

A guided tour of machine learning (theory)

Lorenzo Rosasco
MaLGa, Università degli Studi di Genova, MIT, IIT

RoMaDS Colloquim, University of Rome Tor Vergata, May 9th 2022

AI everywhere (literally...)

AI

iCub >1000 sensors

Data Science

Hep 30Pb/Y

ERC Interview_SLING8.pdf (page 3 of 28)

50 billion microcontrollers in 2019
>100 billions in service

predictive maintenance, connected cars, precision agriculture, personalized fitness and wearables, smart housing, cities, healthcare, etc.

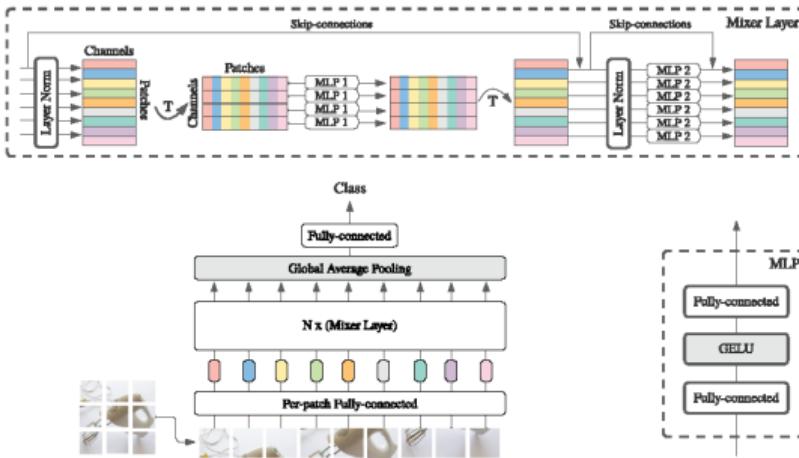
MCU Market History and Forecast

Year	Market [M]	Units [M]	ASP [\$]
12	~90	~15.98	~\$0.55
13	~100	~19.18	~\$0.58
14	~110	~16.68	~\$0.52
15F	~120	~25.48	~\$0.50
16F	~130	~31.88	~\$0.48
17F	~140	~25.78	~\$0.45
18F	~150	~31.28	~\$0.43
19F	~160	~31.28	~\$0.42

UniGe | MoGa

2

The state of affairs



Rethinking machine learning:

- ▶ with statistical mechanics
- ▶ with information theory
- ▶ with tropical geometry
- ▶ ...

Outline

Data driven modeling paradigm

Statistical learning, bias and variance

From statistics to optimization

A theory crisis?

The basic picture

$$(x_i, y_i)_{i=1}^n \quad \mapsto \quad f : X \rightarrow Y$$

Fixing a model

$$w \in \mathbb{R}^p \mapsto f_w$$

Fixing a model

$$w \in \mathbb{R}^p \mapsto f_w$$

$$f_w(x) = \sum_{j=1}^p w^j \phi_j(x)$$

Fixing a model

$$w \in \mathbb{R}^p \mapsto f_w$$

$$f_w(x) = \sum_{j=1}^p w^j \phi_j(x)$$

$$f_w(x) = \sum_{j=1}^p \beta^j \sigma(a_j^\top x + \alpha_j),$$

Fitting (or overfitting??)

$$\min_w \frac{1}{n} \sum_{i=1}^n (f_w(x_i) - y_i)^2$$

Fitting (or overfitting??)

$$\min_w \frac{1}{n} \sum_{i=1}^n (f_w(x_i) - y_i)^2$$

Overparametrization

w has often millions of parameters...data are often (much) less!

Fitting (or overfitting??)

$$\min_w \frac{1}{n} \sum_{i=1}^n (f_w(x_i) - y_i)^2$$

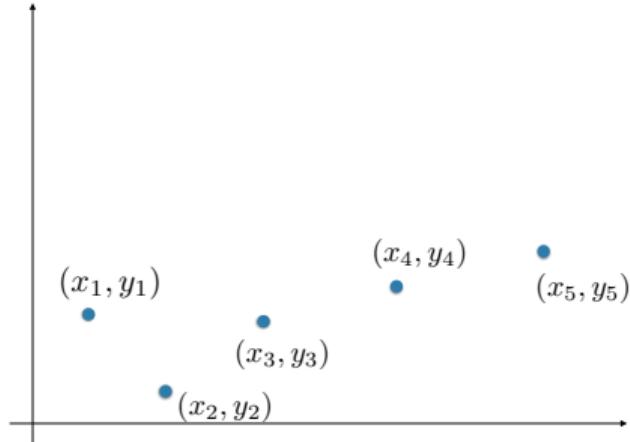
Overparametrization

w has often millions of parameters...data are often (much) less!

