

# HYPERPARAMETER SELECTION FOR HIGH DIMENSIONAL SPARSE LEARNING: APPLICATION TO NEUROIMAGING

## A GENTLE MOTIVATION

**Joseph Salmon**

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



UNIVERSITÉ  
DE  
MONTPELLIER



# JOINT WORKS WITH VARIOUS COLLEAGUES

1

**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



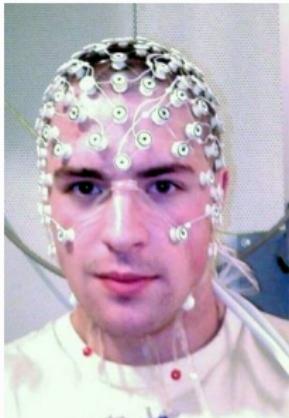
Samuel



Alexandre

# NEUROIMAGING DATA: EEG <sup>(1)</sup> AND MEG <sup>(2)</sup>

2



**(a) EEG**



**(b) MEG=Mag.+Grad.**



**(c) M/EEG**

Photo credit: S. Whitmarsh

<sup>(1)</sup> H. Berger (1929). "Über das elektroenzephalogramm des menschen". In: *Archiv für psychiatrie und nervenkrankheiten*

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

# NEUROIMAGING DATA: EEG <sup>(1)</sup> AND MEG <sup>(2)</sup>

2



**(a) EEG**



**(b) MEG=Mag.+Grad.**



**(c) M/EEG**

Photo credit: S. Whitmarsh

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

---

<sup>(1)</sup> H. Berger (1929). "Über das elektroenzephalogramm des menschen". In: *Archiv für psychiatrie und nervenkrankheiten*

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

# NEUROIMAGING DATA: EEG <sup>(1)</sup> AND MEG <sup>(2)</sup>



**(a) EEG**



**(b) MEG=Mag.+Grad.**



**(c) M/EEG**

Photo credit: S. Whitmarsh

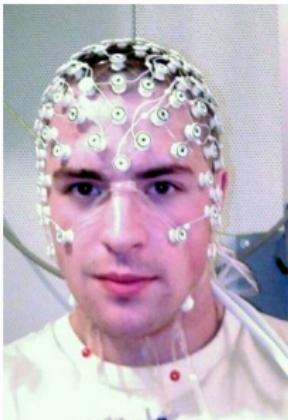
- ▶ **Data Y:** electric and magnetic fields at the head surface
- ▶ **Goal:** which parts of the brain are responsible for the signals?

---

<sup>(1)</sup> H. Berger (1929). "Über das elektroenzephalogramm des menschen". In: *Archiv für psychiatrie und nervenkrankheiten*

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

# NEUROIMAGING DATA: EEG <sup>(1)</sup> AND MEG <sup>(2)</sup>



**(a) EEG**



**(b) MEG=Mag.+Grad.**



**(c) M/EEG**

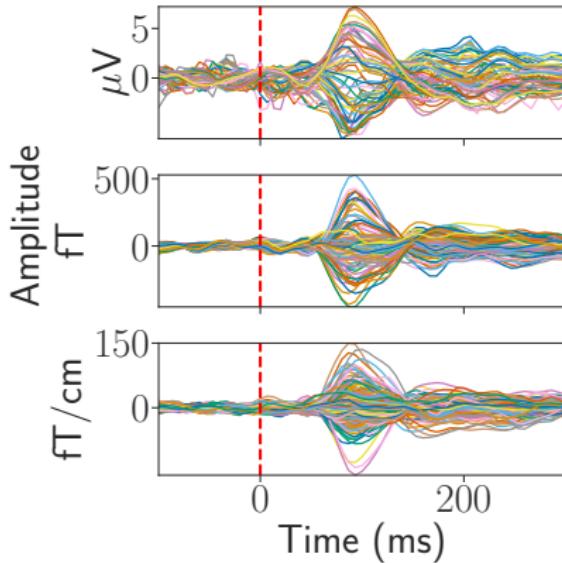
Photo credit: S. Whitmarsh

- ▶ **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

---

<sup>(1)</sup> H. Berger (1929). "Über das elektroenzephalogramm des menschen". In: *Archiv für psychiatrie und nervenkrankheiten*

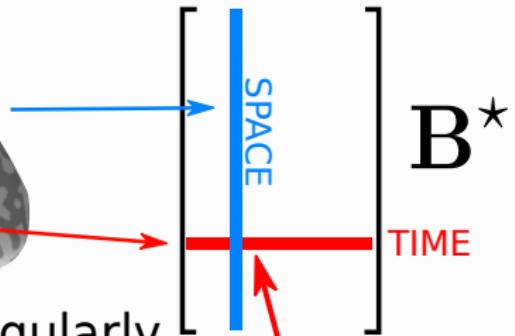
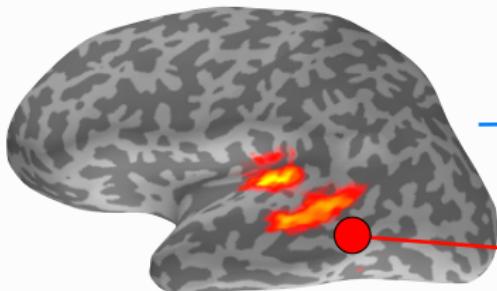
<sup>(2)</sup> D. Cohen (1968). "Magnetoencephalography: evidence of magnetic fields produced by alpha-rhythm currents". In: *Science*



3 modalities:

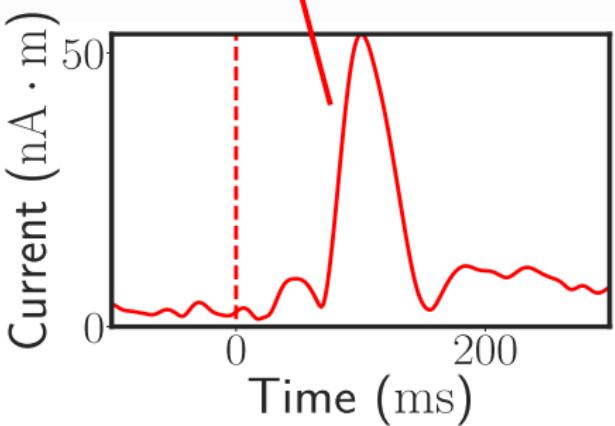
- ▶ EEG
- ▶ MEG: magnometers (amplitude)
- ▶ MEG: gradiometers (gradients)

# SOURCE MODELING

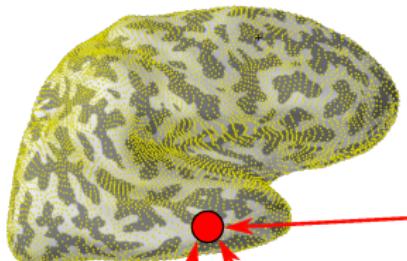


Source candidates regularly spaced in the brain  
(e.g., every 5mm)

