

where  $0 < \eta_t \leq 1$  is the step size, and

$$\mathbf{g}_t = \frac{\partial}{\partial \boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}_t) |_{\boldsymbol{\theta}=\boldsymbol{\theta}_t} \quad (6.166)$$

$$\mathbf{H}_t = \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} Q(\boldsymbol{\theta}, \boldsymbol{\theta}_t) |_{\boldsymbol{\theta}=\boldsymbol{\theta}_t} \quad (6.167)$$

If  $\eta_t = 1$ , [Lan95a] calls this the **gradient EM algorithm**. However, it is possible to use a larger step size to speed up the algorithm, as in the **quasi-Newton EM algorithm** of [Lan95b]. This method also replaces the Hessian in Equation (6.167), which may not be negative definite (for non exponential family models), with a BFGS approximation. This ensures the overall algorithm is an ascent algorithm. Note, however, when the M step cannot be computed in closed form, EM loses some of its appeal over directly optimizing the marginal likelihood with a gradient based solver.

### 6.6.6.5 ECM algorithm

The **ECM** algorithm stands for “expectation conditional maximization”, and refers to optimizing the parameters in the M step sequentially, if they turn out to be dependent. The **ECME** algorithm, which stands for “ECM either” [LR95], is a variant of ECM in which we maximize the expected complete data log likelihood (the  $Q$  function) as usual, or the observed data log likelihood, during one or more of the conditional maximization steps. The latter can be much faster, since it ignores the results of the E step, and directly optimizes the objective of interest. A standard example of this is when fitting the Student T distribution. For fixed  $\nu$ , we can update  $\boldsymbol{\Sigma}$  as usual, but then to update  $\nu$ , we replace the standard update of the form  $\nu^{t+1} = \arg \max_{\nu} Q((\boldsymbol{\mu}^{t+1}, \boldsymbol{\Sigma}^{t+1}, \nu), \boldsymbol{\theta}^t)$  with  $\nu^{t+1} = \arg \max_{\nu} \log p(\mathcal{D} | \boldsymbol{\mu}^{t+1}, \boldsymbol{\Sigma}^{t+1}, \nu)$ . See [MK97] for more information.

### 6.6.6.6 Online EM

When dealing with large or streaming datasets, it is important to be able to learn online, as we discussed in Section 19.7.5. There are two main approaches to **online EM** in the literature. The first approach, known as **incremental EM** [NH98a], optimizes the lower bound  $Q(\boldsymbol{\theta}, q_1, \dots, q_N)$  one  $q_n$  at a time; however, this requires storing the expected sufficient statistics for each data case.

The second approach, known as **stepwise EM** [SI00; LK09; CM09], is based on stochastic gradient descent. This optimizes a local upper bound on  $\ell_n(\boldsymbol{\theta}) = \log p(\mathbf{x}_n | \boldsymbol{\theta})$  at each step. (See [Mai13; Mai15] for a more general discussion of stochastic and incremental bound optimization algorithms.)

## 6.7 The Bayesian learning rule

in this section, we discuss the “**Bayesian learning rule**” [KR21a], which provides a unified framework for deriving many standard (and non-standard) optimization and inference algorithms used in the ML community.

To motivate the BLR, recall the standard **empirical risk minimization** or **ERM** problem, which has the form  $\boldsymbol{\theta}_* = \arg \min_{\boldsymbol{\theta}} \bar{\ell}(\boldsymbol{\theta})$ , where

$$\bar{\ell}(\boldsymbol{\theta}) = \sum_{n=1}^N \ell(\mathbf{y}_n, f_{\boldsymbol{\theta}}(\mathbf{x}_n)) + R(\boldsymbol{\theta}) \quad (6.168)$$

