

HYPERPARAMETER SELECTION FOR HIGH DIMENSIONAL SPARSE LEARNING

WITH NEUROIMAGING MOTIVATIONS IN MIND

Joseph Salmon

IMAG, Univ Montpellier, CNRS
Institut Universitaire de France (IUF)



UNIVERSITÉ DE
MONTPELLIER



JOINT WORKS WITH VARIOUS COLLEAGUES



Quentin Bertrand (MILA)

Quentin Klopfenstein (Université du Luxembourg)

Mathurin Massias (INRIA, OCKHAM them)

Pierre-Antoine Bannier (M2 student, Parietal Team)

Samuel Vaïter (Université Côte d'Azur, CNRS)

Mathieu Blondel (Google Research, Brain team)

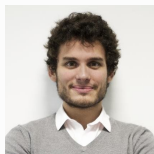
Alexandre Gramfort (INRIA, Parietal Team)



Quentin B.



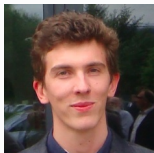
Quentin K.



Mathurin



Pierre-Antoine



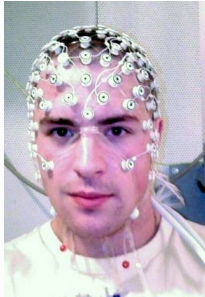
Mathieu



Samuel



Alexandre



(a) EEG



(b) MEG=Mag.+Grad.

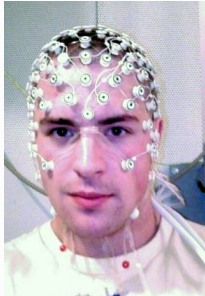


(c) M/EEG

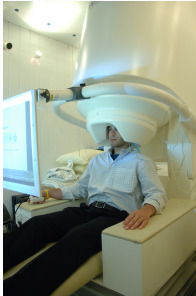
Photo credit: S. Whitmarsh

⁽¹⁾ H. Berger (1929). "Über das elektroenkephalogramm des menschen". In: *Archiv für psychiatrie und nervenkrankheiten*

⁽²⁾ D. Cohen (1968). "Magnetoencephalography: evidence of magnetic fields produced by alpha-rhythm currents". In: *Science*



(a) EEG



(b) MEG=Mag.+Grad.



(c) M/EEG

Photo credit: S. Whitmarsh

- ▶ **Data Y**: electric and magnetic fields at the head surface

⁽¹⁾ H. Berger (1929). "Über das elektroenkephalogramm des menschen". In: *Archiv für psychiatrie und nervenkrankheiten*

⁽²⁾ D. Cohen (1968). "Magnetoencephalography: evidence of magnetic fields produced by alpha-rhythm currents". In: *Science*



(a) EEG



(b) MEG=Mag.+Grad.



(c) M/EEG

Photo credit: S. Whitmarsh

- ▶ **Data** Y : electric and magnetic fields at the head surface
- ▶ **Goal**: which parts of the brain are responsible for the signals?

⁽¹⁾ H. Berger (1929). "Über das elektroenkephalogramm des menschen". In: *Archiv für psychiatrie und nervenkrankheiten*

⁽²⁾ D. Cohen (1968). "Magnetoencephalography: evidence of magnetic fields produced by alpha-rhythm currents". In: *Science*



(a) EEG



(b) MEG=Mag.+Grad.



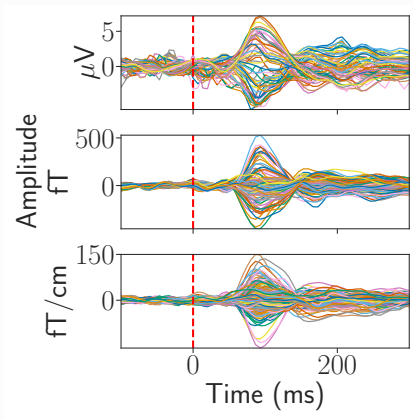
(c) M/EEG

Photo credit: S. Whitmarsh

- ▶ **Data Y :** electric and magnetic fields at the head surface
- ▶ **Goal:** which parts of the brain are responsible for the signals?
- ▶ **Applications:** clinical and cognitive experiments

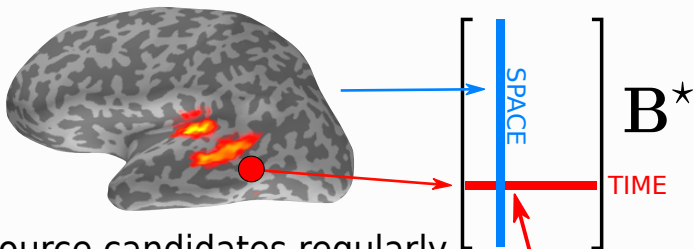
⁽¹⁾ H. Berger (1929). "Über das elektroenkephalogramm des menschen". In: *Archiv für psychiatrie und nervenkrankheiten*

⁽²⁾ D. Cohen (1968). "Magnetoencephalography: evidence of magnetic fields produced by alpha-rhythm currents". In: *Science*



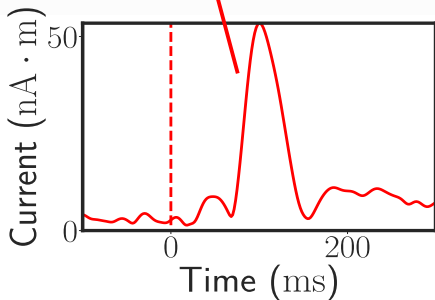
3 modalities:

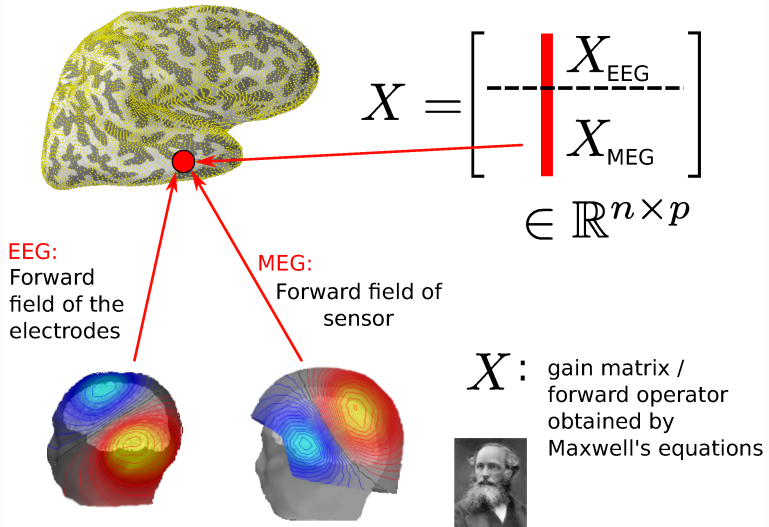
- ▶ EEG
- ▶ MEG: magnetometers (amplitude)
- ▶ MEG: gradiometers (gradients)



Source candidates regularly spaced in the brain (e.g., every 5mm)

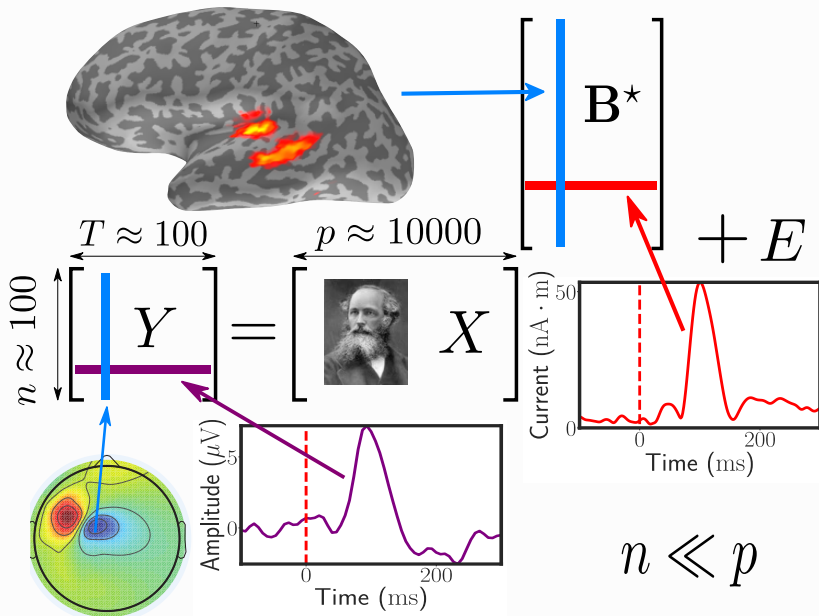
$$B^* \in \mathbb{R}^{p \times T}$$





THE M/EEG INVERSE PROBLEM

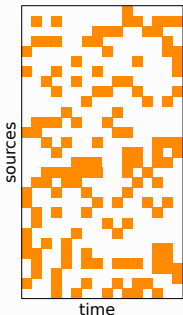
MAXWELL EQUATIONS AND (APPROX.) LINEARITY





Popular convex penalties:

$$\hat{\mathbf{B}} \in \arg \min_{\mathbf{B} \in \mathbb{R}^{p \times T}} \left(\frac{1}{2nT} \|\mathbf{Y} - \mathbf{X}\mathbf{B}\|^2 + \lambda \Omega(\mathbf{B}) \right)$$



Parameter $\hat{\mathbf{B}} \in \mathbb{R}^{p \times T}$

Sparse support: no structure

Penalty: **Lasso**

$$\Omega(\mathbf{B}) = \|\mathbf{B}\|_1 = \sum_{j=1}^p \sum_{k=1}^T |\mathbf{B}_{j,k}|$$

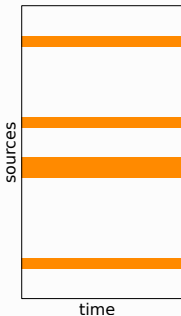
⁽¹⁾ A. Argyriou, T. Evgeniou, and M. Pontil (2008). "Convex multi-task feature learning". In: *Machine Learning* 73.3, pp. 243–272.

⁽²⁾ A. Gramfort, M. Kowalski, and M. Hämmäläinen (2012). "Mixed-norm estimates for the M/EEG inverse problem using accelerated gradient methods". In: *Phys. Med. Biol.* 57.7, pp. 1937–1961



Popular convex penalties: multitask Lasso (MTL)

$$\hat{\mathbf{B}} \in \arg \min_{\mathbf{B} \in \mathbb{R}^{p \times T}} \left(\frac{1}{2nT} \|\mathbf{Y} - \mathbf{X}\mathbf{B}\|^2 + \lambda \Omega(\mathbf{B}) \right)$$



Parameter $\hat{\mathbf{B}} \in \mathbb{R}^{p \times T}$

Sparse support: group structure ✓

Penalty: **Group-Lasso**

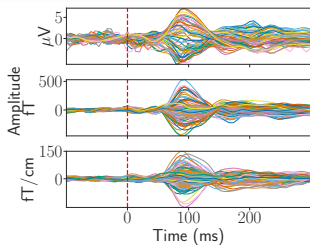
$$\Omega(\mathbf{B}) = \|\mathbf{B}\|_{2,1} = \sum_{j=1}^p \|\mathbf{B}_{j,:}\|_2$$

where $\mathbf{B}_{j,:}$: the j -th row of \mathbf{B}

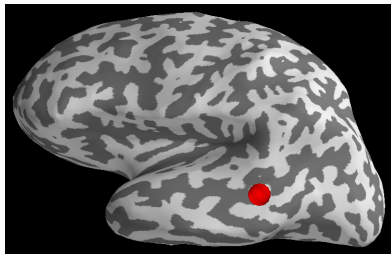
⁽¹⁾ A. Argyriou, T. Evgeniou, and M. Pontil (2008). "Convex multi-task feature learning". In: *Machine Learning* 73.3, pp. 243–272.

