

Machine Learning 2

6. Sparse Linear Models

Lars Schmidt-Thieme

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

Syllabus

			A. Advanced Supervised Learning
Fri.	24.4.	(1)	A.1 Generalized Linear Models
Fri.	1.5.	—	— <i>Labour Day</i> —
Fri.	8.5.	(2)	A.2 Gaussian Processes
Fri.	15.5.	(3)	A.3 Advanced Support Vector Machines
			B. Ensembles
Fri.	22.5.	(4)	B.1 Stacking & B.2 Boosting
Fri.	29.5.	(5)	B.3 Mixtures of Experts
Fri.	5.6.	—	— <i>Pentecoste Break</i> —
			C. Sparse Models
Fri.	12.6.	(6)	C.1 Homotopy and Least Angle Regression
Fri.	19.6.	(7)	C.2 Proximal Gradients
Fri.	26.6.	(8)	C.3 Laplace Priors
Fri.	3.7.	(9)	C.4 Automatic Relevance Determination
			D. Complex Predictors
Fri.	10.7.	(10)	D.1 Latent Dirichlet Allocation (LDA)
Fri.	17.7.	(11)	Q & A

Outline

1. Homotopy Methods: Least Angle Regression
2. Proximal Gradient Methods
3. Laplace Priors (Bayesian Lasso)

Outline

1. Homotopy Methods: Least Angle Regression
2. Proximal Gradient Methods
3. Laplace Priors (Bayesian Lasso)

Sparse Models so far

- ▶ Variable subset selection
 - ▶ forward search, backward search

- ▶ L1 regularization / Lasso
 - ▶ Coordinate descent (shooting algorithm)

L1 Regularization

$$\min_{\hat{\theta} \in \mathbb{R}^P} f(\hat{\theta}) := \ell(y, \hat{y}(\hat{\theta}, X)) + \lambda \|\hat{\theta}\|_1$$

is equivalent to

$$\min_{\hat{\theta} \in \mathbb{R}^P} f(\hat{\theta}) := \ell(y, \hat{y}(\hat{\theta}, X))$$
$$\|\hat{\theta}\|_1 \leq B$$

with

$$B := \|\hat{\theta}^*\|_1$$

L1 Regularization / Equivalence

More generally, given

$$x^* := \arg \min_x f(x) + \lambda g(x), \quad \lambda \geq 0 \quad (1)$$

$$\tilde{x} := \arg \min_{x: g(x) \leq g(x^*)} f(x) \quad (2)$$

then

$$x^* = \tilde{x}$$

because

$$f(\tilde{x}) \stackrel{(2)}{\leq} f(x^*) \stackrel{(1)}{\leq} f(\tilde{x}) + \lambda \underbrace{(g(\tilde{x}) - g(x^*))}_{\leq 0} \leq f(\tilde{x})$$

$$\rightsquigarrow f(\tilde{x}) = f(x^*) \quad \rightsquigarrow \tilde{x} = x^*$$

assuming x^* is unique.

Homotopy Methods

$$\min. f(\hat{\theta}) := \ell(y, \hat{y}(\hat{\theta}, X)) + \lambda \|\hat{\theta}\|_1$$

or equivalently

$$\begin{aligned} \min. f(\hat{\theta}) &:= \ell(y, \hat{y}(\hat{\theta}, X)) \\ &\|\hat{\theta}\|_1 \leq B \end{aligned}$$

- ▶ start with a solution for large $\lambda^{(0)}$ (or equiv. $B^{(0)} := 0$)
 - ▶ then $\hat{\theta}^{(0)} = 0$.

- ▶ stepwise decrease $\lambda^{(t)}$ (or equiv. increase $B^{(t)}$)
 - ▶ learn $\hat{\theta}^{(t)}$ starting from $\hat{\theta}^{(t-1)}$ (**warmstart**).

Homotopy Methods / Prerequisites

For homotopy to work,

1. the parameters as function of λ

$$\hat{\theta}(\lambda) := \arg \min_{\hat{\theta}} \ell(y, \hat{y}(\hat{\theta}, X)) + \lambda \|\hat{\theta}\|_1$$

must be continuous, i.e.,

- ▶ \hat{y} must be continuous in $\hat{\theta}$ and
- ▶ ℓ be continuous in \hat{y} .

2. the steps in $\lambda^{(t)}$ must be small enough.

Homotopy for the L1 Weight of Linear Regression

Most simple model: linear regression

- ▶ model $\hat{y}(\hat{\theta}, X) := X\hat{\theta}$
- ▶ loss $\ell(y, \hat{y}) := \|y - \hat{y}\|_2^2$

Advantage: can find optimal $\lambda^{(t)}$ sequence analytically! (actually $B^{(t)}$)

Homotopy for the L1 Weight of Linear Regression

Most simple model: linear regression

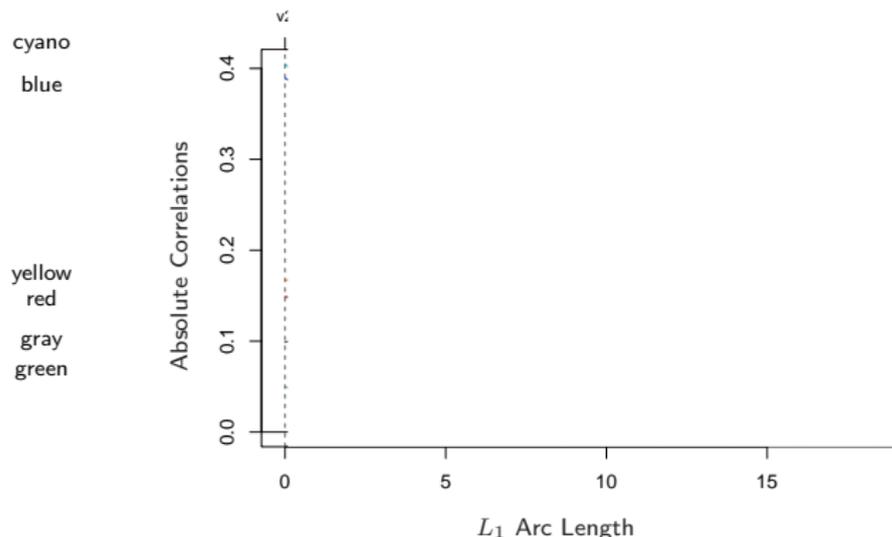
- ▶ model $\hat{y}(\hat{\theta}, X) := X\hat{\theta}$
- ▶ loss $\ell(y, \hat{y}) := \|y - \hat{y}\|_2^2$

Advantage: can find optimal $\lambda^{(t)}$ sequence analytically! (actually $B^{(t)}$)

Imagine the following situation:

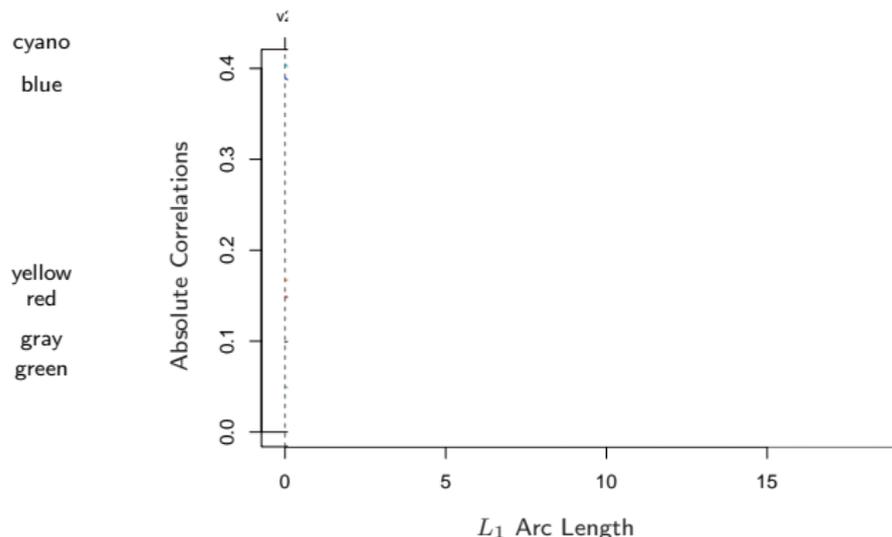
- ▶ initially all parameters $\hat{\theta}_m = 0$.
- ▶ you can add one variable x_m to the model
 - ▶ by setting its parameter $\hat{\theta}_m$ to a small positive or negative value ϵ .
- ▶ the goal is to reduce the error as much as possible.
- ▶ Q: which parameter $\hat{\theta}_m$ would you choose?