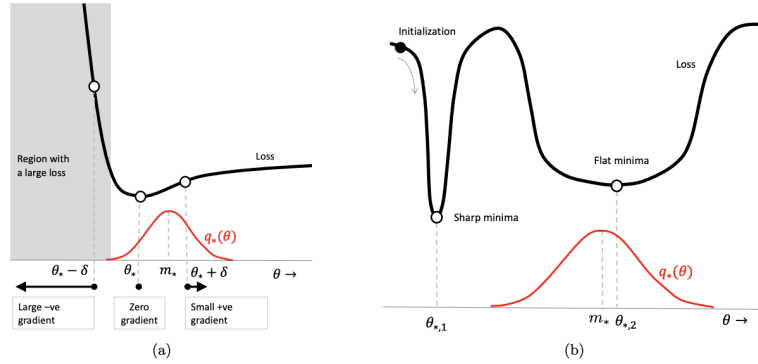


Figure 6.9: Illustration of the robustness obtained by using a Bayesian approach to parameter estimation. (a) When the minimum  $\theta_*$  lies next to a “wall”, the Bayesian solution shifts away from the boundary to avoid large losses due to perturbations of the parameters. (b) The Bayesian solution prefers flat minima over sharp minima, to avoid large losses due to perturbations of the parameters. From Figure 1 of [KR21a]. Used with kind permission of Emtiyaz Khan.

where  $f_{\theta}(\mathbf{x})$  is a prediction function,  $\ell(\mathbf{y}, \hat{\mathbf{y}})$  is a loss function, and  $R(\theta)$  is some kind of regularizer.

Although the regularizer can prevent overfitting, the ERM method can still result in parameter estimates that are not robust. A better approach is to fit a *distribution* over possible parameter values,  $q(\theta)$ . If we minimize the expected loss, we will find parameter settings that will work well even if they are slightly perturbed, as illustrated in Figure 6.9, which helps with robustness and generalization. Of course, if the distribution  $q$  collapses to a single delta function, we will end up with the ERM solution. To prevent this, we add a penalty term, that measures the KL divergence from  $q(\theta)$  to some prior  $\pi_0(\theta) \propto \exp(R(\theta))$ . This gives rise to the following BLR objective:

$$\mathcal{L}(q) = \mathbb{E}_{q(\theta)} \left[ \sum_{n=1}^N \ell(\mathbf{y}_n, f_{\theta}(\mathbf{x}_n)) \right] + D_{\text{KL}}(q(\theta) \parallel \pi_0(\theta)) \quad (6.169)$$

We can rewrite the KL term as

$$D_{\text{KL}}(q(\theta) \parallel \pi_0(\theta)) = \mathbb{E}_{q(\theta)} [R(\theta)] - \mathbb{H}(q(\theta)) \quad (6.170)$$

and hence can rewrite the BLR objective as follows:

$$\mathcal{L}(q) = \mathbb{E}_{q(\theta)} [\bar{\ell}(\theta)] - \mathbb{H}(q(\theta)) \quad (6.171)$$

Below we show that different approximations to this objective recover a variety of different methods in the literature.

### 6.7.1 Deriving inference algorithms from BLR

In this section we show how to derive several different inference algorithms from BLR. (We discuss such algorithms in more detail in Chapter 10.)

### 6.7.1.1 Bayesian inference as optimization

The BLR objective includes standard exact Bayesian inference as a special case, as first shown in [Opt88]. To see this, let us assume the loss function is derived from a log-likelihood:

$$\ell(y, f_{\boldsymbol{\theta}}(\mathbf{x})) = -\log p(\mathbf{y}|\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x})) \quad (6.172)$$

Let  $\mathcal{D} = \{(\mathbf{x}_n, \mathbf{y}_n) : n = 1 : N\}$  be the data we condition on. The Bayesian posterior can be written as

$$p(\boldsymbol{\theta}|\mathcal{D}) = \frac{1}{Z(\mathcal{D})} \pi_0(\boldsymbol{\theta}) \prod_{n=1}^N p(\mathbf{y}_n|\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x}_n)) \quad (6.173)$$

This can be derived by minimizing the BLR, since

$$\mathcal{L}(q) = -\mathbb{E}_{q(\boldsymbol{\theta})} \left[ \sum_{n=1}^N \log p(\mathbf{y}_n|\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x}_n)) \right] + D_{\text{KL}}(q(\boldsymbol{\theta}) \parallel \pi_0(\boldsymbol{\theta})) \quad (6.174)$$

$$= \mathbb{E}_{q(\boldsymbol{\theta})} \left[ \log \frac{q(\boldsymbol{\theta})}{\frac{\pi_0(\boldsymbol{\theta})}{Z(\mathcal{D})} \prod_{n=1}^N p(\mathbf{y}_n|\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x}_n))} \right] - \log Z(\mathcal{D}) \quad (6.175)$$

$$= D_{\text{KL}}(q(\boldsymbol{\theta}) \parallel p(\boldsymbol{\theta}|\mathcal{D})) - \log Z(\mathcal{D}) \quad (6.176)$$

Since  $Z(\mathcal{D})$  is a constant, we can minimize the loss by by setting  $q(\boldsymbol{\theta}) = p(\boldsymbol{\theta}|\mathcal{D})$ .

