

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

もうしん 世間 ふゆうえやく 西マネー

Outline



1. Homotopy Methods: Least Angle Regression

2. Proximal Gradient Methods

3. Laplace Priors

《日》《聞》《臣》《臣》 通言 '오오오

Syllabus



A. Advanced Supervised Learning

Tue. 9.12.	(1)	A.1 Generalized Linear Models
Wed. 10.12.	(2)	A.2 Gaussian Processes
Tue. 16.12.	(3)	A.3 Advanced Support Vector Machines
Wed. 17.12.	(4)	A.4 Neural Networks
Tue. 6.1.	(5)	A.5 Ensembles
Wed. 7.1.	(6)	A.5b Ensembles (ctd.)
Tue. 13.1.	(7)	A.6 Sparse Linear Models
Wed. 14.1.	(8)	
Tue. 20.1.	(9)	
Wed. 21.1.	(10)	
Tue. 27.1.	(11)	
Wed. 28.1.	(12)	
Tue. 3.2.	(13)	
Wed. 4.2.	(14)	

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

<ロト < @ ト < E ト < E ト 三国 のへで</p>

Outline



1. Homotopy Methods: Least Angle Regression

2. Proximal Gradient Methods

3. Laplace Priors

・日本・西本・山田・山田・山田・今日・

Sparse Models so far



- Variable subset selection
 - forward search, backward search
- ► L1 regularization / Lasso
 - Coordinate descent (shooting algorithm)

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

◆□▶ ◆□▶ ◆□▶ ◆□▶ ◆□▶ ◆□

L1 Regularization



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

 $\hat{\theta} \in \mathbb{R}^P$

is equivalent to

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

 $||\hat{\theta}||_1 \leq B$
 $\hat{\theta} \in \mathbb{R}^P$

with

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

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

・ロト・4日×・4日× 4日× 900

Homotopy Methods



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

or equivalently

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

 $||\hat{\theta}||_1 \le B$

- start with a solution for large λ⁽⁰⁾ (or equiv. B := 0)
 then θ̂⁽⁰⁾ = 0.
- stepwise decrease $\lambda^{(t)}$ (or equiv. increase B)
 - learn $\hat{\theta}^{(t)}$ starting from $\hat{\theta}^{(t-1)}$ (warmstart).

・ロト・四ト・モー・ 相下 今々で

Least Angle Regression (LAR)

in step t:



2. regress these predictors on the residuum:

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

$$\hat{\gamma}^{(t)} := \operatorname*{arg\,min}_{\gamma} ||y - \hat{y}^{(t-1)} - X^{(t)}\gamma||_{2}$$

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

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, $\mathbb{P} \times \mathbb{R} \to \mathbb{R} \to \mathbb{R} \to \mathbb{R} \to \mathbb{R}$ Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany



5 / 31



Least Angle Regression (LAR): step length

Residuum correlations after the update

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

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

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. Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

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



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

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

or for negative correlations:

γi

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

$$\alpha = \frac{C_{\max}^{(t-1)} + C_{m}^{(t-1)}}{C_{\max}^{(t-1)} + X_{.,m}^{T} X^{(t)} \hat{\gamma}^{(t)}}$$
elding
$$\alpha := \min \left\{ \left(\frac{C_{\max}^{(t-1)} - C_{m}^{(t-1)}}{C_{\max}^{(t-1)} - X_{.,m}^{T} X^{(t)} \hat{\gamma}^{(t)}} \right)_{0}, \left(\frac{C_{\max}^{(t-1)} + C_{m}^{(t-1)}}{C_{\max}^{(t-1)} + X_{.,m}^{T} X^{(t)} \hat{\gamma}^{(t)}} \right)_{0} | m \in \{1, \dots, M\} \setminus A^{(t)} \}$$





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.

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

8 / 31

[HTFF05, p. 75]

Example





Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

Sac

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 β_m^(t) for an active predictor m becomes 0 or changes sign, drop it from the active set.

