We have:

$$\Pr(\gamma = \text{rainy night} | \mathbf{h} = \text{wet grass}) = \frac{0.9 \cdot 0.5}{0.475} \approx 0.947 \quad (1.10)$$

$$\Pr(\gamma = \text{rainy night} | \mathbf{h} = \text{dry grass}) = \frac{0.1 \cdot 0.5}{0.525} \approx 0.095 \quad (1.11)$$

$$\Pr(\gamma = \text{dry night} | \mathbf{h} = \text{wet grass}) = \frac{0.05 \cdot 0.5}{0.475} \approx 0.053 \quad (1.12)$$

$$\Pr(\gamma = \text{dry night} | \mathbf{h} = \text{dry grass}) = \frac{0.95 \cdot 0.5}{0.525} \approx 0.905 \quad (1.13)$$

We can interpret these numbers as follows. If we step outside, and observe wet grass, it is significantly more likely that it rained overnight, than it didn't. A

reasonable conclusion given the observation of wet grass is then that it rained overnight. Similarly, if we observe dry grass, it is significantly more likely that it did not rain overnight, making this a reasonable conclusion.

Example 1.2 (Predicting the Output of a Linear Channel) Let us consider a signal $\mathbf{h} \in \mathbb{R}^M$, sent over a linear channel parametrized by $w^o \in \mathbb{R}^M$, resulting in a measurement γ :

$$\gamma = \mathbf{h}^\top w^o + \mathbf{v} \quad (1.14)$$

Here, \mathbf{v} denotes measurement noise, that we may assume to be normally distributed $\mathbf{v} \sim \mathcal{N}(0, \sigma_v^2)$ with variance $\sigma_v^2 > 0$. Suppose we would like to predict the most likely measurement γ for a given signal \mathbf{h} . We can then formulate the conditional pdf:

$$f_{\gamma|\mathbf{h}}(\gamma|h) = \frac{1}{\sqrt{2\pi\sigma_v^2}} e^{-\frac{1}{2} \frac{(\gamma - \mathbf{h}^\top w^o)^2}{\sigma_v^2}} \quad (1.15)$$

It then follows directly that:

$$\gamma^* = \arg \max_{\gamma \in \Gamma} f_{\gamma|\mathbf{h}}(\gamma|h) = \mathbf{h}^\top w^o \quad (1.16)$$

The above examples have illustrated how the Bayesian framework (1.1) can be used to formalize inference problems given data, conditional likelihoods and prior information. In its current form, however, the framework has two important limitations.

- First, relation (1.1) and the following examples provided a rule to map a single observation h to an estimate γ^* of γ . In most inference and learning settings, we are provided with multiple realizations $\{h_n\}_{n=1}^N$ of the random variable \mathbf{h} , and seek to reconcile multiple realizations in generating γ^* . We show in Section 1.1.1 how (1.1) can be generalized to such settings.
- Perhaps more importantly, the resulting estimation rules are driven in large part on the conditional likelihoods (e.g., $f_{\gamma|\mathbf{h}}(\gamma|h)$ when estimating γ from \mathbf{h}). This essentially requires full knowledge of the statistical model relating the label γ to its feature \mathbf{h} . For instance, in the case of the linear channel treated in Example 1.2, this translates into knowledge of the channel taps contained in w^o , along with the distribution of the noise term \mathbf{v} . In practice, we do not expect to have access to this knowledge, and need to estimate relevant conditional distributions directly from data. We refer to the process of estimating models from data as “learning” and develop it from (1.1) in 1.2.

1.1.1 Inference using a Batch of Features

Let us now consider a setting where, for a single observation of the random variable γ , we observe multiple related random variables $\{\mathbf{h}_n\}_{n=1}^N$ through their realizations $\{h_n\}_{n=1}^N$. Although this condition can be relaxed significantly, we will assume for the time being that the observations $\{\mathbf{h}_n\}_{n=1}^N$ are distributed identically and independently after conditioning on the state γ . Formally, this translates to:

$$f_{\{\mathbf{h}_n\}_{n=1}^N|\gamma}(\{\mathbf{h}_n\}_{n=1}^N|\gamma) \stackrel{(a)}{=} \prod_{n=1}^N f_{\mathbf{h}_n|\gamma}(h_n|\gamma) \stackrel{(b)}{=} \prod_{n=1}^N f_{\mathbf{h}|\gamma}(h_n|\gamma) \quad (1.17)$$