Of course, we can use other kinds of loss, not just log likelihoods. This results in a framework known as **generalized Bayesian inference** [BHW16; KJD19; KJD21]. See Section 14.1.3 for more discussion.

### 6.7.1.2 Optimization of BLR using natural gradient descent

In general, we cannot compute the exact posterior  $q(\boldsymbol{\theta}) = p(\boldsymbol{\theta}|\mathcal{D})$ , so we seek an approximation. We will assume that  $q(\boldsymbol{\theta})$  is an exponential family distribution, such as a multivariate Gaussian, where the mean represents the standard point estimate of  $\boldsymbol{\theta}$  (as in ERM), and the covariance represents our uncertainty (as in Bayes). Hence  $q$  can be written as follows:

$$q(\boldsymbol{\theta}) = h(\boldsymbol{\theta}) \exp[\boldsymbol{\lambda}^\top \mathcal{T}(\boldsymbol{\theta}) - A(\boldsymbol{\lambda})] \quad (6.177)$$

where  $\boldsymbol{\lambda}$  are the natural parameters,  $\mathcal{T}(\boldsymbol{\theta})$  are the sufficient statistics,  $A(\boldsymbol{\lambda})$  is the log partition function, and  $h(\boldsymbol{\theta})$  is the base measure, which is usually a constant. The BLR loss becomes

$$\mathcal{L}(\boldsymbol{\lambda}) = \mathbb{E}_{q_{\boldsymbol{\lambda}}(\boldsymbol{\theta})} [\bar{\ell}(\boldsymbol{\theta})] - \mathbb{H}(q_{\boldsymbol{\lambda}}(\boldsymbol{\theta})) \quad (6.178)$$

We can optimize this using natural gradient descent (Section 6.4). The update becomes

$$\boldsymbol{\lambda}_{t+1} = \boldsymbol{\lambda}_t - \eta_t \tilde{\nabla}_{\boldsymbol{\lambda}} \left[ \mathbb{E}_{q_{\boldsymbol{\lambda}_t}(\boldsymbol{\theta})} [\bar{\ell}(\boldsymbol{\theta})] - \mathbb{H}(q_{\boldsymbol{\lambda}_t}) \right] \quad (6.179)$$

where  $\tilde{\nabla}_{\boldsymbol{\lambda}}$  denotes the natural gradient. We discuss how to compute these natural gradients in Section 6.4.5. In particular, we can convert it to regular gradients wrt the moment parameters  $\boldsymbol{\mu}_t = \boldsymbol{\mu}(\boldsymbol{\lambda}_t)$ . This gives

$$\boldsymbol{\lambda}_{t+1} = \boldsymbol{\lambda}_t - \eta_t \nabla_{\boldsymbol{\mu}} \mathbb{E}_{q_{\boldsymbol{\mu}_t}(\boldsymbol{\theta})} [\bar{\ell}(\boldsymbol{\theta})] + \eta_t \nabla_{\boldsymbol{\mu}} \mathbb{H}(q_{\boldsymbol{\mu}_t}) \quad (6.180)$$

From Equation (6.84) we have

$$\nabla_{\boldsymbol{\mu}} \mathbb{H}(q) = -\boldsymbol{\lambda} - \nabla_{\boldsymbol{\mu}} \mathbb{E}_{q_{\boldsymbol{\mu}}(\boldsymbol{\theta})} [\log h(\boldsymbol{\theta})] \quad (6.181)$$

Hence the update becomes

$$\boldsymbol{\lambda}_{t+1} = \boldsymbol{\lambda}_t - \eta_t \nabla_{\boldsymbol{\mu}} \mathbb{E}_{q_{\boldsymbol{\mu}_t}(\boldsymbol{\theta})} [\bar{\ell}(\boldsymbol{\theta})] - \eta_t \boldsymbol{\lambda}_t - \eta_t \nabla_{\boldsymbol{\mu}} \mathbb{E}_{q_{\boldsymbol{\mu}}(\boldsymbol{\theta})} [\log h(\boldsymbol{\theta})] \quad (6.182)$$

$$= (1 - \eta_t) \boldsymbol{\lambda}_t - \eta_t \nabla_{\boldsymbol{\mu}} \mathbb{E}_{q_{\boldsymbol{\mu}}(\boldsymbol{\theta})} [\bar{\ell}(\boldsymbol{\theta}) + \log h(\boldsymbol{\theta})] \quad (6.183)$$

