

# Implicit differentiation for hyperparameter optimization of Lasso-type models

**Quentin Bertrand** (Inria)

<https://QB3.github.io>

Joint work with:

**Quentin Klopfenstein** (Univ. Bourgogne Franche-Comté)

**Mathieu Blondel** (Google)

**Samuel Vaïter** (CNRS)

**Alexandre Gramfort** (Inria)

**Joseph Salmon** (IMAG, Univ. Montpellier, CNRS)

Motivation

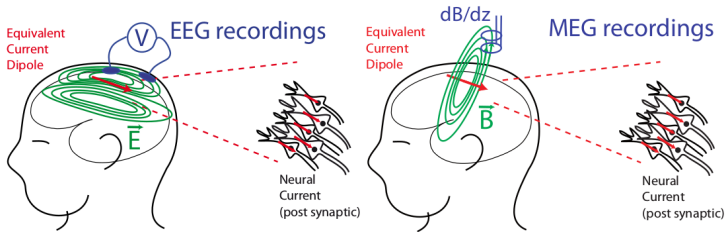
Hyperparameter optimization

Hypergradient computation

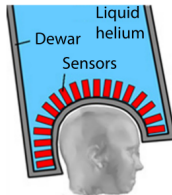
Experiments

# M/EEG inverse problem for brain imaging

- ▶ sensors: electric and magnetic fields during a cognitive task
- ▶ goal: which parts of the brain are responsible for the signals?
- ▶ applications: epilepsy treatment, brain aging, anesthesia risks

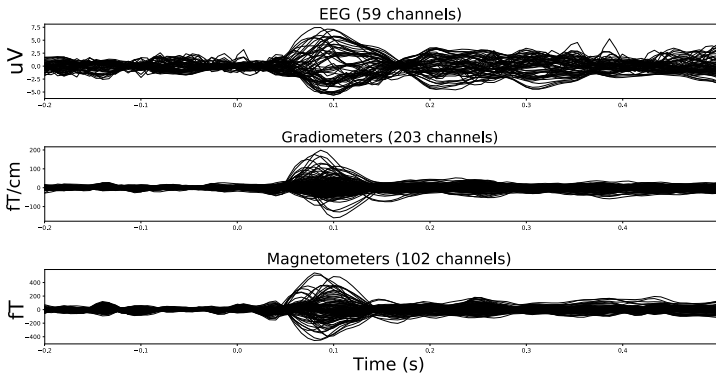


First EEG  
recordings  
in 1929  
by H. Berger



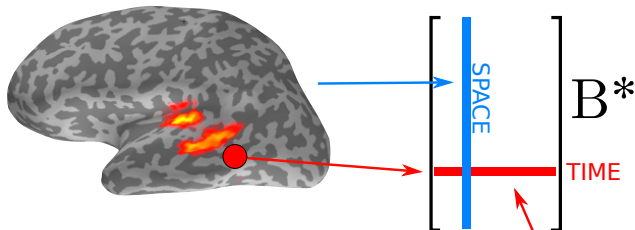
Hôpital La Timone  
Marseille, France

# M/EEG data

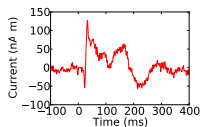


- 3 different types of sensor

# Source modeling (discretization with voxels)

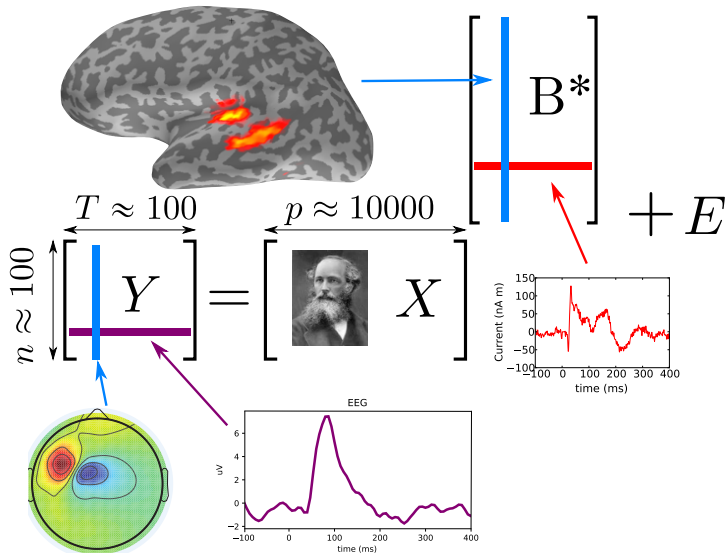


Position a few thousands candidate sources over the brain (e.g., every 5mm)



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

# The M/EEG inverse problem: modeling

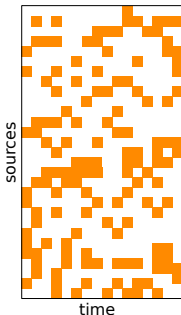


$$n \ll p$$

# Multi-Task penalties<sup>(1)</sup>

Popular convex penalties considered:

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



Sparse support: no structure

Penalty: **Lasso type**

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

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

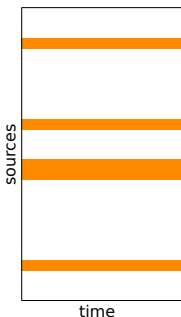
---

<sup>(1)</sup>G. Obozinski, B. Taskar, and M. I. Jordan. "Joint covariate selection and joint subspace selection for multiple classification problems". In: *Statistics and Computing* 20.2 (2010), pp. 231–252.

# Multi-Task penalties<sup>(1)</sup>

Popular convex penalties considered: Multi-Task Lasso (MTL)

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



Sparse support: group structure ✓

Penalty: **Group-Lasso type**

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

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

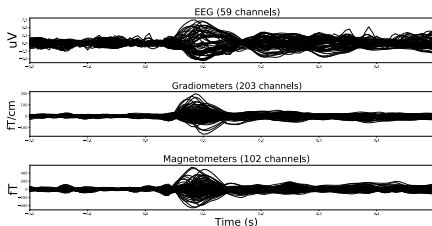
Parameter  $\hat{\mathbf{B}} \in \mathbb{R}^{p \times T}$

---

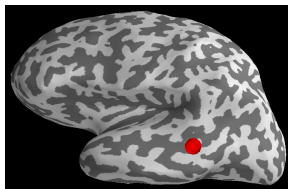
<sup>(1)</sup>G. Obozinski, B. Taskar, and M. I. Jordan. "Joint covariate selection and joint subspace selection for multiple classification problems". In: *Statistics and Computing* 20.2 (2010), pp. 231–252.



