

# Machine Learning 2

## C. Sparse Models / 3. Automatic Relevance Determination

Lars Schmidt-Thieme

Information Systems and Machine Learning Lab (ISMLL)  
Institute for Computer Science  
University of Hildesheim, Germany

# Syllabus

## A. Advanced Supervised Learning

- Fri. 12.4. (1) A.1 Generalized Linear Models
- Fri. 26.4. (2) A.2 Gaussian Processes
- Fri. 3.5. (3) A.2b Gaussian Processes (ctd.)
- Fri. 10.5. (4) A.3 Advanced Support Vector Machines

## B. Ensembles

- Fri. 17.5. (5) B.1 Stacking
- Fri. 24.5. (6) B.2 Boosting
- Fri. 31.5. (7) B.3 Mixtures of Experts

## C. Sparse Models

- Fri. 7.6. (8) C.1 Homotopy and Least Angle Regression
- Fri. 14.6. — — Pentecoste Break —
- Fri. 21.6. (9) C.2 Proximal Gradients
- Fri. 28.6. (10) C.3 Laplace Priors
- Fri. 29.6. (11) C.4 Automatic Relevance Determination

## D. Complex Predictors

- Fri. 6.7. (12) D.1 Latent Dirichlet Allocation (LDA)
- Fri. 12.7. (13) Q & A

# Outline

## 1. Automatic Relevance Determination (ARD)

# Outline

## 1. Automatic Relevance Determination (ARD)

# Linear Regression plus ARD Regularization

Linear Regression plus L2 regularization (Ridge Regression):

$$p(y_n | x_n, \beta, \sigma_y^2) := \mathcal{N}(y_n | \beta^T x_n, \sigma_y^2)$$

$$p(\beta) := \mathcal{N}(\beta | 0, \Sigma_\beta := \sigma_\beta^2 I)$$

Linear Regression plus ARD Regularization:

$$p(y_n | x_n, \beta, \sigma_y^2) := \mathcal{N}(y_n | \beta^T x_n, \sigma_y^2)$$

$$p(\beta) := \mathcal{N}(\beta | 0, \Sigma_\beta := \text{diag}(\sigma_{\beta_1}^2, \dots, \sigma_{\beta_M}^2))$$

# Linear Regression plus ARD Regularization

Linear Regression plus L2 regularization (Ridge Regression):

$$p(y_n | x_n, \beta, \sigma_y^2) := \mathcal{N}(y_n | \beta^T x_n, \sigma_y^2)$$

$$p(\beta) := \mathcal{N}(\beta | 0, \Sigma_\beta := \sigma_\beta^2 I)$$

Linear Regression plus ARD Regularization:

$$p(y_n | x_n, \beta, \sigma_y^2) := \mathcal{N}(y_n | \beta^T x_n, \sigma_y^2)$$

$$p(\beta) := \mathcal{N}(\beta | 0, \Sigma_\beta := \text{diag}(\sigma_{\beta_1}^2, \dots, \sigma_{\beta_M}^2))$$

Idea:

- ▶ use a **different regularization weight** for each predictor  $x_m$ .

# Linear Regression plus ARD Regularization

Linear Regression plus L2 regularization (Ridge Regression):

$$p(y_n | x_n, \beta, \sigma_y^2) := \mathcal{N}(y_n | \beta^T x_n, \sigma_y^2)$$

$$p(\beta) := \mathcal{N}(\beta | 0, \Sigma_\beta := \sigma_\beta^2 I)$$

Linear Regression plus ARD Regularization:

$$p(y_n | x_n, \beta, \sigma_y^2) := \mathcal{N}(y_n | \beta^T x_n, \sigma_y^2)$$

$$p(\beta) := \mathcal{N}(\beta | 0, \Sigma_\beta := \text{diag}(\sigma_{\beta_1}^2, \dots, \sigma_{\beta_M}^2))$$

Idea:

- ▶ use a **different regularization weight** for each predictor  $x_m$ .
- ▶ but  $M$  hyperparameters are too many to learn by grid search.