where (a) holds by conditional independence and (b) holds by identical distribution of the random variables \mathbf{h}_n . We can then adjust the framework (1.1) to include multiple observations as follows:

$$\begin{aligned} \gamma^* &\triangleq \arg \max_{\gamma \in \Gamma} f_{\gamma|\{\mathbf{h}_n\}_{n=1}^N}(\gamma|\{\mathbf{h}_n\}_{n=1}^N) \\ &= \arg \max_{\gamma \in \Gamma} \frac{f_{\{\mathbf{h}_n\}_{n=1}^N|\gamma}(\{\mathbf{h}_n\}_{n=1}^N|\gamma) \cdot f_{\gamma}(\gamma)}{f_{\{\mathbf{h}_n\}_{n=1}^N}(\{\mathbf{h}_n\}_{n=1}^N)} \\ &= \arg \max_{\gamma \in \Gamma} f_{\{\mathbf{h}_n\}_{n=1}^N|\gamma}(\{\mathbf{h}_n\}_{n=1}^N|\gamma) \cdot f_{\gamma}(\gamma) \\ &= \arg \max_{\gamma \in \Gamma} \left(\prod_{n=1}^N f_{\mathbf{h}|\gamma}(h_n|\gamma) \right) \cdot f_{\gamma}(\gamma) \end{aligned} \quad (1.18)$$

As we saw in Example 1.2 there are situations where optimizing the *logarithm* of the posterior distribution is more tractable. The fact that the logarithmic function is monotonically increasing, ensures that:

$$\begin{aligned} \gamma^* &= \arg \max_{\gamma \in \Gamma} \log \left[\left(\prod_{n=1}^N f_{\mathbf{h}|\gamma}(h_n|\gamma) \right) \cdot f_{\gamma}(\gamma) \right] \\ &= \arg \max_{\gamma \in \Gamma} \left\{ \sum_{n=1}^N \log f_{\mathbf{h}|\gamma}(h_n|\gamma) + \log f_{\gamma}(\gamma) \right\} \\ &= \arg \min_{\gamma \in \Gamma} \left\{ - \sum_{n=1}^N \log f_{\mathbf{h}|\gamma}(h_n|\gamma) - \log f_{\gamma}(\gamma) \right\} \end{aligned} \quad (1.19)$$

We can normalize by the sample size N , since the minimizing argument γ^* is unaffected by scaling:

$$\boxed{\gamma^* = \arg \min_{\gamma \in \Gamma} \left\{ -\frac{1}{N} \sum_{n=1}^N \log f_{\mathbf{h}|\gamma}(h_n|\gamma) - \frac{1}{N} \log f_{\gamma}(\gamma) \right\}} \quad (1.20)$$

Example 1.3 (Inverting a Linear Channel) We consider a binary class variable $\gamma \in \{+1, -1\}$. Suppose that the observed feature \mathbf{h} is normally distributed with mean that depends on the class variable γ . Specifically, we model $\mathbf{h}_n \sim \mathcal{N}(w^o, \sigma_v^2 I_M)$ when $\gamma = +1$ and $\mathbf{h}_n \sim \mathcal{N}(-w^o, \sigma_v^2 I_M)$ when $\gamma = -1$. This relation can be summarized through the linear model:

$$\mathbf{h} = \gamma w^o + \mathbf{v} \quad (1.21)$$

with $\mathbf{v} \sim \mathcal{N}(0, \sigma_v^2)$. In this sense, we may view the present example as the inverse problem to Example 1.2. The linear relationship translates to the conditional distribution:

$$f_{\mathbf{h}|\gamma}(\mathbf{h}|\gamma) = \frac{1}{\sqrt{(2\pi)^M \sigma_v^M}} e^{-\frac{1}{2} \frac{\|\mathbf{h} - \gamma w^o\|^2}{\sigma_v^2}} \quad (1.22)$$

