

Learning Hyperparameters via MLE-II, and Introduction to Multiparameter Models

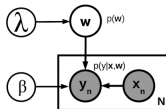
Piyush Rai

Topics in Probabilistic Modeling and Inference (CS698X)

Jan 18, 2018

Learning Hyperparameters

- Let's (re-)consider the Bayesian linear regression model



- Ideally, we would like to **learn** the hyperparams as well (rather than using cross-validation)
- When also learning the hyperparams, the set of unknowns will be $(\mathbf{w}, \beta, \lambda)$
- Suppose we have a prior $p(\mathbf{w}, \beta, \lambda) = p(\mathbf{w} | \lambda) p(\beta) p(\lambda)$
- Ideally, would like to get the **joint posterior** over all the unknowns

$$p(\mathbf{w}, \beta, \lambda | \mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \beta, \lambda) p(\mathbf{w}, \lambda, \beta)}{p(\mathbf{y} | \mathbf{X})} = \frac{p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \beta) p(\mathbf{w} | \lambda) p(\beta) p(\lambda)}{p(\mathbf{y} | \mathbf{X})}$$

where $p(\mathbf{y} | \mathbf{X}) = \int p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \beta) p(\mathbf{w} | \lambda) p(\beta) p(\lambda) d\mathbf{w} d\lambda d\beta$ (ouch! intractable)

- Note: Computing posterior predictive distribution would also require integrating out hyperparams

$$p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \int p(y_* | \mathbf{x}_*, \mathbf{w}, \beta) p(\mathbf{w}, \beta, \lambda | \mathbf{X}, \mathbf{y}) d\mathbf{w} d\beta d\lambda \quad (\text{.. again, intractable})$$

Learning Hyperparameters

- Intractability in hyperparameter estimation can usually be handled via
 - Doing [point estimation for hyperparams](#) and learning full posterior for parameters ([today's lecture](#))
 - This procedure is commonly known as [Type-II MLE](#), or [evidence approximation](#), or [empirical Bayes](#)
 - Note: This won't exactly be 100% Bayesian but nevertheless very widely used
 - Using [general approximate inference methods](#), e.g.,
 - MCMC, Variational Bayesian Inference (we will look at these later)
- In either case, the resulting algorithms work in an alternating fashion
 - Basically, alternate between estimating the parameters and hyperparameters (until convergence)
 - Many algorithms have such a flavor (EM, Gibbs sampling, variational inference, etc.)
- Note: This approach to learning hyperparams is not limited to only learning “hyperparams” (instead of hyperparams, we might have a large number of parameters in a [multiparameter model](#))

Learning Hyperparameters

- For Bayesian linear regression, the posterior over the unknowns can be written as

$$p(\mathbf{w}, \beta, \lambda | \mathbf{X}, \mathbf{y}) = p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta, \lambda) p(\beta, \lambda | \mathbf{X}, \mathbf{y}) \quad (\text{from product rule})$$

- Computing $p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta, \lambda)$ is easy for Bayesian linear regression if β and λ are “known”
- However, computing $p(\beta, \lambda | \mathbf{X}, \mathbf{y})$ is NOT easy. Note that it is given by

$$p(\beta, \lambda | \mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{X}, \beta, \lambda) p(\beta) p(\lambda)}{p(\mathbf{y} | \mathbf{X})}$$

.. where the denominator $p(\mathbf{y} | \mathbf{X}) = \int p(\mathbf{y} | \mathbf{X}, \beta, \lambda) p(\beta) p(\lambda) d\beta d\lambda$ is intractable in general

- One work-around to this intractability is to do point estimation for these hyperparameters

Learning Hyperparameters via Point Estimation

- Let's approximate $p(\beta, \lambda|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{X}, \beta, \lambda)p(\beta)p(\lambda)}{p(\mathbf{y}|\mathbf{X})}$ by a point function δ at its mode

$$p(\beta, \lambda|\mathbf{X}, \mathbf{y}) \approx \delta(\hat{\beta}, \hat{\lambda}) \quad \text{where} \quad \hat{\beta}, \hat{\lambda} = \arg \max_{\beta, \lambda} p(\beta, \lambda|\mathbf{X}, \mathbf{y})$$