⁽²⁾ A. Gramfort, M. Kowalski, and M. Hämmäläinen (2012). "Mixed-norm estimates for the M/EEG inverse problem using accelerated gradient methods". In: *Phys. Med. Biol.* 57.7, pp. 1937–1961

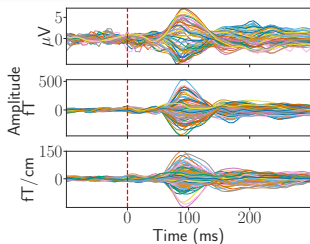
SUMMARY OF THE PROBLEM SETTING



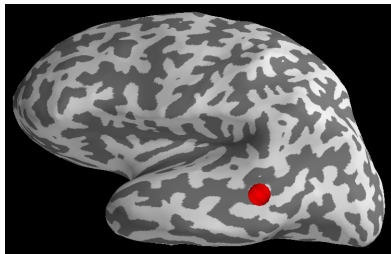
What you have: $Y \in \mathbb{R}^{n \times T}$



What you want: $B \in \mathbb{R}^{p \times T}$



What you have: $Y \in \mathbb{R}^{n \times T}$



What you want: $B \in \mathbb{R}^{p \times T}$

This is typically done using optimization based estimators

$$\hat{B} \in \arg \min_{B \in \mathbb{R}^{p \times T}} \left(\frac{1}{2nT} \|Y - XB\|_F^2 + \lambda \Omega(B) \right)$$



$$\hat{\mathbf{B}} \in \arg \min_{\mathbf{B} \in \mathbb{R}^{p \times T}} \left(\frac{1}{2nT} \|Y - X\mathbf{B}\|_F^2 + \lambda \Omega(\mathbf{B}) \right)$$

Covered in this presentation

- ▶ How to efficiently select the regularization parameter λ ?^{(1), (2)}

Not covered in this presentation

- ▶ How to efficiently solve this optimization problem?⁽³⁾
- ▶ How to handle spatial correlation?⁽⁴⁾

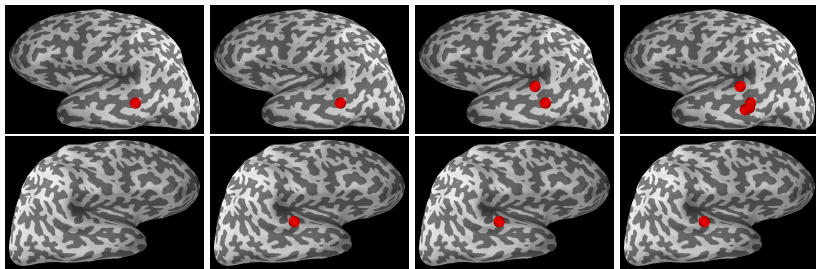
⁽¹⁾ Q Bertrand, Q Klopfenstein, M. Blondel, et al. (2020). "Implicit differentiation of Lasso-type models for hyperparameter optimization". In: *ICML*.

⁽²⁾ Q Bertrand, Q Klopfenstein, M. Massias, et al. (2022). "Implicit differentiation for fast hyperparameter selection in non-smooth convex learning". In: *Submitted to JMLR*.

⁽³⁾ Q Bertrand and M. Massias (2021). "Anderson acceleration of coordinate descent". In: *AISTATS*.

⁽⁴⁾ Q Bertrand, M. Massias, et al. (2019). "Handling correlated and repeated measurements with the smoothed multivariate square-root Lasso". In: *NeurIPS*.

$$\hat{\mathbf{B}} \in \arg \min_{\mathbf{B} \in \mathbb{R}^{p \times T}} \left(\frac{1}{2nT} \|\mathbf{Y} - \mathbf{X}\mathbf{B}\|_F^2 + \lambda \|\mathbf{B}\|_{2,1} \right)$$



$\lambda = 0.85\lambda_{\max}$

$\lambda = 0.82\lambda_{\max}$

$\lambda = 0.80\lambda_{\max}$

$\lambda = 0.75\lambda_{\max}$

Real M/EEG data. Brain source reconstruction using multitask Lasso with multiple λ . Which λ to pick? How to *automatically* select λ ?

- ▶ When $\lambda \geq \lambda_{\max}$, $\hat{\mathbf{B}} = \mathbf{0}$ no sources are recovered



- ▶ Statistical route^{(1), (2)}:
use assumptions on X , provide guarantees but often conservative
- ▶ Bayesian statistics^{(3), (4)}: prior on λ
- ▶ Bayesian optimization,⁽⁵⁾ 0-th order method:

The road today:

- ▶ Hyperparameter optimization⁽⁶⁾: minimize a given criterion $\mathcal{C}(\hat{\beta}(\lambda))$

(1) K. Lounici (2008). "Sup-norm convergence rate and sign concentration property of Lasso and Dantzig estimators". In: *Electron. J. Stat.* 2, pp. 90–102.

(2) K. Lounici, M. Pontil, et al. (2011). "Oracle inequalities and optimal inference under group sparsity". In: *Ann. Statist.* 39.4, pp. 2164–2204.

(3) M. E. Tipping (2001). "Sparse Bayesian learning and the relevance vector machine". In: *J. Mach. Learn. Res.* 1, pp. 211–244.

(4) M. Figueiredo (2001). "Adaptive Sparseness Using Jeffreys Prior". In: *NIPS*, pp. 697–704.

(5) F. Hutter, J. Lücke, and L. Schmidt-Thieme (2015). "Beyond Manual Tuning of Hyperparameters". In: *Künstliche Intell.* 29.4, pp. 329–337.

(6) R. Kohavi and G. H. John (1995). "Automatic parameter selection by minimizing estimated error". In: *ICML*, pp. 304–312.

Possible selection criterion:

- ▶ Good generalization^{(1), (2)} of $\hat{\beta}(\lambda)$
- ▶ AIC/BIC,⁽³⁾ SURE⁽⁴⁾ that controls model complexity

⁽¹⁾ L. R. A. Stone and J.C. Ramer (1965). "Estimating WAIS IQ from Shipley Scale scores: Another cross-validation". In: *Journal of clinical psychology* 21.3, pp. 297–297.

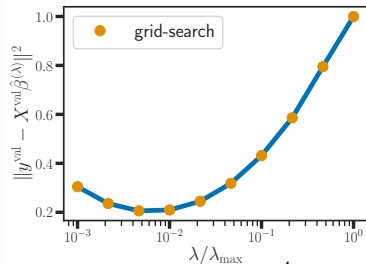
⁽²⁾ K. Lounici, K. Meziani, and B. Riu (2021). "Muddling Labels for Regularization, a novel approach to generalization". In: *arXiv preprint arXiv:2102.08769*.

⁽³⁾ W. Liu and Y. Yang (2011). "Parametric or nonparametric? A parametricness index for model selection". In: *Ann. Statist.* 39.4, pp. 2074–2102.

⁽⁴⁾ C. M. Stein (1981). "Estimation of the mean of a multivariate normal distribution". In: *Ann. Statist.* 9.6, pp. 1135–1151.

Possible selection criterion:

- ▶ Good generalization^{(1), (2)} of $\hat{\beta}(\lambda)$
- ▶ AIC/BIC,⁽³⁾ SURE⁽⁴⁾ that controls model complexity



Real-sim dataset, $n \approx p \approx 10^4$
 Validation loss as a function of λ .

Simplified example ($T = 1$):

Model: Lasso

$$\hat{\beta}(\lambda) \in \arg \min_{\beta \in \mathbb{R}^p} \frac{\|y^{\text{train}} - X^{\text{train}}\beta\|^2}{2n} + \lambda \|\beta\|_1$$

Criterion: held-out loss

$$\arg \min_{\lambda} \|y^{\text{test}} - X^{\text{test}}\hat{\beta}(\lambda)\|^2$$

(1) L. R. A. Stone and J.C. Ramer (1965). "Estimating WAIS IQ from Shipley Scale scores: Another cross-validation". In: *Journal of clinical psychology* 21.3, pp. 297–297.

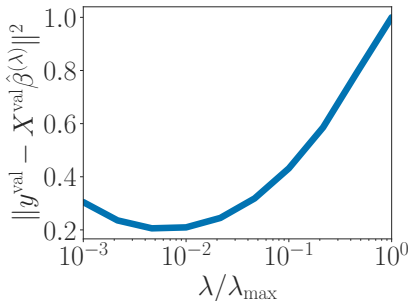
(2) K. Lounici, K. Meziani, and B. Riu (2021). "Muddling Labels for Regularization, a novel approach to generalization". In: *arXiv preprint arXiv:2102.08769*.

(3) W. Liu and Y. Yang (2011). "Parametric or nonparametric? A parametricness index for model selection". In: *Ann. Statist.* 39.4, pp. 2074–2102.

(4) C. M. Stein (1981). "Estimation of the mean of a multivariate normal distribution". In: *Ann. Statist.* 9.6, pp. 1135–1151.



$$\begin{array}{l}
 \text{outer optimization problem} \\
 \arg \min_{\lambda \in \mathbb{R}} \left\{ \mathcal{L}(\lambda) := \|y^{\text{test}} - X^{\text{test}} \hat{\beta}(\lambda)\|^2 \right\} \\
 \text{s.t. } \hat{\beta}(\lambda) \in \underbrace{\arg \min_{\beta \in \mathbb{R}^p} \frac{1}{2n} \|y^{\text{train}} - X^{\text{train}} \beta\|^2 + \lambda \|\beta\|_1}_{\text{inner optimization problem}}
 \end{array}$$



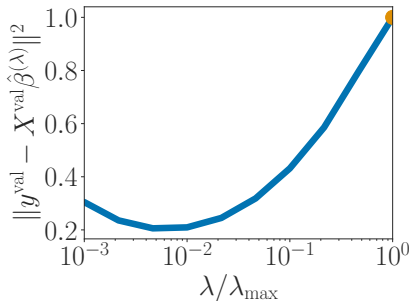
Grid search evaluation
(sequential)

⁽¹⁾ P. Ochs et al. (2015). "Bilevel optimization with nonsmooth lower level problems". In: *SSVM*. vol. 9087, pp. 654–665

⁽²⁾ F. Pedregosa (2016). "Hyperparameter optimization with approximate gradient". In: *ICML*. vol. 48, pp. 737–746



$$\begin{aligned}
 & \text{outer optimization problem} \\
 & \arg \min_{\lambda \in \mathbb{R}} \left\{ \mathcal{L}(\lambda) := \|y^{\text{test}} - X^{\text{test}} \hat{\beta}(\lambda)\|^2 \right\} \\
 & \text{s.t. } \hat{\beta}(\lambda) \in \arg \min_{\beta \in \mathbb{R}^p} \underbrace{\frac{1}{2n} \|y^{\text{train}} - X^{\text{train}} \beta\|^2 + \lambda \|\beta\|_1}_{\text{inner optimization problem}}
 \end{aligned}$$