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

◆□▶ ◆□▶ ★∃▶ ★∃▶ ★目★ 少々で

Example









Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

[HTFF05, p. 75]

@ ▶

Outline



1. Homotopy Methods: Least Angle Regression

2. Proximal Gradient Methods

3. Laplace Priors

うせん 判所 《田》《田》《田》《日》



$$\operatorname{prox}_{f}(x^{0}) := \arg\min_{x} f(x) + \frac{1}{2} ||x - x^{0}||_{2}^{2}$$

- 「日本 本語 * 本語 * 本語 * 人口 * 、



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

$$\operatorname{prox}_{f}(x^{0}) := \arg\min_{x} f(x) + \frac{1}{2} ||x - x^{0}||_{2}^{2}$$

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

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

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

◆□▶ ◆□▶ ◆□▶ ◆□▶ ◆□▶ ◆□



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

$$\operatorname{prox}_{f}(x^{0}) := \arg\min_{x} f(x) + \frac{1}{2} ||x - x^{0}||_{2}^{2}$$

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

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

•
$$f := \lambda ||x||_1$$
 :

$$\begin{aligned} \mathsf{prox}_f(x^0) = & \mathsf{soft}(x^0, \lambda) := (\mathsf{soft}(x^0_n, \lambda))_{n=1,\dots,N} \\ & \mathsf{soft}(z, \lambda) := \mathsf{sign}(z)(|z| - \lambda)_0 \end{aligned}$$

・ 日 ト 《 母 ト 《 臣 ト 《 臣 ト 《 国 ト 《 の ヘ



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

$$\operatorname{prox}_{f}(x^{0}) := \arg\min_{x} f(x) + \frac{1}{2} ||x - x^{0}||_{2}^{2}$$

Can be solved analytically for some typical (possibly non-differentiable) regularization functions: 1 0

- $f := \lambda ||x||_2^2$: $\operatorname{prox}_f(x^0) = \frac{1}{2\lambda + 1} x^0$
- ► $f := \lambda ||x||_1$: $\operatorname{prox}_f(x^0) = \operatorname{soft}(x^0, \lambda) := (\operatorname{soft}(x^0, \lambda))_{n=1,...,N}$ $\operatorname{soft}(z, \lambda) := \operatorname{sign}(z)(|z| - \lambda)_0$

►
$$f := \lambda ||x||_0$$
:
 $\operatorname{prox}_f(x^0) = \operatorname{hard}(x^0, \lambda) := (\operatorname{hard}(x^0_n, \lambda))_{n=1,...,N},$
 $\operatorname{hard}(z, \lambda) := \delta(|z| \ge \lambda) z$



f



More Analytical Solutions for the Proximal Problem

▶ find x with minimal f in a vicinity of a given x⁰:

$$\operatorname{prox}_{f}(x^{0}) := \operatorname{arg\,min}_{x} f(x) + \frac{1}{2} ||x - x^{0}||_{2}^{2}$$
$$I := I_{C} \text{ for a convex set } C \text{ and } I_{C}(x) := \begin{cases} 0, & \text{if } x \in C \\ \infty, & \text{else} \end{cases}$$
$$\operatorname{prox}_{f}(x^{0}) = \operatorname{arg\,min}_{x \in C} ||x - x^{0}||_{2}^{2} =: \operatorname{proj}_{C}(x^{0})$$



- More Analytical Solutions for the Proximal Problem
 - find x with minimal f in a vicinity of a given x^0 :

$$prox_f(x^0) := \arg\min_x f(x) + \frac{1}{2}||x - x^0||_2^2$$
$$:= I_C \text{ for a convex set } C \text{ and } I_C(x) := \begin{cases} 0, & \text{if } x \in C \\ \infty, & \text{else} \end{cases}$$
$$prox_f(x^0) = \arg\min_{x \in C} ||x - x^0||_2^2 =: \operatorname{proj}_C(x^0)$$