# Summary



What you have:  $Y$



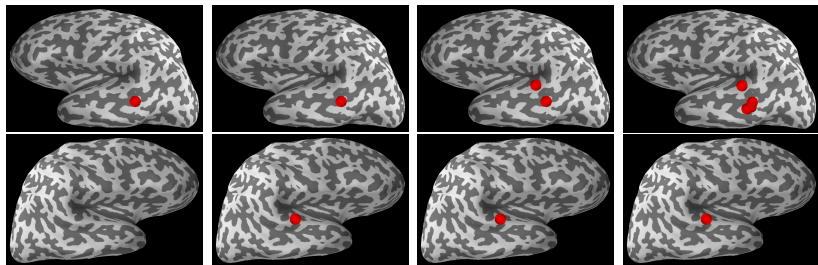
What you want  $B$

This is typically done using optimization based estimators:

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

## Which $\lambda$ to pick?

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



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

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

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

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

**Real MEEG 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

# Which $\lambda$ to pick? A statistical perspective<sup>(2)</sup>

## (i.i.d. case, Single-Task, $y = X\beta + \sigma^*\varepsilon$ )

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

---

---

### Theorem

---

---

- ▶ i.i.d. Gaussian noise
- ▶ +  $X$  satisfying the “Restricted Eigenvalue” property
- ▶ +  $\lambda = 2\sigma_* \sqrt{\frac{2 \log(p/\delta)}{n}}$
- ▶  $\implies$  with probability  $1 - \delta$ :

$$\frac{1}{n} \|X\beta^* - X\hat{\beta}^{(\lambda)}\|^2 \leq \frac{18}{\kappa_{s^*}^2} \frac{\sigma_*^2 s^*}{n} \log\left(\frac{p}{\delta}\right)$$

---

---

BUT  $\sigma_*$  is unknown in practice !

---

<sup>(2)</sup>P. J. Bickel, Y. Ritov, and A. B. Tsybakov. “Simultaneous analysis of Lasso and Dantzig selector”. In: *Ann. Statist.* 37.4 (2009), pp. 1705–1732.

# Which $\lambda$ to pick? A statistical perspective<sup>(2)</sup>

## (i.i.d. case, Single-Task, $y = X\beta + \sigma^*\varepsilon$ )

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

---

---

### Theorem

---

---

- ▶ i.i.d. Gaussian noise
- ▶ +  $X$  satisfying the “Restricted Eigenvalue” property
- ▶ +  $\lambda = 2\sigma_* \sqrt{\frac{2 \log(p/\delta)}{n}}$
- ▶  $\implies$  with probability  $1 - \delta$ :

$$\frac{1}{n} \|X\beta^* - X\hat{\beta}^{(\lambda)}\|^2 \leq \frac{18}{\kappa_{s^*}^2} \frac{\sigma_*^2 s^*}{n} \log\left(\frac{p}{\delta}\right)$$

---

---

**BUT**  $\sigma_*$  is unknown in practice !

---

<sup>(2)</sup>P. J. Bickel, Y. Ritov, and A. B. Tsybakov. “Simultaneous analysis of Lasso and Dantzig selector”. In: *Ann. Statist.* 37.4 (2009), pp. 1705–1732.

## Which $\lambda$ to pick? A statistical perspective II<sup>(3)</sup> (i.i.d. case, Single-Task)

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

---

---

### Theorem

---

---

- ▶ i.i.d. Gaussian noise
- ▶ +  $X$  satisfying the “Restricted Eigenvalue” property
- ▶ +  $\lambda = 2\sqrt{\frac{2\log(p/\delta)}{n}}$
- ▶  $\implies$  with probability  $1 - \delta$ :

$$\frac{1}{n} \|X\beta^* - X\hat{\beta}^{(\lambda)}\|^2 \leq \frac{18}{\kappa_{s^*}^2} \frac{\sigma_{s^*}^2}{n} \log\left(\frac{p}{\delta}\right)$$

---

---

$\lambda$  does not depend on  $\sigma_*$  anymore!

---

<sup>(3)</sup>A. Belloni, V. Chernozhukov, and L. Wang. “Square-root Lasso: pivotal recovery of sparse signals via conic programming”. In: *Biometrika* 98.4 (2011), pp. 791–806.

## Which $\lambda$ to pick? A statistical perspective II<sup>(3)</sup> (i.i.d. case, Single-Task)

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

---

---

### Theorem

---

---

- ▶ i.i.d. Gaussian noise
- ▶ +  $X$  satisfying the “Restricted Eigenvalue” property
- ▶ +  $\lambda = 2\sqrt{\frac{2\log(p/\delta)}{n}}$
- ▶  $\implies$  with probability  $1 - \delta$ :

$$\frac{1}{n} \|X\beta^* - X\hat{\beta}^{(\lambda)}\|^2 \leq \frac{18}{\kappa_{s^*}^2} \frac{\sigma_{s^*}^2}{n} \log\left(\frac{p}{\delta}\right)$$

---

---

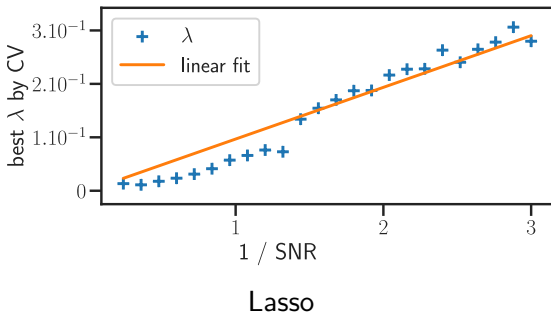
$\lambda$  does not depend on  $\sigma_*$  anymore!

---

<sup>(3)</sup>A. Belloni, V. Chernozhukov, and L. Wang. “Square-root Lasso: pivotal recovery of sparse signals via conic programming”. In: *Biometrika* 98.4 (2011), pp. 791–806.

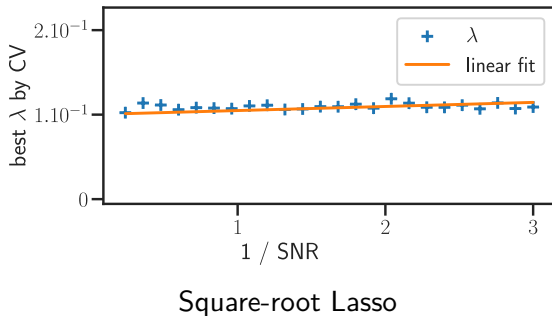
# Which $\lambda$ to pick? A statistical perspective III

$$\hat{\beta}_{\text{Lasso}} \in \arg \min_{\beta \in \mathbb{R}^p} \left( \frac{1}{2n} \|y - X\beta\|^2 + \lambda \|\beta\|_1 \right)$$



# Which $\lambda$ to pick? A statistical perspective III

$$\hat{\beta}_{\sqrt{\text{Lasso}}} \in \arg \min_{\beta \in \mathbb{R}^p} \left( \frac{1}{\sqrt{n}} \|y - X\beta\| + \lambda \|\beta\|_1 \right)$$





## Which $\lambda$ to pick? A statistical perspective III

- ▶  $\lambda \sim \sigma^*$  and  $\lambda$  independent of  $\sigma^*$  confirmed in practice ✓
- ▶ Strong statistical assumptions, not verified in practice ✗
- ▶ Still unknown quantities in the closed-form formula for  $\lambda$ : still needs calibration in practice ✗

# Hyperparameter optimization (HO)

Possible selection criterion:

- ▶ Good generalization<sup>(4)</sup> of  $\hat{\beta}^{(\lambda)}$
- ▶ AIC/BIC,<sup>(5)</sup> SURE<sup>(6)</sup> that control model complexity

---

<sup>(4)</sup>L. R. A. Stone and J.C. Ramer. "Estimating WAIS IQ from Shipley Scale scores: Another cross-validation". In: *Journal of clinical psychology* 21.3 (1965), pp. 297–297.

