

Bayesian Learning

Lecture 6 - Bayesian regularization

Mattias Villani 🧑

Department of Statistics
Stockholm University



Lecture overview

- Non-linear regression
- Regularization priors

Polynomial regression

Polynomial regression

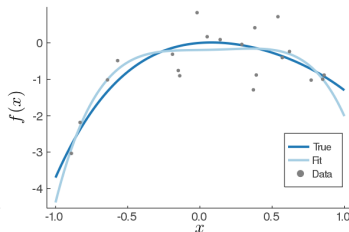
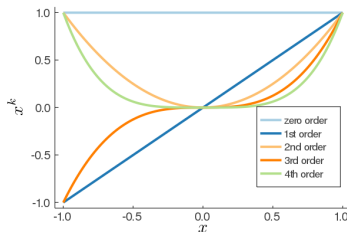
$$f(x_i) = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \dots + \beta_k x_i^k, \quad \text{for } i = 1, \dots, n.$$

$$\mathbf{y} = \mathbf{X}\beta + \varepsilon,$$

where i th row of \mathbf{X} is

$$(1, x_i, x_i^2, \dots, x_i^k).$$

- Still **linear in β** and $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$. Bayes unchanged.



Spline regression

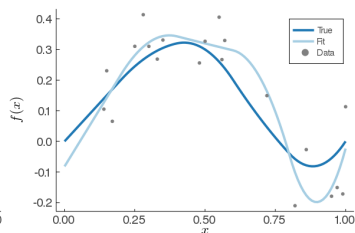
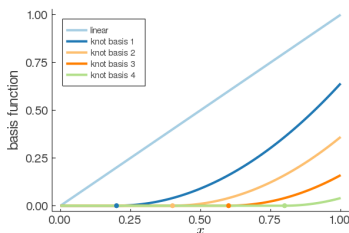
- Polynomials are too global. Need more local basis functions.
- Truncated quadratic splines** with **knot locations** $\kappa_1, \dots, \kappa_m$:

$$b_j(x) = \begin{cases} (x - \kappa_j)^2 & \text{if } x > \kappa_j \\ 0 & \text{otherwise} \end{cases}$$

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \varepsilon,$$

where i th row of \mathbf{X} is

$$(1, x_i, b_1(x_i), \dots, b_m(x_i)).$$



Regularization prior - Ridge

- Too many knots leads to **over-fitting**.
- **Smoothness/shrinkage/regularization prior**

$$\beta_i | \sigma^2 \stackrel{\text{iid}}{\sim} N\left(0, \frac{\sigma^2}{\lambda}\right)$$

- Larger λ gives smoother fit. Note: $\Omega_0 = \lambda I$ in conjugate prior.
- Equivalent to **penalized likelihood**:

$$-2 \cdot \log p(\beta | \sigma^2, \mathbf{y}, \mathbf{X}) \propto (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) + \lambda \beta^T \beta$$

- Posterior mean gives **ridge regression** estimator

$$\tilde{\beta} = (\mathbf{X}^T \mathbf{X} + \lambda I)^{-1} \mathbf{X}^T \mathbf{y}$$

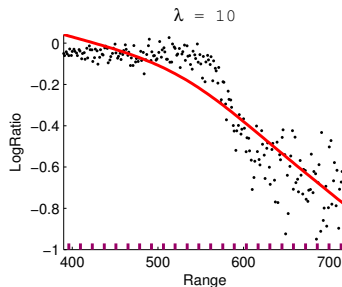
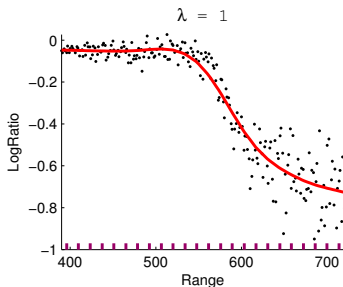
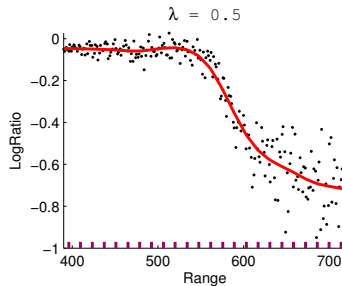
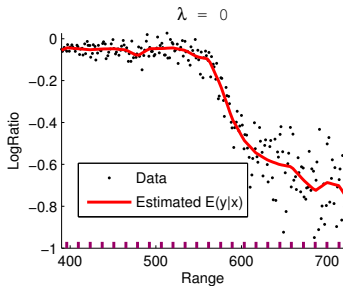
- **Shrinkage** toward zero