Example



Q: which parameter will you pick initially to reduce the loss maximally?

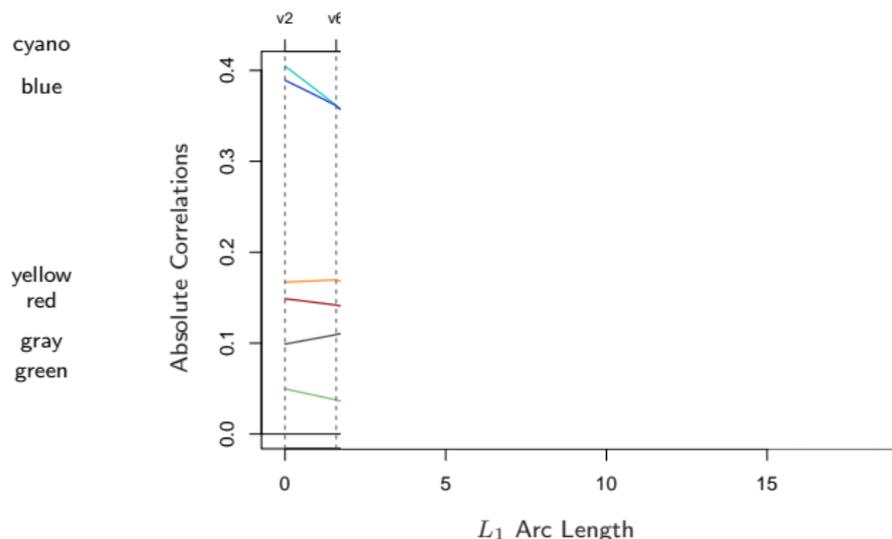
Example



Q: which parameter will you pick initially to reduce the loss maximally?

Q₂: How will the cyano curve go?

Example



Q: which parameter will you pick initially to reduce the loss maximally?

Q₂: How will the cyano curve go?

Least Angle Regression (LAR)

in step t :

1. choose the predictors with largest correlation with the residuum (**active predictors**):

$$C^{(t-1)} := X^T (y - \hat{y}^{(t-1)})$$

$$A^{(t)} := \arg \max_m |C_m^{(t-1)}| \quad (\text{a set!})$$

2. regress these predictors on the residuum:

$$X^{(t)} := X_{\cdot, A^{(t)}}$$

$$\hat{\gamma}^{(t)} := \arg \min_{\gamma} \|y - \hat{y}^{(t-1)} - X^{(t)}\gamma\|_2$$

$$= (X^{(t)T} X^{(t)})^{-1} X^{(t)T} (y - \hat{y}^{(t-1)})$$

3. update parameters in this direction:

$$\hat{\beta}^{(t)} := \hat{\beta}^{(t-1)} + \alpha \Delta^{(t)} \hat{\gamma}^{(t)}$$

Note: $\Delta_{m_k, k}^{(t)} := 1$ for $A^{(t)} := \{m_1, m_2, \dots, m_K\}$, $\Delta_{m, k}^{(t)} := 0$ otherwise.

Least Angle Regression (LAR): step length

Residuum correlations after the update

$$\begin{aligned}
 C^{(t)} &= X^T (y - \hat{y}^{(t)}) = X^T (y - X \hat{\beta}^{(t)}) = X^T (y - X(\hat{\beta}^{(t-1)} + \alpha \Delta^{(t)} \hat{\gamma}^{(t)})) \\
 &= C^{(t-1)} - \alpha X^T X \Delta^{(t)} \hat{\gamma}^{(t)} \\
 &= C^{(t-1)} - \alpha X^{(t)T} X^{(t)} \hat{\gamma}^{(t)}
 \end{aligned}$$

are uniformly reduced for active predictors:

$$C^{(t)}|_{A^{(t)}} = C^{(t-1)}|_{A^{(t)}} - \alpha X^{(t)T} X^{(t)} \hat{\gamma}^{(t)} = (1 - \alpha) C^{(t-1)}|_{A^{(t)}}$$

and may also change for non-active predictors:

$$C_m^{(t)} = C_m^{(t-1)} - \alpha X_{\cdot, m}^T X^{(t)} \hat{\gamma}^{(t)}$$

Note: Maybe a mistake somewhere here. Final formula for α differs from the one in the paper.

Least Angle Regression (LAR): step length (2/2)

Reduce until another predictor has same (max) residuum correlation:

$$C_m^{(t)} = C_m^{(t-1)} - \alpha X_{:,m}^T X^{(t)} \hat{\gamma}^{(t)} \stackrel{!}{=} (1 - \alpha) C_{\max}^{(t-1)}$$

$$\alpha = \frac{C_{\max}^{(t-1)} - C_m^{(t-1)}}{C_{\max}^{(t-1)} - X_{:,m}^T X^{(t)} \hat{\gamma}^{(t)}}$$

or for negative correlations:

$$C_m^{(t)} = C_m^{(t-1)} - \alpha X_{:,m}^T X^{(t)} \hat{\gamma}^{(t)} \stackrel{!}{=} -(1 - \alpha) C_{\max}^{(t-1)}$$

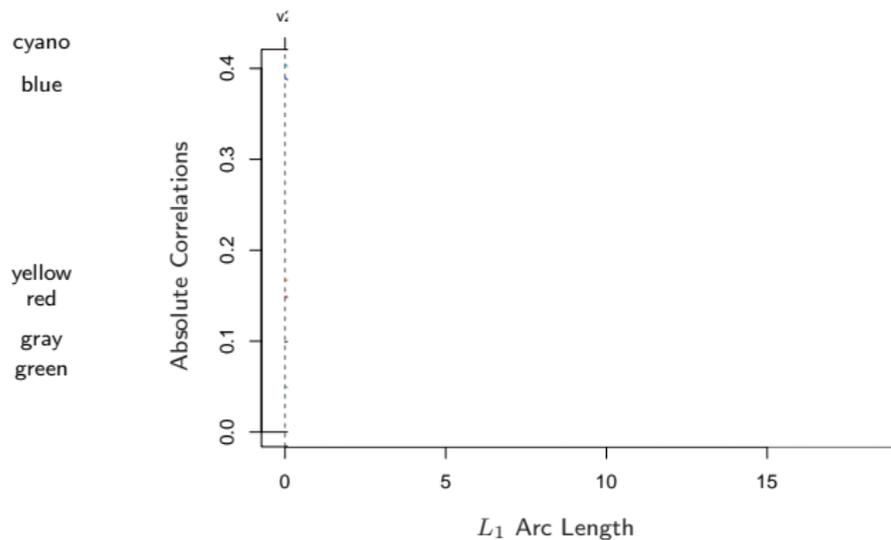
$$\alpha = \frac{C_{\max}^{(t-1)} + C_m^{(t-1)}}{C_{\max}^{(t-1)} + X_{:,m}^T X^{(t)} \hat{\gamma}^{(t)}}$$

yielding

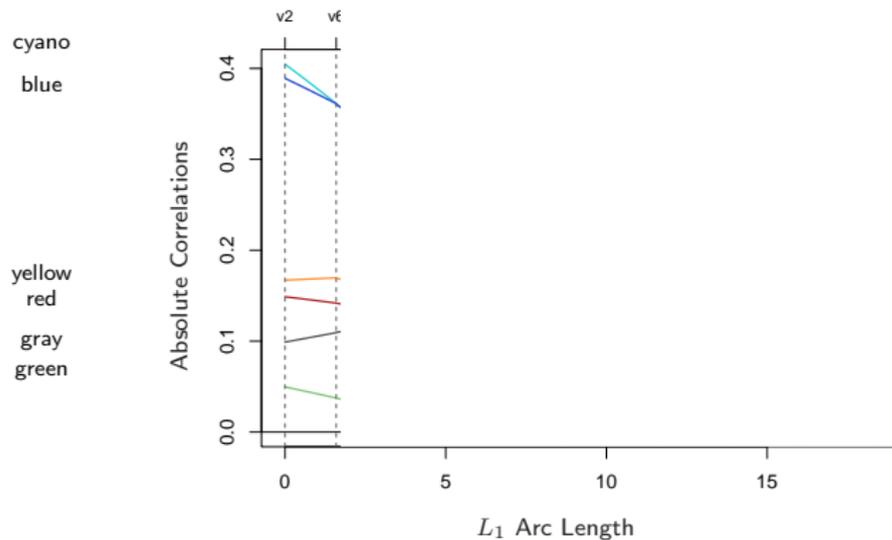
$$\alpha := \minpos \left\{ \frac{C_{\max}^{(t-1)} - C_m^{(t-1)}}{C_{\max}^{(t-1)} - X_{:,m}^T X^{(t)} \hat{\gamma}^{(t)}}, \frac{C_{\max}^{(t-1)} + C_m^{(t-1)}}{C_{\max}^{(t-1)} + X_{:,m}^T X^{(t)} \hat{\gamma}^{(t)}} \right\}$$

$$| m \in \{1, \dots, M\} \setminus A^{(t)} \}, \quad \minpos(X) := \min\{x \in X \mid x > 0\}$$