# Linear Regression plus ARD Regularization

Linear Regression plus L2 regularization (Ridge Regression):

$$p(y_n | x_n, \beta, \sigma_y^2) := \mathcal{N}(y_n | \beta^T x_n, \sigma_y^2)$$

$$p(\beta) := \mathcal{N}(\beta | 0, \Sigma_\beta := \sigma_\beta^2 I)$$

Linear Regression plus ARD Regularization:

$$p(y_n | x_n, \beta, \sigma_y^2) := \mathcal{N}(y_n | \beta^T x_n, \sigma_y^2)$$

$$p(\beta) := \mathcal{N}(\beta | 0, \Sigma_\beta := \text{diag}(\sigma_{\beta_1}^2, \dots, \sigma_{\beta_M}^2))$$

$$p(\sigma_y^2) := \text{InvGamma}(\sigma_y^2 | c, d)$$

$$p(\sigma_{\beta_m}^2) := \text{InvGamma}(\sigma_{\beta_m}^2 | a, b), \quad m = 1, \dots, M$$

Idea:

- ▶ use a **different regularization weight** for each predictor  $x_m$ .
- ▶ but  $M$  hyperparameters are too many to learn by grid search.
- ▶ hence put a **hyperprior** on top.



# Empirical Bayes

Maximum Likelihood (ML):

$$\hat{\theta} := \arg \max_{\theta} p(\mathcal{D} | \theta)$$

Full Bayes:

$$(\hat{\theta}, \hat{\eta}) \sim p(\theta, \eta | \mathcal{D}) \propto p(\mathcal{D} | \theta) p(\theta | \eta) p(\eta)$$

# Empirical Bayes

Maximum Likelihood (ML):  $\hat{\theta} := \arg \max_{\theta} p(\mathcal{D} | \theta)$

Maximum A Posteriori (MAP):  $\hat{\theta} := \arg \max_{\theta} p(\mathcal{D} | \theta) p(\theta | \eta)$

Full Bayes:

$$(\hat{\theta}, \hat{\eta}) \sim p(\theta, \eta | \mathcal{D}) \propto p(\mathcal{D} | \theta) p(\theta | \eta) p(\eta)$$

# Empirical Bayes

Maximum Likelihood (ML):  $\hat{\theta} := \arg \max_{\theta} p(\mathcal{D} | \theta)$

Maximum A posteriori (MAP):  $\hat{\theta} := \arg \max_{\theta} p(\mathcal{D} | \theta) p(\theta | \eta)$

ML-II (Empirical Bayes):  $\hat{\eta} := \arg \max_{\eta} p(\mathcal{D} | \eta)$

$$= \arg \max_{\eta} \int p(\mathcal{D} | \theta) p(\theta | \eta) d\theta$$

MAP-II:  $\hat{\theta} \sim p(\theta | \mathcal{D}, \hat{\eta}) \propto p(\mathcal{D} | \theta) p(\theta | \hat{\eta})$

$\hat{\eta} := \arg \max_{\eta} p(\mathcal{D} | \eta) p(\eta)$

$$= \arg \max_{\eta} \int p(\mathcal{D} | \theta) p(\theta | \eta) p(\eta) d\theta$$

$\hat{\theta} \sim p(\theta | \mathcal{D}, \hat{\eta}) \propto p(\mathcal{D} | \theta) p(\theta | \hat{\eta})$

Full Bayes:  $(\hat{\theta}, \hat{\eta}) \sim p(\theta, \eta | \mathcal{D}) \propto p(\mathcal{D} | \theta) p(\theta | \eta) p(\eta)$

# Marginal Likelihood

Without hyperpriors:

$$\begin{aligned} p(y | X, \sigma_y^2, \Sigma_\beta) &= \int \mathcal{N}(y | X\beta, \sigma_y^2 I) \mathcal{N}(\beta | \mathbf{0}, \Sigma_\beta) d\beta \\ &= \mathcal{N}(y | \mathbf{0}, \sigma_y^2 I + X\Sigma_\beta X^T) \end{aligned}$$

# Marginal Likelihood

Without hyperpriors:

$$\begin{aligned}
 p(y | X, \sigma_y^2, \Sigma_\beta) &= \int \mathcal{N}(y | X\beta, \sigma_y^2 I) \mathcal{N}(\beta | \mathbf{0}, \Sigma_\beta) d\beta \\
 &= \mathcal{N}(y | \mathbf{0}, \underbrace{\sigma_y^2 I + X\Sigma_\beta X^T}_{=: C_y})
 \end{aligned}$$

$$\ell(\sigma_y^2, \Sigma_\beta) := -\log p(y | X, \sigma_y^2, \Sigma_\beta) \propto \log |C_y| + y^T C_y^{-1} y$$