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



# DESIGN MATRIX - FORWARD OPERATOR

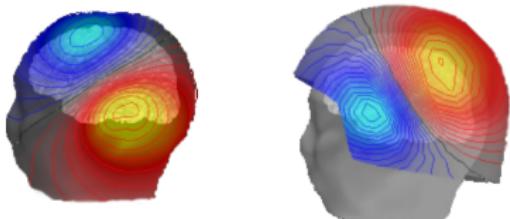


$$X = \begin{bmatrix} X_{\text{EEG}} \\ X_{\text{MEG}} \end{bmatrix}$$

$$\in \mathbb{R}^{n \times p}$$

**EEG:**  
Forward field of the electrodes

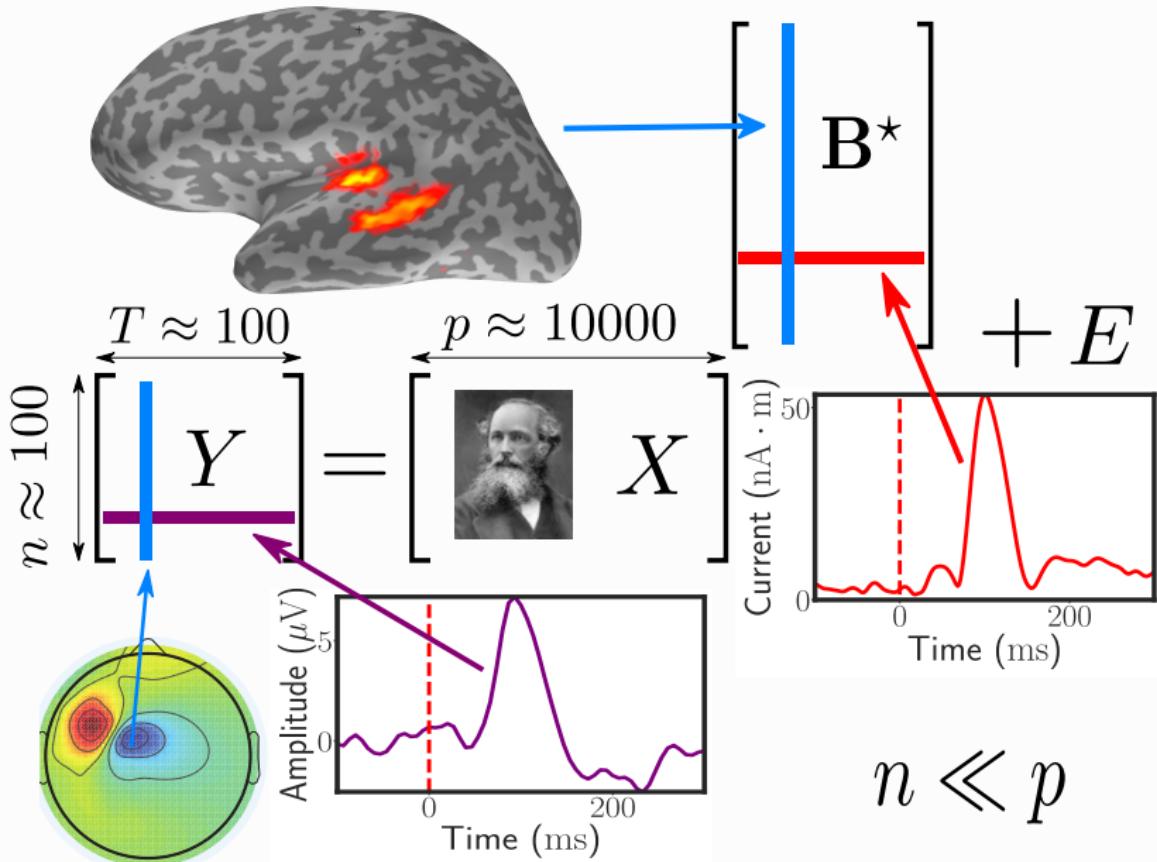
**MEG:**  
Forward field of sensor



$X$ : gain matrix /  
forward operator  
obtained by  
Maxwell's equations

# THE M/EEG INVERSE PROBLEM

MAXWELL EQUATIONS AND (APPROX.) LINEARITY

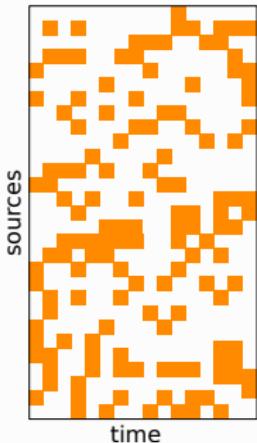


# MULTITASK PENALTIES<sup>(1)</sup> AND MEG<sup>(2)</sup>



Popular convex penalties:

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



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*

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

# MULTITASK PENALTIES<sup>(1)</sup> AND MEG<sup>(2)</sup>



Popular convex penalties: multitask Lasso (MTL)

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



Sparse support: group structure ✓

Penalty: **Group-Lasso**

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

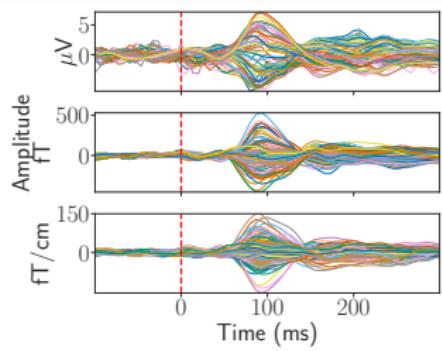
Parameter  $\hat{\mathbf{B}} \in \mathbb{R}^{p \times T}$

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

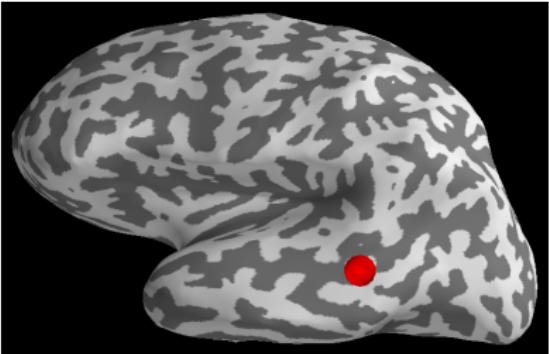
(1) A. Argyriou, T. Evgeniou, and M. Pontil (2008). "Convex multi-task feature learning". In: *Machine Learning*

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

# SUMMARY OF THE PROBLEM SETTING

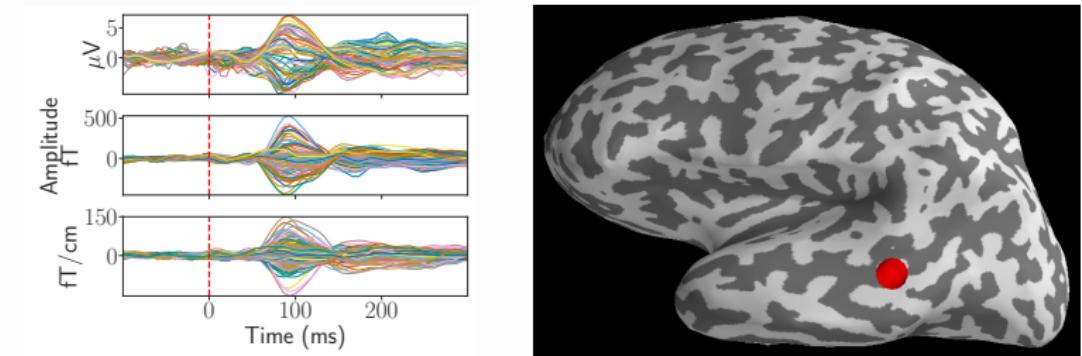


What you have:  $\textcolor{blue}{Y} \in \mathbb{R}^{n \times T}$



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

# SUMMARY OF THE PROBLEM SETTING



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

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

This is typically done using optimization based estimators

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

# SUMMARY OF CONTRIBUTIONS



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

Covered in this presentation

- ▶ How to efficiently select the regularization parameter  $\lambda$ ?<sup>(1), (2)</sup>

Not covered in this presentation

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

---

<sup>(1)</sup> Q. Bertrand, Q. Klopfenstein, M. Blondel, et al. (2020). "Implicit differentiation of Lasso-type models for hyperparameter optimization". In: [ICML](#).

<sup>(2)</sup> Q. Bertrand, Q. Klopfenstein, M. Massias, et al. (2021). "Implicit differentiation for fast hyperparameter selection in non-smooth convex learning". In: [Submitted to JMLR](#).

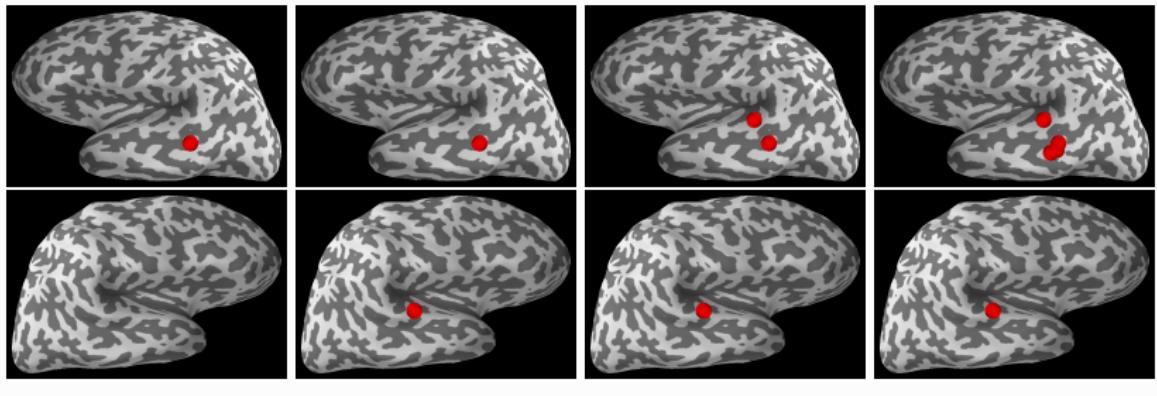
<sup>(3)</sup> Q. Bertrand and M. Massias (2021). "Anderson acceleration of coordinate descent". In: [AISTATS](#).

<sup>(4)</sup> Q. Bertrand, M. Massias, et al. (2019). "Handling correlated and repeated measurements with the smoothed Multivariate square-root Lasso". In: [NeurIPS](#).

# WHICH $\lambda$ TO PICK?



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



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

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

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

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

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

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



- ▶ Statistical route<sup>(1), (2)</sup>: assumptions on the design matrix  $X$ , provide guarantees but are often too
- ▶ Bayesian statistics<sup>(3), (4)</sup>: prior on  $\lambda$
- ▶ Bayesian optimisation,<sup>(5)</sup> 0-th order method:

The road today:

- ▶ Hyperparameter optimization<sup>(6)</sup> : 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) K. Lounici, M. Pontil, et al. (2009). "Taking Advantage of Sparsity in Multi-Task Learning". In: *arXiv preprint arXiv:0903.1468*.

(3) M. E. Tipping (2001). "Sparse Bayesian learning and the relevance vector machine". In: *Journal of Machine Learning Research*.

(4) M. Figueiredo (2001). "Adaptive Sparseness Using Jeffreys Prior.". In: *Advances in Neural Information Processing Systems*.

(5) F. Hutter, J. Lücke, and L. Schmidt-Thieme (2015). "Beyond manual tuning of hyperparameters". In: *KI-Künstliche Intelligenz*.

(6) R. Kohavi and G. H. John (1995). "Automatic parameter selection by minimizing estimated error". In: *Machine Learning Proceedings*.



Possible selection criterion:

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

---

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

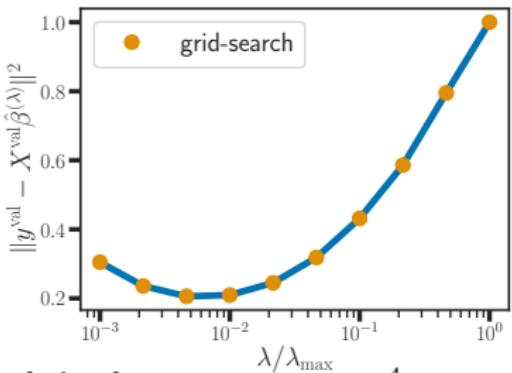
(2) K. Lounici, K. Meziani, and B. Riu (2021). "Muddling Labels for Regularization, a novel approach to generalization". In: *arXiv preprint arXiv:2102.08769*.