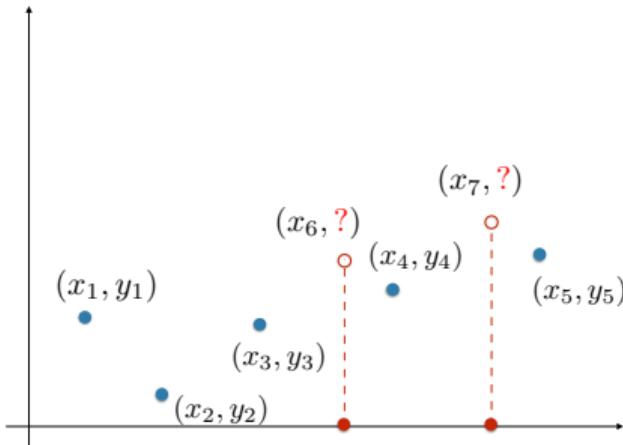
"With four parameters I can fit an elephant, and with five I can make him wiggle his trunk"

von Neumann

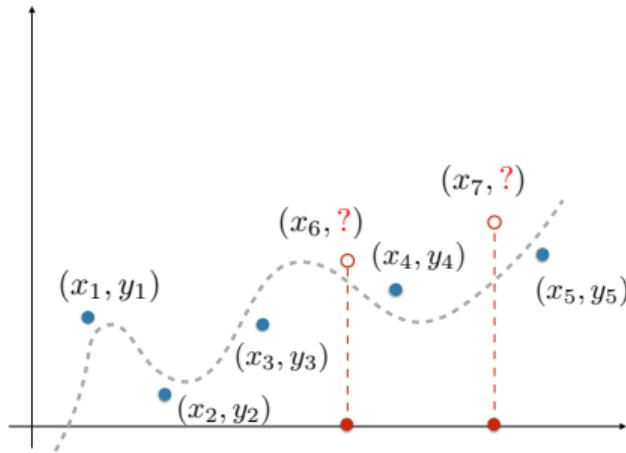
Learning is not (just) fitting, but prediction



Learning is not (just) fitting, but prediction



Learning is not (just) fitting, but prediction



Predictions from random and noisy samples

Learning pipeline

Model fitting (regularized)

$$\widehat{w}_\theta = \underset{\|w\| \leq \theta}{\operatorname{argmin}} \frac{3}{n} \sum_{i=1}^{n/3} (f_w(x_i) - y_i)^2$$

Learning pipeline

Model fitting (regularized)

$$\hat{w}_\theta = \underset{\|w\| \leq \theta}{\operatorname{argmin}} \frac{3}{n} \sum_{i=1}^{n/3} (f_w(x_i) - y_i)^2$$

Model tuning

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \frac{3}{n} \sum_{i=n/3+1}^{2n/3} (f_{\hat{w}_\theta}(x_i) - y_i)^2$$

Learning pipeline

Model fitting (regularized)

$$\hat{w}_\theta = \underset{\|w\| \leq \theta}{\operatorname{argmin}} \frac{3}{n} \sum_{i=1}^{n/3} (f_w(x_i) - y_i)^2$$

Model tuning

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \frac{3}{n} \sum_{i=n/3+1}^{2n/3} (f_{\hat{w}_\theta}(x_i) - y_i)^2$$

Model assessment

$$\frac{3}{n} \sum_{i=2n/3+1}^n (f_{\hat{w}_{\hat{\theta}_\theta}}(x_i) - y_i)^2$$

Classic vs data driven modeling

- ▶ Paradigm shift in modeling, **driven by data** availability.
- ▶ Careful **pipeline** needed.
- ▶ **Theoretical** guidance needed.

ML theory

- ▶ Representation: "Which model?"
- ▶ Generalization: "How accurate is my model?"
- ▶ Optimization: "How can I compute my model?"

Outline

Data driven modeling paradigm

Statistical learning, bias and variance

From statistics to optimization

A theory crisis?

Statistical machine learning

- $(X, Y) \sim P$ random variables in $\mathbb{R}^d \times \mathbb{R}$, and $(x_1, y_1), \dots, (x_n, y_n) \sim P^n$.

Statistical machine learning

- $(X, Y) \sim P$ random variables in $\mathbb{R}^d \times \mathbb{R}$, and $(x_1, y_1), \dots, (x_n, y_n) \sim P^n$.
- $\ell : \mathbb{R} \times \mathbb{R} \rightarrow [0, \infty)$ loss function, e.g. $\ell(f(x), y) = (y - f(x))^2$.

Statistical machine learning

- $(X, Y) \sim P$ random variables in $\mathbb{R}^d \times \mathbb{R}$, and $(x_1, y_1), \dots, (x_n, y_n) \sim P^n$.
- $\ell : \mathbb{R} \times \mathbb{R} \rightarrow [0, \infty)$ loss function, e.g. $\ell(f(x), y) = (y - f(x))^2$.