Applying (1.20), and assuming an equal prior $f_\gamma(+1) = f_\gamma(-1) = \frac{1}{2}$, we obtain:

$$\begin{aligned} \gamma^* &= \arg \min_{\gamma \in \pm 1} \frac{1}{N} \sum_{n=1}^N \|\mathbf{h}_n - \gamma w^o\|^2 \\ &= \arg \min_{\gamma \in \pm 1} \frac{1}{N} \sum_{n=1}^N \left\{ \|\mathbf{h}_n\|^2 - 2\gamma \mathbf{h}_n^\top w^o + \|\gamma w^o\|^2 \right\} \\ &= \arg \min_{\gamma \in \pm 1} \left\{ \frac{1}{N} \sum_{n=1}^N \|\mathbf{h}_n\|^2 - 2\gamma \left(\frac{1}{N} \sum_{n=1}^N \mathbf{h}_n \right)^\top w^o + \|\gamma w^o\|^2 \right\} \\ &\stackrel{(a)}{=} \arg \min_{\gamma \in \pm 1} \left\{ -2\gamma \left(\frac{1}{N} \sum_{n=1}^N \mathbf{h}_n \right)^\top w^o + \gamma^2 \|w^o\|^2 \right\} \\ &\stackrel{(b)}{=} \arg \min_{\gamma \in \pm 1} \left\{ -2\gamma \left(\frac{1}{N} \sum_{n=1}^N \mathbf{h}_n \right)^\top w^o \right\} \\ &\stackrel{(c)}{=} \arg \max_{\gamma \in \pm 1} \left\{ 2\gamma \left(\frac{1}{N} \sum_{n=1}^N \mathbf{h}_n \right)^\top w^o \right\} \\ &\stackrel{(d)}{=} \text{sign} \left\{ \left(\frac{1}{N} \sum_{n=1}^N \mathbf{h}_n \right)^\top w^o \right\} \end{aligned} \quad (1.23)$$

In this sequence of reformulations, (a) follows from the fact that $\frac{1}{N} \sum_{n=1}^N \|\mathbf{h}_n\|^2$ is independent of γ , (b) uses the fact that $\gamma^2 = 1$ with probability one and (c) reverses minimization and maximization. Step (d) uses the fact that

$$\gamma \left(\frac{1}{N} \sum_{n=1}^N \mathbf{h}_n \right)^\top w^o \geq 0 \quad (1.24)$$

if, and only if, $\gamma = \text{sign} \left\{ \left(\frac{1}{N} \sum_{n=1}^N h_n \right)^\top w^o \right\}$.

1.2 FROM INFERENCE TO LEARNING

Our previous discussion, culminating in (1.20), has uncovered how to perform optimal inference of a random quantity of interest γ given observations $\{h_n\}_{n=1}^N$, models $f_{h|\gamma}(\cdot|\gamma)$ and prior information $f_\gamma(\gamma)$. In many practical situations, we are not provided with prior knowledge about the model, and instead need to estimate $f_{h|\gamma}(\cdot|\gamma)$ from data before performing inference. We refer to the process of estimating $f_{h|\gamma}(\cdot|\gamma)$ from data as *learning*. As we will see in this section, the learning of models can be formalized using a Bayesian framework analogous to (1.1). To this end, we will assume that the conditional likelihood $f_{h|\gamma}(\cdot|\gamma)$ is parametrized by a set of learnable parameters $w \in \mathbb{R}^{M_w}$, denoted by $f_{h|\gamma,w}(\cdot|\gamma, w)$. Under this parametrization, learning the conditional distribution of γ given h is equivalent to learning the parameters w that parametrize $f_{h|\gamma,w}(\cdot|\gamma, w)$.

To formulate a procedure for learning w from pairs $\{h, \gamma\}$ we mirror the argument in Section 1.1.1. Suppose we collect labeled data in the form of pairs $\{h_n, \gamma_n\}_{n=1}^N$. If we suppose that feature-label pairs $\{h_n, \gamma_n\}$ are identically and independently distributed for a given set of parameters w , we can factorize:

$$\begin{aligned} f_{\{h_n, \gamma_n\}_{n=1}^N | w} \left(\{h_n, \gamma_n\}_{n=1}^N | w \right) &\stackrel{(a)}{=} \prod_{n=1}^N f_{h_n, \gamma_n | w} (h_n, \gamma_n | w) \\ &\stackrel{(b)}{=} \prod_{n=1}^N f_{h, \gamma | w} (h_n, \gamma_n | w) \end{aligned} \quad (1.25)$$