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



**Real-sim dataset**,  $n \approx p \approx 10^4$

Validation loss as a function of  $\lambda$ .

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**

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

(1) L. R. A. Stone and J. 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*.

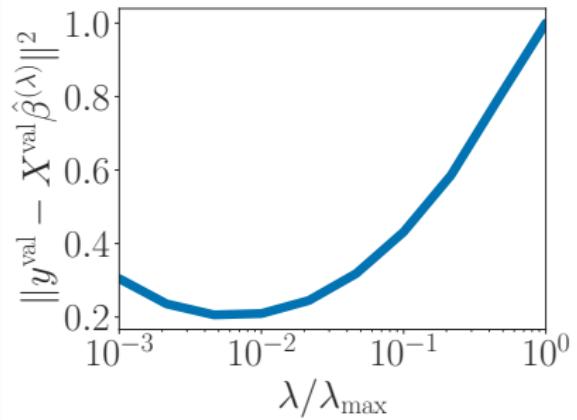
(3) W. Liu and Y. Yang (2011). "Parametric or nonparametric? A parametricness index for model selection". In: *Ann. Statist.* 39.4, pp. 2074–2102.

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

# HO AS A BILEVEL OPTIMIZATION PROBLEM<sup>(1), (2)</sup>



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



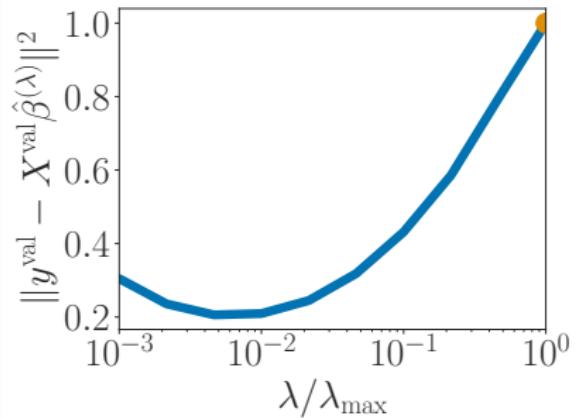
<sup>(1)</sup> P. Ochs et al. (2015). "Bilevel optimization with nonsmooth lower level problems". In: *SSVM*

<sup>(2)</sup> F. Pedregosa (2016). "Hyperparameter optimization with approximate gradient". In: *ICML*

# HO AS A BILEVEL OPTIMIZATION PROBLEM<sup>(1), (2)</sup>



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



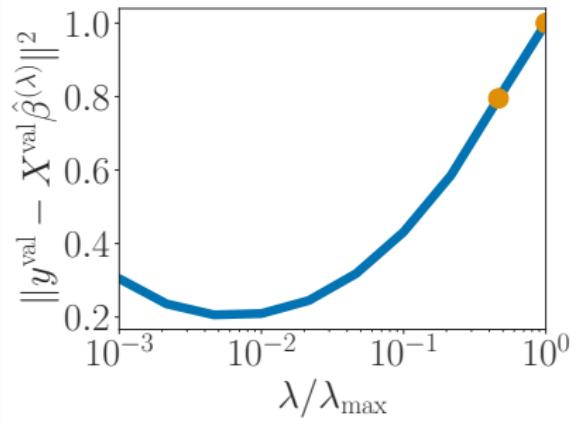
<sup>(1)</sup> P. Ochs et al. (2015). "Bilevel optimization with nonsmooth lower level problems". In: *SSVM*

<sup>(2)</sup> F. Pedregosa (2016). "Hyperparameter optimization with approximate gradient". In: *ICML*

# HO AS A BILEVEL OPTIMIZATION PROBLEM<sup>(1), (2)</sup>



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



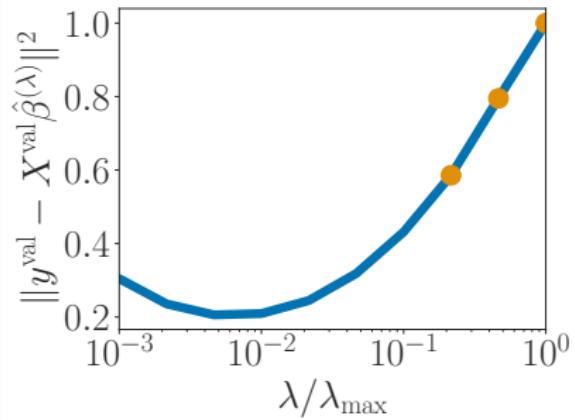
<sup>(1)</sup> P. Ochs et al. (2015). "Bilevel optimization with nonsmooth lower level problems". In: SSVM

<sup>(2)</sup> F. Pedregosa (2016). "Hyperparameter optimization with approximate gradient". In: ICML

# HO AS A BILEVEL OPTIMIZATION PROBLEM<sup>(1), (2)</sup>



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



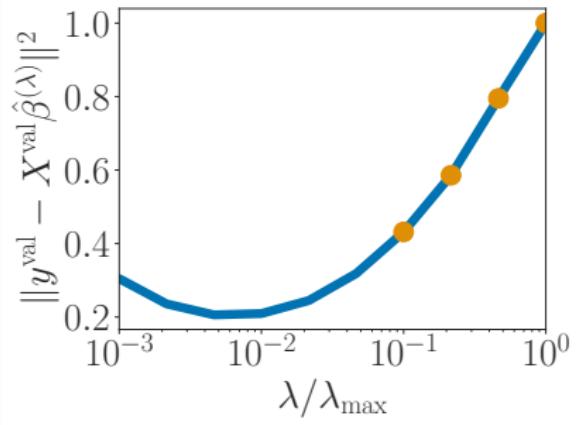
<sup>(1)</sup> P. Ochs et al. (2015). "Bilevel optimization with nonsmooth lower level problems". In: SSVM

<sup>(2)</sup> F. Pedregosa (2016). "Hyperparameter optimization with approximate gradient". In: ICML

# HO AS A BILEVEL OPTIMIZATION PROBLEM<sup>(1), (2)</sup>



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



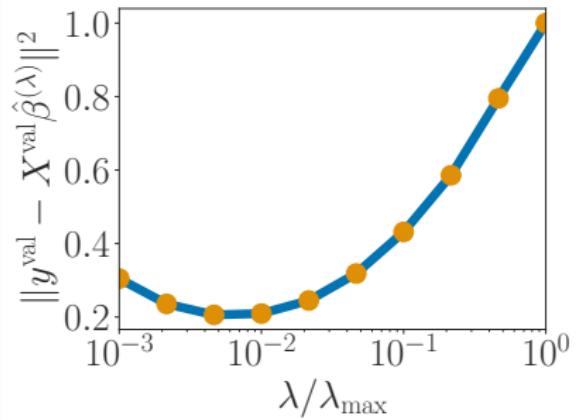
<sup>(1)</sup> P. Ochs et al. (2015). "Bilevel optimization with nonsmooth lower level problems". In: SSVM

<sup>(2)</sup> F. Pedregosa (2016). "Hyperparameter optimization with approximate gradient". In: ICML

# HO AS A BILEVEL OPTIMIZATION PROBLEM<sup>(1), (2)</sup>

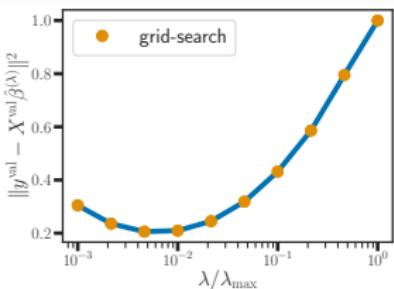


$$\begin{aligned} & \underbrace{\arg \min_{\lambda \in \mathbb{R}} \left\{ \mathcal{L}(\lambda) := \|y^{\text{val}} - X^{\text{val}} \hat{\beta}^{(\lambda)}\|^2 \right\}}_{\text{outer optimization problem}} \\ & \text{s.t. } \hat{\beta}^{(\lambda)} \in \underbrace{\arg \min_{\beta \in \mathbb{R}^p} \frac{1}{2n}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



<sup>(1)</sup> P. Ochs et al. (2015). "Bilevel optimization with nonsmooth lower level problems". In: SSVM

<sup>(2)</sup> F. Pedregosa (2016). "Hyperparameter optimization with approximate gradient". In: ICML



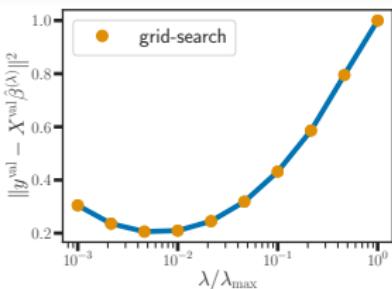
$$\begin{aligned} & \arg \min_{\lambda \in \mathbb{R}} \left\{ \mathcal{L}(\lambda) := \|y^{\text{val}} - X^{\text{val}} \hat{\beta}^{(\lambda)}\|^2 \right\} \\ & \text{s.t. } \hat{\beta}^{(\lambda)} \in \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{2n} \|y^{\text{train}} - X^{\text{train}} \beta\|^2 +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

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

<sup>(1)</sup> J. Bergstra and Y. Bengio (2012). "Random search for hyper-parameter optimization". In: *Journal of Machine Learning Research*.

<sup>(2)</sup> E. Brochu, V. M. Cora, and N. D. Freitas (2010). "A tutorial on Bayesian optimization of expensive cost functions, with application to active user modeling and hierarchical reinforcement learning". In: *arXiv preprint arXiv:1012.2599*.