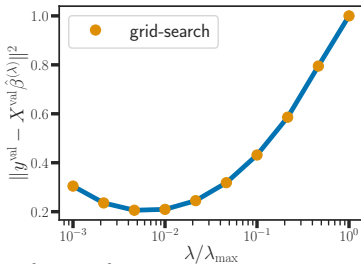
<sup>(5)</sup>W. Liu, Y. Yang, et al. "Parametric or nonparametric? A parametricness index for model selection". In: *Ann. Statist.* 39.4 (2011), pp. 2074–2102.

<sup>(6)</sup>C. M. Stein. "Estimation of the mean of a multivariate normal distribution". In: *Ann. Statist.* 9.6 (1981), pp. 1135–1151.

# Hyperparameter optimization (HO)

Possible selection criterion:

- ▶ Good generalization<sup>(4)</sup> of  $\hat{\beta}(\lambda)$
- ▶ AIC/BIC,<sup>(5)</sup> SURE<sup>(6)</sup> that control model complexity



**Real-sim dataset**

Validation loss as a function of  $\lambda$ .

**Example**

**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**

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

---

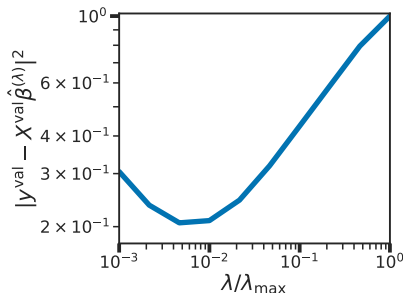
<sup>(4)</sup>L. R. A. Stone and J.C. Ramer. "Estimating WAIS IQ from Shipley Scale scores: Another cross-validation". In: *Journal of clinical psychology* 21.3 (1965), pp. 297–297.

<sup>(5)</sup>W. Liu, Y. Yang, et al. "Parametric or nonparametric? A parametricness index for model selection". In: *Ann. Statist.* 39.4 (2011), pp. 2074–2102.

<sup>(6)</sup>C. M. Stein. "Estimation of the mean of a multivariate normal distribution". In: *Ann. Statist.* 9.6 (1981), pp. 1135–1151.

# HO as a bilevel optimization problem<sup>(7)(8)</sup>

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

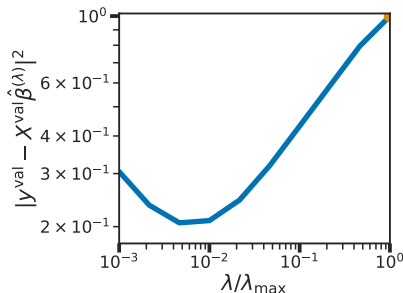


<sup>(7)</sup>P. Ochs et al. "On iteratively reweighted algorithms for nonsmooth nonconvex optimization in computer vision". In: *SIAM Journal on Imaging Sciences* 8.1 (2015), pp. 331–372.

<sup>(8)</sup>F. Pedregosa. "Hyperparameter optimization with approximate gradient". In: *ICML*. 2016.

# HO as a bilevel optimization problem<sup>(7)(8)</sup>

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

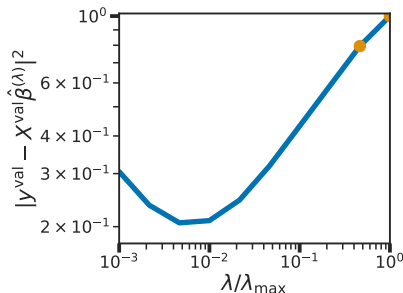


<sup>(7)</sup>P. Ochs et al. "On iteratively reweighted algorithms for nonsmooth nonconvex optimization in computer vision". In: *SIAM Journal on Imaging Sciences* 8.1 (2015), pp. 331–372.

<sup>(8)</sup>F. Pedregosa. "Hyperparameter optimization with approximate gradient". In: *ICML*. 2016.

# HO as a bilevel optimization problem<sup>(7)(8)</sup>

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

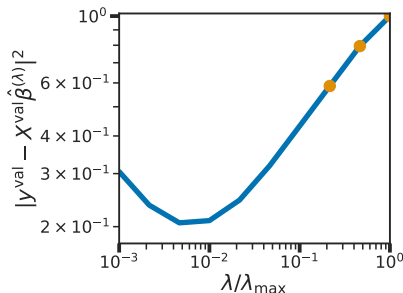


<sup>(7)</sup>P. Ochs et al. "On iteratively reweighted algorithms for nonsmooth nonconvex optimization in computer vision". In: *SIAM Journal on Imaging Sciences* 8.1 (2015), pp. 331–372.

<sup>(8)</sup>F. Pedregosa. "Hyperparameter optimization with approximate gradient". In: *ICML*. 2016.

# HO as a bilevel optimization problem<sup>(7)(8)</sup>

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

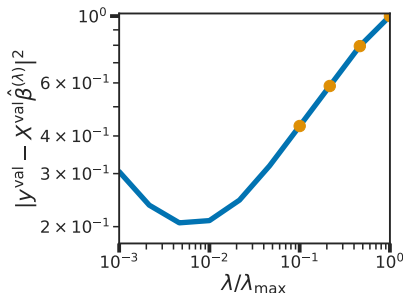


<sup>(7)</sup>P. Ochs et al. "On iteratively reweighted algorithms for nonsmooth nonconvex optimization in computer vision". In: *SIAM Journal on Imaging Sciences* 8.1 (2015), pp. 331–372.

<sup>(8)</sup>F. Pedregosa. "Hyperparameter optimization with approximate gradient". In: *ICML*. 2016.

# HO as a bilevel optimization problem<sup>(7)(8)</sup>

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



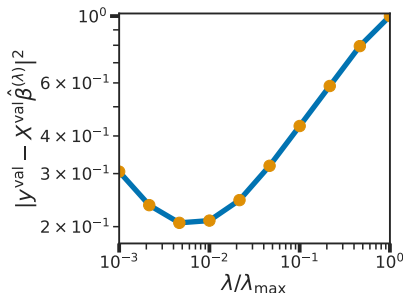
<sup>(7)</sup>P. Ochs et al. "On iteratively reweighted algorithms for nonsmooth nonconvex optimization in computer vision". In: *SIAM Journal on Imaging Sciences* 8.1 (2015), pp. 331–372.

<sup>(8)</sup>F. Pedregosa. "Hyperparameter optimization with approximate gradient". In: *ICML*. 2016.



# HO as a bilevel optimization problem<sup>(7)(8)</sup>

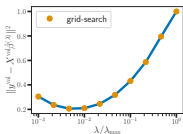
$$\begin{aligned} & \text{outer optimization problem} \\ \arg \min_{\lambda \in \mathbb{R}} & \left\{ \mathcal{L}(\lambda) := \|y^{\text{val}} - X^{\text{val}} \hat{\beta}^{(\lambda)}\|^2 \right\} \\ \text{s.t. } & \underbrace{\hat{\beta}^{(\lambda)} \in \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{2n}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



<sup>(7)</sup>P. Ochs et al. "On iteratively reweighted algorithms for nonsmooth nonconvex optimization in computer vision". In: *SIAM Journal on Imaging Sciences* 8.1 (2015), pp. 331–372.