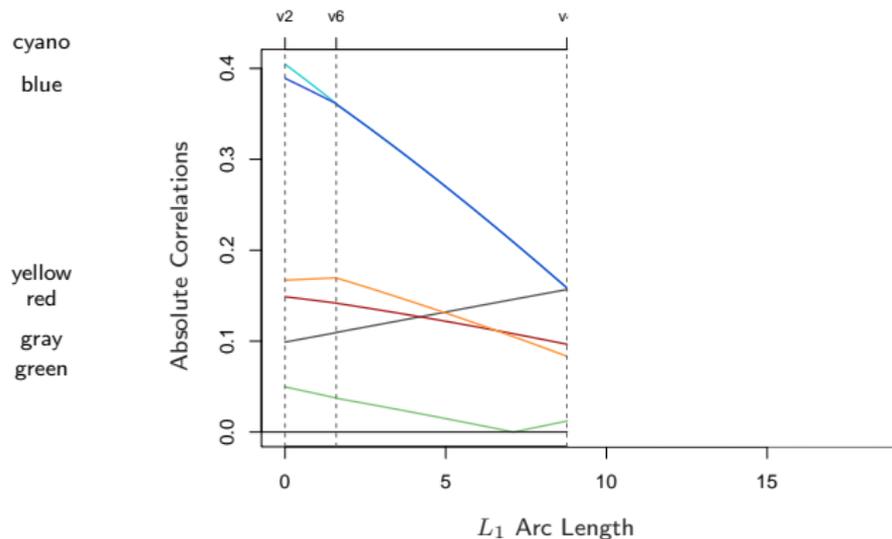
Example



Example



Example



Example

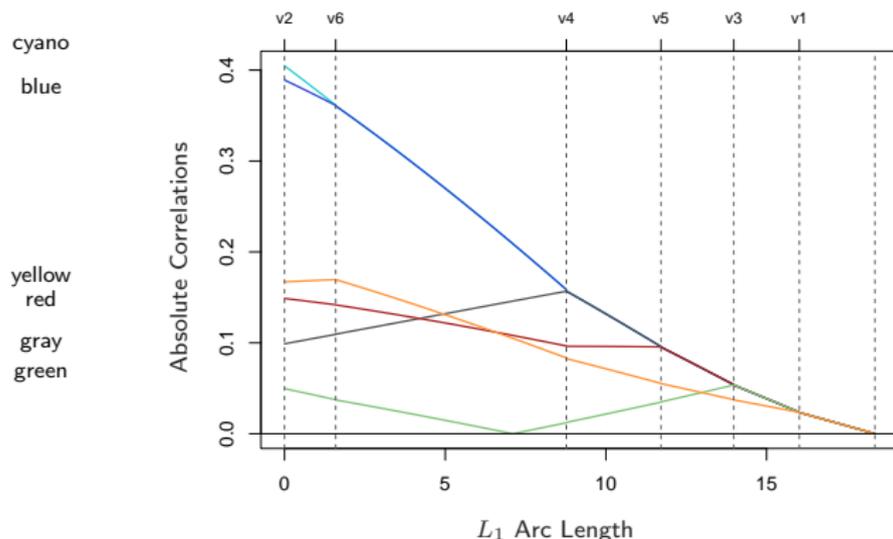
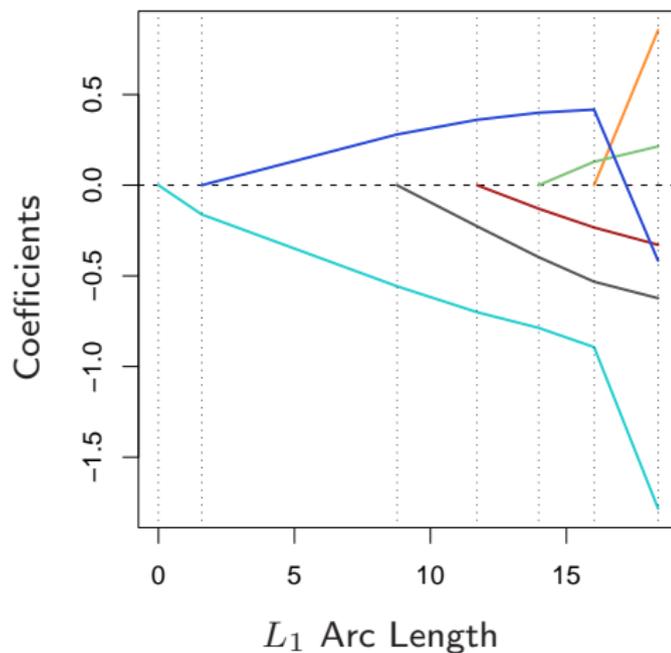


FIGURE 3.14. Progression of the absolute correlations during each step of the LAR procedure, using a simulated data set with six predictors. The labels at the top of the plot indicate which variables enter the active set at each step. The step length are measured in units of L_1 arc length.

Example

Least Angle Regression



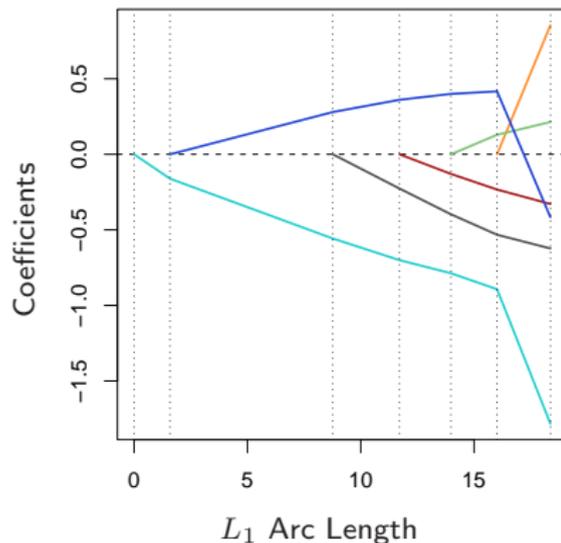
[?, p. 75]

Remarks

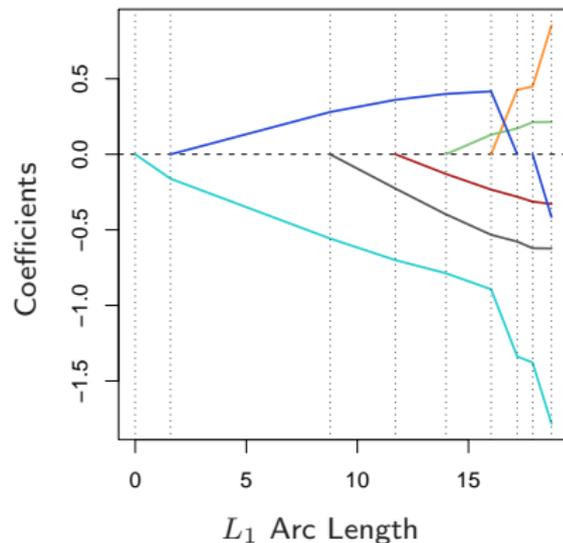
- ▶ algorithm can be used two ways:
 1. Estimate parameters for **all** λ (**regularization path**)
 2. Estimate parameters for **a specific** λ (Homotopy method)
 - ▶ start with large $\lambda^{(0)}$, stop once $\lambda^{(t)} < \lambda$ reached.
- ▶ not straightforward to extend from regression to GLMs
- ▶ LAR can be modified to solve the LASSO:
 - ▶ if the parameter $\beta_m^{(t)}$ for an active predictor m becomes 0 or changes sign, drop it from the active set.
- ▶ also called Least Angle Regression and Shrinkage (LARS)

Example

Least Angle Regression



Lasso



[?, p. 75]

Outline

1. Homotopy Methods: Least Angle Regression
2. Proximal Gradient Methods
3. Laplace Priors (Bayesian Lasso)

Regularized

We want to compute models

$$\theta^* = \arg \min_{\theta} \underbrace{L(\theta)}_{\text{Loss}} + \underbrace{R(\theta)}_{\text{Regularization}}$$

Even when R is not differentiable, e.g.