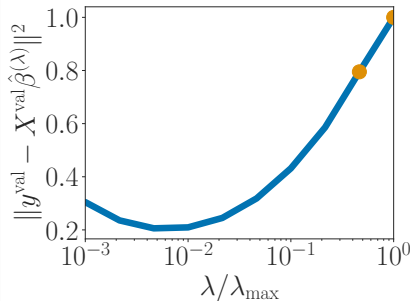
Grid search evaluation
(sequential)

⁽¹⁾ P. Ochs et al. (2015). "Bilevel optimization with nonsmooth lower level problems". In: *SSVM*. vol. 9087, pp. 654–665

⁽²⁾ F. Pedregosa (2016). "Hyperparameter optimization with approximate gradient". In: *ICML*. vol. 48, pp. 737–746



$$\begin{array}{l}
 \text{outer optimization problem} \\
 \arg \min_{\lambda \in \mathbb{R}} \left\{ \mathcal{L}(\lambda) := \|y^{\text{test}} - X^{\text{test}} \hat{\beta}(\lambda)\|^2 \right\} \\
 \text{s.t. } \hat{\beta}(\lambda) \in \arg \min_{\beta \in \mathbb{R}^p} \underbrace{\frac{1}{2n} \|y^{\text{train}} - X^{\text{train}} \beta\|^2 + \lambda \|\beta\|_1}_{\text{inner optimization problem}}
 \end{array}$$



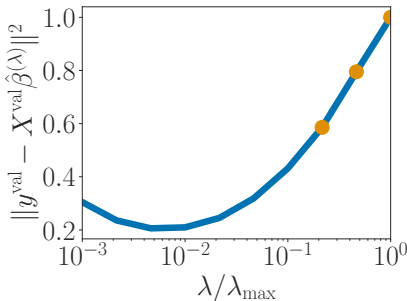
Grid search evaluation
(sequential)

⁽¹⁾ P. Ochs et al. (2015). "Bilevel optimization with nonsmooth lower level problems". In: *SSVM*. vol. 9087, pp. 654–665

⁽²⁾ F. Pedregosa (2016). "Hyperparameter optimization with approximate gradient". In: *ICML*. vol. 48, pp. 737–746



$$\begin{aligned}
 & \text{outer optimization problem} \\
 & \arg \min_{\lambda \in \mathbb{R}} \left\{ \mathcal{L}(\lambda) := \|y^{\text{test}} - X^{\text{test}} \hat{\beta}(\lambda)\|^2 \right\} \\
 & \text{s.t. } \hat{\beta}(\lambda) \in \arg \min_{\beta \in \mathbb{R}^p} \underbrace{\frac{1}{2n} \|y^{\text{train}} - X^{\text{train}} \beta\|^2 + \lambda \|\beta\|_1}_{\text{inner optimization problem}}
 \end{aligned}$$



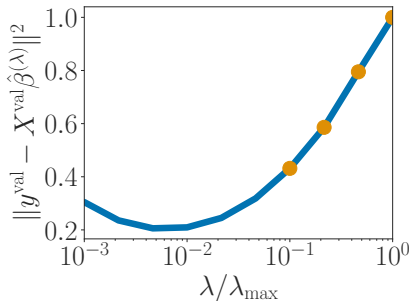
Grid search evaluation
(sequential)

⁽¹⁾ P. Ochs et al. (2015). "Bilevel optimization with nonsmooth lower level problems". In: *SSVM*. vol. 9087, pp. 654–665

⁽²⁾ F. Pedregosa (2016). "Hyperparameter optimization with approximate gradient". In: *ICML*. vol. 48, pp. 737–746



$$\begin{array}{l}
 \text{outer optimization problem} \\
 \arg \min_{\lambda \in \mathbb{R}} \left\{ \mathcal{L}(\lambda) := \|y^{\text{test}} - X^{\text{test}} \hat{\beta}(\lambda)\|^2 \right\} \\
 \text{s.t. } \hat{\beta}(\lambda) \in \underbrace{\arg \min_{\beta \in \mathbb{R}^p} \frac{1}{2n} \|y^{\text{train}} - X^{\text{train}} \beta\|^2 + \lambda \|\beta\|_1}_{\text{inner optimization problem}}
 \end{array}$$



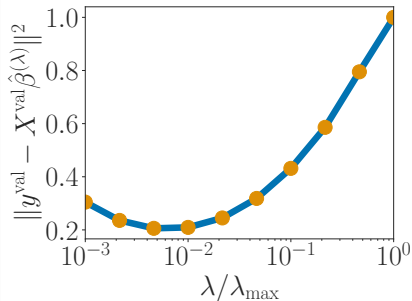
Grid search evaluation
(sequential)

⁽¹⁾ P. Ochs et al. (2015). "Bilevel optimization with nonsmooth lower level problems". In: *SSVM*. vol. 9087, pp. 654–665

⁽²⁾ F. Pedregosa (2016). "Hyperparameter optimization with approximate gradient". In: *ICML*. vol. 48, pp. 737–746



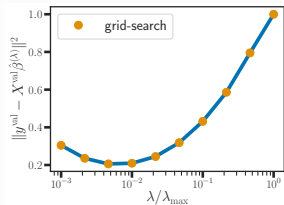
$$\begin{array}{l}
 \text{outer optimization problem} \\
 \arg \min_{\lambda \in \mathbb{R}} \left\{ \mathcal{L}(\lambda) := \|y^{\text{test}} - X^{\text{test}} \hat{\beta}(\lambda)\|^2 \right\} \\
 \text{s.t. } \hat{\beta}(\lambda) \in \arg \min_{\beta \in \mathbb{R}^p} \underbrace{\frac{1}{2n} \|y^{\text{train}} - X^{\text{train}} \beta\|^2 + \lambda \|\beta\|_1}_{\text{inner optimization problem}}
 \end{array}$$



Grid search evaluation
(sequential)

⁽¹⁾ P. Ochs et al. (2015). "Bilevel optimization with nonsmooth lower level problems". In: *SSVM*. vol. 9087, pp. 654–665

⁽²⁾ F. Pedregosa (2016). "Hyperparameter optimization with approximate gradient". In: *ICML*. vol. 48, pp. 737–746



$$\arg \min_{\lambda \in \mathbb{R}} \left\{ \mathcal{L}(\lambda) := \|y^{\text{test}} - X^{\text{test}} \hat{\beta}(\lambda)\|^2 \right\}$$

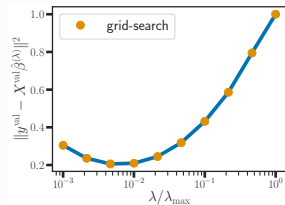
$$\text{s.t. } \hat{\beta}(\lambda) \in \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{2n} \|y^{\text{train}} - X^{\text{train}} \beta\|^2 + \lambda \|\beta\|_1$$

- Grid-search, random-search,⁽¹⁾ SMBO⁽²⁾:
0-order methods to solve bilevel optimization problem

⁽¹⁾ J. Bergstra and Y. Bengio (2012). "Random search for hyper-parameter optimization". In: *J. Mach. Learn. Res.* 13.2.

⁽²⁾ E. Brochu, V. M. Cora, and N. De Freitas (2010). *A tutorial on Bayesian optimization of expensive cost functions, with application to active user modeling and hierarchical reinforcement learning*. Tech. rep.

⁽³⁾ F. Pedregosa (2016). "Hyperparameter optimization with approximate gradient". In: *ICML*. vol. 48, pp. 737–746.



$$\arg \min_{\lambda \in \mathbb{R}} \left\{ \mathcal{L}(\lambda) := \|y^{\text{test}} - X^{\text{test}} \hat{\beta}(\lambda)\|^2 \right\}$$

$$\text{s.t. } \hat{\beta}(\lambda) \in \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{2n} \|y^{\text{train}} - X^{\text{train}} \beta\|^2 + \lambda \|\beta\|_1$$

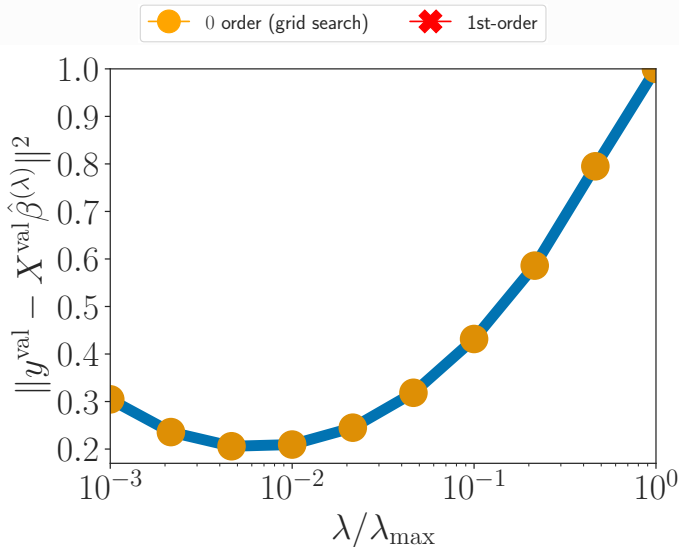
- ▶ Grid-search, random-search,⁽¹⁾ SMBO⁽²⁾:
0-order methods to solve bilevel optimization problem
- ▶ **Idea:** if \mathcal{L} is differentiable, use 1st-order optimization
 - ▶ Compute gradient: $\nabla_{\lambda} \mathcal{L}$
 - ▶ Perform gradient descent step⁽³⁾:

$$\lambda^{(t+1)} = \lambda^{(t)} - \rho \nabla_{\lambda} \mathcal{L}(\lambda^{(t)}) \quad \text{with } \rho > 0$$

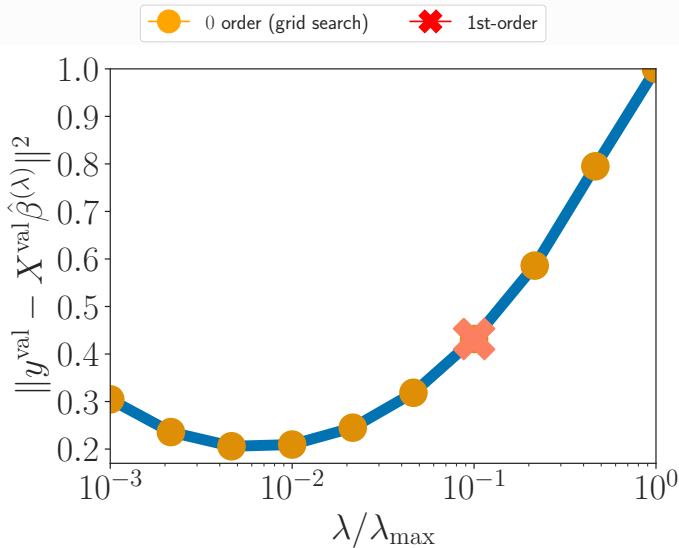
⁽¹⁾ J. Bergstra and Y. Bengio (2012). "Random search for hyper-parameter optimization". In: *J. Mach. Learn. Res.* 13.2.