<sup>(8)</sup>F. Pedregosa. "Hyperparameter optimization with approximate gradient". In: *ICML*. 2016.

# Grid-search as a 0-order optimization method



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

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

- ▶ Grid-search, random-search,<sup>(9)</sup> SMBO<sup>(10)</sup>:  
0-order methods to solve bilevel optimization problem
- ▶ **Idea:** if  $\mathcal{L}$  is differentiable, use first order optimization, *i.e.*, compute  $\nabla_{\lambda} \mathcal{L}$
- ▶ Once  $\nabla_{\lambda} \mathcal{L}(\lambda)$  is computed, use gradient descent<sup>(11)</sup>:  
 $\lambda^{(t+1)} = \lambda^{(t)} - \rho \nabla_{\lambda} \mathcal{L}(\lambda^{(t)})$  with suitable  $\rho > 0$

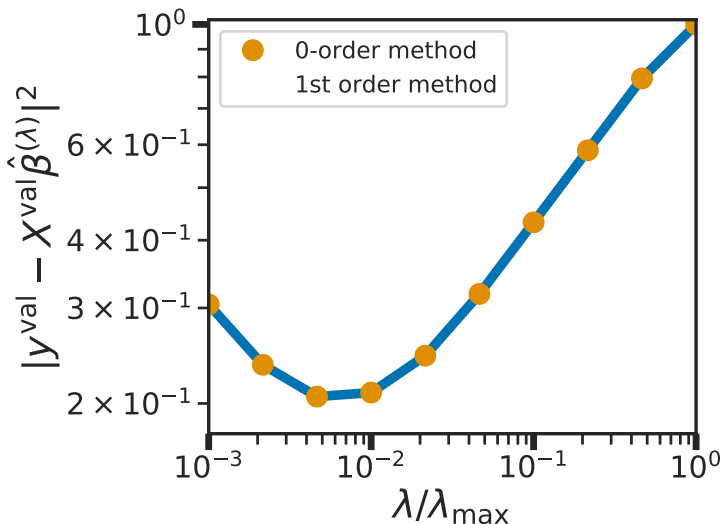
---

<sup>(9)</sup> J. Bergstra and Y. Bengio. “Random search for hyper-parameter optimization”. In: *J. Mach. Learn. Res.* (2012).

<sup>(10)</sup> E. Brochu, V. M. Cora, and N. De Freitas. “A tutorial on Bayesian optimization of expensive cost functions, with application to active user modeling and hierarchical reinforcement learning”. In: (2010).

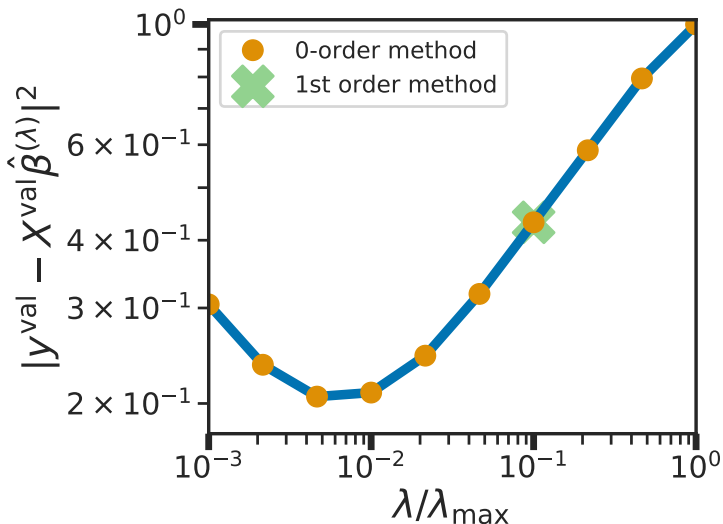
<sup>(11)</sup> F. Pedregosa. “Hyperparameter optimization with approximate gradient”. In: *ICML*. 2016.

## First order optimization in $\lambda$ , Lasso



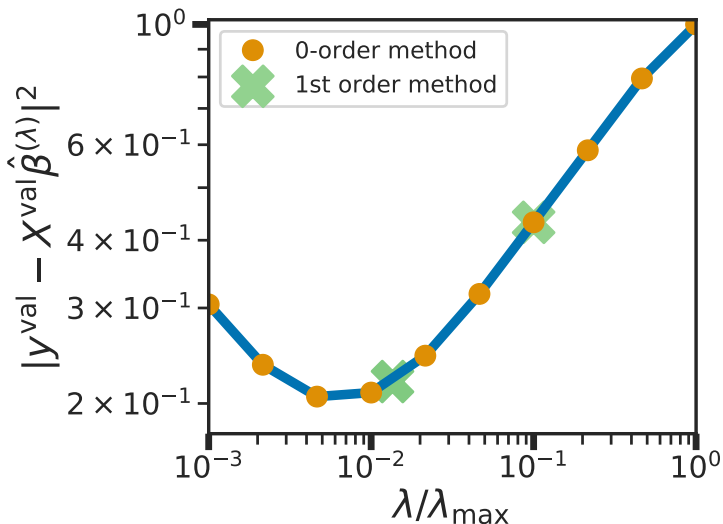
**Real-sim dataset.** Validation loss as a function of  $\lambda$ .

## First order optimization in $\lambda$ , Lasso



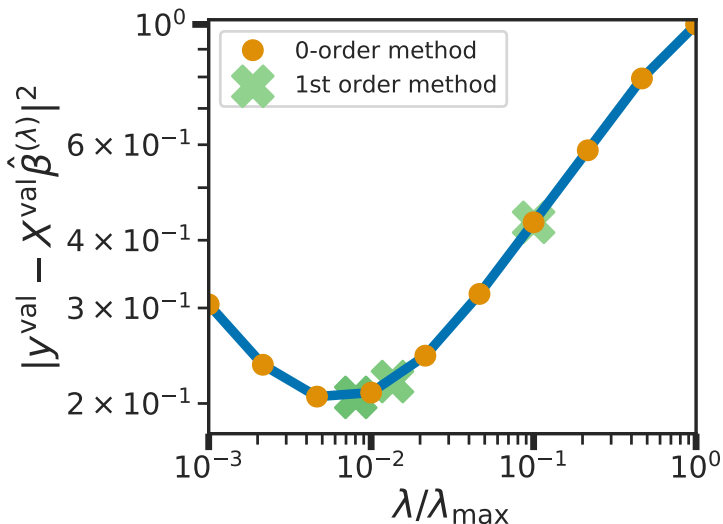
**Real-sim dataset.** Validation loss as a function of  $\lambda$ .