<sup>(3)</sup> F. Pedregosa (2016). "Hyperparameter optimization with approximate gradient". In: *ICML*.



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

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

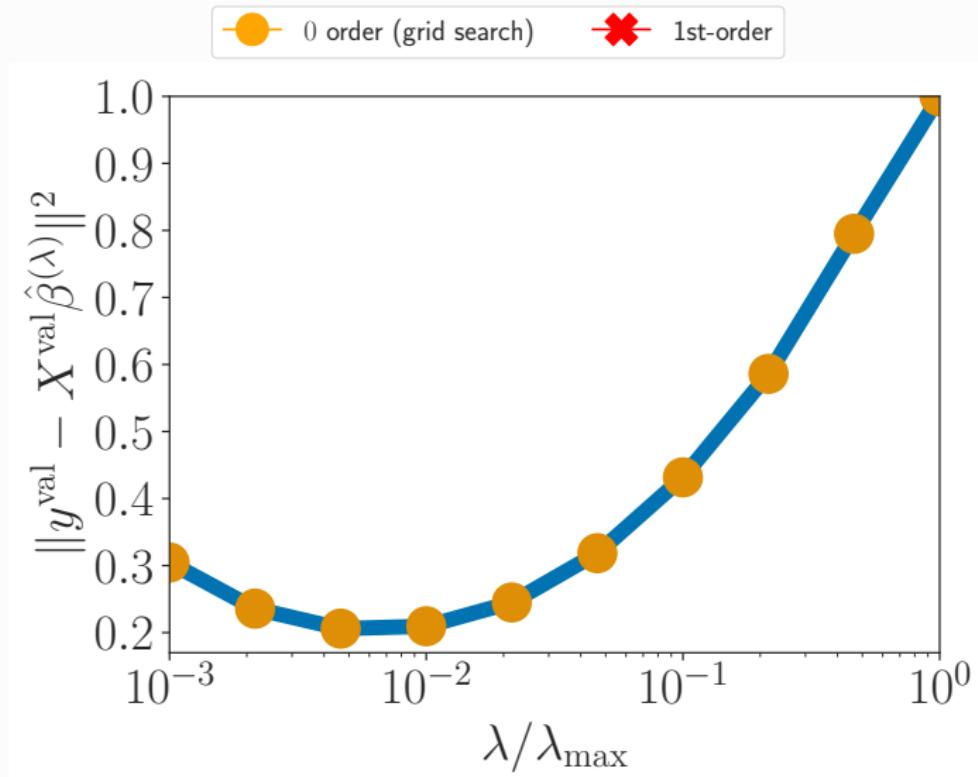
$$\lambda^{(t+1)} = \lambda^{(t)} - \rho \nabla_{\lambda} \mathcal{L}(\lambda^{(t)}) \quad \text{with } \rho > 0$$

<sup>(1)</sup> J. Bergstra and Y. Bengio (2012). "Random search for hyper-parameter optimization". In: *Journal of Machine Learning Research*.

<sup>(2)</sup> E. Brochu, V. M. Cora, and N. D. Freitas (2010). "A tutorial on Bayesian optimization of expensive cost functions, with application to active user modeling and hierarchical reinforcement learning". In: *arXiv preprint arXiv:1012.2599*.

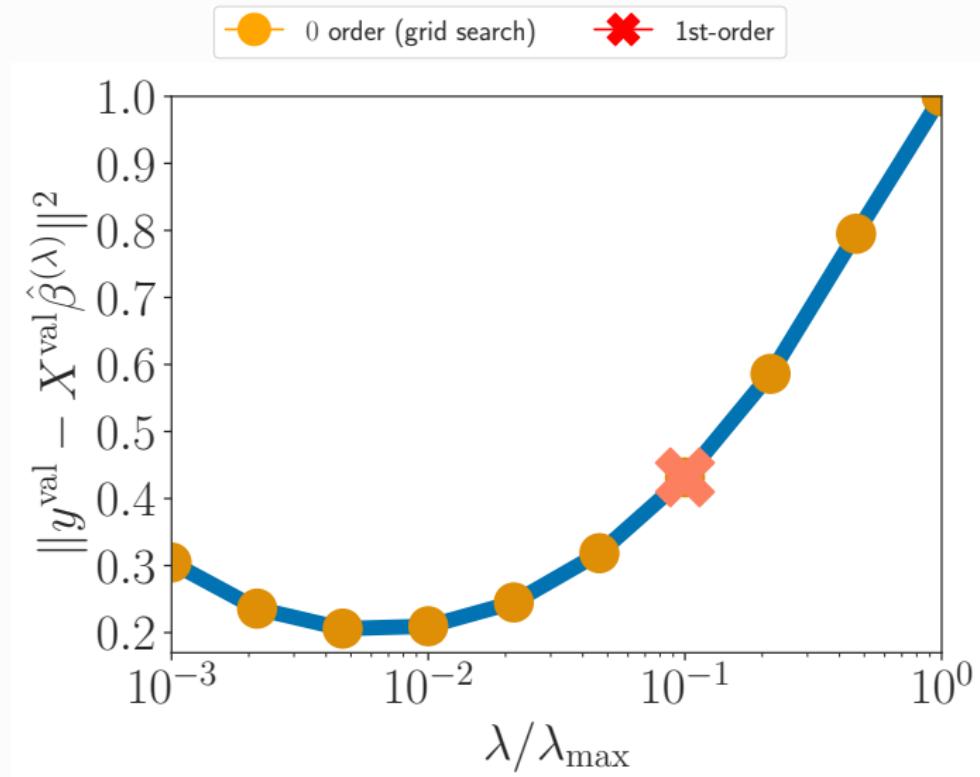
<sup>(3)</sup> F. Pedregosa (2016). "Hyperparameter optimization with approximate gradient". In: *ICML*.

# FIRST-ORDER OPTIMIZATION IN $\lambda$ , LASSO



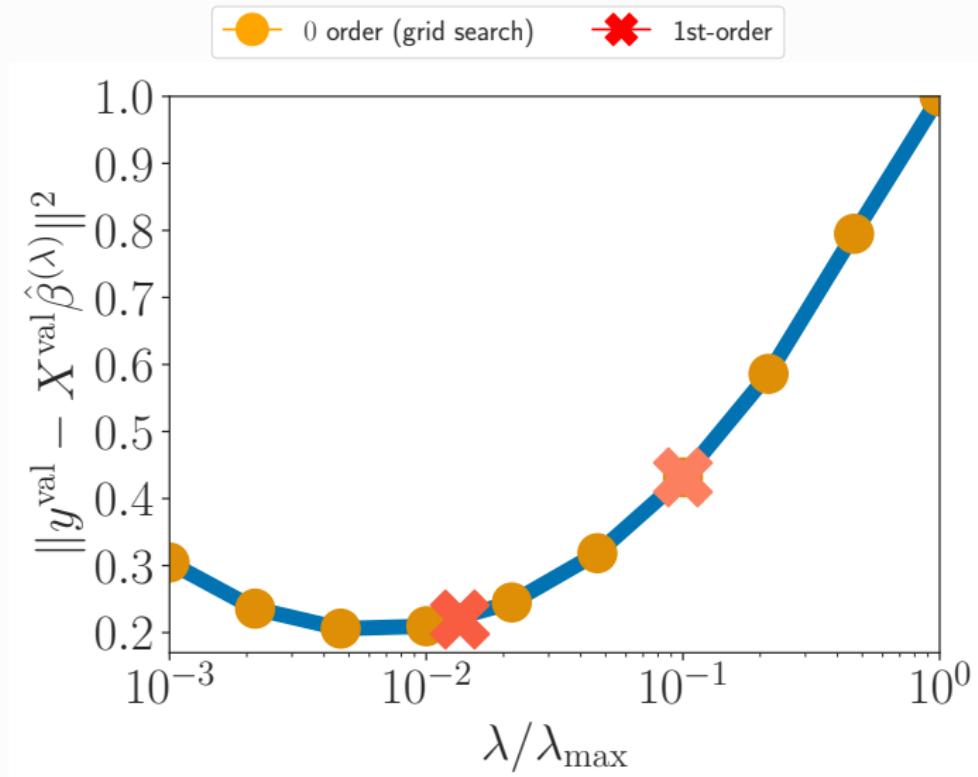
**Real-sim dataset**,  $n \approx p \approx 10^4$ . Validation loss as a function of  $\lambda$ .