# Marginal Likelihood

Without hyperpriors:

$$\begin{aligned}
 p(y | X, \sigma_y^2, \Sigma_\beta) &= \int \mathcal{N}(y | X\beta, \sigma_y^2 I) \mathcal{N}(\beta | \mathbf{0}, \Sigma_\beta) d\beta \\
 &= \mathcal{N}(y | \mathbf{0}, \underbrace{\sigma_y^2 I + X\Sigma_\beta X^T}_{=: C_y})
 \end{aligned}$$

$$\ell(\sigma_y^2, \Sigma_\beta) := -\log p(y | X, \sigma_y^2, \Sigma_\beta) \propto \log |C_y| + y^T C_y^{-1} y$$

With hyperpriors:

$$\begin{aligned}
 \ell(\sigma_y^2, \Sigma_\beta) &:= -\log p(y | X, \sigma_y^2, \Sigma_\beta) p(\sigma_y^2 | c, d) p(\Sigma_\beta | a, b) \\
 &\propto \log |C_y| + y^T C_y^{-1} y + \sum_{m=1}^M (-a \log \sigma_{\beta_m}^2 - b / \sigma_{\beta_m}^2) \\
 &\quad - c \log \sigma_y^2 - d / \sigma_y^2
 \end{aligned}$$

# Inferring Parameters $\beta$

$$\begin{aligned}
 p(\beta \mid X, y, \sigma_y^2, \Sigma_\beta) &= \frac{1}{Z} \mathcal{N}(\beta \mid 0, \Sigma_\beta) \mathcal{N}(y \mid X\beta, \sigma_y^2 I) \\
 &= \mathcal{N}(\beta \mid \mu_\beta := \frac{1}{\sigma_y^2} C_\beta X^T y, C_\beta := (\frac{1}{\sigma_y^2} X^T X + \Sigma_\beta^{-1})^{-1})
 \end{aligned}$$

for  $\Sigma_\beta = \infty I$ : unregularized estimates

$$= \mathcal{N}(\beta \mid (X^T X)^{-1} X^T y, \sigma_y^2 (X^T X)^{-1})$$

for  $\sigma_y^2 = \infty$ : overregularized estimates

$$= \mathcal{N}(\beta \mid 0, \Sigma_\beta)$$

# Equivalent MAP Estimation Problem

$$\ell(\beta) = \frac{1}{\sigma_y^2} \|y - X\beta\|_2^2 + \min_{\sigma_{\beta_1}^2, \dots, \sigma_{\beta_M}^2 \geq 0} \log |\sigma_y^2 I + X \text{diag}(\sigma_{\beta}^2) X^T| + \sum_{m=1}^M \frac{\beta_m^2}{\sigma_{\beta_m}^2}$$



# Learning ARD I: via EM

$$\begin{aligned}
 \ell(\sigma_y^2, \Sigma_\beta, \mu_\beta, C_\beta) &:= E_{\beta \sim \mathcal{N}(\mu_\beta, C_\beta)}(\log p(y | X, \beta, \sigma_y^2, \Sigma_\beta)) \\
 &= E_{\beta \sim \mathcal{N}(\mu_\beta, C_\beta)}(\log \mathcal{N}(y | X\beta, \sigma_y^2) + \log \mathcal{N}(\beta | 0, \Sigma_\beta)) \\
 &\quad + \sum_{m=1}^M \log \text{InvGamma}(\sigma_{\beta_m}^2 | a, b) + \log \text{InvGamma}(\sigma_y^2 | c, d) \\
 &\propto E_{\beta \sim \mathcal{N}(\mu_\beta, C_\beta)}\left(-\frac{N}{2} \log \sigma_y^2 - \frac{1}{2\sigma_y^2} \|y - X\beta\|^2 - \frac{1}{2} \sum_m \log \sigma_{\beta_m}^2 - \frac{1}{2} \text{tr} \Sigma_\beta^{-1} \beta \beta^T\right. \\
 &\quad \left. + \sum_{m=1}^M (-a \log \sigma_{\beta_m}^2 - b/\sigma_{\beta_m}^2) - c \log \sigma_y^2 - d/\sigma_y^2\right) \\
 &= -\frac{N}{2} \log \sigma_y^2 - \frac{1}{2\sigma_y^2} (\|y - X\mu_\beta\|^2 + \text{tr}(X^T X C_\beta)) - \frac{1}{2} \sum_m \log \sigma_{\beta_m}^2 \\
 &\quad - \frac{1}{2} \text{tr} \Sigma_\beta^{-1} (\mu_\beta \mu_\beta^T + C_\beta) + \sum_{m=1}^M (-a \log \sigma_{\beta_m}^2 - b/\sigma_{\beta_m}^2) - c \log \sigma_y^2 - d/\sigma_y^2
 \end{aligned}$$