## First order optimization in $\lambda$ , Lasso



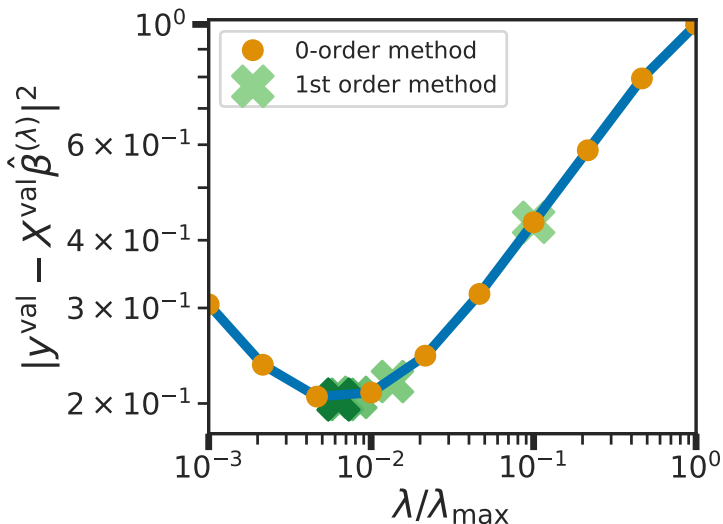
**Real-sim dataset.** Validation loss as a function of  $\lambda$ .

## First order optimization in $\lambda$ , Lasso



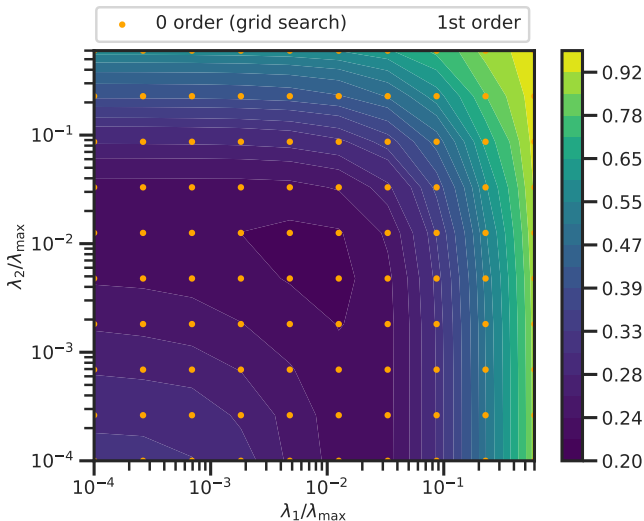
**Real-sim dataset.** Validation loss as a function of  $\lambda$ .

## First order optimization in $\lambda$ , Lasso



**Real-sim dataset.** Validation loss as a function of  $\lambda$ .

# First order optimization in $\lambda$ , Enet

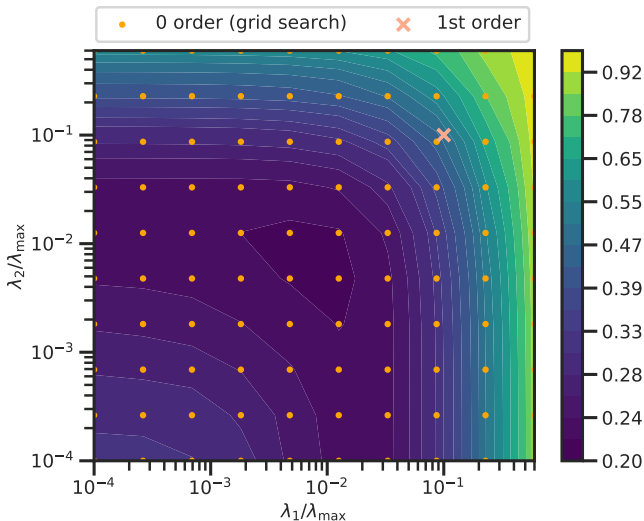


Real-sim dataset, level sets of the validation loss (held out)

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



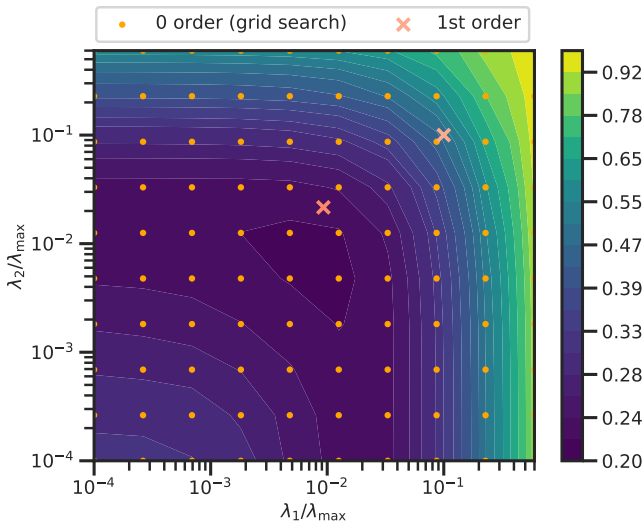
# First order optimization in $\lambda$ , Enet



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

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

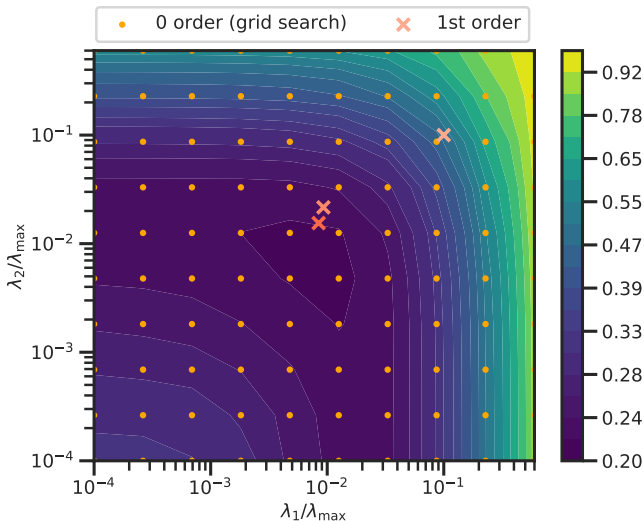
# First order optimization in $\lambda$ , Enet



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

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

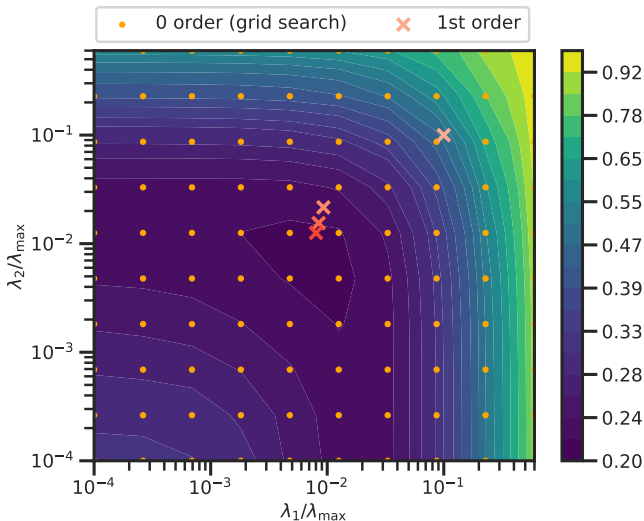
# First order optimization in $\lambda$ , Enet