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

$$\begin{aligned} w^* &\triangleq \arg \max_{w \in \Omega} f_{w | \{h_n, \gamma_n\}_{n=1}^N} \left(w | \{h_n, \gamma_n\}_{n=1}^N \right) \\ &= \arg \max_{w \in \Omega} \frac{f_{\{h_n, \gamma_n\}_{n=1}^N | w} \left(\{h_n, \gamma_n\}_{n=1}^N | w \right) \cdot f_w(w)}{f_{\{h_n, \gamma_n\}_{n=1}^N} \left(\{h_n, \gamma_n\}_{n=1}^N \right)} \\ &= \arg \max_{w \in \Omega} f_{\{h_n, \gamma_n\}_{n=1}^N | w} \left(\{h_n, \gamma_n\}_{n=1}^N | w \right) \cdot f_w(w) \\ &= \arg \max_{w \in \Omega} \left(\prod_{n=1}^N f_{h, \gamma | w} (h_n, \gamma_n | w) \right) \cdot f_w(w) \end{aligned} \quad (1.26)$$

Following the same argument that led to (1.20), we arrive at:

$$w^* = \arg \min_{w \in \Omega} \left\{ -\frac{1}{N} \sum_{n=1}^N \log f_{\mathbf{h}, \gamma | \mathbf{w}}(h_n, \gamma_n | w) - \frac{1}{N} \log f_{\mathbf{w}}(w) \right\} \quad (1.27)$$

Example 1.4 (Ridge Regression for Channel Estimation) Let us consider again the linear channel from Example 1.2:

$$\gamma = \mathbf{h}^\top \mathbf{w} + v \quad (1.28)$$

We assume that the parameters \mathbf{w} are independent of all other random variables, impose a normal prior of the form $\mathbf{w} \sim \mathcal{N}(0, \sigma_w^2 I_M)$. We continue to assume $v \sim \mathcal{N}(0, \sigma_v^2)$ and $\mathbf{h} \sim \mathcal{N}(0, \sigma_h^2)$. It then follows that:

$$f_{\mathbf{w}}(w) = \frac{1}{\sqrt{(2\pi)^M \sigma_w^2}} e^{-\frac{1}{2} \frac{\|\mathbf{w}\|^2}{\sigma_w^2}} \quad (1.29)$$

$$\begin{aligned} f_{\mathbf{h}, \gamma | \mathbf{w}}(h, \gamma | w) &= f_{\gamma | \mathbf{h}, \mathbf{w}}(\gamma | h, w) \cdot f_{\mathbf{h} | \mathbf{w}}(h | w) \\ &\stackrel{(a)}{=} f_{\gamma | \mathbf{h}, \mathbf{w}}(\gamma | h, w) \cdot f_{\mathbf{h}}(h) \\ &= \frac{1}{\sqrt{2\pi \sigma_v^2}} e^{-\frac{1}{2} \frac{(\gamma - \mathbf{h}^\top \mathbf{w})^2}{\sigma_v^2}} \cdot \frac{1}{\sqrt{(2\pi)^{M_h} \sigma_h^{2M_h}}} e^{-\frac{1}{2} \frac{\|\mathbf{h}\|^2}{\sigma_h^2}} \end{aligned} \quad (1.30)$$

Plugging these relations into (1.27), and removing constant terms independent of w , we obtain:

$$w^* = \arg \min_{w \in \Omega} \left\{ \frac{1}{2N} \sum_{n=1}^N (\gamma_n - \mathbf{h}_n^\top \mathbf{w})^2 + \frac{1}{2N} \frac{\sigma_v^2}{\sigma_w^2} \|\mathbf{w}\|^2 \right\} \quad (1.31)$$

1.2.1 Asymptotic Behavior

It is instructive to examine the asymptotic behavior of the optimal estimate (1.27) as the sample size N grows. We have for all w such that $f_{\mathbf{w}}(w) > 0$:

$$\begin{aligned} &\lim_{N \rightarrow \infty} \left\{ -\frac{1}{N} \sum_{n=1}^N \log f_{\mathbf{h}, \gamma | \mathbf{w}}(h_n, \gamma_n | w) - \frac{1}{N} \log f_{\mathbf{w}}(w) \right\} \\ &= \lim_{N \rightarrow \infty} \left\{ -\frac{1}{N} \sum_{n=1}^N \log f_{\mathbf{h}, \gamma | \mathbf{w}}(h_n, \gamma_n | w) \right\} - \lim_{N \rightarrow \infty} \left\{ \frac{1}{N} \log f_{\mathbf{w}}(w) \right\} \\ &\stackrel{(a)}{=} \lim_{N \rightarrow \infty} \left\{ -\frac{1}{N} \sum_{n=1}^N \log f_{\mathbf{h}, \gamma | \mathbf{w}}(h_n, \gamma_n | w) \right\} \\ &\stackrel{(b)}{=} -\mathbb{E}_{\mathbf{h}, \gamma} \{ \log f_{\mathbf{h}, \gamma | \mathbf{w}}(\mathbf{h}, \gamma | w) \} \end{aligned} \quad (1.32)$$