⁽²⁾ E. Brochu, V. M. Cora, and N. De Freitas (2010). *A tutorial on Bayesian optimization of expensive cost functions, with application to active user modeling and hierarchical reinforcement learning*. Tech. rep.

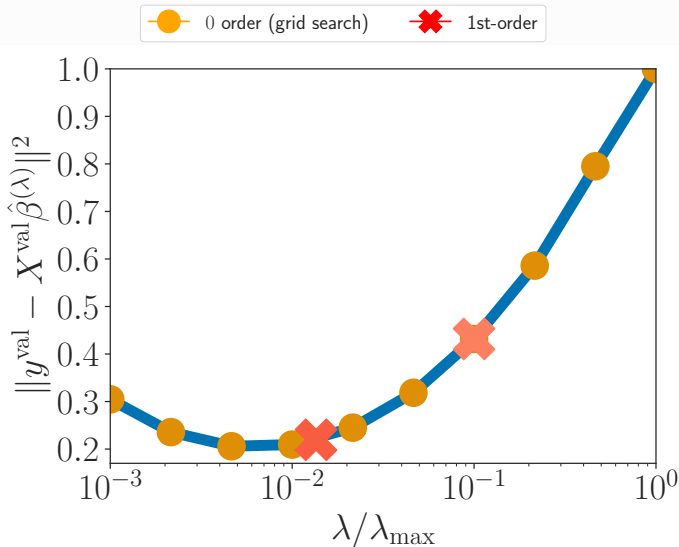
⁽³⁾ F. Pedregosa (2016). "Hyperparameter optimization with approximate gradient". In: *ICML*. vol. 48, pp. 737–746.



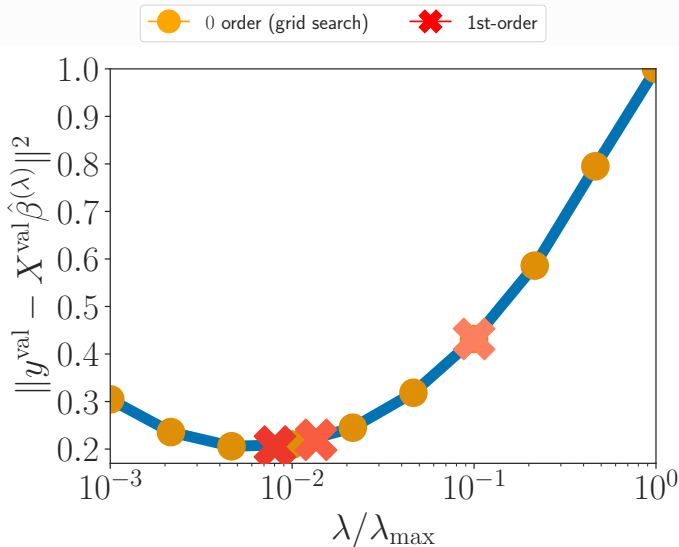
Real-sim dataset, $n \approx p \approx 10^4$. Validation loss as a function of λ .



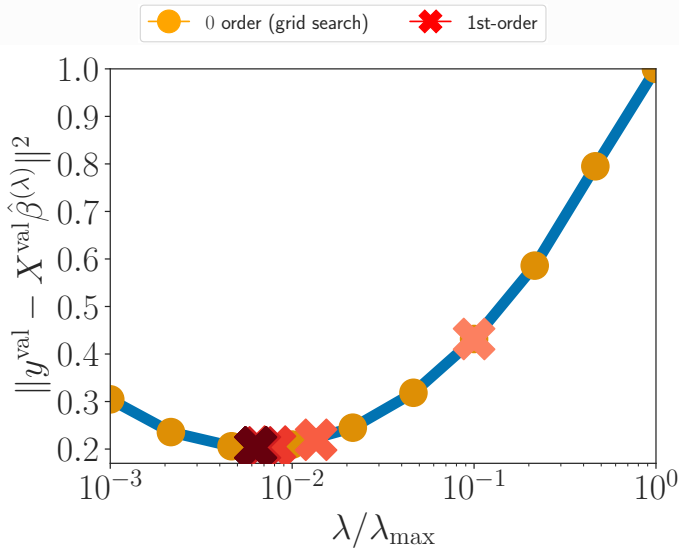
Real-sim dataset, $n \approx p \approx 10^4$. Validation loss as a function of λ .



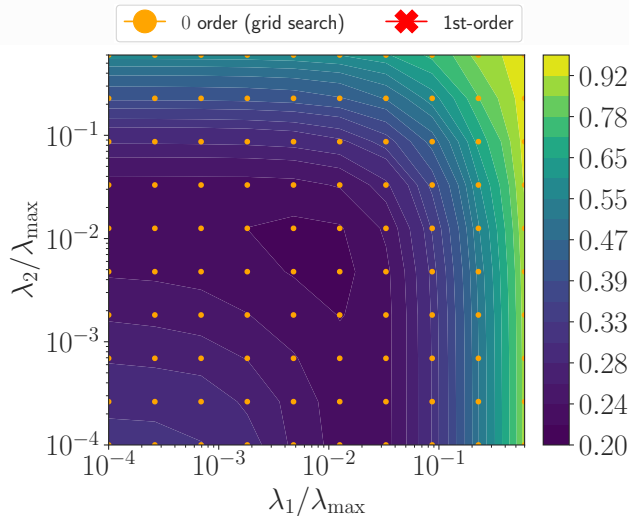
Real-sim dataset, $n \approx p \approx 10^4$. Validation loss as a function of λ .



Real-sim dataset, $n \approx p \approx 10^4$. Validation loss as a function of λ .

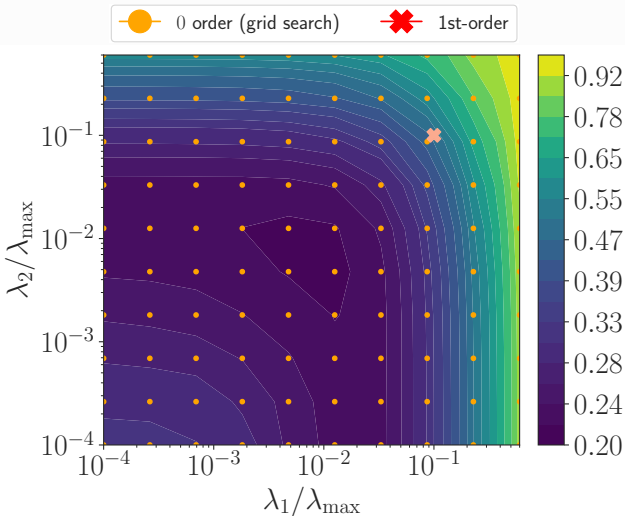


Real-sim dataset, $n \approx p \approx 10^4$. Validation loss as a function of λ .



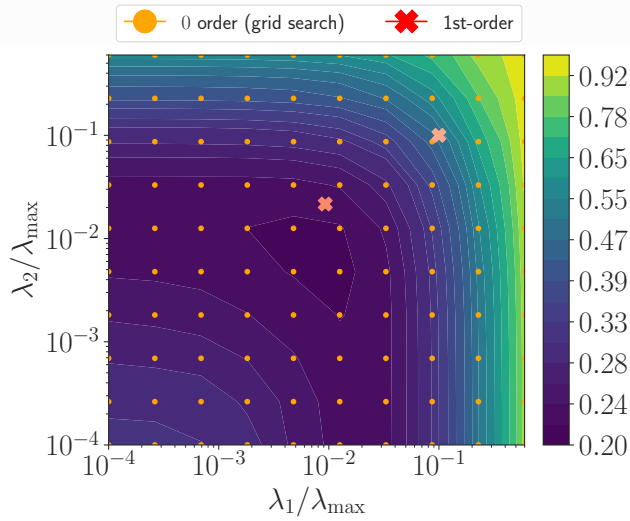
Real-sim dataset, level sets of the validation loss (hold-out)

$$\arg \min_{\beta} \frac{1}{2n} \|y^{\text{train}} - X^{\text{train}} \beta\|^2 + \lambda_1 \|\beta\|_1 + \frac{\lambda_2}{2} \|\beta\|^2$$



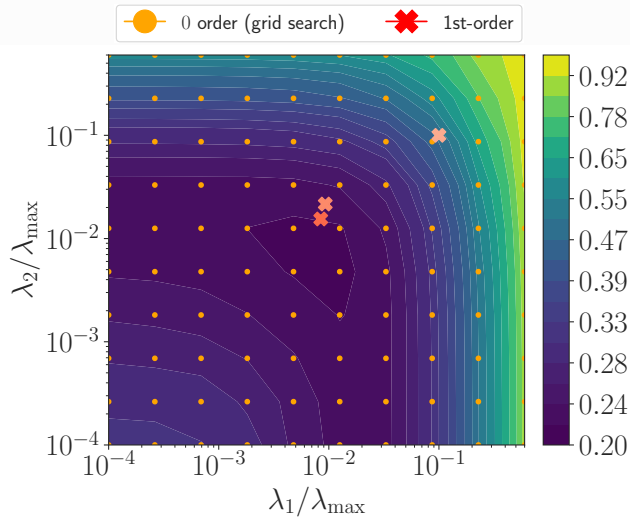
Real-sim dataset, level sets of the validation loss (hold-out)

$$\arg \min_{\beta} \frac{1}{2n} \|y^{\text{train}} - X^{\text{train}} \beta\|^2 + \lambda_1 \|\beta\|_1 + \frac{\lambda_2}{2} \|\beta\|^2$$



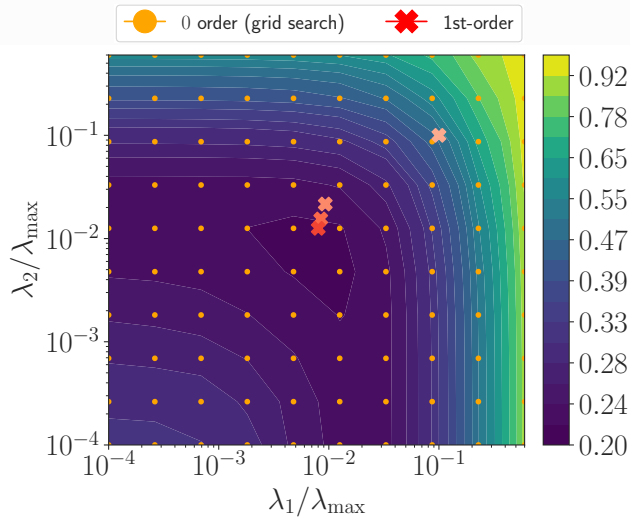
Real-sim dataset, level sets of the validation loss (hold-out)

$$\arg \min_{\beta} \frac{1}{2n} \|y^{\text{train}} - X^{\text{train}} \beta\|^2 + \lambda_1 \|\beta\|_1 + \frac{\lambda_2}{2} \|\beta\|^2$$