For distributions  $q$  with constant base measure  $h(\boldsymbol{\theta})$ , this simplifies to

$$\boldsymbol{\lambda}_{t+1} = (1 - \eta_t) \boldsymbol{\lambda}_t - \eta_t \nabla_{\boldsymbol{\mu}} \mathbb{E}_{q_{\boldsymbol{\mu}}(\boldsymbol{\theta})} [\bar{\ell}(\boldsymbol{\theta})] \quad (6.184)$$

Hence at the fixed point we have

$$\boldsymbol{\lambda}_* = (1 - \eta) \boldsymbol{\lambda}_* - \eta \nabla_{\boldsymbol{\mu}} \mathbb{E}_{q_{\boldsymbol{\mu}}(\boldsymbol{\theta})} [\bar{\ell}(\boldsymbol{\theta})] \quad (6.185)$$

$$\boldsymbol{\lambda}_* = \nabla_{\boldsymbol{\mu}} \mathbb{E}_{q_{\boldsymbol{\mu}}(\boldsymbol{\theta})} [-\bar{\ell}(\boldsymbol{\theta})] = \tilde{\nabla}_{\boldsymbol{\lambda}} \mathbb{E}_{q_{\boldsymbol{\lambda}}(\boldsymbol{\theta})} [-\bar{\ell}(\boldsymbol{\theta})] \quad (6.186)$$

### 6.7.1.3 Conjugate variational inference

In Section 7.3 we showed how to do exact inference in conjugate models. We can derive Equation (7.11) from the BLR by using the fixed point condition in Equation (6.186) to write

$$\boldsymbol{\lambda}_* = \nabla_{\boldsymbol{\mu}} \mathbb{E}_{q_*} [-\bar{\ell}(\boldsymbol{\theta})] = \boldsymbol{\lambda}_0 + \sum_{i=1}^N \underbrace{\nabla_{\boldsymbol{\mu}} \mathbb{E}_{q_*} [\log p(\mathbf{y}_i | \boldsymbol{\theta})]}_{\tilde{\boldsymbol{\lambda}}_i(\mathbf{y}_i)} \quad (6.187)$$

where  $\tilde{\boldsymbol{\lambda}}_i(\mathbf{y}_i)$  are the sufficient statistics for the  $i$ 'th likelihood term.

For models where the joint distribution over the latents factorizes (using a graphical model), we can further decompose this update into a series of local terms. This gives rise to the variational message passing scheme discussed in Section 10.2.7.

### 6.7.1.4 Partially conjugate variational inference

In [Supplementary Section 10.3.1](#), we discuss CVI, which performs variational inference for partially conjugate models, using gradient updates for the non-conjugate parts, and exact Bayesian inference for the conjugate parts.

## 6.7.2 Deriving optimization algorithms from BLR

In this section we show how to derive several different optimization algorithms from BLR. Recall that in BLR, instead of directly minimizing the loss

$$\bar{\ell}(\boldsymbol{\theta}) = \sum_{n=1}^N \ell(\mathbf{y}_n, f_{\boldsymbol{\theta}}(\mathbf{x}_n)) + R(\boldsymbol{\theta}) \quad (6.188)$$

we will instead minimize

$$\mathcal{L}(\boldsymbol{\lambda}) = \mathbb{E}_{q(\boldsymbol{\theta}|\boldsymbol{\lambda})} [\bar{\ell}(\boldsymbol{\theta})] - \mathbb{H}(q(\boldsymbol{\theta}|\boldsymbol{\lambda})) \quad (6.189)$$

Below we show that different approximations to this objective recover a variety of different optimization methods that are used in the literature. (We discuss such algorithms in more detail in [Chapter 6](#).)

