

06-Linear Regression (Bayesian Learning)

Recap: MLE vs MAP parameter estimation

- **MLE:** Maximizing the likelihood of data. Unbiased if data are iid drawn from a linear model with zero-mean noise.
- **MAP:** Expectation over data is weird because data are known. MAP estimation defines the prior distribution of parameters and maximizes the posterior.
- Unfortunately, MAP estimation is deceptive Bayesian in that
 - The goal of machine learning is **NOT** to estimate parameters, but to **make a prediction**.
 - Therefore, we do not really care any particular parameters.

A Bayesian view

- Anything unknown is a random variable
- Anything unrelated to the ultimate task should be marginalized out

The general framework for Bayesian Learning

- We define the prior $p(\theta)$
- We compute the likelihood $p(D|\theta)$
- The posterior is
$$p(\theta|D) = \frac{p(\theta) p(D|\theta)}{p(D)} \propto_p p(\theta) \cdot p(D|\theta)$$
- Note again that we care more about prediction than parameters.
The **predictive distribution/density** for a new data sample x_* is

$$\begin{aligned} P(t_* | x_*, D) &= \int p(t_*, \theta | x_*, D) d\theta && [\text{marginalization}] \\ &= \int p(t_* | \theta, x_*, D) p(\theta | x_*, D) d\theta && [\text{conditional probability}] \\ &= \int p(t_* | \theta, x_*) p(\theta | D) d\theta \\ &\quad \left[\begin{array}{l} p(t_* | \theta, x_*, D) = p(t_* | \theta, x_*) \text{ because } t_* \perp D | \theta \\ p(\theta | x_*, D) = p(\theta | D) \text{ because } \theta \perp x_* | D \end{array} \right] \end{aligned}$$

Bayesian Linear Regression

- We define the prior

$$p(\mathbf{w}) = N(\mathbf{w}; \mathbf{m}_0, \mathbf{S}_0) \propto_{\mathbf{w}} \exp\left\{-\frac{1}{2}(\mathbf{w} - \mathbf{m}_0)^T \mathbf{S}_0^{-1} (\mathbf{w} - \mathbf{m}_0)\right\}$$

- The likelihood of data, given parameters \mathbf{w}

$$p(\mathcal{D}|\mathbf{w}) \propto_{\mathbf{w}} \exp\left\{-\frac{1}{2}(\mathbf{t} - \mathbf{X}\mathbf{w})^T \Lambda^{-1} (\mathbf{t} - \mathbf{X}\mathbf{w})\right\}$$

where $\Lambda = \begin{pmatrix} \sigma^2 & & \\ & \ddots & \\ & & \sigma^2 \end{pmatrix}$ if $\mathbf{t}^{(m)} = \mathbf{w}^T \mathbf{x}^{(m)} + \varepsilon^{(m)}$ and $\varepsilon^{(m)} \text{iid } N(0, \sigma^2)$

- The posterior of \mathbf{w} is

$$p(\mathbf{w}|\mathcal{D}) = \frac{p(\mathbf{w}) p(\mathcal{D}|\mathbf{w})}{p(\mathcal{D})}$$

$$\propto_{\mathbf{w}} p(\mathbf{w}) p(\mathcal{D}|\mathbf{w})$$

$$= \exp\left\{-\frac{1}{2} [\mathbf{w}^T (\mathbf{S}_0^{-1} + \mathbf{X}^T \Lambda^{-1} \mathbf{X}) \mathbf{w} - 2(\mathbf{m}_0^T \mathbf{S}_0^{-1} + \mathbf{y}^T \Lambda^{-1} \mathbf{X}) \mathbf{w} + \text{const}]\right\}$$

[Note, for symmetric matrix like \mathbf{S}, Λ , $\mathbf{a}^T \mathbf{S} \mathbf{b} = \mathbf{b}^T \mathbf{S} \mathbf{a}$]

$$\propto_{\mathbf{w}} \exp\left(-\frac{1}{2}(\mathbf{w} - \mathbf{m}_N)^T \mathbf{S}_N^{-1} (\mathbf{w} - \mathbf{m}_N)\right)$$

$$= N(\mathbf{w}; \mathbf{m}_N, \mathbf{S}_N)$$

where $\mathbf{m}_N = \mathbf{S}_N (\mathbf{S}_0^{-1} \mathbf{m}_0 + \frac{1}{\sigma^2} \mathbf{X}^T \mathbf{t})$

$$\mathbf{S}_N = (\mathbf{S}_0^{-1} + \mathbf{X}^T \Lambda^{-1} \mathbf{X})^{-1} = (\mathbf{S}_0^{-1} + \frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X})^{-1}$$