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

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

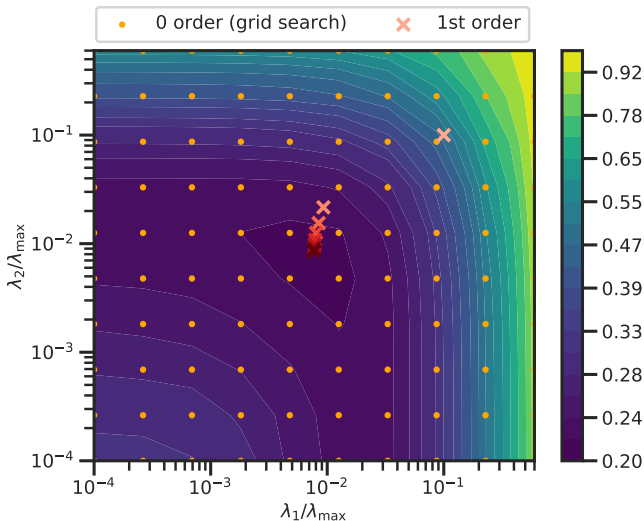
# First order optimization in $\lambda$ , Enet



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

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

# First order optimization in $\lambda$ , Enet



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

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

# What's hard? Computing $\nabla_{\lambda}\mathcal{L}(\lambda)$

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

Once  $\nabla_{\lambda}\mathcal{L}(\lambda)$  is computed life is "easy":

- ▶ Line-search<sup>(12)</sup>
- ▶ LBFGS<sup>(13)</sup>
- ▶ Gradient descent

---

<sup>(12)</sup>J. Nocedal and S. J. Wright. *Numerical optimization*. Second. Springer Series in Operations Research and Financial Engineering. New York: Springer, 2006, Chap. 3.

<sup>(13)</sup>D. Goldfarb. "A family of variable-metric methods derived by variational means". In: *Mathematics of computation* 24.109 (1970), pp. 23–26.

# What's hard? Computing $\nabla_{\lambda}\mathcal{L}(\lambda)$

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

Once  $\nabla_{\lambda}\mathcal{L}(\lambda)$  is computed life is "easy":

- ▶ Line-search<sup>(12)</sup>
- ▶ LBFGS<sup>(13)</sup>
- ▶ Gradient descent

The main challenge is to compute  $\nabla_{\lambda}\mathcal{L}(\lambda)$  for a given  $\lambda$ !

---

<sup>(12)</sup> J. Nocedal and S. J. Wright. *Numerical optimization*. Second. Springer Series in Operations Research and Financial Engineering. New York: Springer, 2006, Chap. 3.

<sup>(13)</sup> D. Goldfarb. "A family of variable-metric methods derived by variational means". In: *Mathematics of computation* 24.109 (1970), pp. 23–26.

# What's hard? Computing $\nabla_{\lambda}\mathcal{L}(\lambda)$

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

Once  $\nabla_{\lambda}\mathcal{L}(\lambda)$  is computed life is "easy":

- ▶ Line-search<sup>(12)</sup>
- ▶ LBFGS<sup>(13)</sup>
- ▶ Gradient descent

The main challenge is to compute  $\nabla_{\lambda}\mathcal{L}(\lambda)$  for a given  $\lambda$ !

---

<sup>(12)</sup> J. Nocedal and S. J. Wright. *Numerical optimization*. Second. Springer Series in Operations Research and Financial Engineering. New York: Springer, 2006, Chap. 3.

<sup>(13)</sup> D. Goldfarb. "A family of variable-metric methods derived by variational means". In: *Mathematics of computation* 24.109 (1970), pp. 23–26.



## 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 and Jacobian:

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

► Boils down to:

how to compute the Jacobian  $\hat{\mathcal{J}}_{(\lambda)}$  efficiently?

## 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 and Jacobian:

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

► Boils down to:

**how to compute the Jacobian  $\hat{\mathcal{J}}_{(\lambda)}$  efficiently?**

# How to compute $\hat{\mathcal{J}}_{(\lambda)} := (\nabla_{\lambda}\hat{\beta}_1^{(\lambda)}, \dots, \nabla_{\lambda}\hat{\beta}_p^{(\lambda)})$ ?

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

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

- ▶ *Implicit differentiation* (**closed-form** formula)<sup>(14)</sup>:  
need to solve a  $p \times p$  linear system ( $p = \# \text{features}$ )
- ▶ *Automatic differentiation*, *forward*<sup>(15)</sup> or *backward*<sup>(16)</sup>

---

<sup>(14)</sup> J. Larsen et al. "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*. 1996; Y. Bengio. "Gradient-based optimization of hyperparameters". In: *Neural computation* 12.8 (2000), pp. 1889–1900.

<sup>(15)</sup> L. Franceschi et al. "Forward and reverse gradient-based hyperparameter optimization". In: *ICML*. 2017, pp. 1165–1173.

<sup>(16)</sup> J. Domke. "Generic methods for optimization-based modeling". In: *AISTATS*. vol. 22. 2012, pp. 318–326.

# How to compute $\hat{\mathcal{J}}_{(\lambda)} := (\nabla_{\lambda}\hat{\beta}_1^{(\lambda)}, \dots, \nabla_{\lambda}\hat{\beta}_p^{(\lambda)})$ ?

$$\underbrace{\hat{\beta}^{(\lambda)} \in \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{2n} \|y^{\text{train}} - X^{\text{train}} \beta\|^2 + \lambda \|\beta\|_1}_{\text{inner optimization problem}}$$

**"Nonsmooth"** inner optimization problems, **scarce literature**:

- ▶ *Smooth the nonsmooth term*<sup>(17)</sup>
- ▶ Use algorithms with differentiable updates<sup>(18)(19)</sup> (Bregman)

Our contributions:

- ▶ Iterative differentiation can be applied on classical proximal algorithms!
- ▶ Key point on the Jacobian:

$\hat{\mathcal{J}}_{(\lambda)} := (\nabla_{\lambda}\hat{\beta}_1^{(\lambda)}, \dots, \nabla_{\lambda}\hat{\beta}_p^{(\lambda)})$  shares  $\hat{\beta}^{(\lambda)}$ 's **sparsity pattern**

---

<sup>(17)</sup>G. Peyré and J. M. Fadili. "Learning analysis sparsity priors". In: *Sampta*. 2011.

<sup>(18)</sup>P. Ochs et al. "Bilevel optimization with nonsmooth lower level problems". In: *International Conference on Scale Space and Variational Methods in Computer Vision*. 2015, pp. 654–665.

<sup>(19)</sup>J. Frecon, S. Salzo, and M. Pontil. "Bilevel learning of the group lasso structure". In: *Advances in Neural Information Processing Systems*. 2018, pp. 8301–8311.