### 6.7.2.1 Gradient descent

In this section, we show how to derive gradient descent as a special case of BLR. We use as our approximate posterior  $q(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}|\mathbf{m}, \mathbf{I})$ . In this case the natural and moment parameters are equal,  $\boldsymbol{\mu} = \boldsymbol{\lambda} = \mathbf{m}$ . The base measure satisfies the following (from Equation (2.180)):

$$2 \log h(\boldsymbol{\theta}) = -D \log(2\pi) - \boldsymbol{\theta}^\top \boldsymbol{\theta} \quad (6.190)$$

Hence

$$\nabla_{\boldsymbol{\mu}} \mathbb{E}_q [\log h(\boldsymbol{\theta})] = \nabla_{\boldsymbol{\mu}} (-D \log(2\pi) - \boldsymbol{\mu}^\top \boldsymbol{\mu} - D) = -\boldsymbol{\mu} = -\boldsymbol{\lambda} = -\mathbf{m} \quad (6.191)$$

Thus from Equation (6.183) the BLR update becomes

$$\mathbf{m}_{t+1} = (1 - \eta_t) \mathbf{m}_t + \eta_t \mathbf{m}_t - \eta_t \nabla_{\mathbf{m}} \mathbb{E}_{q_{\mathbf{m}}(\boldsymbol{\theta})} [\bar{\ell}(\boldsymbol{\theta})] \quad (6.192)$$

We can remove the expectation using the first order delta method (Section 6.5.5):

$$\nabla_{\mathbf{m}} \mathbb{E}_{q_{\mathbf{m}}(\boldsymbol{\theta})} [\bar{\ell}(\boldsymbol{\theta})] \approx \nabla_{\boldsymbol{\theta}} \bar{\ell}(\boldsymbol{\theta})|_{\boldsymbol{\theta}=\mathbf{m}} \quad (6.193)$$

Putting these together gives the gradient descent update:

$$\mathbf{m}_{t+1} = \mathbf{m}_t - \eta_t \nabla_{\boldsymbol{\theta}} \bar{\ell}(\boldsymbol{\theta})|_{\boldsymbol{\theta}=\mathbf{m}_t} \quad (6.194)$$

### 6.7.2.2 Newton's method

In this section, we show how to derive Newton's second order optimization method as a special case of BLR, as first shown in [Kha+18].

Suppose we assume  $q(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}|\mathbf{m}, \mathbf{S}^{-1})$ . The natural parameters are

$$\boldsymbol{\lambda}^{(1)} = \mathbf{S} \mathbf{m}, \quad \boldsymbol{\lambda}^{(2)} = -\frac{1}{2} \mathbf{S} \quad (6.195)$$

The mean (moment) parameters are

$$\boldsymbol{\mu}^{(1)} = \mathbf{m}, \quad \boldsymbol{\mu}^{(2)} = \mathbf{S}^{-1} + \mathbf{m} \mathbf{m}^\top \quad (6.196)$$

Since the base measure is constant (see Equation (2.193)), from Equation (6.184) we have

$$\mathbf{S}_{t+1} \mathbf{m}_{t+1} = (1 - \eta_t) \mathbf{S}_t \mathbf{m}_t - \eta_t \nabla_{\boldsymbol{\mu}^{(1)}} \mathbb{E}_{q_{\boldsymbol{\mu}}(\boldsymbol{\theta})} [\bar{\ell}(\boldsymbol{\theta})] \quad (6.197)$$

$$\mathbf{S}_{t+1} = (1 - \eta_t) \mathbf{S}_t + 2\eta_t \nabla_{\boldsymbol{\mu}^{(2)}} \mathbb{E}_{q_{\boldsymbol{\mu}}(\boldsymbol{\theta})} [\bar{\ell}(\boldsymbol{\theta})] \quad (6.198)$$

In Section 6.4.5.1 we show that

$$\nabla_{\boldsymbol{\mu}^{(1)}} \mathbb{E}_{q(\boldsymbol{\theta})} [\bar{\ell}(\boldsymbol{\theta})] = \mathbb{E}_{q(\boldsymbol{\theta})} [\nabla_{\boldsymbol{\theta}} \bar{\ell}(\boldsymbol{\theta})] - \mathbb{E}_{q(\boldsymbol{\theta})} [\nabla_{\boldsymbol{\theta}}^2 \bar{\ell}(\boldsymbol{\theta})] \mathbf{m} \quad (6.199)$$

$$\nabla_{\boldsymbol{\mu}^{(2)}} \mathbb{E}_{q(\boldsymbol{\theta})} [\bar{\ell}(\boldsymbol{\theta})] = \frac{1}{2} \mathbb{E}_{q(\boldsymbol{\theta})} [\nabla_{\boldsymbol{\theta}}^2 \bar{\ell}(\boldsymbol{\theta})] \quad (6.200)$$