- ▶ $R(\theta) = \|\theta\|_1$ (L^1 regularization, LASSO)
- ▶ $R(\theta) = I_C(\theta) = \begin{cases} 0 & : \theta \in C \\ \infty & : \theta \notin C \end{cases}$ (hard constraint)

Regularized

We want to compute models

$$\theta^* = \arg \min_{\theta} \underbrace{L(\theta)}_{\text{Loss}} + \underbrace{R(\theta)}_{\text{Regularization}}$$

Even when R is not differentiable, e.g.

▶ $R(\theta) = \|\theta\|_1$ (L^1 regularization, LASSO)

▶ $R(\theta) = I_C(\theta) = \begin{cases} 0: & \theta \in C \\ \infty: & \theta \notin C \end{cases}$ (hard constraint)

Observation: For simple loss functions, we can sometimes compute θ^* analytically

$$\arg \min_{\theta} \frac{1}{2} \|\theta - y\|_2^2 + \lambda \|\theta\|_1 = \text{soft}(y, \lambda)$$

Proximal Problem

- ▶ find x with minimal f **in a vicinity of a given** x^0 :

$$\text{prox}_f(x^0) := \arg \min_x f(x) + \frac{1}{2} \|x - x^0\|_2^2$$

Proximal Problem

- ▶ find x with minimal f **in a vicinity of a given x^0** :

$$\text{prox}_f(x^0) := \arg \min_x f(x) + \frac{1}{2} \|x - x^0\|_2^2$$

Can be solved analytically for some typical (possibly non-differentiable) regularization functions:

- ▶ $f := \lambda \|x\|_2^2$:
$$\text{prox}_f(x^0) = \frac{1}{2\lambda + 1} x^0$$

Proximal Problem

- ▶ find x with minimal f **in a vicinity of a given x^0** :

$$\text{prox}_f(x^0) := \arg \min_x f(x) + \frac{1}{2} \|x - x^0\|_2^2$$

Can be solved analytically for some typical (possibly non-differentiable) regularization functions:

- ▶ $f := \lambda \|x\|_2^2$: $\text{prox}_f(x^0) = \frac{1}{2\lambda + 1} x^0$

$$\begin{aligned} 0 &\stackrel{!}{=} \frac{\partial(\lambda x^T x + \frac{1}{2}(x - x^0)^T (x - x^0))}{\partial x} \\ &= 2\lambda x + (x - x^0) = (2\lambda + 1)x - x^0 \\ \rightsquigarrow x &= \frac{1}{2\lambda + 1} x^0 \end{aligned}$$

Proximal Problem

- ▶ find x with minimal f **in a vicinity of a given x^0** :

$$\text{prox}_f(x^0) := \arg \min_x f(x) + \frac{1}{2} \|x - x^0\|_2^2$$

Can be solved analytically for some typical (possibly non-differentiable) regularization functions:

- ▶ $f := \lambda \|x\|_2^2$:
$$\text{prox}_f(x^0) = \frac{1}{2\lambda + 1} x^0$$

- ▶ $f := \lambda \|x\|_1$:

$$\text{prox}_f(x^0) = \text{soft}(x^0, \lambda) := (\text{soft}(x_n^0, \lambda))_{n=1, \dots, N}$$

$$\text{soft}(z, \lambda) := \text{sign}(z)(|z| - \lambda)_0$$

Proximal Problem

- ▶ find x with minimal f **in a vicinity of a given** x^0 :

$$\text{prox}_f(x^0) := \arg \min_x f(x) + \frac{1}{2} \|x - x^0\|_2^2$$

Can be solved analytically for some typical (possibly non-differentiable) regularization functions:

- ▶ $f := \lambda \|x\|_2^2$:
$$\text{prox}_f(x^0) = \frac{1}{2\lambda + 1} x^0$$

- ▶ $f := \lambda \|x\|_1$:

$$\text{prox}_f(x^0) = \text{soft}(x^0, \lambda) := (\text{soft}(x_n^0, \lambda))_{n=1, \dots, N}$$

$$\text{soft}(z, \lambda) := \text{sign}(z)(|z| - \lambda)_0$$

- ▶ $f := \lambda \|x\|_0$:

$$\text{prox}_f(x^0) = \text{hard}(x^0, \lambda) := (\text{hard}(x_n^0, \lambda))_{n=1, \dots, N},$$

$$\text{hard}(z, \lambda) := \delta(|z| \geq \lambda) z$$

More Analytical Solutions for the Proximal Problem

- find x with minimal f **in a vicinity of a given** x^0 :

$$\text{prox}_f(x^0) := \arg \min_x f(x) + \frac{1}{2} \|x - x^0\|_2^2$$

$f := I_C$ for a **convex set** C and $I_C(x) := \begin{cases} 0, & \text{if } x \in C \\ \infty, & \text{else} \end{cases}$

$$\text{prox}_f(x^0) = \arg \min_{x \in C} \|x - x^0\|_2^2 =: \text{proj}_C(x^0)$$

More Analytical Solutions for the Proximal Problem

- ▶ find x with minimal f **in a vicinity of a given** x^0 :

$$\text{prox}_f(x^0) := \arg \min_x f(x) + \frac{1}{2} \|x - x^0\|_2^2$$

$f := I_C$ for a **convex set** C and $I_C(x) := \begin{cases} 0, & \text{if } x \in C \\ \infty, & \text{else} \end{cases}$

$$\text{prox}_f(x^0) = \arg \min_{x \in C} \|x - x^0\|_2^2 =: \text{proj}_C(x^0)$$

- ▶ **rectangles / box constraints** $C := [l_1, u_1] \times [l_2, u_2] \times \dots \times [l_N, u_N]$:
 $\text{prox}_f(x^0) = \text{clip}(x^0, C)$ with $\text{clip}(x^0, C)_n := \min\{\max\{x_n^0, l_n\}, u_n\}$

More Analytical Solutions for the Proximal Problem

- ▶ find x with minimal f **in a vicinity of a given** x^0 :

$$\text{prox}_f(x^0) := \arg \min_x f(x) + \frac{1}{2} \|x - x^0\|_2^2$$

$f := I_C$ for a **convex set** C and $I_C(x) := \begin{cases} 0, & \text{if } x \in C \\ \infty, & \text{else} \end{cases}$

$$\text{prox}_f(x^0) = \arg \min_{x \in C} \|x - x^0\|_2^2 =: \text{proj}_C(x^0)$$

- ▶ **rectangles / box constraints** $C := [l_1, u_1] \times [l_2, u_2] \times \dots \times [l_N, u_N]$:

$$\text{prox}_f(x^0) = \text{clip}(x^0, C) \quad \text{with } \text{clip}(x^0, C)_n := \min\{\max\{x_n^0, l_n\}, u_n\}$$

- ▶ **euclidean balls** $C := \{x \mid \|x\|_2 \leq 1\}$:

$$\text{prox}_f(x^0) = \begin{cases} \frac{x^0}{\|x^0\|_2}, & \text{if } \|x^0\|_2 > 1 \\ x^0, & \text{else} \end{cases}$$

More Analytical Solutions for the Proximal Problem

- find x with minimal f **in a vicinity of a given x^0** :

$$\text{prox}_f(x^0) := \arg \min_x f(x) + \frac{1}{2} \|x - x^0\|_2^2$$

$f := I_C$ for

- **L1 balls** $C := \{x \mid \|x\|_1 \leq 1\}$:

$$\text{prox}_f(x^0) = \begin{cases} \text{soft}(x^0, \lambda), & \text{if } \|x^0\|_1 > 1 \\ x^0, & \text{else} \end{cases}$$

$$\text{for } \lambda \text{ with } \sum_{n=1}^N (|x_n^0| - \lambda)_0 \stackrel{!}{=} 1$$