# Iterative forward differentiation on PGD

$$\hat{\beta}^{(\lambda)} \in \arg \min_{\beta \in \mathbb{R}^p} \overbrace{f}^{\text{smooth}}(\beta) + \lambda \overbrace{g}^{\text{non-smooth}}(\beta) \quad (1)$$

---

**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

**return**  $\beta$

---

# Iterative forward differentiation on PGD

$$\hat{\beta}^{(\lambda)} \in \arg \min_{\beta \in \mathbb{R}^p} \overbrace{f}^{\text{smooth}}(\beta) + \lambda \overbrace{g}^{\text{non-smooth}}(\beta) \quad (1)$$

---

**Algorithm:** Iterative forward diff. (for PGD)

---

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

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

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

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

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

**return**  $\beta$

---

# Iterative forward differentiation on PGD

$$\hat{\beta}^{(\lambda)} \in \arg \min_{\beta \in \mathbb{R}^p} \overbrace{f}^{\text{smooth}}(\beta) + \lambda \overbrace{g}^{\text{non-smooth}}(\beta) \quad (1)$$

---

**Algorithm:** Iterative forward diff. (for PGD)

---

**init** :  $\beta = 0_p, \mathcal{J} = 0_p, L$

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

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

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

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

$\mathcal{J} \leftarrow \partial_z \text{prox}_{\lambda g/L}(z) dz$  // diff w.r.t.  $\lambda$ : chain rule

**return**  $\beta, \mathcal{J}$

---

# Iterative forward differentiation on PGD

$$\hat{\beta}^{(\lambda)} \in \arg \min_{\beta \in \mathbb{R}^p} \overbrace{f}^{\text{smooth}}(\beta) + \lambda \overbrace{g}^{\text{non-smooth}}(\beta) \quad (1)$$

---

**Algorithm:** Forward iterative differentiation (for PGD)

---

**init** :  $\beta = 0_p, \mathcal{J} = 0_p, L$

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

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

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

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

$\mathcal{J} \leftarrow \partial_z \text{prox}_{\lambda g/L}(z) dz$  // diff w.r.t.  $\lambda$ : chain rule  
         $+ \partial_\lambda \text{prox}_{\lambda g/L}(z)$  // do not forget this term!

**return**  $\beta, \mathcal{J}$

---



# Forward iterative differentiation on BCD

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

- ▶ Iterative forward can generalize to **coordinate descent** (BCD, state of art algorithm for the Lasso)
- ▶ Convergence of the Jacobian sequence  $\mathcal{J}$ ?

# Forward iterative differentiation on BCD

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

- ▶ Iterative forward can generalize to **coordinate descent** (BCD, state of art algorithm for the Lasso)
- ▶ **Convergence** of the Jacobian sequence  $\mathcal{J}$ ?

## Contribution

- ▶ Prove Jacobian sequence convergence

# Forward iterative differentiation on BCD

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

- ▶ Iterative forward can generalize to **coordinate descent** (BCD, state of art algorithm for the Lasso)
- ▶ **Convergence** of the Jacobian sequence  $\mathcal{J}$ ?

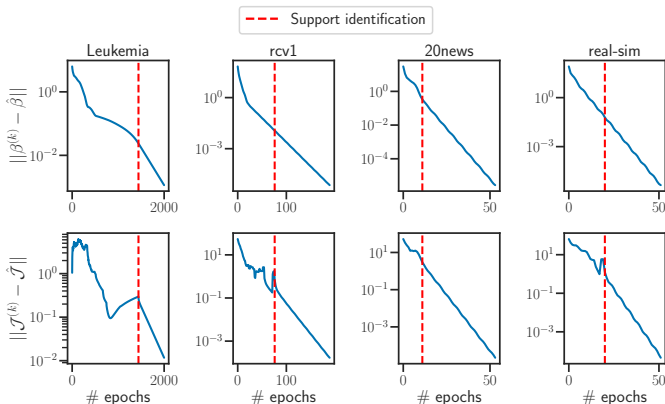
## Contribution

- ▶ Prove Jacobian sequence convergence

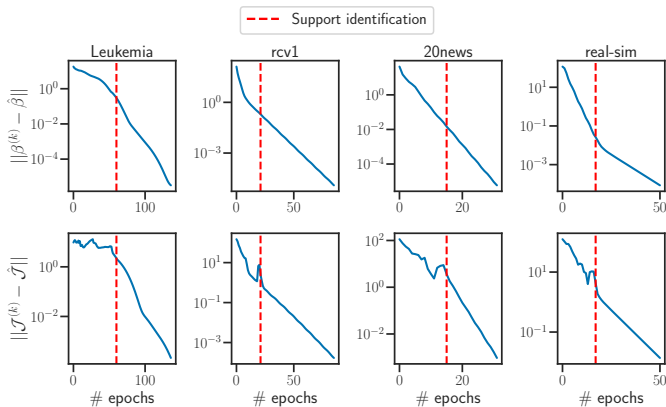
# Local linear convergence of the Jacobian (I)

## Proposition: forward diff. convergence (Lasso)

Assuming that the Lasso inner optimization has a unique minimizer, then the Jacobian sequence based on forward diff. of BCD converges to the true Jacobian. Once the support (*i.e.*, non-zeros coefs.) has been identified, convergence is linear.



# Local linear convergence of the Jacobian (II)



**Exemple:** sparse logistic regression

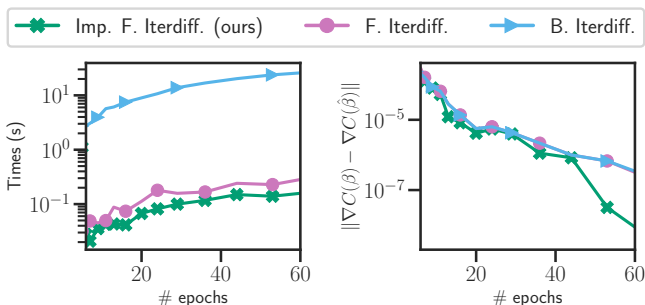
$$\arg \min_{\beta \in \mathbb{R}^p} \frac{1}{n} \sum_i \frac{1}{1 + \exp(-y_i X_{i,:} \beta)} + \lambda \|\beta\|_1$$

# Proposed algorithm: Implicit forward diff.

- ▶ Jacobian  $\hat{\mathcal{J}}^{(\lambda)}$  shares  $\hat{\beta}^{(\lambda)}$  sparsity pattern
- ▶ Leverage sparsity to **speed up computation**

2-step algorithm:

1. Solve the inner Lasso problem to get  $\hat{\beta}^{(\lambda)}$  and its support  $\hat{S}^{(\lambda)}$
2. Compute Jacobian only on the support  $\hat{S}^{(\lambda)}$  using the forward iterations of coordinate descent



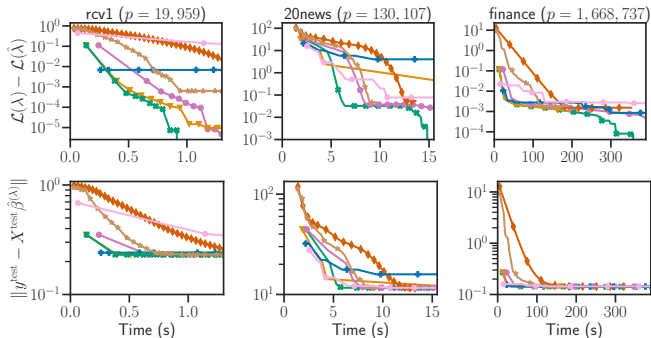
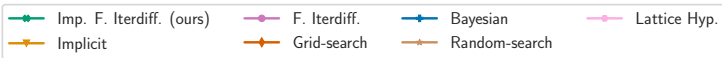
Convergence on synthetic data