Hence the update for the precision matrix becomes

$$\mathbf{S}_{t+1} = (1 - \eta_t) \mathbf{S}_t + \eta_t \mathbb{E}_{q_t} [\nabla_{\boldsymbol{\theta}}^2 \bar{\ell}(\boldsymbol{\theta})] \quad (6.201)$$

For the precision weighted mean, we have

$$\mathbf{S}_{t+1}\mathbf{m}_{t+1} = (1 - \eta_t)\mathbf{S}_t\mathbf{m}_t - \eta_t\mathbb{E}_{q_t} [\nabla_{\boldsymbol{\theta}}\bar{\ell}(\boldsymbol{\theta})] + \eta_t\mathbb{E}_{q_t} [\nabla_{\boldsymbol{\theta}}^2\bar{\ell}(\boldsymbol{\theta})] \mathbf{m}_t \quad (6.202)$$

$$= \mathbf{S}_{t+1}\mathbf{m}_t - \eta_t\mathbb{E}_{q_t} [\nabla_{\boldsymbol{\theta}}\bar{\ell}(\boldsymbol{\theta})] \quad (6.203)$$

Hence

$$\mathbf{m}_{t+1} = \mathbf{m}_t - \eta_t\mathbf{S}_{t+1}^{-1}\mathbb{E}_{q_t} [\nabla_{\boldsymbol{\theta}}\bar{\ell}(\boldsymbol{\theta})] \quad (6.204)$$

We can recover Newton's method in three steps. First set the learning rate to  $\eta_t = 1$ , based on an assumption that the objective is convex. Second, treat the iterate as  $\mathbf{m}_t = \boldsymbol{\theta}_t$ . Third, apply the delta method to get

$$\mathbf{S}_{t+1} = \mathbb{E}_{q_t} [\nabla_{\boldsymbol{\theta}}^2\bar{\ell}(\boldsymbol{\theta})] \approx \nabla_{\boldsymbol{\theta}}^2\bar{\ell}(\boldsymbol{\theta})|_{\boldsymbol{\theta}=\mathbf{m}_t} \quad (6.205)$$

and

$$\mathbb{E}_q [\nabla_{\boldsymbol{\theta}}\bar{\ell}(\boldsymbol{\theta})] \approx \nabla_{\boldsymbol{\theta}}\bar{\ell}(\boldsymbol{\theta})|_{\boldsymbol{\theta}=\mathbf{m}_t} \quad (6.206)$$

This gives Newton's update:

$$\mathbf{m}_{t+1} = \mathbf{m}_t - [\nabla_{\mathbf{m}}^2\bar{\ell}(\mathbf{m}_t)]^{-1}[\nabla_{\mathbf{m}}\bar{\ell}(\mathbf{m}_t)] \quad (6.207)$$

### 6.7.2.3 Variational online Gauss-Newton

In this section, we describe the **Variational Online Gauss-Newton** or **VOGN** method of [Kha+18]. This is an approximate second order optimization method that can be derived from the BLR in several steps, as we show below.

First, we use a diagonal Gaussian approximation to the posterior,  $q_t(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}|\boldsymbol{\theta}_t, \mathbf{S}_t^{-1})$ , where  $\mathbf{S}_t = \text{diag}(\mathbf{s}_t)$  is a vector of precisions. Following Section 6.7.2.2, we get the following updates:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta_t \frac{1}{\mathbf{s}_{t+1}} \odot \mathbb{E}_{q_t} [\nabla_{\boldsymbol{\theta}}\bar{\ell}(\boldsymbol{\theta})] \quad (6.208)$$

$$\mathbf{s}_{t+1} = (1 - \eta_t)\mathbf{s}_t + \eta_t\mathbb{E}_{q_t} [\text{diag}(\nabla_{\boldsymbol{\theta}}^2\bar{\ell}(\boldsymbol{\theta}))] \quad (6.209)$$

where  $\odot$  is elementwise multiplication, and the division by  $\mathbf{s}_{t+1}$  is also elementwise.