- Comments

- Easy-to-remember result:

$$\text{Gaussian} \cdot \text{Gaussian} = \text{Gaussian}$$

Since the prior and the posterior take the same form (Gaussian distribution), it is known as a **conjugate prior**. Choosing a conjugate prior, unfortunately, is purely for computational convenience.

- For an unbounded random variable like Gaussian, we cannot define uniform prior properly. Suppose we have no real prior information for the weights in linear regression, it's intuitive to assume

$$\mathbf{m}_0 = \mathbf{0}, \quad \mathbf{S}_0 = \lambda^2 \mathbf{I}$$

In this case,

$$\mathbf{m}_N = \frac{1}{\sigma^2} \mathbf{S}_N \mathbf{X}^T \mathbf{t}$$

$$\mathbf{S}_N = (\frac{1}{\lambda^2} \mathbf{I} + \frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X})^{-1}$$

$$\propto_{\mathbf{w}} p(\mathbf{w})$$

If, for the prior, $\lambda \rightarrow 0$, i.e., $\cancel{\frac{1}{\lambda} w}$, then $m_N \rightarrow v - m_0$
If, for the prior, $\lambda \rightarrow \infty$, i.e., $\cancel{\frac{1}{\lambda} w}$, then $m_N \rightarrow \hat{w}_{MLE}$

If $\sigma \rightarrow 0$, i.e., data are less noisy, $m_N \rightarrow \hat{w}_{MLE}$
If $\sigma \rightarrow \infty$, i.e., data are very noisy, $m_N \rightarrow 0 = m_0$

No surprise,
as was clear
in MAP learning

- The predictive density is

$$\begin{aligned}
p(t_* | \mathbf{X}_*, \mathcal{D}) &= \int p(t_* | \mathbf{x}_*, \mathbf{w}) p(\mathbf{w} | \mathcal{D}) d\mathbf{w} \\
&= \int N(t_* | \mathbf{w}^T \mathbf{x}_*, \sigma^2) \cdot N(\mathbf{w} | m_N, S_N) d\mathbf{w} \\
&\propto \int \exp \left\{ -\frac{1}{2\sigma^2} (t_* - \mathbf{w}^T \mathbf{x}_*)^T (t_* - \mathbf{w}^T \mathbf{x}_*) \right\} \cdot \exp \left\{ -\frac{1}{2} (\mathbf{w} - m_N)^T S_N^{-1} (\mathbf{w} - m_N) \right\} d\mathbf{w} \\
&= N(t_* | \mu_*, \sigma_*^2)
\end{aligned}$$

The observation of t_* following a normal distribution is due to the quadratic form in both \mathbf{w} and y_* . By completing the square for \mathbf{w} , which integrates to 1, we still get a quadratic for t_*

where $\mu_* = m_N^T \mathbf{x}_* = \mathbf{x}_*^T (\mathbf{X}^T \mathbf{X} + \frac{\lambda^2}{\sigma^2} \mathbf{I})^{-1} \mathbf{X}^T \mathbf{t}$
 $\sigma_*^2 = \sigma^2 + \mathbf{x}_*^T S_N \mathbf{x}_*$

- Bayesian linear regression
 - With zero-mean Gaussian prior, the mean prediction is the same as MAP
 - In addition, Bayesian linear regression also estimates a standard deviation, which reflects the model's degree of uncertainty.
 - This is useful for **active learning**, where the model could suggest data samples (in Bayesian linear regression, those with higher variance) for humans to label.

