# Bayesian inference and prediction in finite regression models

Carl Edward Rasmussen

October 10th, 2023

# Key concepts

Bayesian inference in finite, parametric models

- we contrast maximum likelihood with Bayesian inference
- when both prior and likelihood are Gaussian, all calculations are tractable
  - the posterior on the parameters is Gaussian
  - the predictive distribution is Gaussian
  - the marginal likelihood is tractable
- we observe the contrast
  - in maximum likelihood the data fit gets better with larger models (overfitting)
  - the marginal likelihood prefers an intermediate model size (Occam's Razor)

# Maximum likelihood, parametric model

Supervised parametric learning:

- data: $x, y$
- model $\mathcal{M}$: $y = f_w(x) + \varepsilon$

Gaussian likelihood:

$$p(y|x, w, \mathcal{M}) \propto \prod_{n=1}^{N} \exp(-\tfrac{1}{2}(y_n - f_w(x_n))^2/\sigma_{\text{noise}}^2).$$

Maximize the likelihood:

$$w_{\text{ML}} = \underset{w}{\operatorname{argmax}} \, p(y|x, w, \mathcal{M}).$$

Make predictions, by plugging in the ML estimate:

$$p(y_*|x_*, w_{\text{ML}}, \mathcal{M})$$

# Bayesian inference, parametric model

Posterior parameter distribution by Bayes rule ($p(a|b)p(b) = p(a)p(b|a)$):

$$p(w|x, y, \mathcal{M})p(y|x, \mathcal{M}) = p(w|\mathcal{M})p(y|x, w, \mathcal{M})$$

Making predictions (marginalizing out the parameters):

$$
\begin{aligned}
p(y_*|x_*, x, y, \mathcal{M}) &= \int p(y_*, w|x, y, x_*, \mathcal{M})dw \\
&= \int p(y_*|w, x_*, \mathcal{M})p(w|x, y, \mathcal{M})dw.
\end{aligned}
$$

Marginal likelihood:

$$p(y|x, \mathcal{M}) = \int p(w|\mathcal{M})p(y|x, w, \mathcal{M})dw.$$

# Posterior and predictive distribution in detail

For a linear-in-the-parameters model with Gaussian priors and Gaussian noise:

- Gaussian *prior* on the weights: $p(\boldsymbol{w}|\mathcal{M}) = \mathcal{N}(\boldsymbol{w}; \, 0, \, \sigma_w^2 \, \mathbf{I})$
- Gaussian *likelihood* of the weights: $p(\boldsymbol{y}|\boldsymbol{x}, \boldsymbol{w}, \mathcal{M}) = \mathcal{N}(\boldsymbol{y}; \, \boldsymbol{\Phi} \, \boldsymbol{w}, \, \sigma_{\text{noise}}^2 \, \mathbf{I})$

Posterior parameter distribution by Bayes rule $p(\mathfrak{a}|\mathfrak{b}) = p(\mathfrak{a})p(\mathfrak{b}|\mathfrak{a})/p(\mathfrak{b})$:

$$p(\boldsymbol{w}|\boldsymbol{x}, \boldsymbol{y}, \mathcal{M}) \; = \; \frac{p(\boldsymbol{w}|\mathcal{M})p(\boldsymbol{y}|\boldsymbol{x}, \boldsymbol{w}, \mathcal{M})}{p(\boldsymbol{y}|\boldsymbol{x}, \mathcal{M})} \; = \; \mathcal{N}(\boldsymbol{w}; \, \boldsymbol{\mu}, \, \boldsymbol{\Sigma})$$

$$\boldsymbol{\Sigma} \; = \; \left(\sigma_{\text{noise}}^{-2} \, \boldsymbol{\Phi}^\top \boldsymbol{\Phi} + \sigma_w^{-2} \, \mathbf{I}\right)^{-1} \quad \text{and} \quad \boldsymbol{\mu} \; = \; \left(\boldsymbol{\Phi}^\top \boldsymbol{\Phi} + \frac{\sigma_{\text{noise}}^2}{\sigma_w^2} \, \mathbf{I}\right)^{-1} \boldsymbol{\Phi}^\top \boldsymbol{y}$$

The predictive distribution is given by:

$$p(y_*|x_*, \boldsymbol{x}, \boldsymbol{y}, \mathcal{M}) \; = \; \int p(y_*|\boldsymbol{w}, x_*, \mathcal{M})p(\boldsymbol{w}|\boldsymbol{x}, \boldsymbol{y}, \mathcal{M})d\boldsymbol{w}$$

$$= \; \mathcal{N}(y_*; \, \boldsymbol{\phi}(x_*)^\top \boldsymbol{\mu}, \, \boldsymbol{\phi}(x_*)^\top \boldsymbol{\Sigma} \boldsymbol{\phi}(x_*) + \sigma_{\text{noise}}^2).$$

# Multiple explanations of the data



Remember that a finite linear model $f(x_n) = \boldsymbol{\phi}(x_n)^\top \boldsymbol{w}$ with prior on the weights $p(\boldsymbol{w}) = \mathcal{N}(\boldsymbol{w}; \, 0, \sigma_w^2 \mathbf{I})$ has a posterior distribution

$$p(\boldsymbol{w}|\mathbf{x}, \mathbf{y}, \mathcal{M}) \, = \, \mathcal{N}(\boldsymbol{w}; \, \boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad \text{with} \quad \begin{aligned} \boldsymbol{\Sigma} &= \left(\sigma_{\text{noise}}^{-2} \boldsymbol{\Phi}^\top \boldsymbol{\Phi} + \sigma_w^{-2}\right)^{-1} \\ \boldsymbol{\mu} &= \left(\boldsymbol{\Phi}^\top \boldsymbol{\Phi} + \frac{\sigma_{\text{noise}}^2}{\sigma_w^2} \mathbf{I}\right)^{-1} \boldsymbol{\Phi}^\top \mathbf{y} \end{aligned}$$
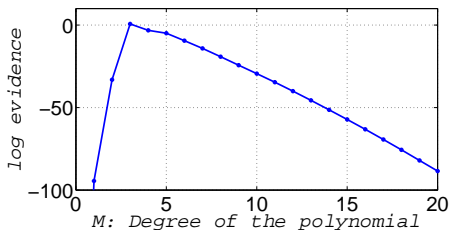
and predictive distribution

$$p(y_*|x_*, \mathbf{x}, \mathbf{y}, \mathcal{M}) \, = \, \mathcal{N}(y_*; \, \boldsymbol{\phi}(x_*)^\top \boldsymbol{\mu}, \, \boldsymbol{\phi}(x_*)^\top \boldsymbol{\Sigma} \boldsymbol{\phi}(x_*) + \sigma_{\text{noise}}^2 \mathbf{I})$$

# Marginal likelihood (Evidence) of our polynomials

Marginal likelihood, or "evidence" of a finite linear model:

$$p(\mathbf{y}|\mathbf{x}, \mathcal{M}) = \int p(\mathbf{w}|\mathcal{M}) p(\mathbf{y}|\mathbf{x}, \mathbf{w}, \mathcal{M}) d\mathbf{w}$$

$$= \mathcal{N}(\mathbf{y};\, 0,\, \sigma_{\mathbf{w}}^2 \, \boldsymbol{\Phi} \, \boldsymbol{\Phi}^\top + \sigma_{\text{noise}}^2 \, \mathbf{I}).$$

Luckily for Gaussian noise there is a closed-form analytical solution!



- The evidence prefers $M = 3$, not simpler, not more complex.
- Too simple models consistently miss most data.
- Too complex models frequently miss some data.