# Learning ARD I: via EM

$$\begin{aligned}
 \ell(\dots) &= -\frac{N}{2} \log \sigma_y^2 - \frac{1}{2\sigma_y^2} (\|y - X\mu_\beta\|^2 + \text{tr}(X^T X C_\beta)) - \frac{1}{2} \sum_m \log \sigma_{\beta_m}^2 \\
 &\quad - \frac{1}{2} \text{tr} \Sigma_\beta^{-1} (\mu_\beta \mu_\beta^T + C_\beta) + \sum_{m=1}^M (-a \log \sigma_{\beta_m}^2 - b/\sigma_{\beta_m}^2) - c \log \sigma_y^2 - d/\sigma_y^2 \\
 &\propto -(2c + N) \log \sigma_y^2 - (2d + \|y - X\mu_\beta\|^2 + \text{tr}(X^T X C_\beta)) \frac{1}{\sigma_y^2} \\
 &\quad - \sum_m (2a + 1) \log \sigma_{\beta_m}^2 + 2b/\sigma_{\beta_m}^2 - \text{tr} \Sigma_\beta^{-1} (\mu_\beta \mu_\beta^T + C_\beta)
 \end{aligned}$$

# Learning ARD I: via EM

$$\ell(\dots) = - (2c + N) \log \sigma_y^2 - (2d + \|y - X\mu_\beta\|^2 + \text{tr}(X^T X C_\beta)) \frac{1}{\sigma_y^2} \\ - \sum_m (2a + 1) \log \sigma_{\beta_m}^2 + 2b/\sigma_{\beta_m}^2 - \text{tr} \Sigma_\beta^{-1} (\mu_\beta \mu_\beta^T + C_\beta)$$

$$0 \stackrel{!}{=} \frac{\partial \ell}{\partial \sigma_{\beta_m}^2} = - (2a + 1) \frac{1}{\sigma_{\beta_m}^2} + (2b + (\mu_\beta)_m^2 + (C_\beta)_{m,m}) / (\sigma_{\beta_m}^2)^2$$

$$\sigma_{\beta_m}^2 = \frac{2b + (\mu_\beta)_m^2 + (C_\beta)_{m,m}}{2a + 1}$$

$$0 \stackrel{!}{=} \frac{\partial \ell}{\partial \sigma_y^2}$$

$$\rightsquigarrow \sigma_y^2 = \frac{2d + \|y - X\mu_\beta\|^2 + \text{tr}(X^T X C_\beta)}{2c + N}$$

which can be accelerated using

$$C_\beta X^T X = \sigma_y^{2\text{old}} (I - C_\beta \Sigma_\beta^{-1}), \quad \text{tr}(\dots) = \sigma_y^{2\text{old}} \sum \left( 1 - \frac{(C_\beta)_{m,m}}{\sigma_\alpha^2} \right)$$

# Learning ARD I: via EM

iteratively fit:

$$\sigma_{\beta_m}^2 := \frac{2b + (\mu_\beta)_m^2 + (C_\beta)_{m,m}}{2a + 1}$$

$$\sigma_y^2 := \frac{2d + \|y - X\mu_\beta\|^2 + \text{tr}(X^T X C_\beta)}{2c + N}$$

$$\mu_\beta := \frac{1}{\sigma_y^2} C_\beta X^T y$$

$$C_\beta := \left( \frac{1}{\sigma_y^2} X^T X + \Sigma_\beta^{-1} \right)^{-1}, \quad \Sigma_\beta := \text{diag}(\sigma_\beta^2)$$

finally yielding:

$$\beta \sim \mathcal{N}(\mu_\beta, C_\beta)$$

# Learning ARD II: Fixed Point Algorithm

iteratively fit:

$$\sigma_{\beta_m}^2 := \frac{2b + (\mu_{\beta})_m^2}{2a + \gamma_m}$$

$$\sigma_y^2 := \frac{2d + \|y - X\mu_{\beta}\|^2}{2c + N - \sum_m \gamma_m}$$

$$C_{\beta} := \left( \frac{1}{\sigma_y^2} X^T X + \Sigma_{\beta}^{-1} \right)^{-1}$$

$$\mu_{\beta} := \frac{1}{\sigma_y^2} C_{\beta} X^T y$$

$$\gamma_m := 1 - \frac{(C_{\beta})_{m,m}}{\sigma_{\beta_m}^2}, \quad m := 1, \dots, M$$