$$\text{As } \lambda \rightarrow \infty, \tilde{\beta} \rightarrow 0$$

- When $\mathbf{X}^T \mathbf{X} = I$

$$\tilde{\beta} = \frac{1}{1 + \lambda} \hat{\beta}$$

Bayesian spline with regularization prior



Regularization prior - Lasso

- **Lasso** is equivalent to posterior mode under Laplace prior

$$\beta_i | \sigma^2 \stackrel{\text{iid}}{\sim} \text{Laplace} \left(0, \frac{\sigma^2}{\lambda} \right)$$

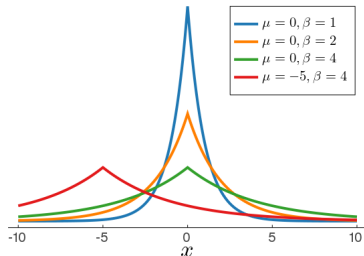
Laplace distribution

$X \sim \text{Laplace}(\mu, \beta)$ for $X \in \mathbb{R}$.

$$p(x) = \frac{1}{2\beta} \exp \left(-\frac{|x - \mu|}{\beta} \right)$$

$$\mathbb{E}(X) = \mu$$

$$\mathbb{V}(X) = 2\beta^2$$



- The **Bayesian shrinkage** prior is **interpretable**. **Not ad hoc**.
- Laplace distribution have heavy tails.
- **Laplace prior**: many β_i close to zero, but some β_i very large.
- Normal distribution have light tails.

Learning the shrinkage

- **Cross-validation** used to determine degree of smoothness, λ .
- Bayesian: λ is **unknown** \Rightarrow **use a prior** for λ !
- $\lambda \sim \text{Inv-}\chi^2(\eta_0, \lambda_0)$.
- **Hierarchical** setup:

$$\begin{aligned} \mathbf{y} | \boldsymbol{\beta}, \sigma^2, \mathbf{X} &\sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2 I_n) \\ \boldsymbol{\beta} | \sigma^2, \lambda &\sim N(0, \sigma^2 \lambda^{-1} I_m) \\ \sigma^2 &\sim \text{Inv} - \chi^2(\nu_0, \sigma_0^2) \\ \lambda &\sim \text{Inv} - \chi^2(\eta_0, \lambda_0) \end{aligned}$$

$$\text{so } \boldsymbol{\Omega}_0 = \lambda I_m.$$

Regression with learned shrinkage

- The **joint posterior** of β , σ^2 and λ is

$$\beta | \sigma^2, \lambda, \mathbf{y} \sim N(\mu_n, \Omega_n^{-1})$$

$$\sigma^2 | \lambda, \mathbf{y} \sim \text{Inv} - \chi^2(\nu_n, \sigma_n^2)$$

$$p(\lambda | \mathbf{y}) \propto \sqrt{\frac{|\Omega_0|}{|\mathbf{X}^T \mathbf{X} + \Omega_0|}} \left(\frac{\nu_n \sigma_n^2}{2} \right)^{-\nu_n/2} \cdot p(\lambda)$$

where $\Omega_0 = \lambda I_m$, and $p(\lambda)$ is the prior for λ , and

$$\mu_n = (\mathbf{X}^T \mathbf{X} + \Omega_0)^{-1} \mathbf{X}^T \mathbf{y}$$

$$\Omega_n = \mathbf{X}^T \mathbf{X} + \Omega_0$$

$$\nu_n = \nu_0 + n$$

$$\nu_n \sigma_n^2 = \nu_0 \sigma_0^2 + \mathbf{y}^T \mathbf{y} - \mu_n^T \Omega_n \mu_n$$

More complexity

- The **location of the knots** can be unknown. Joint posterior:

$$p(\beta, \sigma^2, \lambda, \kappa_1, \dots, \kappa_m | \mathbf{y}, \mathbf{X})$$

- The marginal posterior for $\kappa_1, \dots, \kappa_m$ is a nightmare.
- Simulate from joint posterior by MCMC. Li and Villani (2013).
- The basic spline model can be extended with:
 - ▶ **Heteroscedastic errors** (also modelled with a spline)
 - ▶ **Non-normal errors** (student-t or mixture distributions)
 - ▶ **Autocorrelated/dependent errors** (AR process for the errors)