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

シック 三回 エヨットボット 白マ



- More Analytical Solutions for the Proximal Problem
 - find x with minimal f in a vicinity of a given x^0 :

$$prox_f(x^0) := \arg\min_x f(x) + \frac{1}{2}||x - x^0||_2^2$$
$$:= I_C \text{ for a convex set } C \text{ and } I_C(x) := \begin{cases} 0, & \text{if } x \in C \\ \infty, & \text{else} \end{cases}$$
$$prox_f(x^0) = \arg\min_{x \in C} ||x - x^0||_2^2 =: \operatorname{proj}_C(x^0)$$

- ► rectangles / box constraints $C := [l_1, u_1] \times [l_2, u_2] \times \cdots \times [l_N, u_N]$: $\operatorname{prox}_f(x^0) = \operatorname{clip}(x^0, C)$ with $\operatorname{clip}(x^0, C)_n := \min\{\max\{x_n^0, l_n\}, u_n\}$
- euclidean balls $C := \{x \mid ||x||_2 \le 1\}$:

$$\operatorname{prox}_{f}(x^{0}) = \begin{cases} \frac{x^{0}}{||x^{0}||_{2}}, & \text{if } ||x^{0}||_{2} > 1\\ x^{0}, & \text{else} \end{cases}$$



More Analytical Solutions for the Proximal Problem

▶ find x with minimal f in a vicinity of a given x⁰:

$$\operatorname{prox}_{f}(x^{0}) := \arg\min_{x} f(x) + \frac{1}{2} ||x - x^{0}||_{2}^{2}$$

 $f := I_C$ for

• L1 balls $C := \{x \mid ||x||_1 \le 1\}$:

$$\operatorname{prox}_{f}(x^{0}) = \begin{cases} \operatorname{soft}(x^{0}, \lambda), & \text{if } ||x^{0}||_{1} > 1\\ x^{0}, & \text{else} \end{cases}$$
$$\operatorname{for } \lambda \text{ with } \sum_{n=1}^{N} (|x_{n}^{0} - \lambda|)_{0} \stackrel{!}{=} 1 \end{cases}$$

シック 単則 《川々 《川々 《山々

Generalized Gradient Descent

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

Generalized Gradient Descent:

$$\begin{aligned} x^{(t+1)} &:= \operatorname{prox}_{\alpha^{(t)}h}(x^{(t)} - \alpha^{(t)}\nabla g(x^{(t)})) \\ \text{with } \operatorname{prox}_{f}(x^{0}) &:= \arg\min_{x} f(x) + \frac{1}{2}||x - x^{0}||^{2} \end{aligned}$$

- ► 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).

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

▲帰▶ ▲ヨ▶ ▲ヨ▶ ヨヨ ののの

Machine Learning 2 2. Proximal Gradient Methods

Application to Regularized Loss Minimization



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

- $\blacktriangleright~\ell$ loss, convex and differentiable
 - ▶ e.g., RSS.
- ► *R* regularization, convex, but possibly not differentiable

► e.g.,
$$||\theta||_1$$
 or $I_C(\theta) := \begin{cases} 0, & \theta \in C \\ \infty, & \text{else} \end{cases}$

《日》《圖》《臣》《臣》 副言 のへの

Machine Learning 2 2. Proximal Gradient Methods



Application to Regularized Loss Minimization Minimizing

$$\theta^{(t+1)} := \operatorname*{arg\,min}_{\theta} R(\theta) + \ell(\theta)$$

using a Taylor expansion around previous estimate $\theta^{(t)}$:

$$\ell(\theta^{(t+1)}) \approx \ell(\theta^{(t)}) + \nabla \ell(\theta^{(t)})^{\mathsf{T}}(\theta - \theta^{(t)}) + (\theta - \theta^{(t)})^{\mathsf{T}} H(\theta - \theta^{(t)})$$

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

$$\approx \ell(\theta^{(t)}) + \nabla \ell(\theta^{(t)})^{\mathsf{T}}(\theta - \theta^{(t)}) + \alpha^{(t)} ||\theta - \theta^{(t)}||_2^2$$

$$\propto \alpha^{(t)} ||\theta - (\theta^{(t)} - \frac{1}{\alpha^{(t)}} \nabla \ell(\theta^{(t)}))||_2^2$$

yields a proximal problem

$$\arg\min_{\theta} \frac{1}{2\alpha^{(t)}} R(\theta) + \frac{1}{2} ||\theta - (\theta^{(t)} - \frac{1}{\alpha^{(t)}} \nabla \ell(\theta^{(t)}))||_{2}^{2}$$
$$= \operatorname{prox}_{\frac{1}{2\alpha^{(t)}} R} (\theta^{(t)} - \frac{1}{\alpha^{(t)}} \nabla \ell(\theta^{(t)}))$$