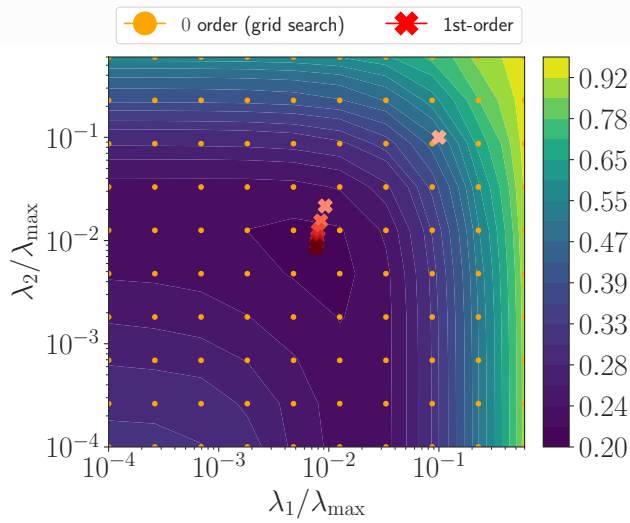
Real-sim dataset, level sets of the validation loss (hold-out)

$$\arg \min_{\beta} \frac{1}{2n} \|y^{\text{train}} - X^{\text{train}} \beta\|^2 + \lambda_1 \|\beta\|_1 + \frac{\lambda_2}{2} \|\beta\|^2$$



Real-sim dataset, level sets of the validation loss (hold-out)

$$\arg \min_{\beta} \frac{1}{2n} \|y^{\text{train}} - X^{\text{train}} \beta\|^2 + \lambda_1 \|\beta\|_1 + \frac{\lambda_2}{2} \|\beta\|^2$$



Real-sim dataset, level sets of the validation loss (hold-out)

$$\arg \min_{\beta} \frac{1}{2n} \|y^{\text{train}} - X^{\text{train}} \beta\|^2 + \lambda_1 \|\beta\|_1 + \frac{\lambda_2}{2} \|\beta\|^2$$



$$\arg \min_{\lambda \in \mathbb{R}} \left\{ \mathcal{L}(\lambda) := C(\hat{\beta}(\lambda)) := \|y^{\text{test}} - X^{\text{test}} \hat{\beta}(\lambda)\|^2 \right\}$$
$$\text{s.t. } \hat{\beta}(\lambda) \in \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{2n} \|y^{\text{train}} - X^{\text{train}} \beta\|^2 + \lambda \|\beta\|_1$$

(1) J. Nocedal and S. J. Wright (2006). *Numerical optimization*. Second. Springer Series in Operations Research. Springer.

(2) D. C. Liu and J. Nocedal (1989). "On the limited memory BFGS method for large scale optimization". In: *Math. Program.*



$$\arg \min_{\lambda \in \mathbb{R}} \left\{ \mathcal{L}(\lambda) := C(\hat{\beta}(\lambda)) := \|y^{\text{test}} - X^{\text{test}} \hat{\beta}(\lambda)\|^2 \right\}$$
$$\text{s.t. } \hat{\beta}(\lambda) \in \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{2n} \|y^{\text{train}} - X^{\text{train}} \beta\|^2 + \lambda \|\beta\|_1$$

Once $\nabla_{\lambda} \mathcal{L}(\lambda)$ is computed, one can use standard first-order methods:

- ▶ Line-search⁽¹⁾
- ▶ L-BFGS⁽²⁾
- ▶ Gradient descent

⁽¹⁾ J. Nocedal and S. J. Wright (2006). *Numerical optimization*. Second. Springer Series in Operations Research. Springer.

⁽²⁾ D. C. Liu and J. Nocedal (1989). "On the limited memory BFGS method for large scale optimization". In: *Math. Program.*



$$\arg \min_{\lambda \in \mathbb{R}} \left\{ \mathcal{L}(\lambda) := C(\hat{\beta}(\lambda)) := \|y^{\text{test}} - X^{\text{test}} \hat{\beta}(\lambda)\|^2 \right\}$$
$$\text{s.t. } \hat{\beta}(\lambda) \in \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{2n} \|y^{\text{train}} - X^{\text{train}} \beta\|^2 + \lambda \|\beta\|_1$$

Once $\nabla_{\lambda} \mathcal{L}(\lambda)$ is computed, one can use standard first-order methods:

- ▶ Line-search⁽¹⁾
- ▶ L-BFGS⁽²⁾
- ▶ Gradient descent

Main contribution here: compute $\nabla_{\lambda} \mathcal{L}(\lambda)$ for a given $\lambda = \begin{pmatrix} \lambda_1 \\ \vdots \\ \lambda_r \end{pmatrix} \in \mathbb{R}^r$

⁽¹⁾ J. Nocedal and S. J. Wright (2006). *Numerical optimization*. Second. Springer Series in Operations Research. Springer.

⁽²⁾ D. C. Liu and J. Nocedal (1989). "On the limited memory BFGS method for large scale optimization". In: *Math. Program.*



$$\begin{aligned} \arg \min_{\lambda \in \mathbb{R}} \left\{ \mathcal{L}(\lambda) := C(\hat{\beta}(\lambda)) := \|y^{\text{test}} - X^{\text{test}} \hat{\beta}(\lambda)\|^2 \right\} \\ \text{s.t. } \hat{\beta}(\lambda) \in \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{2n} \|y^{\text{train}} - X^{\text{train}} \beta\|^2 + \lambda \|\beta\|_1 \end{aligned}$$

$$\arg \min_{\lambda \in \mathbb{R}} \left\{ \mathcal{L}(\lambda) := C(\hat{\beta}(\lambda)) := \|y^{\text{test}} - X^{\text{test}} \hat{\beta}(\lambda)\|^2 \right\}$$

$$\text{s.t. } \hat{\beta}(\lambda) \in \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{2n} \|y^{\text{train}} - X^{\text{train}} \beta\|^2 + \lambda \|\beta\|_1$$

Chain rule:

$$\nabla_{\lambda} \mathcal{L}(\lambda) = \underbrace{\hat{\mathcal{J}}_{(\lambda)}^{\top}}_{\substack{:= (\nabla_{\lambda} \hat{\beta}_1^{(\lambda)}, \dots, \nabla_{\lambda} \hat{\beta}_p^{(\lambda)}) \\ \rightarrow \text{main technical challenge}}} \nabla_{\beta} C(\hat{\beta}(\lambda))$$

► Boils down to:

how to compute the Jacobian $\hat{\mathcal{J}}_{(\lambda)} \in \mathbb{R}^{p \times r}$ efficiently?



"Smooth" inner optimization problems, **well studied**:

- ▶ *Implicit differentiation* (**closed-form** formula)^{(1), (2)}:
need to solve a $p \times p$ linear system ($p = \#$ features)
- ▶ *Automatic differentiation*, *forward*⁽³⁾ or *reverse*⁽⁴⁾ mode

Example:

$$\underbrace{\hat{\beta}^{(\lambda)} \in \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{2n} \|y^{\text{train}} - X^{\text{train}} \beta\|^2 + \frac{\lambda}{2} \|\beta\|^2}_{\text{inner optimization problem}}$$

(1) J. Larsen et al. (1996). "Design and regularization of neural networks: the optimal use of a validation set". In: *Neural Networks for Signal Processing VI. Proceedings of the 1996 IEEE Signal Processing Society Workshop*.

(2) Y. Bengio (2000). "Gradient-based optimization of hyperparameters". In: *Neural comput.* 12.8, pp. 1889–1900.

(3) L. Franceschi et al. (2017). "Forward and reverse gradient-based hyperparameter optimization". In: *ICML*, pp. 1165–1173.

(4) J. Domke (2012). "Generic methods for optimization-based modeling". In: *AISTATS*. vol. 22, pp. 318–326.



Optimization problem:

$$\min_{\beta \in \mathbb{R}^p} f(\beta) + \lambda g(\beta)$$

Properties:

- ▶ f convex, gradient L -Lipschitz
- ▶ g convex but non necessarily smooth (can have kinks)

Example: $f(\beta) = \frac{1}{2n} \|X\beta - y\|^2$, $g(\beta) = \lambda \|\beta\|_1$



Optimization problem:

$$\min_{\beta \in \mathbb{R}^p} f(\beta) + \lambda g(\beta)$$

Properties:

- ▶ f convex, gradient L -Lipschitz
- ▶ g convex but non necessarily smooth (can have kinks)

Example: $f(\beta) = \frac{1}{2n} \|X\beta - y\|^2$, $g(\beta) = \lambda \|\beta\|_1$

Rem: fix step size (sub-)gradient descent does not converge: take $f = 0$, $g = |\cdot|$ and use $\beta_0 = 1/2$, $\alpha = 1$ (ping-pong!)

Properties:

- ▶ g (convex)
- ▶ its **proximal** operator⁽¹⁾ has a closed-form:

$$\text{prox}_{\lambda g}(\beta^0) := \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{2} \|\beta - \beta^0\|^2 + \lambda g(\beta)$$



J.-J. Moreau

⁽¹⁾ Jean-Jacques Moreau (1962). "Fonctions convexes duales et points proximaux dans un espace hilbertien". In: *C. R. Acad. Sci. Paris* 255, pp. 2897–2899.

⁽²⁾ N. Parikh and S. Boyd (2014). "Proximal Algorithms". In: *Foundations and Trends in Machine Learning* 1.3, pp. 127–239.

⁽³⁾ H. H. Bauschke and P. L. Combettes (2011). *Convex analysis and monotone operator theory in Hilbert spaces*. Springer, pp. xvi+468.



Properties:

- ▶ g (convex)
- ▶ its **proximal** operator⁽¹⁾ has a closed-form:

$$\text{prox}_{\lambda g}(\beta^0) := \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{2} \|\beta - \beta^0\|^2 + \lambda g(\beta)$$



J.-J. Moreau

Majorizer (at step t): $\beta \mapsto f(\beta^t) + \langle \nabla f(\beta^t), \beta - \beta^t \rangle + \frac{L \|\beta^t - \beta\|^2}{2} + \lambda g(\beta)$

⁽¹⁾ Jean-Jacques Moreau (1962). "Fonctions convexes duales et points proximaux dans un espace hilbertien". In: *C. R. Acad. Sci. Paris* 255, pp. 2897–2899.

⁽²⁾ N. Parikh and S. Boyd (2014). "Proximal Algorithms". In: *Foundations and Trends in Machine Learning* 1.3, pp. 127–239.

⁽³⁾ H. H. Bauschke and P. L. Combettes (2011). *Convex analysis and monotone operator theory in Hilbert spaces*. Springer, pp. xvi+468.

Properties:

- ▶ g (convex)
- ▶ its **proximal** operator⁽¹⁾ has a closed-form:

$$\text{prox}_{\lambda g}(\beta^0) := \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{2} \|\beta - \beta^0\|^2 + \lambda g(\beta)$$



J.-J. Moreau

Majorizer (at step t): $\beta \mapsto f(\beta^t) + \langle \nabla f(\beta^t), \beta - \beta^t \rangle + \frac{L \|\beta^t - \beta\|^2}{2} + \lambda g(\beta)$

Update rule (minimize the majorizer) :

$$\beta^{t+1} = \text{prox}_{\frac{\lambda g}{L}}\left(\beta^t - \frac{1}{L} \nabla f(\beta^t)\right)$$

⁽¹⁾ Jean-Jacques Moreau (1962). "Fonctions convexes duales et points proximaux dans un espace hilbertien". In: *C. R. Acad. Sci. Paris* 255, pp. 2897–2899.