Second, we use the delta approximation to replace expectations by plugging in the mean. Third we use a minibatch approximation to the gradient and diagonal Hessian:

$$\hat{\nabla}_{\boldsymbol{\theta}}\bar{\ell}(\boldsymbol{\theta}) = \frac{N}{M} \sum_{i \in \mathcal{M}} \nabla_{\boldsymbol{\theta}}\ell(y_i, f_{\boldsymbol{\theta}}(\mathbf{x}_i)) + \nabla_{\boldsymbol{\theta}}R(\boldsymbol{\theta}) \quad (6.210)$$

$$\hat{\nabla}_{\boldsymbol{\theta}_j}^2\bar{\ell}(\boldsymbol{\theta}) = \frac{N}{M} \sum_{i \in \mathcal{M}} \nabla_{\boldsymbol{\theta}_j}^2\ell(y_i, f_{\boldsymbol{\theta}}(\mathbf{x}_i)) + \nabla_{\boldsymbol{\theta}_j}^2R(\boldsymbol{\theta}) \quad (6.211)$$

where  $M$  is the minibatch size.

For some non-convex problems, such as DNNs, the Hessian may be not be positive definite, so we can get better results using a Gauss-Newton approximation, based on the squared gradients instead of the Hessian:

$$\hat{\nabla}_{\theta_j}^2 \bar{\ell}(\theta) \approx \frac{N}{M} \sum_{i \in \mathcal{M}} [\nabla_{\theta_j} \ell(y_i, f_{\theta}(\mathbf{x}_i))]^2 + \nabla_{\theta_j}^2 R(\theta) \quad (6.212)$$

This is also faster to compute.

Putting all this together gives rise to the **Online Gauss-Newton** or **OGN** method of [Osa+19a].

If we drop the delta approximation, and work with expectations, we get the **Variational Online Gauss-Newton** or **VOGN** method of [Kha+18]. We can approximate the expectations by sampling. In particular, VOGN uses the following **weight perturbation** method

$$\mathbb{E}_{q_t} [\hat{\nabla}_{\theta} \bar{\ell}(\theta)] \approx \hat{\nabla}_{\theta} \bar{\ell}(\theta_t + \epsilon_t) \quad (6.213)$$

where  $\epsilon_t \sim \mathcal{N}(\mathbf{0}, \text{diag}(\mathbf{s}_t))$ . It is also possible to approximate the Fisher information matrix directly; this results in the **Variational Online Generalized Gauss-Newton** or **VOGGN** method of [Osa+19a].

#### 6.7.2.4 Adaptive learning rate SGD

In this section, we show how to derive an update rule which is very similar to the **RMSprop** [Hin14] method, which is widely used in deep learning. The approach we take is similar to that VOGN in Section 6.7.2.3. We use the same diagonal Gaussian approximation,  $q_t(\theta) = \mathcal{N}(\theta | \theta_t, \mathbf{S}_t^{-1})$ , where  $\mathbf{S}_t = \text{diag}(\mathbf{s}_t)$  is a vector of precisions. We then use the delta method to eliminate expectations:

$$\theta_{t+1} = \theta_t - \eta_t \frac{1}{\mathbf{s}_{t+1}} \odot \nabla_{\theta} \bar{\ell}(\theta_t) \quad (6.214)$$

$$\mathbf{s}_{t+1} = (1 - \eta_t) \mathbf{s}_t + \eta_t \text{diag}(\nabla_{\theta}^2 \bar{\ell}(\theta_t)) \quad (6.215)$$

where  $\odot$  is elementwise multiplication. If we allow for different learning rates we get

$$\theta_{t+1} = \theta_t - \alpha_t \frac{1}{\mathbf{s}_{t+1}} \odot \nabla_{\theta} \bar{\ell}(\theta_t) \quad (6.216)$$

$$\mathbf{s}_{t+1} = (1 - \beta_t) \mathbf{s}_t + \beta_t \text{diag}(\hat{\nabla}_{\theta}^2 \bar{\ell}(\theta_t)) \quad (6.217)$$

Now suppose we replace the diagonal Hessian approximation with the sum of the squares per-sample gradients:

$$\text{diag}(\nabla_{\theta}^2 \bar{\ell}(\theta_t)) \approx \hat{\nabla} \bar{\ell}(\theta_t) \odot \hat{\nabla} \bar{\ell}(\theta_t) \quad (6.218)$$