- An alternative way to compute the integral

- Monte Carlo (MC) sampling**

$$\begin{aligned}
\mathbb{E}_{x \sim p(x)} [f(x)] &= \int p(x) f(x) dx \\
&\approx \frac{1}{K} \sum_{k=1}^K f(x^{(k)}) \quad \text{where } x^{(k)} \stackrel{iid}{\sim} p(x)
\end{aligned}$$

The approximation error converges to zero if $M \rightarrow 0$ due to the Law of Large Numbers (LLN).

Note:

- The approximation holds for subjective probabilities because LLN is derived from probability axioms (which does not depend on how we interpret probabilities), and because intuitively, we can indeed sample infinitely many x 's following our subjective probability, in which case, such sampling becomes a repeatable trial.
- The expectation can be thought of as weighted average, where the weights are probabilities. When we do MC sampling, the weights are implicit in how we obtain $x^{(1)}, \dots, x^{(K)}$ and thus, MC sampling uses a direct average.
- The expectation can also be thought of as an "analytic" average, where we consider the precise probability of x taking a particular value. Thus, expectation always gives an exact number (with "=" sign). MC sampling is an approximation with finite sampled values, and thus, always having the " \approx " sign.
- o Applying MC sampling to Bayesian learning

$$p(t_* | \mathbf{x}_*, \mathcal{D}) = \frac{1}{K} \sum_{k=1}^K p(t_* | \mathbf{x}_*, \mathbf{w}^{(k)})$$

where $\mathbf{w}^{(k)}$ iid $p(\mathbf{w}^{(k)} | \mathcal{D})$

Choosing the prior

- For orthodox Bayesian, the prior is what a human believes. The prior is not subject to optimize.
 - o Bayesian learning never overfits or underfits, because conditional Bayesian analysis, i.e., minimizing the Bayesian expected loss (over the posterior of parameters) is the only fundamentally correct decision principle. However, we may encounter the mis-fitting problem (similar to overfitting or underfitting) due to a bad prior.
- In practice, we may choose the prior by validation, or prior is learnable by
- Maximum likelihood estimation

Suppose $\theta \sim p(\theta | \eta)$
 $t \sim p(t | \theta)$ [conditioning on \mathbf{x} is implicit]

η can be learned by MLE

$$\hat{\eta} = \underset{\eta}{\operatorname{argmax}} p(t | \eta) = \underset{\eta}{\operatorname{argmax}} \int p(t | \theta) p(\theta | \eta) d\theta$$

This is known as Type-II maximum likelihood (Type-II ML, or ML-II), or empirical Bayes.

- **MAP-II**

Likewise, if we have a prior for $\boldsymbol{\eta}$, we could perform type-II max a posterior

$$\hat{\boldsymbol{\eta}} = \arg \max_{\boldsymbol{\eta}} p(\boldsymbol{\eta} | \mathbf{t}) = \arg \max_{\boldsymbol{\eta}} \int p(\mathbf{t} | \boldsymbol{\theta}) p(\boldsymbol{\theta} | \boldsymbol{\eta}) p(\boldsymbol{\eta}) d\boldsymbol{\theta}$$

- **Full-Bayes** (marginalization over $\boldsymbol{\eta}$ and $\boldsymbol{\theta}$)

$$p(t_* | \mathbf{x}_*, \mathcal{D}) \propto \iint p(\boldsymbol{\eta}) \cdot p(\boldsymbol{\theta} | \boldsymbol{\eta}) \cdot p(\mathcal{D} | \boldsymbol{\theta}) \cdot p(t_* | \mathbf{x}_*, \boldsymbol{\theta}) d\boldsymbol{\theta} d\boldsymbol{\eta}$$

- Final notes:

- Theoretically, we may define a prior for the parameters of $p(\boldsymbol{\eta})$, and this forms a deeper hierarchy of Bayes. In practice, we are happy with two-level prior.
- Usually, we set the prior for $\boldsymbol{\eta}$ as the simplest or most intuitive one, e.g., standard norm or non-informative. Sometimes, a hierarchical prior may be more robust than one-level prior.
- Technically, a hierarchical prior is merely a convenient vehicle for humans to express prior $p(\boldsymbol{\theta})$. It can be shown (how?) that a hierarchical prior can be collapsed to a usually prior, just like the usually parameter $\boldsymbol{\theta}$ collapses in ML-II and MAP-II.