Deriving Generalized Gradient Descent (1/2)

$\min_x f(x) := g(x) + h(x)$, g, h convex, g differentiable, h possibly not

using a **Taylor expansion of g around previous solution $x^{(t)}$** :

$$g(x) \approx g(x^{(t)}) + \nabla g(x^{(t)})(x - x^{(t)}) + \frac{1}{2}(x - x^{(t)})^T H(x - x^{(t)})$$

and **diagonal approximation of the Hessian $H \approx \frac{1}{\alpha^{(t)}} I$**

$$\begin{aligned} &\approx g(x^{(t)}) + \nabla g(x^{(t)})(x - x^{(t)}) + \frac{1}{2\alpha^{(t)}}(x - x^{(t)})^T (x - x^{(t)}) \\ &= \frac{1}{2\alpha^{(t)}}(x - x^{(t)} + 2\alpha^{(t)}\nabla g(x^{(t)}))^T (x - x^{(t)}) + \text{const} \\ &= \frac{1}{2\alpha^{(t)}}(x - (x^{(t)} - \alpha^{(t)}\nabla g(x^{(t)})))^T (x - (x^{(t)} - \alpha^{(t)}\nabla g(x^{(t)}))) \\ &\quad + \text{const} \\ &= \frac{1}{2\alpha^{(t)}}\|x - (x^{(t)} - \alpha^{(t)}\nabla g(x^{(t)}))\|^2 + \text{const} \end{aligned}$$

Deriving Generalized Gradient Descent (2/2)

$$\min_x f(x) := g(x) + h(x), \quad g, h \text{ convex, } g \text{ differentiable, } h \text{ possibly}$$

$$g(x) = \frac{1}{2\alpha^{(t)}} \|x - (x^{(t)} - \alpha^{(t)} \nabla g(x^{(t)}))\|^2 + \text{const}$$

yields a proximal problem

$$\begin{aligned} \min_x f(x) &= \frac{1}{2\alpha^{(t)}} \|x - (x^{(t)} - \alpha^{(t)} \nabla g(x^{(t)}))\|^2 + h(x) \\ &\propto \frac{1}{2} \|x - (x^{(t)} - \alpha^{(t)} \nabla g(x^{(t)}))\|^2 + \alpha^{(t)} h(x) \\ &= \text{prox}_{\alpha^{(t)} h}(x^{(t)} - \alpha^{(t)} \nabla g(x^{(t)})) \end{aligned}$$

$$\text{with } \text{prox}_q(x^0) := \arg \min_x q(x) + \frac{1}{2} \|x - x^0\|^2$$

Generalized Gradient Descent

$$\min_x g(x) + h(x), \quad g, h \text{ convex, } g \text{ differentiable}$$

Generalized Gradient Descent:

$$x^{(t+1)} := \text{prox}_{\alpha^{(t)}h}(x^{(t)} - \alpha^{(t)}\nabla g(x^{(t)}))$$

$$\text{with } \text{prox}_q(x^0) := \arg \min_x q(x) + \frac{1}{2}\|x - x^0\|^2$$

- ▶ two-step approach:
 1. minimize component g via gradient descent
 2. minimize component h via prox operator
- ▶ requires control of step size $\alpha^{(t)}$
- ▶ generalizes gradient descent to objective functions with non-differentiable additive components
- ▶ convergence rate $O(1/t)$.

Application to Regularized Loss Minimization

$$\min \quad f(\theta) := \ell(\theta) + R(\theta)$$

- ▶ ℓ loss, convex and differentiable
 - ▶ e.g., RSS.
- ▶ R regularization, convex, but possibly not differentiable
 - ▶ e.g., $\|\theta\|_1$ or $I_C(\theta) := \begin{cases} 0, & \theta \in C \\ \infty, & \text{else} \end{cases}$

Special Cases

$$\begin{aligned}\theta^{(t+1)} &:= \text{prox}_{\alpha^{(t)}R}(\theta^{(t)} - \alpha^{(t)}\nabla\ell(\theta^{(t)})) \\ &= \arg \min_{\theta} \alpha^{(t)}R(\theta) + \frac{1}{2}\|\theta - (\theta^{(t)} - \alpha^{(t)}\nabla\ell(\theta^{(t)}))\|_2^2\end{aligned}$$

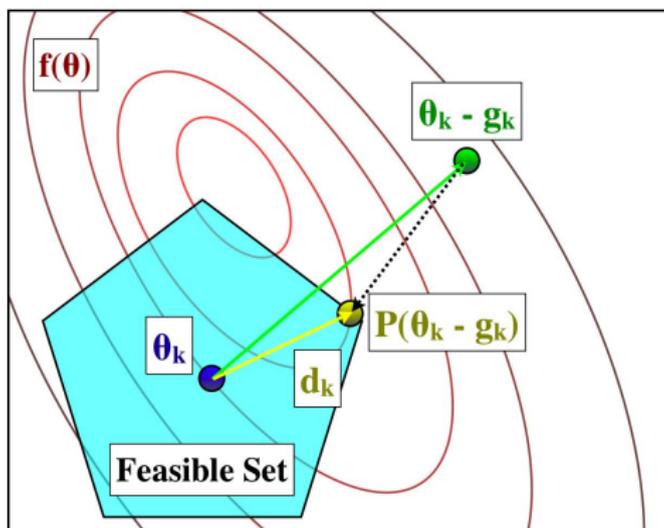
1. $R = 0$ yields **gradient descent**:

$$\theta^{(t+1)} = \theta^{(t)} - \alpha^{(t)}\nabla\ell(\theta^{(t)})$$

2. $R = I_C$ yields **projected gradient descent**:

$$\theta^{(t+1)} = \text{proj}_C(\theta^{(t)} - \alpha^{(t)}\nabla\ell(\theta^{(t)}))$$

Special Cases: Projected Gradient Descent



- ▶ Instead of taking a gradient step and then project, we could compute the smallest stepsize that does not leave the feasible area (“guarded gradient descent”).
- ▶ Q: Which next iterate would “guarded gradient descent” find instead?
- ▶ Now assume the current iterate θ_t is on the upper right border of the feasible area.
- ▶ Q: Which next iterate would “guarded gradient descent” find now?
- ▶ Q: How about projected gradient descent?

[?, fig. 13.11]

Special Cases

$$\begin{aligned}\theta^{(t+1)} &:= \text{prox}_{\alpha^{(t)}R}(\theta^{(t)} - \alpha^{(t)}\nabla\ell(\theta^{(t)})) \\ &= \arg \min_{\theta} \alpha^{(t)}R(\theta) + \frac{1}{2}\|\theta - (\theta^{(t)} - \alpha^{(t)}\nabla\ell(\theta^{(t)}))\|_2^2\end{aligned}$$

3. $R = \lambda\|\theta\|_1$ yields **iterative soft thresholding**:

$$\theta^{(t+1)} = \text{soft}(\theta^{(t)} - \alpha^{(t)}\nabla\ell(\theta^{(t)}), \lambda\alpha^{(t)})$$

Stepsizes $\alpha^{(t)}$

Taylor expansion of the Gradient:

$$\begin{aligned}\nabla\ell(\theta) &\approx \nabla\ell(\theta^{(t)}) + \nabla^2\ell(\theta^{(t)})(\theta - \theta^{(t)}) \approx \nabla\ell(\theta^{(t)}) + \frac{1}{\alpha^{(t)}}(\theta - \theta^{(t)}) \\ &\rightsquigarrow \alpha^{(t)}\nabla\ell(\theta^{(t)}) - \nabla\ell(\theta^{(t-1)}) \approx (\theta^{(t)} - \theta^{(t-1)})\end{aligned}$$

Idea:

$$\begin{aligned}\alpha^{(t)} &:= \arg \min_{\alpha} \|(\theta^{(t)} - \theta^{(t-1)}) - \alpha(\nabla\ell(\theta^{(t)}) - \nabla\ell(\theta^{(t-1)}))\|_2^2 \\ &= \frac{(\theta^{(t)} - \theta^{(t-1)})^T (\theta^{(t)} - \theta^{(t-1)})}{(\theta^{(t)} - \theta^{(t-1)})^T (\nabla\ell(\theta^{(t)}) - \nabla\ell(\theta^{(t-1)}))}\end{aligned}$$

called **Barzilai-Borwein stepsize** or **spectral stepsize**.