Special Cases



$$\theta^{(t+1)} := \operatorname{prox}_{\frac{1}{2\alpha^{(t)}}R}(\theta^{(t)} - \frac{1}{\alpha^{(t)}}\nabla\ell(\theta^{(t)}))$$
$$= \operatorname{arg\,min}_{\theta} \frac{1}{2\alpha^{(t)}}R(\theta) + \frac{1}{2}||\theta - (\theta^{(t)} - \frac{1}{\alpha^{(t)}}\nabla\ell(\theta^{(t)}))||_{2}^{2}$$

1. R = 0 yields gradient descent:

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

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

$$\theta^{(t+1)} = \operatorname{proj}_{\mathcal{C}}(\theta^{(t)} - \frac{1}{\alpha^{(t)}} \nabla \ell(\theta^{(t)}))$$

- 《日》 《聞》 《臣》 《臣》 (四) 《 (1)

Machine Learning 2 2. Proximal Gradient Methods



Special Cases: Projected Gradient Descent



Special Cases



$$\begin{split} \theta^{(t+1)} &:= \operatorname{prox}_{\frac{1}{2\alpha^{(t)}}R}(\theta^{(t)} - \frac{1}{\alpha^{(t)}}\nabla\ell(\theta^{(t)})) \\ &= \operatorname{arg\,min}_{\theta} \frac{1}{2\alpha^{(t)}}R(\theta) + \frac{1}{2}||\theta - (\theta^{(t)} - \frac{1}{\alpha^{(t)}}\nabla\ell(\theta^{(t)}))||_{2}^{2} \end{split}$$

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

$$heta^{(t+1)} = \operatorname{soft}(heta^{(t)} - rac{1}{lpha^{(t)}}
abla \ell(heta^{(t)}), rac{\lambda}{2lpha^{(t)}})$$

もうてい 正則 ふかく ふやく (型を) とう

Stepsizes $\alpha^{(t)}$



$$\alpha^{(t)}I(\theta^{(t)} - \theta^{(t-1)}) \stackrel{!}{\approx} H(\theta^{(t)} - \theta^{(t-1)}) = \nabla \ell(\theta^{(t)}) - \nabla \ell(\theta^{(t-1)})$$
$$\alpha^{(t)} := \arg\min_{\alpha} ||\alpha^{(t)}(\theta^{(t)} - \theta^{(t-1)}) - (\nabla \ell(\theta^{(t)}) - \nabla \ell(\theta^{(t-1)}))$$
$$= \frac{(\theta^{(t)} - \theta^{(t-1)})^T (\nabla \ell(\theta^{(t)}) - \nabla \ell(\theta^{(t-1)}))}{(\theta^{(t)} - \theta^{(t-1)})^T (\theta^{(t)} - \theta^{(t-1)})}$$

called Barzilai-Borwein stepsize or spectral stepsize.

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

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

◆□▶ ◆□▶ ◆□▶ ◆□▶ ◆□▶ ◆□



Iterative Shrinkage and Thresholding Algorithm(ISTA)