⁽²⁾ N. Parikh and S. Boyd (2014). "Proximal Algorithms". In: *Foundations and Trends in Machine Learning* 1.3, pp. 127–239.

⁽³⁾ H. H. Bauschke and P. L. Combettes (2011). *Convex analysis and monotone operator theory in Hilbert spaces*. Springer, pp. xvi+468.

Properties:

- ▶ g (convex)
- ▶ its **proximal** operator⁽¹⁾ has a closed-form:

$$\text{prox}_{\lambda g}(\beta^0) := \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{2} \|\beta - \beta^0\|^2 + \lambda g(\beta)$$



J.-J. Moreau

Majorizer (at step t): $\beta \mapsto f(\beta^t) + \langle \nabla f(\beta^t), \beta - \beta^t \rangle + \frac{L \|\beta^t - \beta\|^2}{2} + \lambda g(\beta)$

Update rule (minimize the majorizer) :

$$\beta^{t+1} = \text{prox}_{\frac{\lambda g}{L}}\left(\beta^t - \frac{1}{L} \nabla f(\beta^t)\right)$$

- ▶ Proximal algorithms / recipes⁽²⁾
- ▶ Associated theory / analysis⁽³⁾

⁽¹⁾ Jean-Jacques Moreau (1962). "Fonctions convexes duales et points proximaux dans un espace hilbertien". In: *C. R. Acad. Sci. Paris* 255, pp. 2897–2899.

⁽²⁾ N. Parikh and S. Boyd (2014). "Proximal Algorithms". In: *Foundations and Trends in Machine Learning* 1.3, pp. 127–239.

⁽³⁾ H. H. Bauschke and P. L. Combettes (2011). *Convex analysis and monotone operator theory in Hilbert spaces*. Springer, pp. xvi+468.



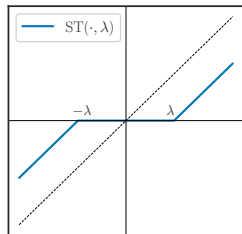
► **Soft-Thresholding**

Closed form solution for the prox of $|\cdot|$:

$$\text{prox}_{\lambda|\cdot|}(\beta^0) = \text{ST}(\beta^0, \lambda)$$

$$:= \arg \min_{\beta \in \mathbb{R}} \left(\frac{1}{2}(\beta^0 - \beta)^2 + \lambda|\beta| \right)$$

$$= \text{sign}(\beta^0) \cdot \max(0, |\beta^0| - \lambda)$$



► **Soft-Thresholding**

Closed form solution for the prox of | · |:

$$\text{prox}_{\lambda|\cdot|}(\beta^0) = \text{ST}(\beta^0, \lambda)$$

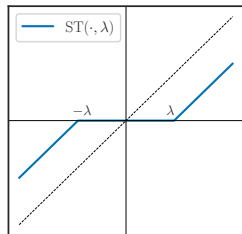
$$:= \arg \min_{\beta \in \mathbb{R}} \left(\frac{1}{2}(\beta^0 - \beta)^2 + \lambda|\beta| \right)$$

$$= \text{sign}(\beta^0) \cdot \max(0, |\beta^0| - \lambda)$$

► Componentwise **Soft-Thresholding**

Closed form solution for the || · ||₁:

$$\left[\text{prox}_{\lambda\|\cdot\|_1}(\beta^0) \right]_j = \left[\text{ST}(\beta^0, \lambda) \right]_j, \quad \text{for all } j \in [p]$$



FORWARD-MODE DIFFERENTIATION ⁽¹⁾ OF PGD ⁽²⁾ PROXIMAL GRADIENT DESCENT (PGD) ⁽³⁾



$$\hat{\beta}(\lambda) \in \arg \min_{\beta \in \mathbb{R}^p} \overbrace{f}^{\text{smooth}}(\beta) + \lambda \overbrace{g}^{\text{non-smooth}}(\beta)$$

Algorithm: Proximal gradient descent PGD

```
init :  $\beta = 0_p, \quad , L$   
for iter = 1, ..., do  
     $z \leftarrow \beta - \frac{1}{L} \nabla f(\beta)$  // gradient step  
  
     $\beta \leftarrow \text{prox}_{\lambda g/L}(z)$  // proximal step: thresholding for us  
  
return  $\beta$ 
```

⁽¹⁾ R. E. Wengert (1964). "A simple automatic derivative evaluation program". In: *Communications of the ACM* 7.8, pp. 463–464

⁽²⁾ C.-A. Deledalle et al. (2014). "Stein Unbiased GrAdient estimator of the Risk (SUGAR) for multiple parameter selection". In: *SIAM J. Imaging Sci.* 7.4

⁽³⁾ B. Martinet (1970). "Brève communication. Régularisation d'inéquations variationnelles par approximations successives". In: *Revue française d'informatique et de recherche opérationnelle. Série rouge* 4.R3, pp. 154–158

FORWARD-MODE DIFFERENTIATION ⁽¹⁾ OF PGD ⁽²⁾ PROXIMAL GRADIENT DESCENT (PGD) ⁽³⁾



$$\hat{\beta}^{(\lambda)} \in \arg \min_{\beta \in \mathbb{R}^p} \overbrace{f}^{\text{smooth}}(\beta) + \lambda \overbrace{g}^{\text{non-smooth}}(\beta)$$

Algorithm: Forward-mode differentiation of PGD

init : $\beta = 0_p, \mathcal{J} = 0_p, L$ (Lipschitz-constant for f)

for iter = 1, ..., **do**

$z \leftarrow \beta - \frac{1}{L} \nabla f(\beta)$ // gradient step

$dz \leftarrow (\text{Id}_p - \frac{1}{L} \nabla^2 f(\beta)) \mathcal{J}$ // diff w.r.t. λ : chain rule

$\beta \leftarrow \text{prox}_{\lambda g/L}(z)$ // proximal step: thresholding for us

return β

⁽¹⁾ R. E. Wengert (1964). "A simple automatic derivative evaluation program". In: *Communications of the ACM* 7.8, pp. 463–464

⁽²⁾ C.-A. Deledalle et al. (2014). "Stein Unbiased GrAdient estimator of the Risk (SUGAR) for multiple parameter selection". In: *SIAM J. Imaging Sci.* 7.4

⁽³⁾ B. Martinet (1970). "Brève communication. Régularisation d'inéquations variationnelles par approximations successives". In: *Revue française d'informatique et de recherche opérationnelle. Série rouge* 4.R3, pp. 154–158

FORWARD-MODE DIFFERENTIATION ⁽¹⁾ OF PGD ⁽²⁾ PROXIMAL GRADIENT DESCENT (PGD) ⁽³⁾



$$\hat{\beta}^{(\lambda)} \in \arg \min_{\beta \in \mathbb{R}^p} \overbrace{f}^{\text{smooth}}(\beta) + \lambda \overbrace{g}^{\text{non-smooth}}(\beta)$$

Algorithm: Forward-mode differentiation of PGD

init : $\beta = 0_p, \mathcal{J} = 0_p, L$ (Lipschitz-constant for f)

for iter = 1, ..., **do**

```
     $z \leftarrow \beta - \frac{1}{L} \nabla f(\beta)$  // gradient step
     $dz \leftarrow (\text{Id}_p - \frac{1}{L} \nabla^2 f(\beta)) \mathcal{J}$  // diff w.r.t.  $\lambda$ : chain rule
     $\beta \leftarrow \text{prox}_{\lambda g/L}(z)$  // proximal step: thresholding for us
     $\mathcal{J} \leftarrow \partial_z \text{prox}_{\lambda g/L}(z) dz$  // diff w.r.t.  $\lambda$ : chain rule
     $\quad + \partial_\lambda \text{prox}_{\lambda g/L}(z)$  // do not forget this term!
```

return β, \mathcal{J}

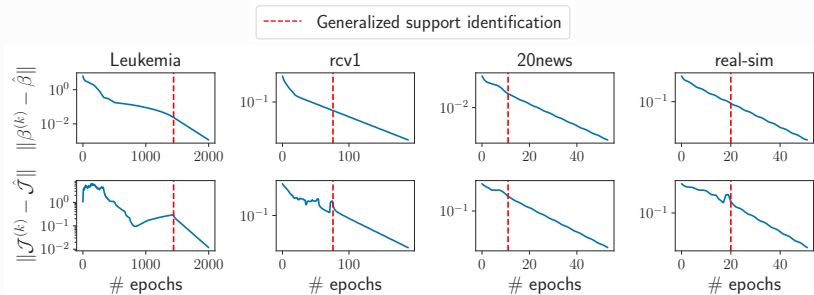
⁽¹⁾ R. E. Wengert (1964). "A simple automatic derivative evaluation program". In: *Communications of the ACM* 7.8, pp. 463–464

⁽²⁾ C.-A. Deledalle et al. (2014). "Stein Unbiased GrAdient estimator of the Risk (SUGAR) for multiple parameter selection". In: *SIAM J. Imaging Sci.* 7.4

⁽³⁾ B. Martinet (1970). "Brève communication. Régularisation d'inéquations variationnelles par approximations successives". In: *Revue française d'informatique et de recherche opérationnelle. Série rouge* 4.R3, pp. 154–158

Forward diff. PCD convergence, Lasso

Provided that X (the design matrix) is not pathological, the sequence generated by PCD is converging to $\hat{\beta}$, and the Jacobian sequence based on forward differentiation converges to the true Jacobian. Moreover, once the support (the non-zero coefs.) has been identified, convergence is linear.⁽¹⁾



⁽¹⁾ Q. Bertrand, Q. Klopfenstein, M. Blondel, et al. (2020). "Implicit differentiation of Lasso-type models for hyperparameter optimization". In: *ICML*



$$\hat{\beta}(\lambda) \in \arg \min_{\beta \in \mathbb{R}^p} f(\beta, \lambda)$$

⁽¹⁾ Y. Bengio (2000). "Gradient-based optimization of hyperparameters". In: *Neural comput.* 12.8, pp. 1889–1900



$$\hat{\beta}^{(\lambda)} \in \arg \min_{\beta \in \mathbb{R}^p} f(\beta, \lambda)$$

$$\nabla_{\beta} f(\hat{\beta}^{(\lambda)}, \lambda) = 0$$

⁽¹⁾ Y. Bengio (2000). "Gradient-based optimization of hyperparameters". In: *Neural comput.* 12.8, pp. 1889–1900



$$\hat{\beta}^{(\lambda)} \in \arg \min_{\beta \in \mathbb{R}^p} f(\beta, \lambda)$$

$$\nabla_{\beta} f(\hat{\beta}^{(\lambda)}, \lambda) = 0$$

$$\nabla_{\beta, \lambda}^2 f(\hat{\beta}^{(\lambda)}, \lambda) + \hat{\mathcal{J}}_{(\lambda)}^{\top} \nabla_{\beta}^2 f(\hat{\beta}^{(\lambda)}, \lambda) = 0$$