- ▶ does not guarantee decreasing objective values.
- ▶ can be used with any gradient descent method.

Iterative Shrinkage and Thresholding Algorithm (ISTA)

- ▶ proximal gradient descent for L1 regularization
 - ▶ iterative soft thresholding
- ▶ Barzilai-Borwein stepsize
- ▶ in outer loop, homotopy on λ
 - ▶ i.e., gradually reducing $\lambda^{(t)}$ to λ

Note: This algorithm is called Sparse Reconstruction by Separable Approximation (SpaRSA) in the literature.

ISTA Algorithm

```

1 learn-l1reg-ista( $X \in \mathbb{R}^{N \times M}$ ,  $y \in \mathbb{R}^N$ ,  $\lambda > 0$ ,  $s \in (0, 1)$ ,  $M$ ) :
2    $\theta := 0$ ,  $r := y$ ,  $\tilde{\lambda} := \infty$ ,  $\alpha := 1$ 
3   for  $t := 1, 2, 3, \dots$  while  $\tilde{\lambda} \neq \lambda$ :
4      $\tilde{\lambda} := \max(\lambda, s \|X^T r\|_\infty)$ 
5     while  $\ell(\theta) + \lambda \|\theta\|_1$  did not increase too much in the last  $M$  steps:
6        $\theta^{\text{old}} := \theta$ 
7        $\tilde{\theta} := \theta - \alpha \nabla \ell(\theta)$ 
8        $\theta := \text{soft}(\tilde{\theta}, \tilde{\lambda} \alpha)$ 
9        $\alpha := \frac{(\theta - \theta^{\text{old}})^T (\theta - \theta^{\text{old}})}{(\theta - \theta^{\text{old}})^T (\nabla \ell(\theta) - \nabla \ell(\theta^{\text{old}}))}$ 
10       $r := y - X\theta$ 
11  return  $\theta$ 
  
```

Nesterov's Accelerated Generalized Gradient Descent

$$\min_x g(x) + h(x), \quad g, h \text{ convex, } g \text{ differentiable}$$

Generalized Gradient Descent:

$$x^{(t+1)} := \text{prox}_{\alpha^{(t)}h}(x^{(t)} + \frac{t-1}{t+2}(x^{(t)} - x^{(t-1)}) - \alpha^{(t)}\nabla g(x^{(t)}))$$

$$\text{with } \text{prox}_f(x^0) := \arg \min_x f(x) + \frac{1}{2}\|x - x^0\|^2$$

- ▶ added **momentum term**
- ▶ works also for vanilla gradient descent ($h = 0$)
- ▶ convergence rate $O(1/t^2)$!
- ▶ beware, there are at least 3 versions of **Nesterov's method**.

Fast Iterative Shrinkage and Thresholding Alg. (FISTA)

$$\theta^{(t+1)} := \text{prox}_{\alpha^{(t)}R} \left(\theta^{(t)} + \frac{t-1}{t+2} (\theta^{(t)} - \theta^{(t-1)}) - \alpha^{(t)} \nabla \ell(\theta^{(t)}) \right)$$

for $R = \lambda \|\cdot\|_1$ yields iterative soft thresholding:

$$\theta^{(t+1)} = \text{soft} \left(\theta^{(t)} + \frac{t-1}{t+2} (\theta^{(t)} - \theta^{(t-1)}) - \alpha^{(t)} \nabla \ell(\theta^{(t)}), \lambda \alpha^{(t)} \right)$$

using **Nesterov's Accelerated Generalized Gradient Descent**.

FISTA vs ISTA

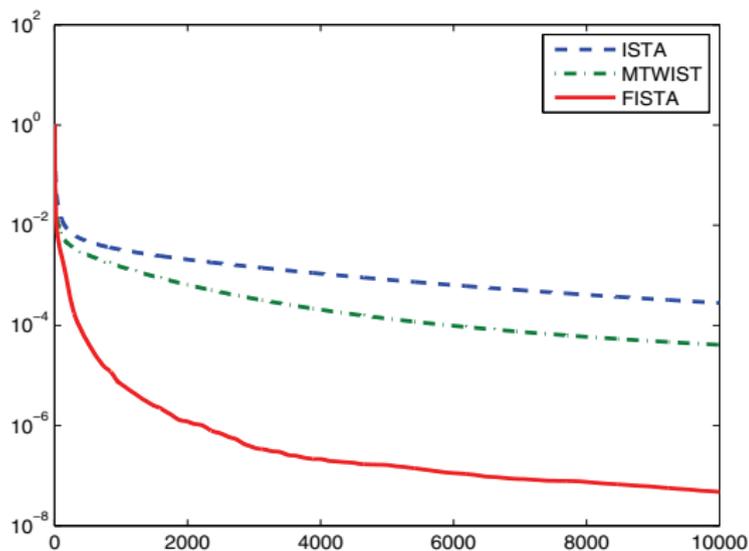


Figure 5. Comparison of function value errors $F(\mathbf{x}_k) - F(\mathbf{x}^*)$ of ISTA, MTWIST, and FISTA.

[?, p. 19]

Outline

1. Homotopy Methods: Least Angle Regression
2. Proximal Gradient Methods
3. Laplace Priors (Bayesian Lasso)

Bayesian Regression

$$\hat{\beta} = \arg \min_{\beta} \underbrace{L(\beta)}_{\text{Loss}} + \lambda \underbrace{R(\beta)}_{\text{Regularization}}$$

↓ "Bayesianize"

$$\hat{\beta} = \arg \max_{\beta} p(\beta | X, y)$$

Bayesian Regression

$$\hat{\beta} = \arg \min_{\beta} \underbrace{L(\beta)}_{\text{Loss}} + \lambda \underbrace{R(\beta)}_{\text{Regularization}}$$

↓ "Bayesianize"

$$\hat{\beta} = \arg \max_{\beta} p(\beta | X, y)$$

$$\underbrace{p(\beta | X, y)}_{\text{posterior}} \propto \underbrace{p(y | X, \beta)}_{\text{likelihood}} \cdot \underbrace{p(\beta)}_{\text{prior}}$$

Bayesian Regression

$$\hat{\beta} = \arg \min_{\beta} \underbrace{L(\beta)}_{\text{Loss}} + \lambda \underbrace{R(\beta)}_{\text{Regularization}}$$

↓ "Bayesianize"

$$\hat{\beta} = \arg \max_{\beta} p(\beta | X, y)$$

$$\underbrace{p(\beta | X, y)}_{\text{posterior}} \propto \underbrace{p(y | X, \beta)}_{\text{likelihood}} \cdot \underbrace{p(\beta)}_{\text{prior}}$$

► $p(y | X, \beta) = \mathcal{N}(y | X\beta, \sigma^2 I) \iff$ Bayesian Linear Regression

Laplace Priors correspond to L1 regularization

L2 regularization:

$$f(\beta) := \|y - X\beta\|_2^2 + \lambda\|\beta\|_2^2$$

Gaussian priors:

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

$$\begin{aligned} p(\beta) &:= \mathcal{N}(\beta | 0, \frac{1}{\lambda} I) \\ &= (2\pi\lambda)^{-M/2} e^{-\frac{1}{2}\lambda\|\beta\|_2^2} \end{aligned}$$

using negative logposterior as objective function:

$$f(\beta; X, y, \sigma^2 \text{ or } \lambda) := -\log p(y | X, \beta) p(\beta)$$

Laplace Priors correspond to L1 regularization

L2 regularization:

$$f(\beta) := \|y - X\beta\|_2^2 + \lambda \|\beta\|_2^2$$

Gaussian priors:

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

$$\begin{aligned} p(\beta) &:= \mathcal{N}(\beta | 0, \frac{1}{\lambda} I) \\ &= (2\pi\lambda)^{-M/2} e^{-\frac{1}{2}\lambda \|\beta\|_2^2} \end{aligned}$$

L1 regularization:

$$f(\beta) := \|y - X\beta\|_2^2 + \lambda \|\beta\|_1$$

Laplace priors:

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

$$\begin{aligned} p(\beta_m) &:= \text{Lap}(\beta_m | 0, \frac{1}{\lambda}) \\ &= \frac{1}{2} \lambda e^{-\lambda |\beta_m|} \end{aligned}$$

using negative logposterior as objective function:

$$f(\beta; X, y, \sigma^2 \text{ or } \lambda) := -\log p(y | X, \beta) p(\beta)$$

Laplace as Gaussian Scale Mixture

- ▶ problem: MAP cannot be found analytically.
- ▶ idea: rewrite the Laplace as a **Gaussian-Scale-Mixture** with Exponential priors:

$$\text{Lap}(\beta_i | 0, \frac{1}{\lambda}) = \int \mathcal{N}(\beta_i | 0, \tau_i^2) \text{Exp}(\tau_i^2 | \frac{1}{2}\lambda^2) d\tau^2$$

i.e. each parameter is distributed as $\beta_i \sim \mathcal{N}(0, \tau_i^2)$ with $\tau_i^2 \sim \text{Exp}(\frac{1}{2}\lambda^2)$

Laplace as Gaussian Scale Mixture

- ▶ problem: MAP cannot be found analytically.
- ▶ idea: rewrite the Laplace as a **Gaussian-Scale-Mixture** with Exponential priors:

$$\text{Lap}(\beta_i | 0, \frac{1}{\lambda}) = \int \mathcal{N}(\beta_i | 0, \tau_i^2) \text{Exp}(\tau_i^2 | \frac{1}{2}\lambda^2) d\tau^2$$

i.e. each parameter is distributed as $\beta_i \sim \mathcal{N}(0, \tau_i^2)$ with $\tau_i^2 \sim \text{Exp}(\frac{1}{2}\lambda^2)$

↪ posterior distribution:

$$p(\beta, \sigma^2 | X, y, \tau^2) \propto \underbrace{p(y | X, \beta, \sigma^2)}_{=\mathcal{N}(y|X\beta, \sigma^2 I)} \cdot \underbrace{p(\beta | \tau^2)}_{=\mathcal{N}(\beta|0, \text{diag}(\tau^2))} \cdot \underbrace{p(\tau^2 | \lambda)}_{=\text{Exp}(\tau^2|\frac{1}{2}\lambda^2)} \cdot \underbrace{p(\sigma^2)}$$

with $(\tau_i)_{m=1\dots M}$ **latent variables**,
 λ the regularization strength hyperparameter,

Laplace as Gaussian Scale Mixture

- ▶ problem: MAP cannot be found analytically.
- ▶ idea: rewrite the Laplace as a **Gaussian-Scale-Mixture** with Exponential priors:

$$\text{Lap}(\beta_i | 0, \frac{1}{\lambda}) = \int \mathcal{N}(\beta_i | 0, \tau_i^2) \text{Exp}(\tau_i^2 | \frac{1}{2}\lambda^2) d\tau^2$$

i.e. each parameter is distributed as $\beta_i \sim \mathcal{N}(0, \tau_i^2)$ with $\tau_i^2 \sim \text{Exp}(\frac{1}{2}\lambda^2)$

↪ posterior distribution:

$$p(\beta, \sigma^2 | X, y, \tau^2) \propto \underbrace{p(y | X, \beta, \sigma^2)}_{=\mathcal{N}(y|X\beta, \sigma^2 I)} \cdot \underbrace{p(\beta | \tau^2)}_{=\mathcal{N}(\beta|0, \text{diag}(\tau^2))} \cdot \underbrace{p(\tau^2 | \lambda)}_{=\text{Exp}(\tau^2|\frac{1}{2}\lambda^2)} \cdot \underbrace{p(\sigma^2)}_{=\text{IG}(\sigma^2|a, b)}$$

with $(\tau_i)_{m=1\dots M}$ **latent variables**,
 λ the regularization strength hyperparameter,
 $\text{IG}(\sigma^2 | a, b)$ an **Inverse-Gamma** prior on the variance.

Laplace as Gaussian Scale Mixture

- ▶ problem: MAP cannot be found analytically.
- ▶ idea: rewrite the Laplace as a **Gaussian-Scale-Mixture** with Exponential priors:

$$\text{Lap}(\beta_i | 0, \frac{1}{\lambda}) = \int \mathcal{N}(\beta_i | 0, \tau_i^2) \text{Exp}(\tau_i^2 | \frac{1}{2}\lambda^2) d\tau^2$$

i.e. each parameter is distributed as $\beta_i \sim \mathcal{N}(0, \tau_i^2)$ with $\tau_i^2 \sim \text{Exp}(\frac{1}{2}\lambda^2)$

↪ posterior distribution:

$$p(\beta, \sigma^2 | X, y, \tau^2) \propto \underbrace{p(y | X, \beta, \sigma^2)}_{=\mathcal{N}(y|X\beta, \sigma^2 I)} \cdot \underbrace{p(\beta | \tau^2)}_{=\mathcal{N}(\beta|0, \text{diag}(\tau^2))} \cdot \underbrace{p(\tau^2 | \lambda)}_{=\text{Exp}(\tau^2|\frac{1}{2}\lambda^2)} \cdot \underbrace{p(\sigma^2)}_{=\text{IG}(\sigma^2|a, b)}$$

with $(\tau_i)_{m=1\dots M}$ **latent variables**,

λ the regularization strength hyperparameter,

$\text{IG}(\sigma^2 | a, b)$ an **Inverse-Gamma** prior on the variance.

- ▶ p is now smooth in all parameters! We can apply EM-algorithm!

Inverse Gamma Distribution

- ▶ Gamma distribution:

$$\Gamma(x \mid a, b) := \frac{b^a}{\Gamma(a)} x^{a-1} e^{-bx}$$

$$\mathbb{E}(x) = \frac{a}{b}$$

- ▶ **Inverse Gamma** distribution:

$$\text{IG}(x \mid a, b) := \frac{b^a}{\Gamma(a)} x^{-a-1} e^{-\frac{b}{x}}$$

- ▶ $X \sim \Gamma(a, b) \iff X^{-1} \sim \text{IG}(a, b)$

$$\rightsquigarrow \mathbb{E}\left(\frac{1}{x}\right) = \frac{a}{b}$$

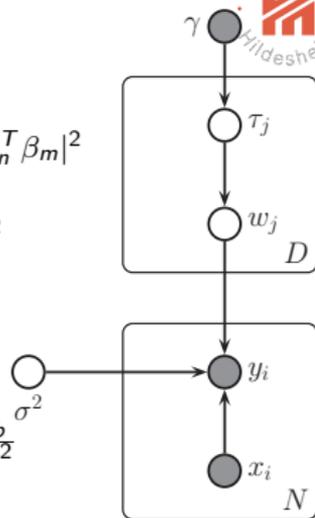
Laplace Prior as Gaussian Scale Mixture

$$p(y_n | x_n, \beta, \sigma^2) := \mathcal{N}(y_n | x_n^T \beta, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2} |y_n - x_n^T \beta|^2}$$

$$p(\beta_m | \tau_m^2) := \mathcal{N}(\beta_m | 0, \tau_m^2) = \frac{1}{\sqrt{2\pi\tau_m^2}} e^{-\frac{1}{2\tau_m^2} |\beta_m|^2}$$

$$p(\tau_m^2) := \text{Exp}(\tau_m^2 | \frac{1}{2}\lambda^2) = \frac{1}{2}\lambda^2 e^{-\frac{1}{2}\lambda^2\tau_m^2}$$

$$p(\sigma^2) := \text{IG}(\sigma^2 | a, b) = \frac{b^a}{\Gamma(a)} \sigma^{-2(1+a)} e^{-\frac{b}{\sigma^2}}$$



Note: $T := \text{diag}(\tau_1^2, \tau_2^2, \dots, \tau_M^2)$

[?, p. 446]