If we also change some scaling factors we can get the RMSprop updates:

$$\theta_{t+1} = \theta_t - \alpha \frac{1}{\sqrt{\mathbf{v}_{t+1}} + c\mathbf{1}} \odot \hat{\nabla}_{\theta} \bar{\ell}(\theta_t) \quad (6.219)$$

$$\mathbf{v}_{t+1} = (1 - \beta) \mathbf{v}_t + \beta [\hat{\nabla} \bar{\ell}(\theta_t) \odot \hat{\nabla} \bar{\ell}(\theta_t)] \quad (6.220)$$

This allows us to use standard deep learning optimizers to get a Gaussian approximation to the posterior for the parameters [Osa+19a].

It is also possible to derive the Adam optimizer [KB15] from BLR by adding a momentum term to RMSprop. See [KR21a; Ait18] for details.

### 6.7.3 Variational optimization

Consider an objective defined in terms of discrete variables. Such objectives are not differentiable and so are hard to optimize. One advantage of BLR is that it optimizes the parameters of a probability distribution, and such expected loss objectives are usually differentiable and smooth. This is called “**variational optimization**” [Bar17], since we are optimizing over a probability distribution.

For example, consider the case of a **binary neural network** where  $\theta_d \in \{0, 1\}$  indicates if weight  $d$  is used or not, we can optimize over the parameters of a Bernoulli distribution,  $q(\boldsymbol{\theta}|\boldsymbol{\lambda}) = \prod_{d=1}^D \text{Ber}(\theta_d|p_d)$ , where  $p_d \in [0, 1]$  and  $\lambda_d = 1/2 \log(p_d/(1 - p_d))$  is the log odds. This is the basis of the **BayesBiNN** approach [MBK20].

If we ignore the entropy and regularizer term, we get the following simplified objective:

$$\mathcal{L}(\boldsymbol{\lambda}) = \int \bar{\ell}(\boldsymbol{\theta}) q(\boldsymbol{\theta}|\boldsymbol{\lambda}) d\boldsymbol{\theta} \quad (6.221)$$

This method has various names: **stochastic relaxation** [SB12; SB13; MMP13], **stochastic approximation** [HHC12; Hu+12], etc. It is closely related to **evolutionary strategies**, which we discuss in Section 6.9.6.

In the case of functions with continuous domains, we can use a Gaussian for  $q(\boldsymbol{\theta}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ . The resulting integral in Equation (6.221) can then sometimes be solved in closed form, as explained in [Mob16]. By starting with a broad variance, and gradually reducing it, we hope the method can avoid poor local optima, similar to simulated annealing (Section 12.9.1). However, we generally get better results by including the entropy term, because then we can automatically learn to adapt the variance. In addition, we can often work with natural gradients, which results in faster convergence.

## 6.8 Bayesian optimization

In this section, we discuss **Bayesian optimization** or **BayesOpt**, which is a model-based approach to black-box optimization, designed for the case where the objective function  $f: \mathcal{X} \rightarrow \mathbb{R}$  is expensive to evaluate (e.g., if it requires running a simulation, or training and testing a particular neural net architecture).

Since the true function  $f$  is expensive to evaluate, we want to make as few function calls (i.e., make as few **queries**  $\mathbf{x}$  to the **oracle**  $f$ ) as possible. This suggests that we should build a **surrogate function** (also called a **response surface model**) based on the data collected so far,  $\mathcal{D}_n = \{(\mathbf{x}_i, y_i) : i = 1 : n\}$ , which we can use to decide which point to query next. There is an inherent tradeoff between picking the point  $\mathbf{x}$  where we think  $f(\mathbf{x})$  is large (we follow the convention in the literature and assume we are trying to maximize  $f$ ), and picking points where we are uncertain about  $f(\mathbf{x})$  but where observing the function value might help us improve the surrogate model. This is another instance of the exploration-exploitation dilemma.

In the special case where the domain we are optimizing over is finite, so  $\mathcal{X} = \{1, \dots, A\}$ , the BayesOpt problem becomes similar to the **best arm identification** problem in the bandit literature (Section 34.4). An important difference is that in bandits, we care about the cost of every action we take, whereas in optimization, we usually only care about the cost of the final solution we find. In other words, in bandits, we want to minimize cumulative regret, whereas in optimization we want to minimize simple or final regret.