# FIRST-ORDER OPTIMIZATION IN $\lambda$ , LASSO



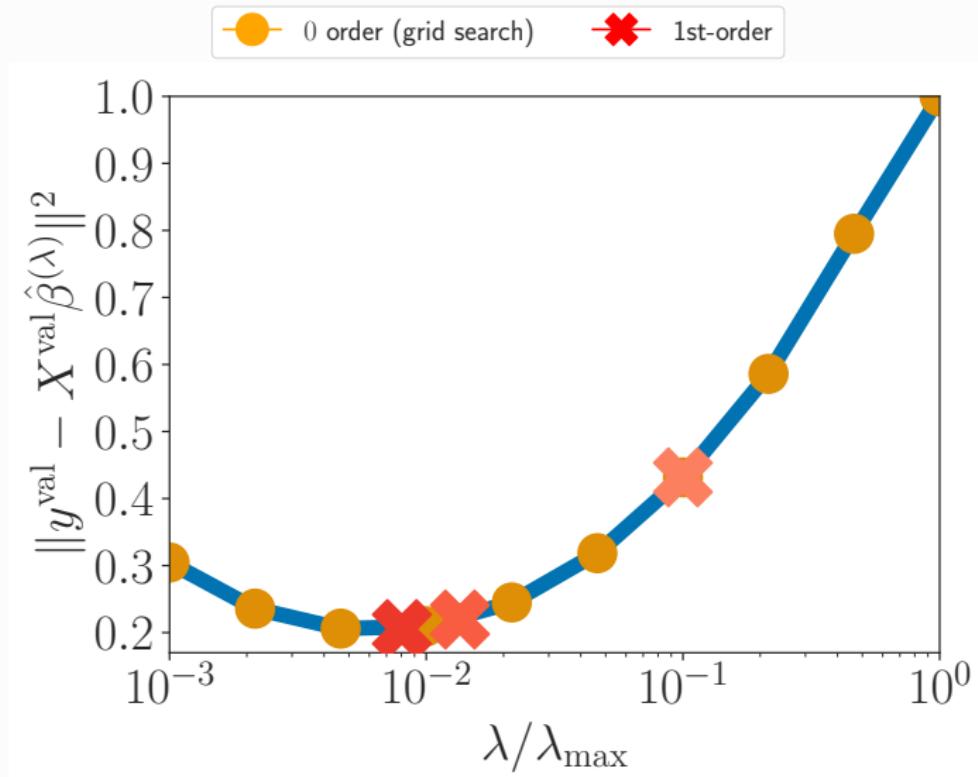
**Real-sim dataset**,  $n \approx p \approx 10^4$ . Validation loss as a function of  $\lambda$ .

# FIRST-ORDER OPTIMIZATION IN $\lambda$ , LASSO



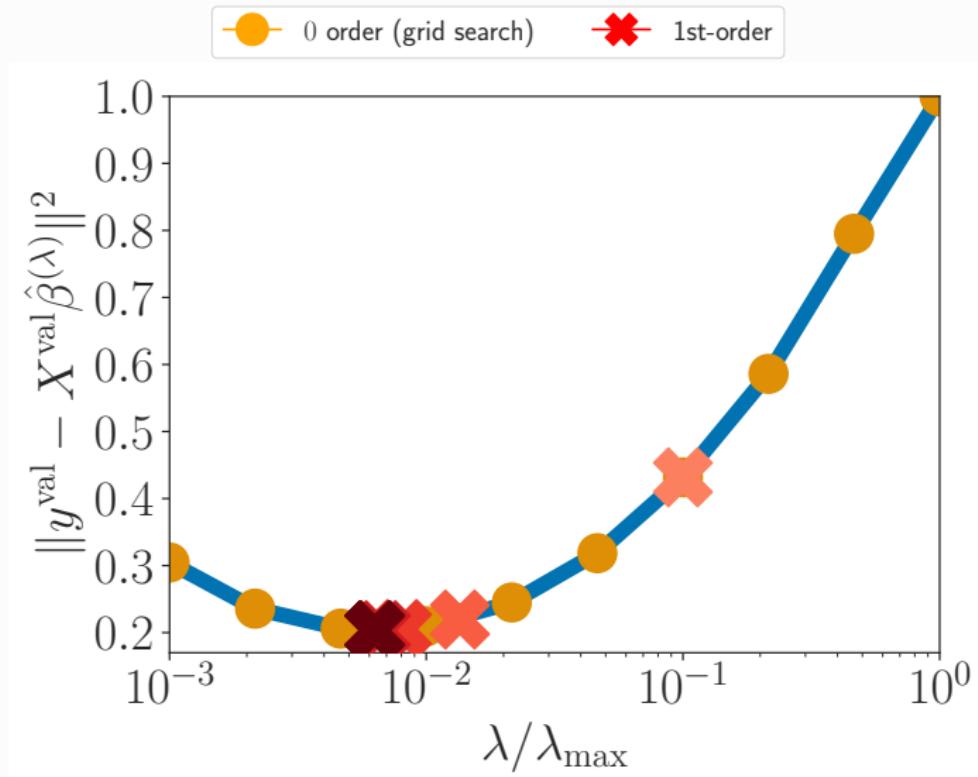
**Real-sim dataset**,  $n \approx p \approx 10^4$ . Validation loss as a function of  $\lambda$ .

# FIRST-ORDER OPTIMIZATION IN $\lambda$ , LASSO



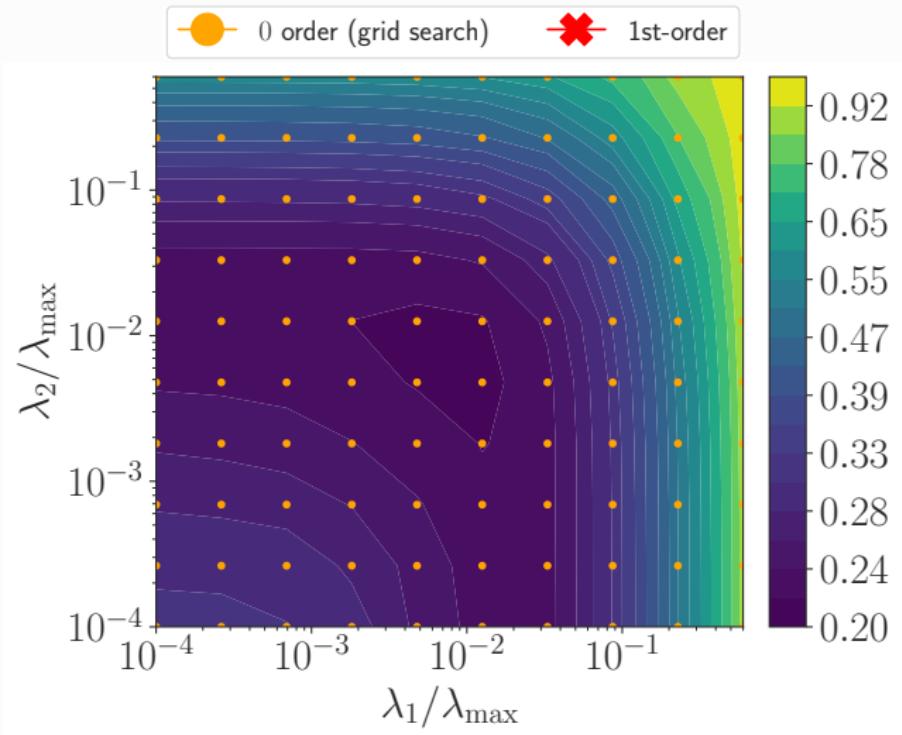
**Real-sim dataset**,  $n \approx p \approx 10^4$ . Validation loss as a function of  $\lambda$ .

# FIRST-ORDER OPTIMIZATION IN $\lambda$ , LASSO



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

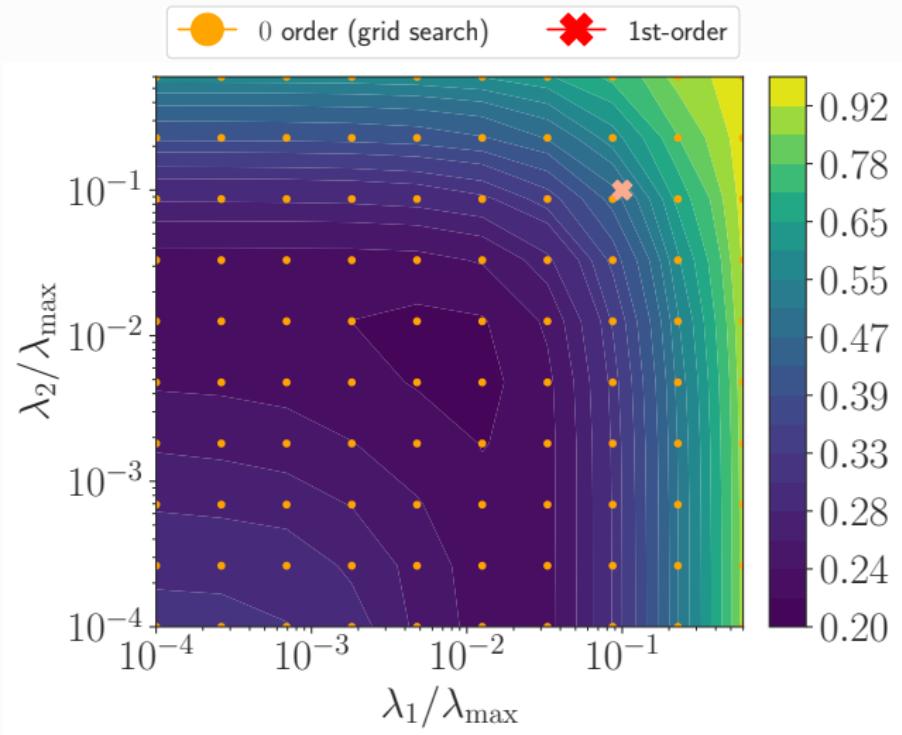
# FIRST-ORDER OPTIMIZATION IN $\lambda$ , ENET