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

Algorithm



Algorithm 13.2: Iterative Shrinkage-Thresholding Algorithm (ISTA) 1 Input: $\mathbf{X} \in \mathbb{R}^{N \times D}$, $\mathbf{y} \in \mathbb{R}^N$, parameters $\lambda \ge 0$, $M \ge 1$, 0 < s < 1; Initialize $\theta_0 = 0$, $\alpha = 1$, $\mathbf{r} = \mathbf{v}$, $\lambda_0 = \infty$; 3 repeat $\lambda_t = \max(s || \mathbf{X}^T \mathbf{r} ||_{\infty}, \lambda) // \text{Adapt the regularizer};$ 4 repeat 5 $\mathbf{g} = \nabla L(\boldsymbol{\theta});$ 6 $\mathbf{u} = \boldsymbol{\theta} - \frac{1}{2}\mathbf{g};$ 7 $\boldsymbol{\theta} = \operatorname{soft}(\mathbf{u}, \frac{\lambda_t}{\alpha});$ 8 Update α using BB stepsize in Equation 13.82 ; 9 **until** $f(\theta)$ increased too much within the past M steps; 10 $\mathbf{r} = \mathbf{y} - \mathbf{X}\boldsymbol{\theta}$ // Update residual ; 11 12 until $\lambda_t = \lambda$;

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

[Mur12, p. 446]



Nesterov's Accelerated Generalized Gradient Descent

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

Generalized Gradient Descent:

$$\begin{aligned} x^{(t+1)} &:= \operatorname{prox}_{\alpha^{(t)}h} \left(x^{(t)} + \frac{t-2}{t+1} (x^{(t)} - x^{(t-1)}) - \alpha^{(t)} \nabla g(x^{(t)}) \right) \\ \text{with } \operatorname{prox}_{f}(x^{0}) &:= \arg\min_{x} f(x) + \frac{1}{2} ||x - x^{0}||^{2} \end{aligned}$$

- 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.



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

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

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

using Nesterov's Accelerated Generalized Gradient Descent.

シック 正則 エル・エット きゃくしゃ

FISTA vs ISTA





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

[BT09, p. 19] <□▷ <률▷ < 토▷ < 토▷ 토▷ 도 ♡ < ♡

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

26 / 31

Outline



1. Homotopy Methods: Least Angle Regression

2. Proximal Gradient Methods

3. Laplace Priors

うせん 判所 《田》《田》《田》《日》

Laplace Priors correspond to L1 regularization

L2 regularization:

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

Gaussian priors:

$$p(y_n \mid x_n, \beta, \sigma^2) := \mathcal{N}(y_n \mid x_n^T \beta, \sigma^2)$$
$$p(\beta) := \mathcal{N}(\beta \mid 0, \tau^2)$$
$$\propto (\tau^2)^{-M/2} e^{-\frac{1}{2\tau^2}\beta^T \beta}$$

using negative loglikelihood as objective function:

$$f(\beta) := -\log p(y \mid X, \beta)$$

・ 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 日 > ・ 4 0 = ・ 4 0 = ・ 4 0 = ・ 4 0 = ・ 4 0 = ・ 4 0 = ・ 4 0 = \cdot 0 = \cdot

Laplace Priors correspond to L1 regularization

L2 regularization:

$$f(heta):=||y-Xeta||_2^2+\lambda||eta||_2^2 \qquad \qquad f(heta):=||y-Xeta||_2^2+\lambda||eta|$$

Gaussian priors:

Laplace priors:

L1 regularization:

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

$$p(\beta) := \mathcal{N}(\beta \mid 0, \tau^2) \qquad p(\beta_m) := \mathsf{Lap}(\beta_m \mid 0, 1/\gamma)$$

$$\propto (\tau^2)^{-M/2} e^{-\frac{1}{2\tau^2}\beta^T \beta} \qquad = \frac{\gamma}{2} e^{-\gamma \mid \beta_m \mid}$$

using negative loglikelihood as objective function:

$$f(\beta) := -\log p(y \mid X, \beta)$$



Machine Learning 2 3. Laplace Priors

Laplace as Gaussian Scale Mixture



$$\mathsf{Lap}(\beta_m \mid 0, 1/\gamma) = \int \mathcal{N}(\beta_m \mid 0, \tau_m^2) \mathsf{Exp}(\tau_m^2 \mid \frac{\gamma^2}{2}) \, d\tau_m^2$$

with exponential distribution

$$\mathsf{Exp}(x \mid \lambda) := \lambda e^{-\lambda x}$$

シック 비판 《파》《파》《曰》

Universität - Hildeshein

Laplace Prior as Gaussian Scale Mixture

$$p(y_n \mid x_n, \beta, \sigma^2) := \mathcal{N}(y_n \mid x_n^T \beta, \sigma^2)$$
$$p(\beta_m \mid \tau_m^2) := \mathcal{N}(\beta_m \mid 0, \tau_m^2)$$
$$p(\tau_m^2) := \mathsf{Exp}(\tau_m^2 \mid \frac{\gamma^2}{2})$$
$$p(\sigma^2) := \mathsf{IG}(\sigma^2 \mid \mathsf{a}_\sigma, \mathsf{b}_\sigma)$$