⁽¹⁾ Y. Bengio (2000). "Gradient-based optimization of hyperparameters". In: *Neural comput.* 12.8, pp. 1889–1900

$$\hat{\beta}^{(\lambda)} \in \arg \min_{\beta \in \mathbb{R}^p} f(\beta, \lambda)$$

$$\nabla_{\beta} f(\hat{\beta}^{(\lambda)}, \lambda) = 0$$

$$\nabla_{\beta, \lambda}^2 f(\hat{\beta}^{(\lambda)}, \lambda) + \hat{\mathcal{J}}_{(\lambda)}^{\top} \nabla_{\beta}^2 f(\hat{\beta}^{(\lambda)}, \lambda) = 0$$

$$\hat{\mathcal{J}}_{(\lambda)}^{\top} = -\nabla_{\beta, \lambda}^2 f(\hat{\beta}^{(\lambda)}, \lambda) \underbrace{\left(\nabla_{\beta}^2 f(\hat{\beta}^{(\lambda)}, \lambda) \right)^{-1}}_{p \times p}$$

- ▶ Need to solve a linear **system of size p**
- ▶ Prohibitive for large p

⁽¹⁾ Y. Bengio (2000). "Gradient-based optimization of hyperparameters". In: *Neural comput.* 12.8, pp.1889–1900



General formulation:

- ▶ Solve a linear **system of size p**
- ▶ Prohibitive for large p

With a sparsity inducing penalty:

- ▶ Solve a linear **system of size S** (sparsity degree of the estimator)
- ▶ $S \ll p$ very often



$$\hat{\beta}^{(\lambda)} = \text{ST} \left(\hat{\beta}^{(\lambda)} - \frac{1}{L} \nabla f(\hat{\beta}^{(\lambda)}), \frac{\lambda}{L} \right)$$

⁽¹⁾ Q. Bertrand, Q. Klopfenstein, M. Massias, et al. (2022). "Implicit differentiation for fast hyperparameter selection in non-smooth convex learning".
In: *Submitted to JMLR*



$$\hat{\beta}^{(\lambda)} = \text{ST} \left(\hat{\beta}^{(\lambda)} - \frac{1}{L} \nabla f(\hat{\beta}^{(\lambda)}), \frac{\lambda}{L} \right)$$

$$\begin{aligned} \hat{\mathcal{J}} &= \partial_{\beta} \text{ST} \left(\hat{\beta}^{(\lambda)} - \frac{1}{L} \nabla f(\hat{\beta}^{(\lambda)}), \frac{\lambda}{L} \right) \left(\text{Id} - \frac{\nabla^2 f}{L} \right) \hat{\mathcal{J}} \\ &+ \partial_{\lambda} \text{ST} \left(\hat{\beta}^{(\lambda)} - \frac{1}{L} \nabla f(\hat{\beta}^{(\lambda)}), \frac{\lambda}{L} \right) \end{aligned}$$



$$\hat{\beta}^{(\lambda)} = \text{ST} \left(\hat{\beta}^{(\lambda)} - \frac{1}{L} \nabla f(\hat{\beta}^{(\lambda)}), \frac{\lambda}{L} \right)$$

$$\begin{aligned} \hat{\mathcal{J}} &= \partial_{\beta} \text{ST} \left(\hat{\beta}^{(\lambda)} - \frac{1}{L} \nabla f(\hat{\beta}^{(\lambda)}), \frac{\lambda}{L} \right) \left(\text{Id} - \frac{\nabla^2 f}{L} \right) \hat{\mathcal{J}} \\ &\quad + \partial_{\lambda} \text{ST} \left(\hat{\beta}^{(\lambda)} - \frac{1}{L} \nabla f(\hat{\beta}^{(\lambda)}), \frac{\lambda}{L} \right) \end{aligned}$$

Key observation, if $\hat{\beta}_j^{(\lambda)} = 0$:

$$\partial_{\beta} \text{ST} \left(\hat{\beta}_j^{(\lambda)} - \frac{1}{L} \nabla_j f(\hat{\beta}^{(\lambda)}), \frac{\lambda}{L} \right) = 0; 0 = \partial_{\lambda} \text{ST} \left(\hat{\beta}_j^{(\lambda)} - \frac{1}{L} \nabla_j f(\hat{\beta}^{(\lambda)}), \frac{\lambda}{L} \right)$$



$$\hat{\beta}^{(\lambda)} = \text{ST} \left(\hat{\beta}^{(\lambda)} - \frac{1}{L} \nabla f(\hat{\beta}^{(\lambda)}), \frac{\lambda}{L} \right)$$

$$\begin{aligned} \hat{\mathcal{J}} &= \partial_{\beta} \text{ST} \left(\hat{\beta}^{(\lambda)} - \frac{1}{L} \nabla f(\hat{\beta}^{(\lambda)}), \frac{\lambda}{L} \right) \left(\text{Id} - \frac{\nabla^2 f}{L} \right) \hat{\mathcal{J}} \\ &\quad + \partial_{\lambda} \text{ST} \left(\hat{\beta}^{(\lambda)} - \frac{1}{L} \nabla f(\hat{\beta}^{(\lambda)}), \frac{\lambda}{L} \right) \end{aligned}$$

Key observation, if $\hat{\beta}_j^{(\lambda)} = 0$:

$$\partial_{\beta} \text{ST} \left(\hat{\beta}_j^{(\lambda)} - \frac{1}{L} \nabla_j f(\hat{\beta}^{(\lambda)}), \frac{\lambda}{L} \right) = 0; 0 = \partial_{\lambda} \text{ST} \left(\hat{\beta}_j^{(\lambda)} - \frac{1}{L} \nabla_j f(\hat{\beta}^{(\lambda)}), \frac{\lambda}{L} \right)$$

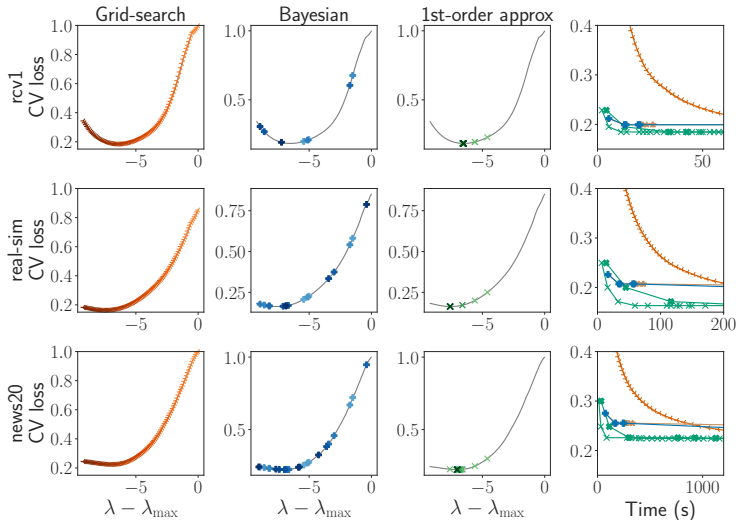
With $\mathcal{S} = \{j \in [p] : \hat{\beta}_j^{(\lambda)} \neq 0\}$ we have $\hat{\mathcal{J}}_{\mathcal{S}^c} = 0$

$$\hat{\mathcal{J}}_{\mathcal{S}} = \partial_{\beta} \text{ST} \left(\hat{\beta}^{(\lambda)} - \frac{1}{L} \nabla f(\hat{\beta}^{(\lambda)}), \frac{\lambda}{L} \right)_{\mathcal{S}} \hat{\mathcal{J}}_{\mathcal{S}} + \partial_{\lambda} \text{ST} \left(\hat{\beta}_j^{(\lambda)} - \frac{1}{L} \nabla_j f(\hat{\beta}^{(\lambda)}), \frac{\lambda}{L} \right)_{\mathcal{S}}$$

EXPERIMENTS I - LASSO CROSS-VALIDATION



◆ 1st-order
 ✕ 1st-order approx
 — Grid-search
 ★ Random-search
 ◆ Bayesian

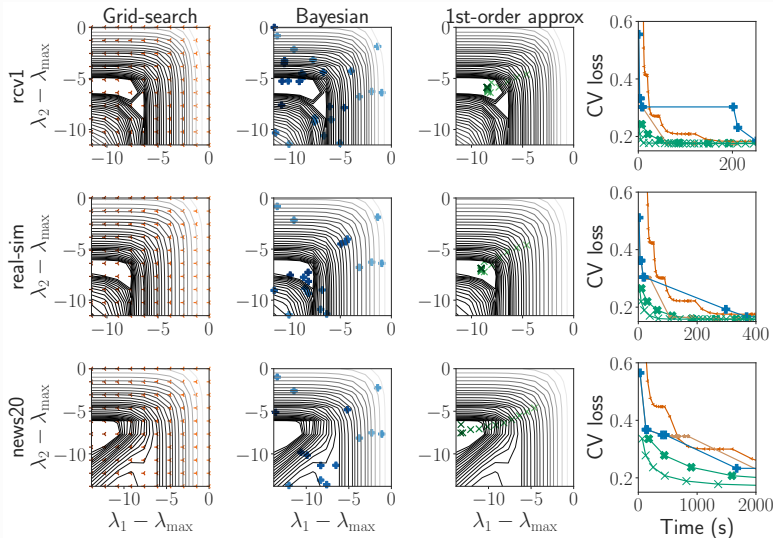


$$\arg \min_{\beta} \frac{1}{2n} \|y^{\text{train}} - X^{\text{train}} \beta\|^2 + e^{\lambda} \|\beta\|_1$$

EXPERIMENTS II - ENET CROSS-VALIDATION

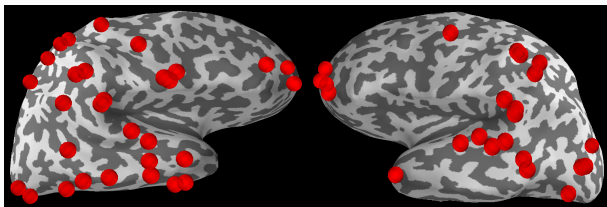


◆ 1st-order
 ✕ 1st-order approx
 — Grid-search
 ★ Random-search
 ◆ Bayesian



$$\arg \min_{\beta} \frac{1}{2n} \|y^{\text{train}} - X^{\text{train}}\beta\|^2 + e^{\lambda_1} \|\beta\|_1 + \frac{e^{\lambda_2}}{2} \|\beta\|^2$$

- ▶ **Outer criterion:** FDMC SURE⁽¹⁾
- ▶ **Inner problems:** vanilla Lasso

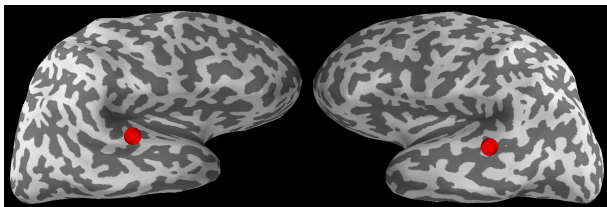


Real M/EEG data, vanilla Lasso (1 hyperparameter λ)

$$\arg \min_{\beta \in \mathbb{R}^p} \frac{1}{2n} \|y - X\beta\|_2^2 + e^\lambda \|\beta\|_1$$

⁽¹⁾ C.-A. Deledalle et al. (2014). "Stein Unbiased Gradient estimator of the Risk (SUGAR) for multiple parameter selection". In: *SIAM J. Imaging Sci.* 7.4.