# Experiments I - Real datasets

- **Outer criterion:** held-out loss. **Inner problems:** the Lasso

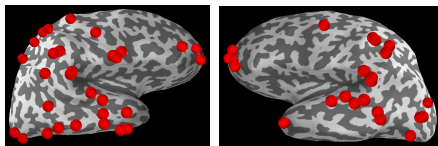
$$\arg \min_{\lambda \in \mathbb{R}} \mathcal{L}(\lambda) := \|y^{\text{val}} - X^{\text{val}} \hat{\beta}^{(\lambda)}\|^2$$

$$\text{s.t. } \hat{\beta}^{(\lambda)} \in \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{2n} \|y^{\text{train}} - X^{\text{train}} \beta\|_2^2 + \lambda \|\beta\|_1$$



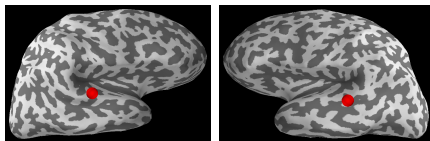
# Experiments III - Real MEEG data

- **Outer criterion:** SURE
- **Inner problems:** the Lasso and weighted Lasso



**Vanilla Lasso (1 parameter)**

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



**Weighted Lasso ( $p$  parameters)**

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



# Limitations

- ▶ May require specific parametrization  $e^\lambda$
- ▶ 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<sup>(20)</sup> nor reweighted Lasso<sup>(21)</sup>
- ▶ Optimized function often **non-convex**: possibly multiple local minima
- ▶ Rely on **line-search**: hidden hyperparameters control the convergence speed

---

<sup>(20)</sup>P. Breheny and J. Huang. "Coordinate descent algorithms for nonconvex penalized regression, with applications to biological feature selection". In: *Ann. Appl. Stat.* 5.1 (2011), p. 232.

<sup>(21)</sup>E. J. Candès, M. B. Wakin, and S. P. Boyd. "Enhancing Sparsity by Reweighted  $l_1$  Minimization". In: *J. Fourier Anal. Applicat.* 14.5-6 (2008), pp. 877–905.

# Conclusion

Hyperparameter optimization cast as a **bilevel optimization problem**, on which we applied 1st order optimization:

- ▶ Proved **locally linear convergence** of the Jacobian
- ▶ **Leverage sparsity** to speed up Jacobian computation

Future work:

- ▶ Extension to other **sparse models**  
(group Lasso, sparse multiclass logistic regression)
- ▶ Extend work on **inexact gradient** to non-smooth inner pb

# Conclusion

Hyperparameter optimization cast as a **bilevel optimization problem**, on which we applied 1st order optimization:

- ▶ Proved **locally linear convergence** of the Jacobian
- ▶ **Leverage sparsity** to speed up Jacobian computation

Future work:

- ▶ Extension to other **sparse models**  
(group Lasso, sparse multiclass logistic regression)
- ▶ Extend work on **inexact gradient** to non-smooth inner pb
- ▶ Paper <https://proceedings.icml.cc/paper/2020/file/e0ab531ec312161511493b002f9be2ee-Paper.pdf>
- ▶ Open source package <https://github.com/QB3/sparse-ho>

# Conclusion

Hyperparameter optimization cast as a **bilevel optimization problem**, on which we applied 1st order optimization:

- ▶ Proved **locally linear convergence** of the Jacobian
- ▶ **Leverage sparsity** to speed up Jacobian computation

## Future work:

- ▶ Extension to other **sparse models**  
(group Lasso, sparse multiclass logistic regression)
- ▶ Extend work on **inexact gradient** to non-smooth inner pb
- ▶ Paper <https://proceedings.icml.cc/paper/2020/file/e0ab531ec312161511493b002f9be2ee-Paper.pdf>
- ▶ Open source package <https://github.com/QB3/sparse-ho>

- ▶ Belloni, A., V. Chernozhukov, and L. Wang. “Square-root Lasso: pivotal recovery of sparse signals via conic programming”. In: *Biometrika* 98.4 (2011), pp. 791–806.
- ▶ Bengio, Y. “Gradient-based optimization of hyperparameters”. In: *Neural computation* 12.8 (2000), pp. 1889–1900.
- ▶ Bergstra, J. and Y. Bengio. “Random search for hyper-parameter optimization”. In: *J. Mach. Learn. Res.* (2012).
- ▶ Bickel, P. J., Y. Ritov, and A. B. Tsybakov. “Simultaneous analysis of Lasso and Dantzig selector”. In: *Ann. Statist.* 37.4 (2009), pp. 1705–1732.
- ▶ Breheny, P. and J. Huang. “Coordinate descent algorithms for nonconvex penalized regression, with applications to biological feature selection”. In: *Ann. Appl. Stat.* 5.1 (2011), p. 232.
- ▶ Brochu, E., V. M. Cora, and N. De Freitas. “A tutorial on Bayesian optimization of expensive cost functions, with application to active user modeling and hierarchical reinforcement learning”. In: (2010).

- ▶ Candès, E. J., M. B. Wakin, and S. P. Boyd. “Enhancing Sparsity by Reweighted  $l_1$  Minimization”. In: *J. Fourier Anal. Applicat.* 14.5-6 (2008), pp. 877–905.
- ▶ Domke, J. “Generic methods for optimization-based modeling”. In: *AISTATS*. Vol. 22. 2012, pp. 318–326.
- ▶ Franceschi, L. et al. “Forward and reverse gradient-based hyperparameter optimization”. In: *ICML*. 2017, pp. 1165–1173.
- ▶ Frecon, J., S. Salzo, and M. Pontil. “Bilevel learning of the group lasso structure”. In: *Advances in Neural Information Processing Systems*. 2018, pp. 8301–8311.
- ▶ Goldfarb, D. “A family of variable-metric methods derived by variational means”. In: *Mathematics of computation* 24.109 (1970), pp. 23–26.
- ▶ Larsen, J. et al. “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*. 1996.

- ▶ Liu, W., Y. Yang, et al. “Parametric or nonparametric? A parametricness index for model selection”. In: *Ann. Statist.* 39.4 (2011), pp. 2074–2102.
- ▶ Nocedal, J. and S. J. Wright. *Numerical optimization*. Second. Springer Series in Operations Research and Financial Engineering. New York: Springer, 2006.
- ▶ Obozinski, G., B. Taskar, and M. I. Jordan. “Joint covariate selection and joint subspace selection for multiple classification problems”. In: *Statistics and Computing* 20.2 (2010), pp. 231–252.
- ▶ Ochs, P. et al. “Bilevel optimization with nonsmooth lower level problems”. In: *International Conference on Scale Space and Variational Methods in Computer Vision*. 2015, pp. 654–665.
- ▶ Ochs, P. et al. “On iteratively reweighted algorithms for nonsmooth nonconvex optimization in computer vision”. In: *SIAM Journal on Imaging Sciences* 8.1 (2015), pp. 331–372.
- ▶ Pedregosa, F. “Hyperparameter optimization with approximate gradient”. In: *ICML*. 2016.

- ▶ Peyré, G. and J. M. Fadili. “Learning analysis sparsity priors”. In: *Sampta*. 2011.
- ▶ Stein, C. M. “Estimation of the mean of a multivariate normal distribution”. In: *Ann. Statist.* 9.6 (1981), pp. 1135–1151.
- ▶ Stone, L. R. A. and J.C. Ramer. “Estimating WAIS IQ from Shipley Scale scores: Another cross-validation”. In: *Journal of clinical psychology* 21.3 (1965), pp. 297–297.