

Parametric Empirical Bayes with ARD for Linear Regression

Dr Alexander D. Shaw
 $a.d.shaw@exeter.ac.uk$
<https://cpnslab.com>

1 Parametric Empirical Bayes with ARD for Linear Regression

Model (single output). Given N observations and p predictors, let $X \in \mathbb{R}^{N \times p}$ and $y \in \mathbb{R}^N$. We assume a Gaussian likelihood

$$p(y | \beta, \sigma^2) = \mathcal{N}(y; X\beta, \sigma^2 I_N), \quad (1)$$

and an ARD Gaussian prior over coefficients

$$p(\beta | \Lambda) = \mathcal{N}(\beta; 0, \Lambda^{-1}), \quad \Lambda = \text{diag}(\lambda_1, \dots, \lambda_p), \quad \lambda_j > 0. \quad (2)$$

Posterior over coefficients. The posterior is Gaussian

$$\Sigma = \left(\Lambda + \frac{1}{\sigma^2} X^\top X \right)^{-1}, \quad (3)$$

$$\mu = \frac{1}{\sigma^2} \Sigma X^\top y. \quad (4)$$

Type-II Maximum Likelihood / Evidence maximisation. Define the marginal likelihood (evidence)

$$\log p(y | \Lambda, \sigma^2) = -\frac{1}{2} \left(N \log 2\pi + \log |C| + y^\top C^{-1} y \right), \quad C = \sigma^2 I_N + X \Lambda^{-1} X^\top. \quad (5)$$

Equivalently, combining posterior terms gives the standard free-energy form:

$$\mathcal{F}(\Lambda, \sigma^2) = \frac{1}{2} \left(\log |\Lambda| + N \log \frac{1}{\sigma^2} - \log |A| \right) - \frac{1}{2} \left(\frac{1}{\sigma^2} \|y - X\mu\|^2 + \mu^\top \Lambda \mu + N \log 2\pi \right), \quad A = \Lambda + \frac{1}{\sigma^2} X^\top X, \quad (6)$$

which is equal to $\log p(y | \Lambda, \sigma^2)$ up to constants. The standard EM/Type-II updates use the *effective degrees of freedom*

$$\gamma_j = 1 - \lambda_j \Sigma_{jj}, \quad j = 1, \dots, p. \quad (7)$$

Then

$$\lambda_j^{\text{new}} = \frac{\gamma_j}{\mu_j^2 + \varepsilon}, \quad \varepsilon > 0 \text{ small (stability)}, \quad (8)$$

$$\sigma^2 \text{new} = \frac{\|y - X\mu\|^2}{N - \sum_{j=1}^p \gamma_j}. \quad (9)$$

Equations (3)–(9) are iterated until convergence.

Predictive distribution. For a new design row $x_\star \in \mathbb{R}^p$,

$$p(y_\star | x_\star, \mathcal{D}) = \mathcal{N}(y_\star; x_\star^\top \mu, \sigma^2 + x_\star^\top \Sigma x_\star). \quad (10)$$

Multi-output extension (shared sparsity). For d outputs, $Y \in \mathbb{R}^{N \times d}$, collect coefficients in $B \in \mathbb{R}^{p \times d}$ with columns $b_{(k)}$. Use independent Gaussian likelihoods with per-output noise σ_k^2 and a shared ARD prior

$$p(B | \Lambda) = \prod_{k=1}^d \mathcal{N}(b_{(k)}; 0, \Lambda^{-1}), \quad \Lambda = \text{diag}(\lambda_1, \dots, \lambda_p). \quad (11)$$

Posteriors (for each k):

$$\Sigma = \left(\Lambda + \frac{1}{\sigma_k^2} X^\top X \right)^{-1}, \quad \mu_{(k)} = \frac{1}{\sigma_k^2} \Sigma X^\top y_{(k)}. \quad (12)$$

Share the ARD precisions via the common Σ (or its diagonal):

$$\gamma_j = 1 - \lambda_j \Sigma_{jj}, \quad (13)$$

$$\lambda_j^{\text{new}} = \frac{\gamma_j}{\frac{1}{d} \sum_{k=1}^d \mu_{j(k)}^2 + \varepsilon}, \quad (14)$$

$$\sigma_k^2 \text{new} = \frac{\|y_{(k)} - X\mu_{(k)}\|^2}{N - \sum_{j=1}^p \gamma_j}, \quad k = 1, \dots, d. \quad (15)$$

Predictive covariance for output k :

$$\text{Var}[y_{\star(k)} | x_\star] = \sigma_k^2 + x_\star^\top \Sigma x_\star. \quad (16)$$