- ▶ **Outer criterion:** FDMC SURE⁽¹⁾
- ▶ **Inner problems:** weighted Lasso



Real M/EEG data, weighted Lasso (p hyperparameters)

$$\arg \min_{\beta \in \mathbb{R}^p} \frac{1}{2n} \|y - X\beta\|_2^2 + \sum_{j=1}^p e^{\lambda_j} |\beta_j|$$

⁽¹⁾ C.-A. Deledalle et al. (2014). "Stein Unbiased Gradient estimator of the Risk (SUGAR) for multiple parameter selection". In: *SIAM J. Imaging Sci.* 7.4.



- ▶ **Local linear convergence** of the Jacobian
- ▶ **Leverage sparsity** to speed up hypergradient computation
- ▶ Open source package

<https://github.com/QB3/sparse-ho>

The screenshot shows the GitHub repository page for 'sparse-ho'. At the top, there is a navigation bar with links for 'sparse-ho', '0.1.dev', 'Examples', 'API', 'GitHub', and 'Site'. Below the navigation bar is the repository name 'sparse-ho' in a large font. Underneath the name are two status badges: 'build passing' (green) and 'codecov 79%' (orange). The main content area contains a paragraph describing the package: 'sparse-ho stands for "sparse hyperparameter optimization". This package implements efficient hyperparameter tuning for sparse machine learning models. It supports models such as the Lasso, the Weighted Lasso, multiclass sparse Logistic regression, SVM, etc.' This is followed by another paragraph: 'Relying on a first order algorithm for bilevel optimization, sparse-ho's performances scales gracefully with the number of hyperparameters to tune.' Below that is a paragraph: 'In order to benchmark performances, the package also implements alternatives such as forward or backward differentiation.' The 'Documentation' section is titled in a large font, followed by the text: 'Please visit [`https://qb3.github.io/sparse-ho`](https://qb3.github.io/sparse-ho) for the latest version of the documentation.' The 'Install' section is also titled in a large font, followed by the text: 'To run the code you first need to clone the repository, and then run, in the folder containing the `setup.py` file (root folder):' Below this text is a code block containing the command: `pip install -e .`

- ▶ Specific parametrization ($\lambda \leftarrow e^\lambda$ for simplicity)
- ▶ Need a **differentiable criterion**: cannot use 0/1-loss
- ▶ Need a **continuous estimator** *w.r.t.* data and hyperparameters: does not apply yet to **non-convex** penalties⁽¹⁾ nor reweighted Lasso⁽²⁾
- ▶ Optimized function often **non-convex**: possibly multiple local minima, use multi-starts
- ▶ Potentially slow and handy **outer solver**

⁽¹⁾ P. Breheny and J. Huang (2011). "Coordinate descent algorithms for nonconvex penalized regression, with applications to biological feature selection". In: *Ann. Appl. Stat.* 5.1, p. 232.

⁽²⁾ E. J. Candès, M. B. Wakin, and S. P. Boyd (2008). "Enhancing Sparsity by Reweighted l_1 Minimization". In: *J. Fourier Anal. Applicat.* 14.5-6, pp. 877-905.



Contributions:

- ▶ 1st-order optimization with nonsmooth inner problem
- ▶ **Local linear convergence** of the Jacobian
- ▶ **Leverage sparsity** to speed up hypergradient computation



Contributions:

- ▶ 1st-order optimization with nonsmooth inner problem
- ▶ **Local linear convergence** of the Jacobian
- ▶ **Leverage sparsity** to speed up hypergradient computation

Future work:

- ▶ Convergence of the bilevel procedure
- ▶ Smarter outer solver



Contributions:

- ▶ 1st-order optimization with nonsmooth inner problem
- ▶ **Local linear convergence** of the Jacobian
- ▶ **Leverage sparsity** to speed up hypergradient computation

Future work:

- ▶ Convergence of the bilevel procedure
- ▶ Smarter outer solver

*"An article about computational science in a scientific publication is **not** the scholarship itself, it is merely **advertising** of the scholarship. The actual scholarship is the complete software development environment and the complete set of instructions which generated the figures."*

J. B. Buckheit and D. L. Donoho⁽¹⁾

"All models are wrong, but some come with good open source implementation and good documentation: so use these."

A. Gramfort (circa 2015)

Contact:

✉ joseph.salmon@umontpellier.fr

🌐 <http://josephsalmon.eu>








Github: @josephsalmon

















Twitter: @salmonjsph
















⁽¹⁾ J. B. Buckheit and D. L. Donoho (1995). "WaveLab and Reproducible Research". In: *Wavelets and Statistics*. Lect. Notes Statist. 103. Springer-Verlag, pp. 55–81.




- 
- Argyriou, A., T. Evgeniou, and M. Pontil (2008). “Convex multi-task feature learning”. In: *Machine Learning* 73.3, pp. 243–272.
- 
- Bauschke, H. H. and P. L. Combettes (2011). *Convex analysis and monotone operator theory in Hilbert spaces*. New York: Springer, pp. xvi+468.
- 
- Bengio, Y. (2000). “Gradient-based optimization of hyperparameters”. In: *Neural comput.* 12.8, pp. 1889–1900.
- 
- Berger, H. (1929). “Über das elektroenkephalogramm des menschen”. In: *Archiv für psychiatrie und nervenkrankheiten*.
- 
- Bergstra, J. and Y. Bengio (2012). “Random search for hyper-parameter optimization”. In: *J. Mach. Learn. Res.* 13.2.
- 
- Bertrand, Q., Q. Klopfenstein, M. Blondel, et al. (2020). “Implicit differentiation of Lasso-type models for hyperparameter optimization”. In: *ICML*.
- 
- Bertrand, Q., Q. Klopfenstein, M. Massias, et al. (2022). “Implicit differentiation for fast hyperparameter selection in non-smooth convex learning”. In: *Submitted to JMLR*.

-  Bertrand, Q. and M. Massias (2021). “Anderson acceleration of coordinate descent”. In: *AISTATS*.
-  Bertrand, Q., M. Massias, et al. (2019). “Handling correlated and repeated measurements with the smoothed multivariate square-root Lasso”. In: *NeurIPS*.
-  Breheny, P. and J. Huang (2011). “Coordinate descent algorithms for nonconvex penalized regression, with applications to biological feature selection”. In: *Ann. Appl. Stat.* 5.1, p. 232.
-  Brochu, E., V. M. Cora, and N. De Freitas (2010). *A tutorial on Bayesian optimization of expensive cost functions, with application to active user modeling and hierarchical reinforcement learning*. Tech. rep.
-  Buckheit, J. B. and D. L. Donoho (1995). “WaveLab and Reproducible Research”. In: *Wavelets and Statistics*. Lect. Notes Statist. 103. New York: Springer-Verlag, pp. 55–81.
-  Candès, E. J., M. B. Wakin, and S. P. Boyd (2008). “Enhancing Sparsity by Reweighted l_1 Minimization”. In: *J. Fourier Anal. Applicat.* 14.5-6, pp. 877–905.

-  Cohen, D. (1968). “Magnetoencephalography: evidence of magnetic fields produced by alpha-rhythm currents”. In: *Science*.
-  Deledalle, C.-A. et al. (2014). “Stein Unbiased GrADient estimator of the Risk (SUGAR) for multiple parameter selection”. In: *SIAM J. Imaging Sci.* 7.4.
-  Domke, J. (2012). “Generic methods for optimization-based modeling”. In: *AISTATS*. Vol. 22, pp. 318–326.
-  Figueiredo, M. (2001). “Adaptive Sparseness Using Jeffreys Prior”. In: *NIPS*, pp. 697–704.
-  Franceschi, L. et al. (2017). “Forward and reverse gradient-based hyperparameter optimization”. In: *ICML*, pp. 1165–1173.
-  Gramfort, A., M. Kowalski, and M. Hämmäläinen (2012). “Mixed-norm estimates for the M/EEG inverse problem using accelerated gradient methods”. In: *Phys. Med. Biol.* 57.7, pp. 1937–1961.
-  Hutter, F., J. Lücke, and L. Schmidt-Thieme (2015). “Beyond Manual Tuning of Hyperparameters”. In: *Künstliche Intell.* 29.4, pp. 329–337.
-  Kohavi, R. and G. H. John (1995). “Automatic parameter selection by minimizing estimated error”. In: *ICML*, pp. 304–312.

-  Larsen, J. et al. (1996). “Design and regularization of neural networks: the optimal use of a validation set”. In: *Neural Networks for Signal Processing VI. Proceedings of the 1996 IEEE Signal Processing Society Workshop*.
-  Liu, D. C. and J. Nocedal (1989). “On the limited memory BFGS method for large scale optimization”. In: *Math. Program.*
-  Liu, W. and Y. Yang (2011). “Parametric or nonparametric? A parametricness index for model selection”. In: *Ann. Statist.* 39.4, pp. 2074–2102.
-  Lounici, K. (2008). “Sup-norm convergence rate and sign concentration property of Lasso and Dantzig estimators”. In: *Electron. J. Stat.* 2, pp. 90–102.
-  Lounici, K., K. Meziani, and B. Riu (2021). “Muddling Labels for Regularization, a novel approach to generalization”. In: *arXiv preprint arXiv:2102.08769*.
-  Lounici, K., M. Pontil, et al. (2011). “Oracle inequalities and optimal inference under group sparsity”. In: *Ann. Statist.* 39.4, pp. 2164–2204.

-  Martinet, B. (1970). “Brève communication. Régularisation d’inéquations variationnelles par approximations successives”. In: *Revue française d’informatique et de recherche opérationnelle. Série rouge* 4.R3, pp. 154–158.
-  Moreau, Jean-Jacques (1962). “Fonctions convexes duales et points proximaux dans un espace hilbertien”. In: *C. R. Acad. Sci. Paris* 255, pp. 2897–2899.
-  Nocedal, J. and S. J. Wright (2006). *Numerical optimization*. Second. Springer Series in Operations Research. New York: Springer.
-  Ochs, P. et al. (2015). “Bilevel optimization with nonsmooth lower level problems”. In: *SSVM*. Vol. 9087, pp. 654–665.
-  Parikh, N. and S. Boyd (2014). “Proximal Algorithms”. In: *Foundations and Trends in Machine Learning* 1.3, pp. 127–239.
-  Pedregosa, F. (2016). “Hyperparameter optimization with approximate gradient”. In: *ICML*. Vol. 48, pp. 737–746.
-  Stein, C. M. (1981). “Estimation of the mean of a multivariate normal distribution”. In: *Ann. Statist.* 9.6, pp. 1135–1151.

-  Stone, L. R. A. and J.C. Ramer (1965). “Estimating WAIS IQ from Shipley Scale scores: Another cross-validation”. In: *Journal of clinical psychology* 21.3, pp. 297–297.
-  Tipping, M. E. (2001). “Sparse Bayesian learning and the relevance vector machine”. In: *J. Mach. Learn. Res.* 1, pp. 211–244.
-  Wengert, R. E. (1964). “A simple automatic derivative evaluation program”. In: *Communications of the ACM* 7.8, pp. 463–464.