シック 三川州 エル・エル・エート



Laplace Prior as Gaussian Scale Mixture



Machine Learning 2 3. Laplace Priors



Laplace Prior as Gaussian Scale Mixture

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

$$p(\beta_m \mid \tau_m^2) := \mathcal{N}(\beta_m \mid 0, \tau_m^2)$$

$$p(\tau_m^2) := \mathsf{Exp}(\tau_m^2 \mid \frac{\gamma^2}{2})$$

$$p(\sigma^2) := \mathsf{IG}(\sigma^2 \mid a_\sigma, b_\sigma)$$
NLL:
$$f(\beta, \tau^2) = -\frac{1}{2\sigma^2} ||y - X\beta||_2^2 - \frac{1}{2}\beta^T \Lambda\beta + \frac{\gamma^2}{2} \sum_{m=1}^M \tau_m^2 \sum_{N=1}^N \lambda \beta = \mathsf{diag}(\frac{1}{\tau_1^2}, \frac{1}{\tau_2^2}, \dots, \frac{1}{\tau_M^2})$$

$$(\mathsf{Murl}_2, \mathsf{p}, \mathsf{446})$$

EM for Laplace Prior NLL:

$$\begin{split} f(\beta, \tau^2) &= -\frac{1}{2\sigma^2} ||y - X\beta||_2^2 - \frac{1}{2}\beta^T \Lambda\beta + \frac{\gamma^2}{2} \sum_{m=1}^M \tau_m^2 \\ \Lambda &:= \mathsf{diag}(\frac{1}{\tau_1^2}, \frac{1}{\tau_2^2}, \dots, \frac{1}{\tau_M^2}) \end{split}$$

1. Expectaction of τ^2 :

$$p(1/\tau_m^2 \mid \beta) = \text{InverseGaussian}(\sqrt{\frac{\gamma^2}{\beta_m^2}}, \gamma^2)$$
$$E(1/\tau_m^2 \mid \beta) = \frac{\gamma}{|\beta_m|}$$



・ロト < 団ト < ヨト < ヨト < ロト



EM for Laplace Prior NLL:



2. Expectaction of σ^2 :

$$p(\sigma^2 \mid \beta) = \mathsf{IG}(\mathsf{a}_{\sigma} + \frac{N}{2}, \mathsf{b}_{\sigma} + \frac{1}{2}(y - X\beta)^T(y - X\beta))$$
$$E(1/\sigma^2 \mid \beta) = \frac{\mathsf{a}_{\sigma} + \frac{N}{2}}{\mathsf{b}_{\sigma} + \frac{1}{2}(y - X\beta)^T(y - X\beta)}$$

・ロット (四マ・山田) (日) (日)



Juniversiter Hildeshein

EM for Laplace Prior

NLL:

$$\begin{split} f(\beta,\tau^2) &= -\frac{1}{2\sigma^2} ||y - X\beta||_2^2 - \frac{1}{2}\beta^T \Lambda\beta + \frac{\gamma^2}{2} \sum_{m=1}^M \tau_m^2 \\ \Lambda &:= \mathsf{diag}(\frac{1}{\tau_1^2}, \frac{1}{\tau_2^2}, \dots, \frac{1}{\tau_M^2}) \end{split}$$

3. Maximizing β : ridge regression

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

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

・ロト < 団ト < ヨト < ヨト < ロト

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.

・日・《四・《四・《四・《日・

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].

シック 正則 《田》《田》《日》

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.



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.

うせん 正則 ふばやえばや きしゃく