Here, (a) follows since $f_{\mathbf{w}}(w) > 0 \implies |\log f_{\mathbf{w}}(w)| < \infty$ and hence:

$$\lim_{N \rightarrow \infty} \left\{ \frac{1}{N} \log f_{\mathbf{w}}(w) \right\} = 0 \quad (1.33)$$

Step (b) follows from the law of large numbers under suitable regularity conditions on the distributions of the involved random variables. The precise nature and rate of convergence as well as limiting distribution of the estimator can be quantified more finely. To avoid a digression, we refer the reader to the appendix for a more detailed discussion on asymptotic properties of Bayesian estimates. For the purpose of our discussion, the key take-away from (1.32) is that it motivates the introduction of the loss:

$$Q(w; \mathbf{h}, \gamma) \triangleq -\log f_{\mathbf{h}, \gamma | \mathbf{w}}(\mathbf{h}, \gamma | w) \quad (1.34)$$

as well as the risk:

$$J(w) \triangleq \mathbb{E}_{\mathbf{h}, \gamma} Q(w; \mathbf{h}, \gamma) = -\mathbb{E}_{\mathbf{h}, \gamma} \log f_{\mathbf{h}, \gamma | \mathbf{w}}(\mathbf{h}, \gamma | w) \quad (1.35)$$

We may then define the asymptotically optimal model:

$$w^o \triangleq \arg \min_{w \in \Omega} J(w) = \arg \min_{w \in \Omega} \left\{ -\mathbb{E}_{\mathbf{h}, \gamma} \log f_{\mathbf{h}, \gamma | \mathbf{w}}(\mathbf{h}, \gamma | w) \right\} \quad (1.36)$$

Example 1.5 (Least Mean-Square (LMS) Estimation) Considering the linear model from Examples 1.4, we find in the limit as $N \rightarrow \infty$ from (1.31), we obtain the negative log-likelihood as:

$$\begin{aligned} & \log f_{\mathbf{h}, \gamma | \mathbf{w}}(\mathbf{h}, \gamma | w) \\ \stackrel{(1.30)}{=} & -\log \left(\frac{1}{\sqrt{2\pi}\sigma_v^2} e^{-\frac{1}{2} \frac{(\gamma - \mathbf{h}^\top \mathbf{w})^2}{\sigma_v^2}} \cdot \frac{1}{\sqrt{(2\pi)^{M_h} \sigma_h^{2M_h}}} e^{-\frac{1}{2} \frac{\|\mathbf{h}\|^2}{\sigma_h^2}} \right) \\ = & -\log \left(e^{-\frac{1}{2} \frac{(\gamma - \mathbf{h}^\top \mathbf{w})^2}{\sigma_v^2}} \right) - \underbrace{\log \left(\frac{1}{\sqrt{2\pi}\sigma_v^2} \cdot \frac{1}{\sqrt{(2\pi)^{M_h} \sigma_h^{2M_h}}} e^{-\frac{1}{2} \frac{\|\mathbf{h}\|^2}{\sigma_h^2}} \right)}_{\triangleq C} \\ = & \frac{1}{2\sigma_v^2} \|\gamma - \mathbf{h}^\top \mathbf{w}\|^2 + C \end{aligned} \quad (1.37)$$

Since $\frac{1}{\sigma_v^2}$ and C are independent of w , we conclude after shifting and scaling:

$$w^o = \arg \min_w \frac{1}{2} \mathbb{E} \|\gamma - \mathbf{h}^\top w\|^2 \quad (1.38)$$

Since the objective function in the case of the least mean-square estimation problem takes the form of a quadratic, we can develop a closed-form expression

for w^o . Indeed, after differentiating, we find:

$$\nabla_w \mathbb{E} \|\gamma - \mathbf{h}^\top w\|^2 = -\underbrace{\mathbb{E} \gamma \mathbf{h}}_{\triangleq r_{\gamma h}} + \underbrace{\mathbb{E} \mathbf{h} \mathbf{h}^\top}_{\triangleq R_h} w = -r_{\gamma h} + R_h w \quad (1.39)$$