- Moreover, if $p(\beta)$, $p(\lambda)$ are uniform/uninformative priors then

$$\hat{\beta}, \hat{\lambda} = \arg \max_{\beta, \lambda} p(\mathbf{y}|\mathbf{X}, \beta, \lambda)$$

- This is point estimation of hyperparams by doing **MLE on the marginal likelihood**

- Also called **Type-II MLE** and **Evidence based method** for hyperparam estimation

- With the above point approximation $p(\beta, \lambda|\mathbf{X}, \mathbf{y}) \approx \delta(\hat{\beta}, \hat{\lambda})$

$$p(\mathbf{w}, \beta, \lambda|\mathbf{X}, \mathbf{y}) = p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda)p(\beta, \lambda|\mathbf{X}, \mathbf{y}) \approx p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \hat{\beta}, \hat{\lambda})$$

- Note: $p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \hat{\beta}, \hat{\lambda})$ is also called the **conditional posterior** of \mathbf{w} given $\lambda = \hat{\lambda}$ and $\beta = \hat{\beta}$ (and for Bayesian linear regression, we know what the form of this conditional posterior is!)

An Iterative Scheme

We can now follow an iterative scheme for estimating \mathbf{w} , λ , and β

- Initialize $\hat{\lambda}$ and $\hat{\beta}$
- Repeat until convergence
 - Estimate the posterior over \mathbf{w} as

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \hat{\lambda}, \hat{\beta}) = \mathcal{N}(\mathbf{w}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

where $\boldsymbol{\Sigma} = (\hat{\beta} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top + \hat{\lambda} \mathbf{I}_D)^{-1}$ and $\boldsymbol{\mu} = \boldsymbol{\Sigma}(\hat{\beta} \sum_{n=1}^N y_n \mathbf{x}_n)$

- Re-estimate $\hat{\lambda}$ as

$$\hat{\lambda} = \frac{\gamma}{\boldsymbol{\mu}^\top \boldsymbol{\mu}} \quad (\text{note: } \gamma \text{ depends on } \boldsymbol{\Sigma})$$

- Re-estimate $\hat{\beta}$ as

$$\hat{\beta} = \frac{N - \gamma}{\sum_{n=1}^N (y_n - \boldsymbol{\mu}^\top \mathbf{x}_n)^2}$$

- Note: This is akin to an Expectation-Maximization (EM) scheme. More on this later.

Using Type-II MLE Approximations

- With the Type-II MLE approximation $p(\beta, \lambda | \mathbf{X}, \mathbf{y}) \approx \delta(\hat{\beta}, \hat{\lambda})$, the posterior over unknowns

$$p(\mathbf{w}, \beta, \lambda | \mathbf{X}, \mathbf{y}) = p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta, \lambda) p(\beta, \lambda | \mathbf{X}, \mathbf{y}) \approx p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \hat{\beta}, \hat{\lambda})$$

- Therefore the overall posterior over the unknowns $\mathbf{w}, \beta, \lambda$ is the same as the conditional posterior of \mathbf{w} treating the hyperparams $\hat{\beta}, \hat{\lambda}$ fixed at their Type-II MLE based estimates
- The posterior predictive distribution can also be approximated as