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

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

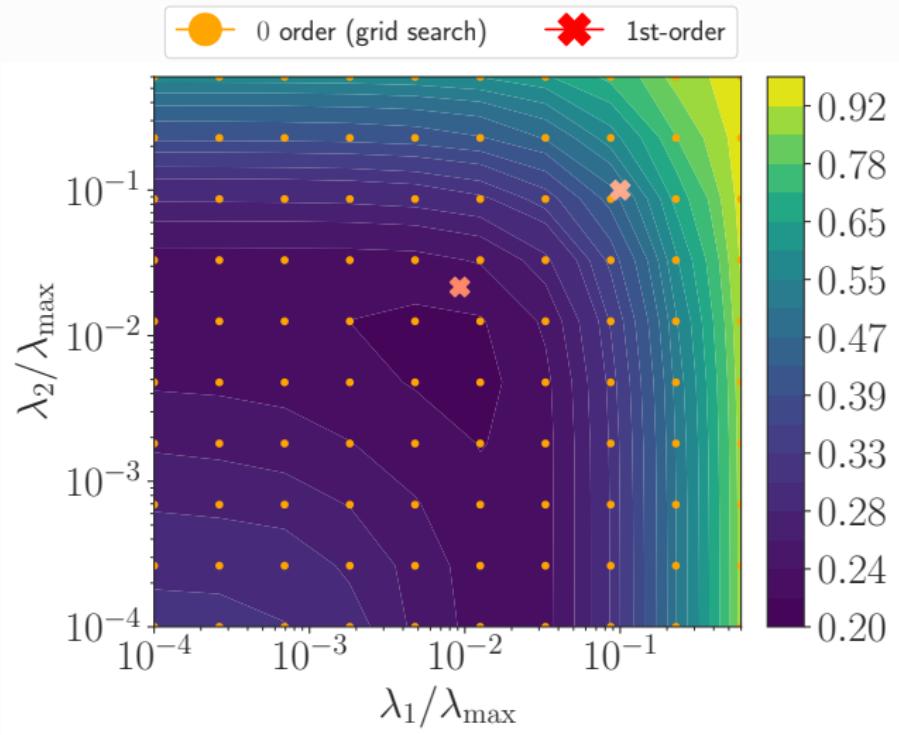
# FIRST-ORDER OPTIMIZATION IN $\lambda$ , ENET



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

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

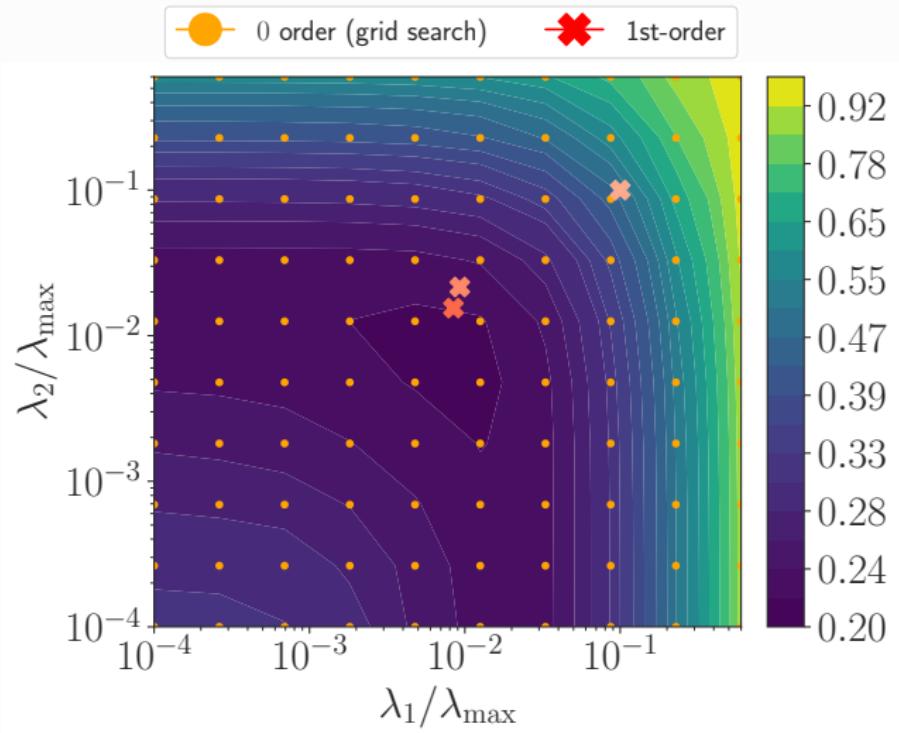
# FIRST-ORDER OPTIMIZATION IN $\lambda$ , ENET



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

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

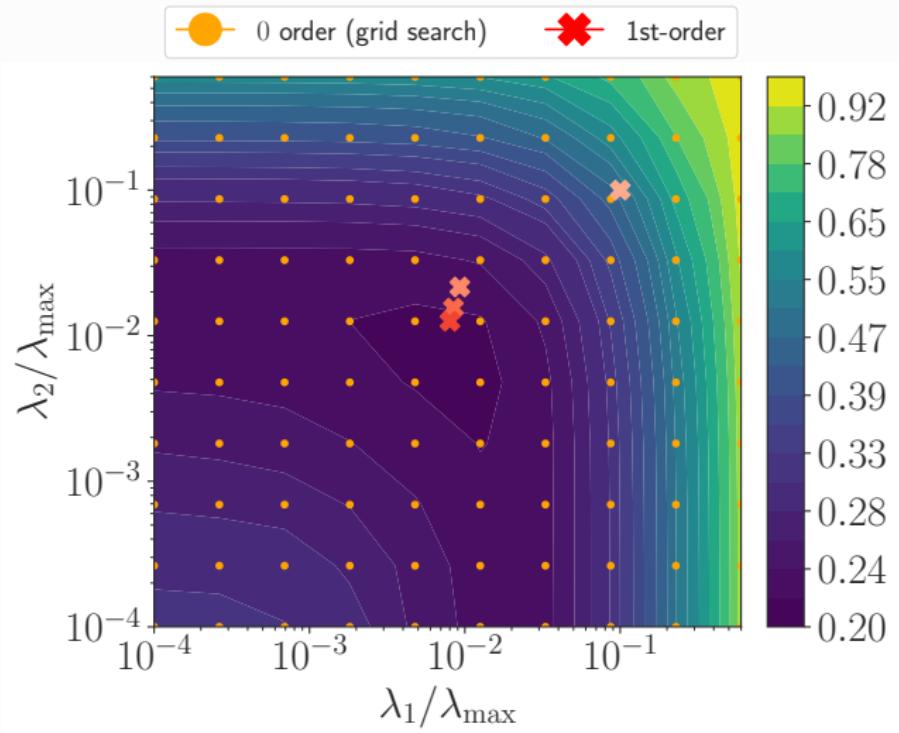
# FIRST-ORDER OPTIMIZATION IN $\lambda$ , ENET



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

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

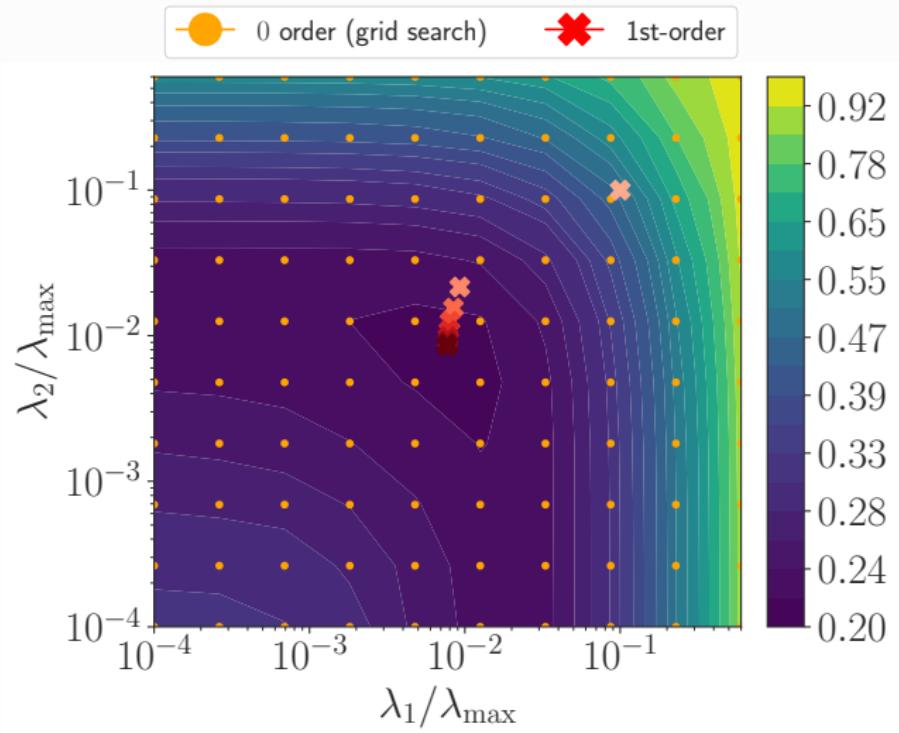
# FIRST-ORDER OPTIMIZATION IN $\lambda$ , ENET



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

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

# FIRST-ORDER OPTIMIZATION IN $\lambda$ , ENET



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

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

# Q: WHAT'S HARD?

A: COMPUTING  $\nabla_{\lambda} \mathcal{L}(\lambda)$



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

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

---

(1) J. Nocedal and S. J. Wright (2006). *Numerical optimization*. Second. Springer Series in Operations Research and Financial Engineering. Springer, Chap. 3.

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

# Q: WHAT'S HARD?

## A: COMPUTING $\nabla_\lambda \mathcal{L}(\lambda)$



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

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

- ▶ Line-search<sup>(1)</sup>
- ▶ L-BFGS<sup>(2)</sup>
- ▶ Gradient descent

---

(1) J. Nocedal and S. J. Wright (2006). *Numerical optimization*. Second. Springer Series in Operations Research and Financial Engineering. Springer, Chap. 3.