At w^o , this derivative must be equal to zero, and hence after rearranging:

$$w^o = R_h^{-1} r_{\gamma h} \quad (1.40)$$

The perceptive reader will notice that we chose the notation w^o to denote the optimal solution to the least-mean square error cost (1.38), while in case of a deterministic linear channel in Example 1.2, w^o corresponds to the *true* channel weights in (1.14). This choice is deliberate. Indeed, if we condition on $\mathbf{w} = w^o$, we have almost surely:

$$\gamma = \mathbf{h}^\top w^o + \mathbf{v} \quad (1.41)$$

If we multiply by \mathbf{h} from the left and take expectations, we find:

$$r_{\gamma h} = R_h w^o \quad (1.42)$$

Rearranging again yields $w^o = R_h^{-1} r_{\gamma h}$, and we can conclude that in the case of a linear channel, the optimal solution w^o of the least-squares cost (1.38) corresponds to the true channel.

1.3 CASE STUDY: DATA FUSION

In the previous sections we introduced the Bayesian framework as a useful methodology to reason about optimal inference and learning. We concluded that in many cases, the answer boils down to solving an optimization problem of the form (1.27) or (1.36), which depends the data and the associated log likelihoods. We will now proceed to illustrate how this same methodology can be applied to develop optimal techniques for fusing multiple data sources.

To this end, we consider a collection of K agents. Each agent has sensing capabilities, and collects a local data set $\{\mathbf{h}_{k,n}, \gamma_{k,n}\}_{n=1}^N$. For simplicity, we assume that the sample size N is common to all agents. Each pair of samples is sampled independently and identically from the distribution:

$$f_{\mathbf{h}_{k,n}, \gamma_{k,n} | \mathbf{w}_k}(h, \gamma | w) = f_{\mathbf{h}, \gamma | \mathbf{w}}(h, \gamma | w) \quad (1.43)$$

Now suppose different agents are completely disconnected, with no way to exchange information. In that case, it is reasonable to formulate independent local learning problems:

$$w_k^* \triangleq \arg \max_{w \in \Omega} f_{\mathbf{w} | \{\mathbf{h}_{k,n}, \gamma_{k,n}\}_{n=1}^N} \left(w | \{h_{k,n}, \gamma_{k,n}\}_{n=1}^N \right) \quad (1.44)$$

Following the argument laid out in Sec. [1.2](#) we arrive at the equivalent formulation:

$$\begin{aligned} w_k^* &= \arg \min_{w \in \Omega} \left\{ -\frac{1}{N} \sum_{n=1}^N \log f_{\mathbf{h}, \gamma | \mathbf{w}}(h_{k,n}, \gamma_{k,n} | w) - \frac{1}{N} \log f_{\mathbf{w}}(w) \right\} \\ &= \arg \min_{w \in \Omega} J_k(w) + R(w) \end{aligned} \quad (1.45)$$

where we defined:

$$J_k(w) \triangleq -\frac{1}{N} \sum_{n=1}^N \log f_{\mathbf{h}, \gamma | \mathbf{w}}(h_{k,n}, \gamma_{k,n} | w) \quad (1.46)$$

$$R(w) \triangleq -\frac{1}{N} \log f_{\mathbf{w}}(w) \quad (1.47)$$

Through the subscript k we emphasize here that because the local samples $\{\mathbf{h}_{k,n}, \gamma_{k,n}\}_{n=1}^N$ are distinct, this implies that $J_k(\cdot)$ as well as the resulting optimal models w_k^o are distinct as well. While we can claim local optimality of each model w_k^o in the sense of [\(1.45\)](#), this model is making use only of the data available to agent k .