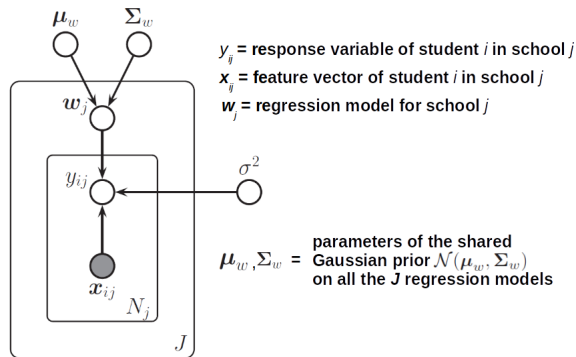
$$\begin{aligned} p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) &= \int p(y_* | \mathbf{x}_*, \mathbf{w}, \beta) p(\mathbf{w}, \beta, \lambda | \mathbf{X}, \mathbf{y}) d\mathbf{w} d\beta d\lambda \\ &= \int p(y_* | \mathbf{x}_*, \mathbf{w}, \beta) p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta, \lambda) p(\beta, \lambda | \mathbf{X}, \mathbf{y}) d\beta d\lambda d\mathbf{w} \\ &\approx \int p(y_* | \mathbf{x}_*, \mathbf{w}, \beta) p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \hat{\beta}, \hat{\lambda}) d\mathbf{w} \end{aligned}$$

- This is also the same as the usual posterior predictive distribution we have seen earlier, except we are treating the hyperparams $\hat{\beta}, \hat{\lambda}$ fixed at their Type-II MLE based estimates

Multiparameter Models

Multiparameter Models

- Multiparameter models arise in many situations, e.g.,
 - Probabilistic models with unknown hyperparameters (e.g., Bayesian linear regression we just saw)
 - **Joint analysis** of data from multiple (and possibly related) groups



- .. and in fact, pretty much in any non-toy example of probabilistic model :)

Multiparameter Models

- A simple two-parameter model: Consider a model with two unknowns θ_1 and θ_2 , e.g.,
 - Gaussian $\mathcal{N}(\mu, \tau^{-1})$ observation model with both unknown mean ($\mu = \theta_1$) and precision ($\tau = \theta_2$)
- Assume a likelihood model $p(\mathbf{y}|\theta_1, \theta_2)$ for observed data \mathbf{y}
- Assume a joint prior distribution $p(\theta_1, \theta_2)$ over both unknowns
- The **joint posterior** distribution for (θ_1, θ_2) would be

$$p(\theta_1, \theta_2|\mathbf{y}) \propto p(\mathbf{y}|\theta_1, \theta_2)p(\theta_1, \theta_2)$$

- If the joint prior $p(\theta_1, \theta_2)$ is conjugate to $p(\mathbf{y}|\theta_1, \theta_2)$, inference would be direct and easy!
 - The joint posterior $p(\theta_1, \theta_2|\mathbf{y})$ will have the same form as the joint prior $p(\theta_1, \theta_2)$
 - Such cases are more of exception than a norm!
 - Coming with with jointly conjugate priors is not always possible/easy

Inference in Multiparameter Models

- There exist several methods to perform inference in multiparameter models
- Some of the inference methods (which we will study) include
 - Specialized inference methods
 - Apply only for some specific types of multiparameter models (will see some examples)
 - Several general-purpose inference methods
 - Markov Chain Monte Carlo (MCMC) methods such as [Gibbs Sampling](#)
 - Variational Bayesian Inference (abbreviated as VB or VI) methods
 - Expectation-Maximization (is actually a special case of VB/VI)
 - Bayesian inference with MLE-II for hyperparameters (we already saw this)
- Note: Some of these methods perform hybrid-style inference
 - Full posteriors for some unknowns, point estimates for others (e.g., EM, MLE-II based methods)
- Most (in fact, all) of these methods perform (usually approximate) inference in an [iterative fashion](#)
 - Basically, [iteratively](#) infer one unknown given all others [fixed](#) at the current value (until convergence)

Marginal Posterior

- Often, instead of (or in addition to) [joint posterior](#), we may be interested in [marginal posterior](#)
- [Marginal posterior](#) is the posterior of one unknown given the data
- In the two-param case, given joint post. $p(\theta_1, \theta_2 | \mathbf{y})$, we can get the marginal posteriors as

$$p(\theta_1 | \mathbf{y}) = \int p(\theta_1, \theta_2 | \mathbf{y}) d\theta_2 \quad \text{and} \quad p(\theta_2 | \mathbf{y}) = \int p(\theta_1, \theta_2 | \mathbf{y}) d\theta_1$$