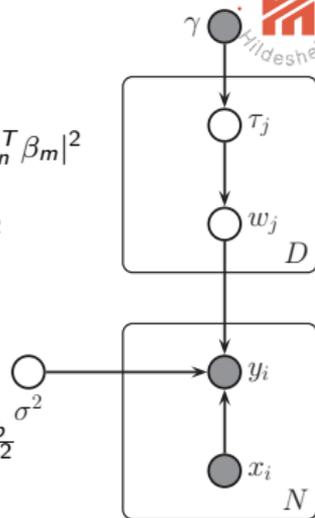
Laplace Prior as Gaussian Scale Mixture

$$p(y_n | x_n, \beta, \sigma^2) := \mathcal{N}(y_n | x_n^T \beta, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2} |y_n - x_n^T \beta|^2}$$

$$p(\beta_m | \tau_m^2) := \mathcal{N}(\beta_m | 0, \tau_m^2) = \frac{1}{\sqrt{2\pi\tau_m^2}} e^{-\frac{1}{2\tau_m^2} |\beta_m|^2}$$

$$p(\tau_m^2) := \text{Exp}(\tau_m^2 | \frac{1}{2}\lambda^2) = \frac{1}{2}\lambda^2 e^{-\frac{1}{2}\lambda^2 \tau_m^2}$$

$$p(\sigma^2) := \text{IG}(\sigma^2 | a, b) = \frac{b^a}{\Gamma(a)} \sigma^{-2(1+a)} e^{-\frac{b}{\sigma^2}}$$



negative logposterior:

$$\begin{aligned} \ell(\beta, \sigma^2 | X, y, \tau^2) &= \frac{1}{2} N \log \sigma^2 + \frac{1}{2\sigma^2} \|y - X\beta\|_2^2 \\ &+ \sum_{m=1}^M \log \tau_m^2 + \frac{1}{2} \beta^T T^{-1} \beta + \frac{1}{2} \lambda^2 \sum_{m=1}^M \tau_m^2 + (1+a) \log \sigma^2 + \frac{b}{\sigma^2} \end{aligned}$$

Note: $T := \text{diag}(\tau_1^2, \tau_2^2, \dots, \tau_M^2)$

[?, p. 446]

E-step for τ^2

We need to compute the expectation of

$$p(\tau^2 | X, y, \beta, \sigma^2) \propto p(\beta | \tau^2)p(\tau^2)$$

where $p(\beta_m | \tau_m^2) = \mathcal{N}(\beta_m | 0, \tau_m^2)$ and $p(\tau_m^2) = \text{Exp}(\tau_m^2 | \frac{1}{2}\lambda^2)$

E-step for τ^2

We need to compute the expectation of

$$p(\tau^2 | X, y, \beta, \sigma^2) \propto p(\beta | \tau^2) p(\tau^2)$$

where $p(\beta_m | \tau_m^2) = \mathcal{N}(\beta_m | 0, \tau_m^2)$ and $p(\tau_m^2) = \text{Exp}(\tau_m^2 | \frac{1}{2}\lambda^2)$

It turns out simpler to estimate $\frac{1}{\tau^2}$: One can show that (tutorial)

$$\frac{1}{\tau^2} | \beta \sim \text{InvGauss}(\sqrt{\frac{\lambda^2}{\beta^2}}, \lambda^2)$$

Where the **Inverse Gaussian distribution** is given by

$$\text{InvGauss}(x | \mu, \nu) = \sqrt{\frac{\nu}{2\pi x^3}} e^{-\frac{\nu}{2\mu^2 x}(x-\mu)^2}$$

with mean $\mathbb{E}[x] = \mu$ and variance $\text{Var}[x] = \mu^3/\nu \implies \boxed{\mathbb{E}\left[\frac{1}{\tau_m^2}\right] = \frac{\lambda}{|\beta_m|}}$

E-step for σ^2

We need to compute the expectation of

$$\begin{aligned} p(\sigma^2 | X, y, \beta, \tau^2) &\propto p(y | X, \beta, \sigma^2) p(\sigma^2) \\ &= \mathcal{N}(y | X\beta, \sigma^2 I) \text{IG}(\sigma^2 | a, b) \end{aligned}$$

E-step for σ^2

We need to compute the expectation of

$$\begin{aligned} p(\sigma^2 | X, y, \beta, \tau^2) &\propto p(y | X, \beta, \sigma^2) p(\sigma^2) \\ &= \mathcal{N}(y | X\beta, \sigma^2 I) \text{IG}(\sigma^2 | a, b) \end{aligned}$$

One can show that (tutorial)

$$\begin{aligned} p(\sigma^2 | X, y, \beta, \tau^2) &= \text{IG}(\sigma^2, a', b') \\ \text{with } a' &:= a + \frac{1}{2}N, \quad b' := b + \frac{1}{2}\|y - X\beta\|_2^2 \end{aligned}$$

Remark on Conjugate Prior

Note that the posterior of σ^2 is again an Inverse Gamma distribution!

$$\underbrace{p(\sigma^2 \mid X, y, \beta)}_{=IG(a', b')} \propto \underbrace{p(y \mid X, \beta, \sigma^2)}_{\mathcal{N}(\mu, \nu)} \underbrace{p(\sigma^2)}_{=IG(a, b)}$$

This is because the IG is a **conjugate prior** to the normal distribution. Conjugate priors let you interpret how the data changes the believe about the parameters. \rightarrow Main reason for choosing this prior!

Remark on Conjugate Prior

Note that the posterior of σ^2 is again an Inverse Gamma distribution!

$$\underbrace{p(\sigma^2 \mid X, y, \beta)}_{=IG(a', b')} \propto \underbrace{p(y \mid X, \beta, \sigma^2)}_{\mathcal{N}(\mu, \nu)} \underbrace{p(\sigma^2)}_{=IG(a, b)}$$

This is because the IG is a **conjugate prior** to the normal distribution. Conjugate priors let you interpret how the data changes the believe about the parameters. \rightarrow Main reason for choosing this prior!

Remark: inverse distributions

Note that the Inverse Gamma distribution is called Inverse Gamma because

$$X \sim \Gamma(a, b) \iff X^{-1} \sim IG(a, b) \quad (1)$$

However, despite the name, the same is **not true** for the Inverse Gaussian!

M-step for β

We need to compute

$$\hat{\beta} = \arg \min_{\beta} \ell(\beta, \sigma^2, \tau^2) = \arg \min_{\beta} \frac{1}{2\sigma^2} \|y - X\beta\|_2^2 + \frac{1}{2}\beta^T T^{-1}\beta$$

where we dropped all terms independent of β . Then

$$\nabla_{\beta} \ell = 0 \iff \left(\frac{1}{\sigma^2} X^T X + T^{-1}\right) \hat{\beta} = \frac{1}{\sigma^2} X^T y$$

So $\boxed{\hat{\beta} = (X^T X + (\frac{1}{\sigma^2} T)^{-1})^{-1} X^T y}$ which is a ridge regression objective!

EM summary

1. Expectation of τ^2 :

$$p\left(\frac{1}{\tau_m^2} \mid \beta\right) = \text{Inv-Gauss}\left(\sqrt{\frac{\lambda^2}{\beta_m^2}}, \lambda^2\right)$$

$$\mathbb{E}\left[\frac{1}{\tau_m^2}\right] = \frac{\lambda}{|\beta_m|}$$

2. Expectation of σ^2 :

$$p(\sigma^2 \mid X, y, \beta) = \text{IG}(\sigma^2 \mid a', b')$$

$$a' := a + \frac{1}{2}N, \quad b' := b + \frac{1}{2}\|y - X\beta\|_2^2$$

$$\mathbb{E}\left[\frac{1}{\sigma^2}\right] = \frac{a'}{b'}$$

3. Maximization w.r.t. β :

$$\ell(\beta) = \frac{1}{2\sigma^2}\|y - X\beta\|_2^2 + \frac{1}{2}\beta^T T^{-1}\beta$$

$$\hat{\beta} = (X^T X + (\frac{1}{\sigma^2} T)^{-1})^{-1} X^T y$$

Why Laplace Prior?

- ▶ Bayesian Lasso
 - ▶ provides posterior distribution, not just point estimates
- ▶ Can be generalized to other models / losses
- ▶ Motivates to experiment with other types of priors, too
- ▶ Less scalable than the other methods, though.

Further Readings

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

References



Christopher M. Bishop.
Pattern recognition and machine learning, volume 1.
springer New York, 2006.



Amir Beck and Marc Teboulle.
A fast iterative shrinkage-thresholding algorithm for linear inverse problems.
SIAM Journal on Imaging Sciences, 2(1):183–202, 2009.



Torbjørn Eltoft, Taesu Kim, and Te-Won Lee.
On the multivariate laplace distribution.
IEEE Signal Processing Letters, 13(5):300–303, 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.