As an alternative to [\(1.45\)](#), we can instead define the *global* inference problem:

$$w^* \triangleq \arg \max_{w \in \Omega} f_{\mathbf{w} | \{\{\mathbf{h}_{k,n}, \gamma_{k,n}\}_{n=1}^N\}_{k=1}^K} \left(w | \left\{ \{\mathbf{h}_{k,n}, \gamma_{k,n}\}_{n=1}^N \right\}_{k=1}^K \right) \quad (1.48)$$

As long as the samples $\{\mathbf{h}_{k,n}, \gamma_{k,n}\}_{n=1}^N$ are also mutually independent between agents after conditioning on \mathbf{w} , we can factorize the conditional likelihoods:

$$\begin{aligned} w^* &\triangleq \arg \max_{w \in \Omega} f_{\mathbf{w} | \{\{\mathbf{h}_{k,n}, \gamma_{k,n}\}_{n=1}^N\}_{k=1}^K} \left(w | \left\{ \{\mathbf{h}_{k,n}, \gamma_{k,n}\}_{n=1}^N \right\}_{k=1}^K \right) \\ &= \arg \max_{w \in \Omega} \frac{f_{\{\{\mathbf{h}_{k,n}, \gamma_{k,n}\}_{n=1}^N\}_{k=1}^K | \mathbf{w}} \left(\left\{ \{\mathbf{h}_{k,n}, \gamma_{k,n}\}_{n=1}^N \right\}_{k=1}^K | w \right) \cdot f_{\mathbf{w}}(w)}{f_{\{\{\mathbf{h}_{k,n}, \gamma_{k,n}\}_{n=1}^N\}_{k=1}^K} \left(\left\{ \{\mathbf{h}_{k,n}, \gamma_{k,n}\}_{n=1}^N \right\}_{k=1}^K \right)} \\ &= \arg \max_{w \in \Omega} f_{\{\{\mathbf{h}_{k,n}, \gamma_{k,n}\}_{n=1}^N\}_{k=1}^K | \mathbf{w}} \left(\left\{ \{\mathbf{h}_{k,n}, \gamma_{k,n}\}_{n=1}^N \right\}_{k=1}^K | w \right) \cdot f_{\mathbf{w}}(w) \\ &= \arg \max_{w \in \Omega} \left(\prod_{k=1}^K f_{\{\mathbf{h}_{k,n}, \gamma_{k,n}\}_{n=1}^N | \mathbf{w}} \left(\{\mathbf{h}_{k,n}, \gamma_{k,n}\}_{n=1}^N | w \right) \right) \cdot f_{\mathbf{w}}(w) \end{aligned} \quad (1.49)$$

After taking logarithms, normalizing and making the substitutions [\(1.46\)](#)–[\(1.47\)](#), we find:

$$\begin{aligned} w^* &= \arg \min_w \left\{ -\frac{1}{KN} \sum_{k=1}^K \sum_{n=1}^N \log f_{\mathbf{h}, \gamma | \mathbf{w}}(h_{k,n}, \gamma_{k,n} | w) - \frac{1}{KN} \log f_{\mathbf{w}}(w) \right\} \\ &= \arg \min_w \frac{1}{K} \sum_{k=1}^K J_k(w) + \frac{1}{K} R(w) \end{aligned} \quad (1.50)$$

We refer to [\(1.50\)](#) as the *aggregate* optimization problem of the collection of

agents, since it is obtained by aggregating data $x_{k,n}$ from across the network. The formulation provides a recipe for optimal learning in a collection of agents.

1.3.1 Example: Least-Squares Fusion

Let us specialize our discussion to the least-squares problem of Example 1.4. For the local models we then have:

$$\begin{aligned} w_k^* &= \arg \min_{w \in \Omega} \left\{ \frac{1}{2N} \sum_{n=1}^N (\gamma_{k,n} - h_{k,n}^\top w)^2 + \frac{1}{2N} \frac{\sigma_v^2}{\sigma_w^2} \|w\|^2 \right\} \\ &= \arg \min_{w \in \Omega} \left\{ \frac{1}{2N} \sum_{n=1}^N (\gamma_{k,n} - h_{k,n}^\top w)^2 + \frac{\rho}{2} \|w\|^2 \right\} \end{aligned} \quad (1.51)$$

where we defined $\rho = \frac{\sigma_v^2}{\sigma_w^2}$ for convenience. Since (1.51) is a quadratic, we can determine a closed-form expression of w_k^* by differentiating:

$$-\frac{1}{N} \sum_{n=1}^N \gamma_{k,n} h_{k,n} + \left(\frac{1}{N} \sum_{n=1}^N h_{k,n} h_{k,n}^\top \right) w_k^* + \frac{\rho}{N} w_k^* = 0 \quad (1.52)$$

or after inverting:

$$\begin{aligned} w_k^* &= \left(\sum_{n=1}^N h_{k,n} h_{k,n}^\top + \rho I \right)^{-1} \left(\sum_{n=1}^N \gamma_{k,n} h_{k,n} \right) \\ &= (H_k + \rho I)^{-1} d_k \end{aligned} \quad (1.53)$$