- In some cases, the marginal posterior can be computed directly using other ways too, e.g.,

$$p(\theta_1 | \mathbf{y}) \propto p(\mathbf{y} | \theta_1) p(\theta_1)$$

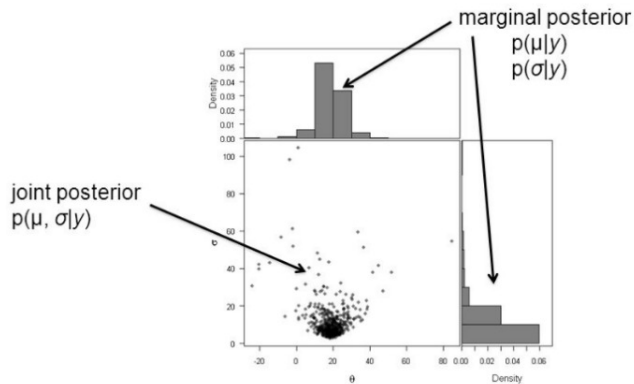
.. but this requires the marginal likelihood $p(\mathbf{y} | \theta_1)$ and the marginal prior $p(\theta_1)$ where

$$\begin{aligned} p(\mathbf{y} | \theta_1) &= \int p(\mathbf{y} | \theta_1, \theta_2) p(\theta_2) d\theta_2 \\ p(\theta_1) &= \int p(\theta_1, \theta_2) p(\theta_2) d\theta_2 \end{aligned}$$

.. getting the marginal likelihood/prior may or may not be easy (depends on distributions involved)

Marginal Posterior: An Illustration

- Often the joint posterior may be hard to visualize/interpret
- Each marginal posterior may provide a better view of the individual parameter's posterior



Conditional Posterior

- **Conditional posterior** is another very useful quantity when doing Bayesian inference
- E.g.: For a joint posterior $p(\theta_1, \theta_2 | \mathbf{y})$, can define conditional posteriors $p(\theta_1 | \theta_2, \mathbf{y})$ or $p(\theta_2 | \theta_1, \mathbf{y})$
- Note that $p(\theta_1 | \theta_2, \mathbf{y})$ assumes that θ_2 is fixed at some value (likewise for $p(\theta_2 | \theta_1, \mathbf{y})$)
- Conditional posteriors are usually easier to derive as compared to the joint $p(\theta_1, \theta_2 | \mathbf{y})$
 - Recall our example of Gaussian mean inference with variance known (and vice-versa)
- Many inference algorithms for joint posterior are based on inferring the conditional posteriors, e.g.,
 - e.g., Gibbs Sampling (an MCMC based sampling algorithm)

Gibbs Sampling

- Gibbs sampler iteratively draws random samples from conditional posteriors
- When run long enough, the sampler produces samples from the joint posterior
- For the simple two-parameter case $\theta = (\theta_1, \theta_2)$, the Gibbs sampler looks like this
 - Initialize $\theta_2^{(0)}$
 - For $s = 1, \dots, S$
 - Draw a random sample for θ_1 as $\theta_1^{(s)} \sim p(\theta_1 | \theta_2^{(s-1)}, \mathbf{y})$
 - Draw a random sample for θ_2 as $\theta_2^{(s)} \sim p(\theta_2 | \theta_1^{(s)}, \mathbf{y})$
- The set of S random samples $\{\theta_1^{(s)}, \theta_2^{(s)}\}_{s=1}^S$ represent the joint posterior distribution $p(\theta_1, \theta_2 | \mathbf{y})$
- More on Gibbs sampling when we discuss MCMC sampling algorithm (above is the high-level idea)
- Many other inference algorithms for inferring the joint posterior (e.g., variational inference) also use conditional posteriors in a similar manner