# Learning ARD III: Iteratively Reweighted L1

The ARD regularization term

$$R(\sigma_\beta^2) := \log |C_Y(\sigma_\beta^2)| = \log |\sigma_Y^2 I + X \Sigma_\beta X^T|, \quad \Sigma_\beta := \text{diag}(\sigma_\beta^2)$$

is concave in  $\sigma_\beta^2$  and thus can be written as

$$R(\sigma_\beta^2) = \min_{\lambda} \lambda^T \sigma_\beta^2 - R^*(\lambda)$$

$$R^*(\lambda) = \min_{\tilde{\sigma}_\beta^2} \lambda^T \tilde{\sigma}_\beta^2 - \log |C_Y(\tilde{\sigma}_\beta^2)|$$

The relaxed function

$$R(\sigma_\beta^2, \lambda) := \lambda^T \sigma_\beta^2 - R^*(\lambda) = \lambda^T \sigma_\beta^2 - \min_{\tilde{\sigma}_\beta^2} \lambda^T \tilde{\sigma}_\beta^2 - \log |C_Y(\tilde{\sigma}_\beta^2)|$$

for fixed  $\sigma_\beta^2$  is minimized by

$$\lambda = \nabla_{\sigma_\beta^2} \log |C_Y(\sigma_\beta^2)|$$

# Learning ARD III: Iteratively Reweighted L1

Instead of  $\sigma_{\beta}^2$  Wipf/Nagarajan 2008 use

$$\sigma_{\beta_m}^2 \xrightarrow{??} \lambda_m^{\frac{1}{2}} |\beta_m|$$

finally yielding the iterative procedure:

$$\beta^{(t+1)} := \arg \min_{\beta} \ell(\beta) + \sum_{m=1}^M \lambda_m^{(t)} |\beta_m|$$

and to find  $\lambda^{(t)}$ :

$$\lambda_m^{(0)} := 1$$

$$\lambda_m^{(t+1)} := (X_{.,m}(\sigma_y^2 I + X \text{diag}(\frac{1}{\lambda_1^{(t)}}, \dots, \frac{1}{\lambda_M^{(t)}}) \text{diag}(|\beta_1^{(t)}|, \dots, |\beta_M^{(t)}|))^{-1} X_{.,m})^{\frac{1}{2}}$$

# ARD for Classification

- ▶ so far, everything was developed for linear regression.
- ▶ for logistic regression, for EM the E-step cannot be done analytically.
  - ▶ possibly use variational approximation
  - ▶ use Gaussian approximation (Laplace approximation)
- ▶ the iteratively reweighted learning algorithm still works.



# Remarks

- ▶ ARD is a good example for a (arguably simple) hierarchical Bayesian model.
- ▶ ARD has to be diligently evaluated against simple baselines such as normalizing the data with a vanilla L1/L2 regularized model.

## Further Readings

- ▶ L1 regularization: [Mur12, chapter 13.3–5], [HTFF05, chapter 3.4, 3.8, 4.4.4], [Bis06, chapter 3.1.4].
  - ▶ LAR, LARS: [HTFF05, chapter 3.4.4], [Mur12, chapter 13.4.2],
- ▶ Non-convex regularizers: [Mur12, chapter 13.6].
- ▶ Automatic Relevance Determination (ARD): [Mur12, chapter 13.7], [HTFF05, chapter 11.9.1], [Bis06, chapter 7.2.2].
- ▶ Sparse Coding: [Mur12, chapter 13.8].

# References



Christopher M. Bishop.

*Pattern recognition and machine learning*, volume 1.

springer New York, 2006.



Trevor Hastie, Robert Tibshirani, Jerome Friedman, and James Franklin.

*The elements of statistical learning: data mining, inference and prediction*, volume 27.

Springer, 2005.



Kevin P. Murphy.

*Machine learning: a probabilistic perspective*.

The MIT Press, 2012.