where we defined for compactness:

$$H_k = \sum_{n=1}^N h_{k,n} h_{k,n}^\top \quad (1.54)$$

$$d_k = \sum_{n=1}^N \gamma_{k,n} h_{k,n} \quad (1.55)$$

For the aggregate problem (1.50) we have analogously:

$$\begin{aligned} w^* &= \left(\sum_{k=1}^K \sum_{n=1}^N h_{k,n} h_{k,n}^\top + \rho I \right)^{-1} \left(\sum_{k=1}^K \sum_{n=1}^N \gamma_{k,n} h_{k,n} \right) \\ &= (H + \rho I)^{-1} d \end{aligned} \quad (1.56)$$

where we defined:

$$H = \sum_{k=1}^K \sum_{n=1}^N h_{k,n} h_{k,n}^\top \quad (1.57)$$

$$d = \sum_{k=1}^K \sum_{n=1}^N \gamma_{k,n} h_{k,n} \quad (1.58)$$

These relationships show how the measurements $h_{k,n}$ ought to be processed to compute an optimal estimate of the weight vector w^o . As it stands, however, evaluating (1.56) requires *central aggregation* of the raw data $h_{k,n}, \gamma_{k,n}$ across the entire network. A natural question is then whether we can learn w^* by sharing only a processed and compressed version of the raw data. In the case of a least-squares problem, as we show in this case study, this is rather straightforward. For general, and more complex learning problems, this will turn out to be substantially more involved, and we will learn the relevant tools in future chapters. Returning to the quadratic problem, note that:

$$(H + \rho I)w^* = d = \sum_{k=1}^K d_k = \sum_{k=1}^K (H_k + \rho I)w_k^* \quad (1.59)$$

It then follows that w^* can be evaluated equivalently as:

$$w^* = \left(\sum_{k=1}^K H_k + \rho I \right)^{-1} \left(\sum_{k=1}^K (H_k + \rho I)w_k^* \right) \quad (1.60)$$

In other words, the globally optimal estimate w^* is fully determined by the locally optimal estimates w_k^* and the local matrices H_k . This means that rather than communicate and collect raw data $h_{k,n}$ at fusion center, and subsequently computing w^* via (1.56), we can instead locally compute H_k and w_k^* , only transmit those processed quantities, and subsequently find w^* via (1.60). The first approach requires the communication of $K(M+1)N$ scalar quantities, where K is the number of agents, M is the dimension of the feature vector, and N is the local sample size. The second approach requires K times the exchange of a matrix of size M^2 and a vector of size M , hence a cost of $K(M+1)M$. Sample sizes are generally significantly larger than the dimension, and hence we benefit from reduced communication cost. We have developed our first distributed algorithm! We will further examine its properties in the problems at the end of this chapter.

1.4 PROBLEMS

1.1 In this problem we implement the least-squares fusion technique in (1.60) in code to verify empirically verify the benefit of data fusion. You are free to choose any unspecified parameters.

- Generate local data sets $\{h_{k,n}, \gamma_{k,n}\}_{n=1}^N$ following the statistical model of Example 1.4.
- Compute w_k^* for each agent along with the globally optimal model w^* .
- Generate a new independent test set $\{\tilde{h}_n, \tilde{\gamma}_n\}_{n=1}^{\tilde{N}}$ and evaluate prediction perfor-

mance of local models:

$$\frac{1}{\tilde{N}} \sum_{n=1}^{\tilde{N}} (\tilde{\gamma}_n - \tilde{h}_n^\top w_k^*)^2 \quad (1.61)$$

as well as of the performance of the global model:

$$\frac{1}{\tilde{N}} \sum_{n=1}^{\tilde{N}} (\tilde{\gamma}_n - \tilde{h}_n^\top w^*)^2 \quad (1.62)$$

Compare the performance for varying choices of N and K . What do you conclude?

1.2 We consider again the setting of Section 1.3 but generalize to the setting where local data sets $\{\mathbf{h}_{k,n}, \gamma_{k,n}\}_{n=1}^{N_k}$ have differing sizes N_k . Adjust the derivation leading up to (1.60) as necessary to find an optimal fusion protocol in this generalized setting. How do the local data sizes N_k affect the fusion?