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

# Q: WHAT'S HARD?

## A: COMPUTING $\nabla_{\lambda} \mathcal{L}(\lambda)$



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

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

- ▶ Line-search<sup>(1)</sup>
- ▶ L-BFGS<sup>(2)</sup>
- ▶ Gradient descent

Main contribution here: compute  $\nabla_{\lambda} \mathcal{L}(\lambda)$  for a given  $\lambda$

---

(1) J. Nocedal and S. J. Wright (2006). *Numerical optimization*. Second. Springer Series in Operations Research and Financial Engineering. Springer, Chap. 3.

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

# How TO COMPUTE $\nabla_{\lambda} \mathcal{L}(\lambda)$ ?

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

# How TO COMPUTE $\nabla_{\lambda} \mathcal{L}(\lambda)$ ?

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

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

Chain rule:

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

- Boils down to:

**how to compute the Jacobian  $\hat{\mathcal{J}}_{(\lambda)} \in \mathbb{R}^{p \times 1}$  efficiently?**

# FORWARD-MODE DIFFERENTIATION <sup>(1), (2)</sup> OF PGD PROXIMAL GRADIENT DESCENT <sup>(3)</sup>

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

---

## Algorithm: Proximal gradient descent PGD

---

```
init :  $\beta = 0_p$ , ,  $L$ 
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  $\beta$ 
```

---

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

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

<sup>(3)</sup> 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 <sup>(1), (2)</sup> OF PGD

## PROXIMAL GRADIENT DESCENT <sup>(3)</sup>

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

---

### 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
return  $\beta$ 
```

---

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

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

<sup>(3)</sup> 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 <sup>(1), (2)</sup> OF PGD

## PROXIMAL GRADIENT DESCENT <sup>(3)</sup>



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

---

### 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
return  $\beta, \mathcal{J}$ 
```

---

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

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

<sup>(3)</sup> 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 <sup>(1), (2)</sup> OF PGD

## PROXIMAL GRADIENT DESCENT <sup>(3)</sup>

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

---

### 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
    +  $\partial_\lambda \text{prox}_{\lambda g / L}(z)$                                 // do not forget this term!
return  $\beta, \mathcal{J}$ 
```

---

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

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

<sup>(3)</sup> 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

# A BRIEF DÉTOUR

## PROXIMAL OPERATORS<sup>(2)</sup>



Motivation for proximal operators:

smooth non-smooth functions (sic),  
*e.g.*, to perform gradient descent

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

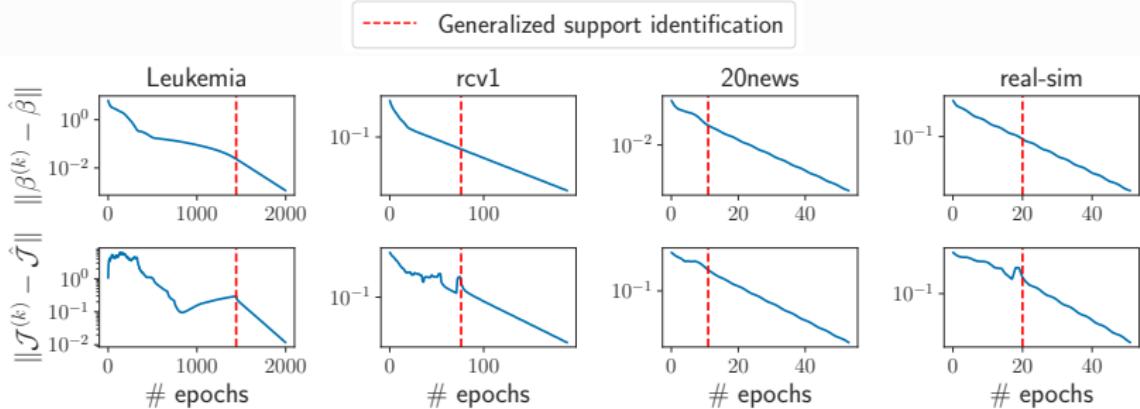


Jean-Jacques Moreau

<sup>(3)</sup> J.-J. Moreau (1962). "Fonctions convexes duales et points proximaux dans un espace hilbertien". In: *C. R. Acad. Sci. Paris* 255, pp. 2897–2899

### 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-zeros coefs.) has been identified, convergence is linear.<sup>(1)</sup>



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

# IMPLICIT DIFFERENTIATION (SMOOTH $\psi$ )<sup>(1)</sup>

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

---

<sup>(1)</sup> Y. Bengio (2000). "Gradient-based optimization of hyperparameters". In: *Neural computation*

# IMPLICIT DIFFERENTIATION (SMOOTH $\psi$ )<sup>(1)</sup>

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

$$\nabla_{\beta} \psi\left(\hat{\beta}^{(\lambda)}, \lambda\right) = 0$$

---

<sup>(1)</sup> Y. Bengio (2000). "Gradient-based optimization of hyperparameters". In: *Neural computation*

# IMPLICIT DIFFERENTIATION (SMOOTH $\psi$ )<sup>(1)</sup>



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

$$\nabla_{\beta} \psi\left(\hat{\beta}^{(\lambda)}, \lambda\right) = 0$$

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

---

<sup>(1)</sup> Y. Bengio (2000). "Gradient-based optimization of hyperparameters". In: *Neural computation*

# IMPLICIT DIFFERENTIATION (SMOOTH $\psi$ )<sup>(1)</sup>



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

$$\nabla_{\beta} \psi\left(\hat{\beta}^{(\lambda)}, \lambda\right) = 0$$

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

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

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

---

<sup>(1)</sup> Y. Bengio (2000). "Gradient-based optimization of hyperparameters". In: *Neural computation*

# IMPLICIT DIFFERENTIATION: BLESSING OF SPARSITY

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

**IMPLICIT DIFF.**<sup>(1)</sup>  $\arg \min_{\beta \in \mathbb{R}^p} f(\beta) + \lambda \sum_j |\beta_j|$



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

---

<sup>(1)</sup> Q. Bertrand, Q. Klopfenstein, M. Massias, et al. (2021). "Implicit differentiation for fast hyperparameter selection in non-smooth convex learning".  
In: *Submitted to JMLR*

**IMPLICIT DIFF.**<sup>(1)</sup>  $\arg \min_{\beta \in \mathbb{R}^p} f(\beta) + \lambda \sum_j |\beta_j|$



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

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

---

<sup>(1)</sup> Q. Bertrand, Q. Klopfenstein, M. Massias, et al. (2021). "Implicit differentiation for fast hyperparameter selection in non-smooth convex learning".  
In: *Submitted to JMLR*

**IMPLICIT DIFF.**<sup>(1)</sup>  $\arg \min_{\beta \in \mathbb{R}^p} f(\beta) + \lambda \sum_j |\beta_j|$



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

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

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

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

---

<sup>(1)</sup> Q. Bertrand, Q. Klopfenstein, M. Massias, et al. (2021). "Implicit differentiation for fast hyperparameter selection in non-smooth convex learning".  
In: *Submitted to JMLR*

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

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

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

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

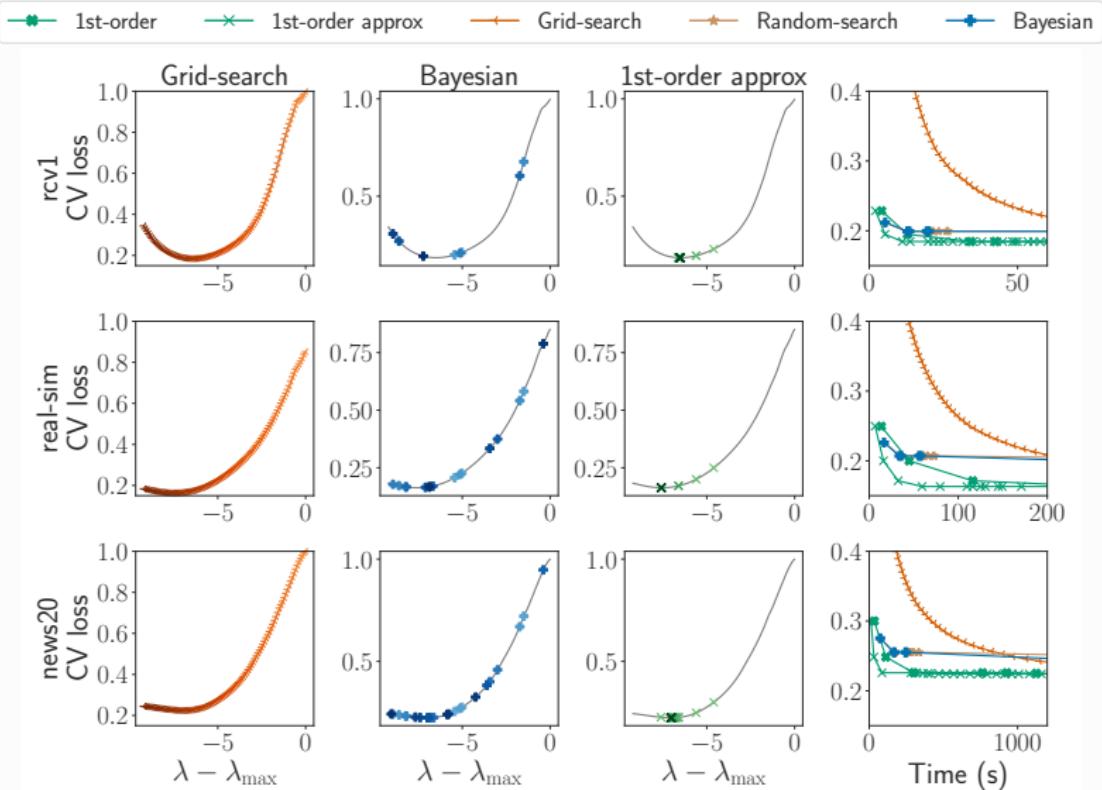
With  $\mathcal{S} = \left\{ j \in [p] : \hat{\beta}_j^{(\lambda)} = 0 \right\}$  we have  $\hat{\mathcal{J}}_{\mathcal{S}^c} = 0$

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

---

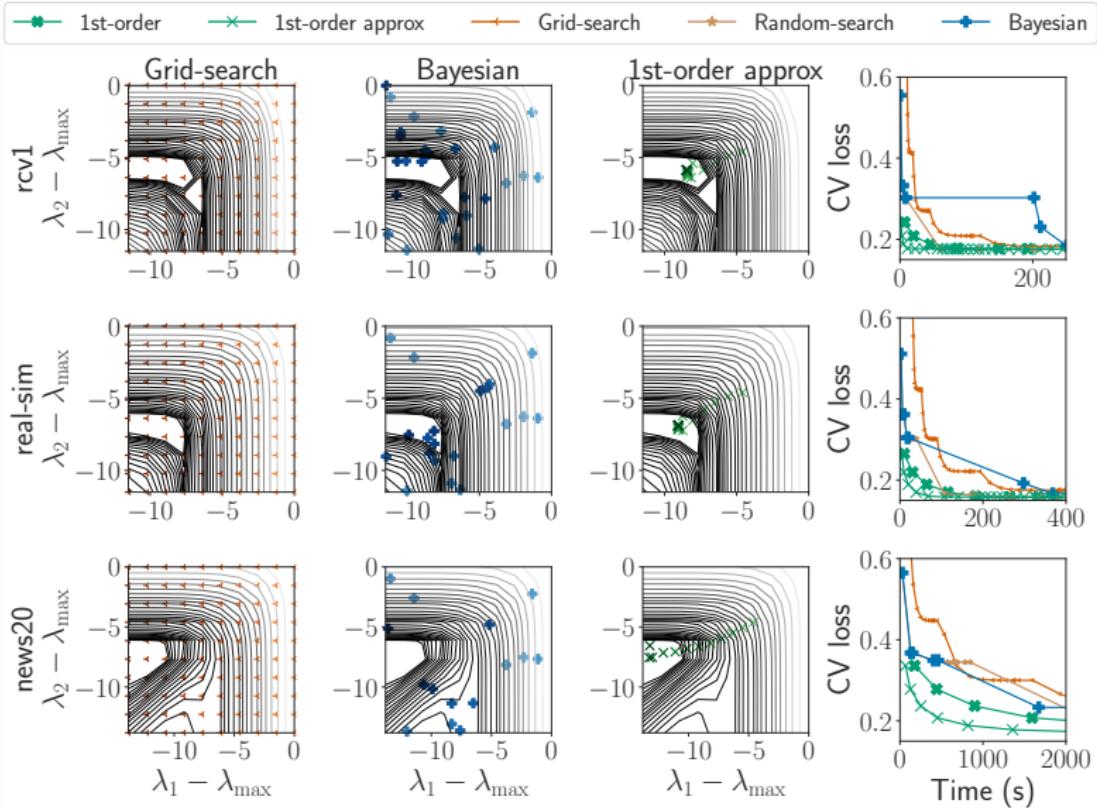
<sup>(1)</sup> Q. Bertrand, Q. Klopfenstein, M. Massias, et al. (2021). "Implicit differentiation for fast hyperparameter selection in non-smooth convex learning". In: *Submitted to JMLR*

# EXPERIMENTS I - LASSO CROSS-VALIDATION



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

# EXPERIMENTS II - ENET CROSS-VALIDATION

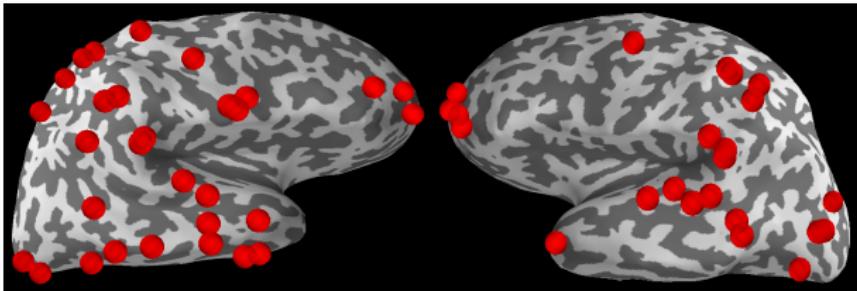


$$\arg \min_{\beta} \frac{1}{2n} \|y^{\text{train}} - X^{\text{train}}\beta\|^2 + e^{\lambda_1} \|\beta\|_1 + \frac{e^{\lambda_2}}{2} \|\beta\|^2$$

# EXPERIMENTS III - REAL MEEG DATA



- Outer criterion: FDMC SURE<sup>(1)</sup>
- Inner problems: vanilla Lasso



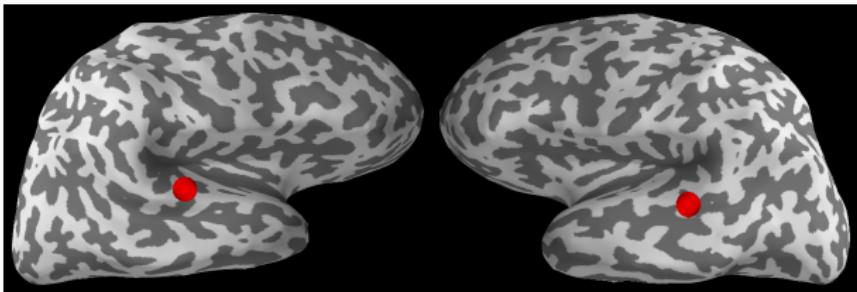
Real M/EEG data, vanilla Lasso (1 hyperparameter  $\lambda$ )

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

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

# EXPERIMENTS III - REAL MEEG DATA

- **Outer criterion:** FDMC SURE<sup>(1)</sup>
- **Inner problems:** weighted Lasso



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

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

<sup>(1)</sup> C.-A. Deledalle et al. (2014). "Stein Unbiased GrAdient estimator of the Risk (SUGAR) for multiple parameter selection". In: SIAM J. Imaging Sci.



- ▶ **Local linear convergence** of the Jacobian
- ▶ **Leverage sparsity** to speed up hypergradient computation
- ▶ Open source package  
<https://github.com/QB3/sparse-ho>

sparse-ho    0.1.dev    Examples    API    GitHub    Site ▾

## sparse-ho

build passing codecov 79%

*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.

Relying on a first order algorithm for bilevel optimization, *sparse-ho*'s performances scales gracefully with the number of hyperparameters to tune.

In order to benchmark performances, the package also implements alternatives such as forward or backward differentiation.

## Documentation

Please visit '<https://qb3.github.io/sparse-ho>' for the latest version of the documentation.

## Install

To run the code you first need to clone the repository, and then run, in the folder containing the `setup.py` file (root folder):