Intercept and standardisation (as in code). If an intercept column is included in X , one may leave it unpenalised by fixing a very small precision (or excluding it from ARD). For numerical stability, we standardise:

$$\tilde{X}_{ij} = \frac{X_{ij} - \bar{X}_j}{s_j}, \quad \tilde{y} = \frac{y - \bar{y}}{s_y},$$

fit in (\tilde{X}, \tilde{y}) to obtain $(\tilde{\mu}, \tilde{\Sigma})$, then map back:

$$\mu_j = \frac{s_y}{s_j} \tilde{\mu}_j, \quad \Sigma_{jj'} = \frac{s_y^2}{s_j s_{j'}} \tilde{\Sigma}_{jj'}.$$

Algorithm 1 PEB–ARD (shared sparsity, multi-output)

- 1: **Input:** X, Y ; initialise $\lambda_j \leftarrow 1$, $\sigma_k^2 \leftarrow \text{Var}(y_{(k)})$.
 - 2: **repeat**
 - 3: Form $A = \Lambda + \frac{1}{\sigma^2} X^\top X$ with $\sigma^2 = \frac{1}{d} \sum_k \sigma_k^2$.
 - 4: Compute $\Sigma = A^{-1}$, and $\mu_{(k)} = \frac{1}{\sigma_k^2} \Sigma X^\top y_{(k)}$ for all k .
 - 5: $\gamma_j \leftarrow 1 - \lambda_j \Sigma_{jj}$.
 - 6: $\lambda_j \leftarrow \frac{\gamma_j}{\frac{1}{d} \sum_k \mu_{j(k)}^2 + \varepsilon}$ (clip to $[\lambda_{\min}, \lambda_{\max}]$).
 - 7: $\sigma_k^2 \leftarrow \frac{\|y_{(k)} - X\mu_{(k)}\|^2}{N - \sum_j \gamma_j}$ (floor at σ_{\min}^2).
 - 8: **until** converged
-

Algorithm (Type-II EM). In practice we compute Σ and solves with Cholesky plus jitter; if needed we apply an SVD-based SPD repair. We also prune rank-deficient columns of X via QR before fitting (always retaining a constant column if present).

Numerical implementation. To ensure numerical stability when $X^\top X$ or A is nearly singular, the implementation performs Cholesky factorisation with progressively increased *jitter* ϵI until positive definiteness is achieved. If Cholesky fails, the symmetric matrix $(A + A^\top)/2$ is repaired via singular value decomposition (SVD) and eigenvalue flooring. Prior to fitting, the design matrix is orthogonalised with a rank-revealing QR decomposition to remove collinear or zero-variance columns, while always retaining any constant (intercept) column. All matrix inversions are avoided in favour of triangular solves.

Connection to empirical Bayes and ridge regression. This procedure maximises the marginal likelihood $p(Y | X, \Lambda, \sigma^2)$, treating (Λ, σ^2) as hyperparameters estimated at the group or “second” level—hence the term *parametric empirical Bayes (PEB)*. When all λ_j are constrained equal,

$$\lambda_j = \lambda \quad \forall j,$$

the update reduces to standard Bayesian ridge regression with precision λ and noise variance σ^2 estimated by evidence optimisation. Automatic relevance determination (ARD) generalises this by learning one precision per feature, leading to sparsity as irrelevant predictors acquire large λ_j .

Interpretation in variational-free-energy form. The negative free energy (evidence lower bound) for this Gaussian model is

$$\mathcal{F} = \mathbb{E}_q[\log p(y, \beta)] - \mathbb{E}_q[\log q(\beta)],$$

where $q(\beta) = \mathcal{N}(\mu, \Sigma)$ is the variational posterior. Under the conjugate Gaussian model, \mathcal{F} is maximised exactly by the updates above, and therefore coincides with the log-evidence. In this view, the algorithm performs coordinate ascent on \mathcal{F} with respect to $\{\mu, \Sigma, \Lambda, \sigma^2\}$.

Computational complexity. For p predictors and N samples, each iteration requires forming $X^\top X$ and $X^\top Y$ once ($\mathcal{O}(Np^2)$) and solving a $p \times p$ linear system ($\mathcal{O}(p^3)$) per update. In typical applications $p \ll N$, so runtime is dominated by the matrix multiplications. The QR pre-pruning reduces p to the numerical rank of X , ensuring stable and efficient computation.

Summary. The resulting posterior over coefficients

$$\beta | Y, X, \hat{\Lambda}, \hat{\sigma}^2 \sim \mathcal{N}(\mu, \Sigma)$$

provides both point estimates and credible intervals (from Σ) for inference, as well as predictive uncertainty

$$y_\star | x_\star, \mathcal{D} \sim \mathcal{N}(x_\star^\top \mu, \hat{\sigma}^2 + x_\star^\top \Sigma x_\star).$$

This compact form directly underlies the accompanying MATLAB routine `peb_ard_novar.m`, which implements the updates in Eqs. (3)–(9) using Cholesky/SVD solves, automatic relevance determination, and optional shared sparsity across multiple outputs.