```
pip install -e .
```



*"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<sup>(1)</sup>

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

A. Gramfort (circa 2015)

---

<sup>(1)</sup> J. B. Buckheit and D. L. Donoho (1995). "Wavelab and reproducible research". In: *Wavelets and statistics*. Springer, pp. 55–81.

# How TO REACH ME



## Contact:

✉️ [joseph.salmon@umontpellier.fr](mailto:joseph.salmon@umontpellier.fr)  
🌐 <http://josephsalmon.eu>

**Github:** [@josephsalmon](#)



**Twitter:** [@salmonjsph](#)



-  Argyriou, A., T. Evgeniou, and M. Pontil (2008). "Convex multi-task feature learning". In: *Machine Learning*.
-  Bengio, Y. (2000). "Gradient-based optimization of hyperparameters". In: *Neural computation*.
-  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: *Journal of Machine Learning Research*.
-  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. (2021). "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*.
-  Brochu, E., V. M. Cora, and N. D. Freitas (2010). "A tutorial on Bayesian optimization of expensive cost functions, with application to active user modeling and hierarchical reinforcement learning". In: *arXiv preprint arXiv:1012.2599*.
-  Buckheit, J. B. and D. L. Donoho (1995). "Wavelab and reproducible research". In: *Wavelets and statistics*. Springer, pp. 55–81.
-  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.*
-  Figueiredo, M. (2001). "Adaptive Sparseness Using Jeffreys Prior". In: *Advances in Neural Information Processing Systems*.
-  Gramfort, A., M. Kowalski, and M. Hämäläinen (2012). "Mixed-norm estimates for the M/EEG inverse problem using accelerated gradient methods". In: *Phys. Med. Biol.*

-  Hutter, F., J. Lücke, and L. Schmidt-Thieme (2015). "Beyond manual tuning of hyperparameters". In: *KI-Künstliche Intelligenz*.
-  Kohavi, R. and G. H. John (1995). "Automatic parameter selection by minimizing estimated error". In: *Machine Learning Proceedings*.
-  Liu, D. C. and J. Nocedal (1989). "On the limited memory BFGS method for large scale optimization". In: *Mathematical programming*.
-  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.*
-  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. (2009). "Taking Advantage of Sparsity in Multi-Task Learning". In: *arXiv preprint arXiv:0903.1468*.

-  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, J.-J. (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 and Financial Engineering. New York: Springer.
-  Ochs, P. et al. (2015). "Bilevel optimization with nonsmooth lower level problems". In: *SSVM*.
-  Pedregosa, F. (2016). "Hyperparameter optimization with approximate gradient". In: *ICML*.
-  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. 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: *Journal of Machine Learning Research*.
-  Wengert, R. E. (1964). "A simple automatic derivative evaluation program". In: *Communications of the ACM* 7.8, pp. 463–464.