Problem: minimize

$$L(f) = \mathbb{E}[\ell(f(X), Y)],$$

given only $(x_1, y_1), \dots, (x_n, y_n) \sim P^n$.

ERM and its excess risk

$$\widehat{w}_\theta = \underset{\|w\| \leq \theta}{\operatorname{argmin}} \widehat{L}(f_w), \quad \widehat{L}(f) = \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i)$$

$$\widehat{f}_\theta = f_{\widehat{w}_\theta}$$

ERM and its excess risk

$$\widehat{w}_\theta = \underset{\|w\| \leq \theta}{\operatorname{argmin}} \widehat{L}(f_w), \quad \widehat{L}(f) = \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i)$$
$$\widehat{f}_\theta = f_{\widehat{w}_\theta}$$

Excess risk

$$L(\widehat{f}_\theta) - \min L(f)$$

Error decomposition

Population algorithm

$$f_\theta = f_{w_\theta}, \quad w_\theta = \underset{\|w\| \leq \theta}{\operatorname{argmin}} L(f_w)$$

Error decomposition

Population algorithm

$$f_\theta = f_{w_\theta}, \quad w_\theta = \underset{\|w\| \leq \theta}{\operatorname{argmin}} L(f_w)$$

$$L(\hat{f}_\theta) - \min L(f) = \underbrace{L(\hat{f}_\theta) - L(f_\theta)}_{\text{Estimation error}} + \underbrace{L(f_\theta) - \min L(f)}_{\text{Approximation error}}$$

Approximation error

Assume

$$|\ell(y, f(x)) - \ell(y, f(x'))| \leq C_\ell |f(x) - f'(x)|$$

Lemma

Let $L(f_*) = \min L(f)$, then

$$L(f_\theta) - \min L(f) \leq C_\ell \min_{\|w\| \leq \theta} \|f_\theta - f_*\|_{L^1(\mathbb{P})}$$

Approximation error

Assume

$$|\ell(y, f(x)) - \ell(y, f_\theta(x))| \leq C_\ell |f(x) - f_\theta(x)|$$

Lemma

Let $L(f_*) = \min L(f)$, then

$$L(f_\theta) - \min L(f) \leq C_\ell \min_{\|w\| \leq \theta} \|f_\theta - f_*\|_{L^1(P)}$$

Proof.

$$L(f_\theta) - L(f_*) = \min_{\|w\| \leq \theta} \int (\ell(f(x), y) - \ell(f_*(x), y)) dP(x, y) \leq C_\ell \min_{\|w\| \leq \theta} \int |f(x) - f_*(x)| dP(x, y)$$

□

Universality

A model is universal if for all f_*

$$\lim_{\theta \rightarrow \infty} \|f_\theta - f_*\|_{L^1(P)} = 0.$$

e.g. Kernel methods and neural nets.

[DeVore, Lorentz '93, Pinkus '99]

Smoothness conditions

Assume

$$f_* \in \mathcal{H}_s,$$

for some smoothness class \mathcal{H}_s . e.g. the Sobolev space $W^{s,2}$.

Smoothness conditions

Assume

$$f_* \in \mathcal{H}_s,$$

for some smoothness class \mathcal{H}_s . e.g. the Sobolev space $W^{s,2}$.

Approximation results ensure that

$$\min_{\|w\| \leq \theta} \|f_\theta - f_*\|_{L^1(P)} \lesssim a(\theta, s)$$

where $a(\theta, s)$ decays for θ increasing and rate depending on s , e.g. $\theta^{-s/d}$

[DeVore, Lorentz, '93]

Estimation error

Lemma

By definition of ERM, it holds

$$L(\hat{f}_\theta) - L(f_\theta) \leq C_\ell \sup_{\|w\| \leq \theta} |\hat{L}(f_w) - L(f_w)|$$

Estimation error

Lemma

By definition of ERM, it holds

$$L(\hat{f}_\theta) - L(f_\theta) \leq C_\ell \sup_{\|w\| \leq \theta} |\hat{L}(f_w) - L(f_w)|$$

Proof.

$$L(\hat{f}_\theta) - L(f_\theta) = L(\hat{f}_\theta) - \hat{L}(\hat{f}_\theta) + \underbrace{\hat{L}(\hat{f}_\theta) - \hat{L}(f_\theta)}_{\leq 0} + \hat{L}(f_\theta) - L(f_\theta)$$

□

[Vapnik, Chervonenkis, '77, Gyorfi, Devroye, Lugosi, '96]

Capacity measures

Empirical process

$$\sup_{\|w\| \leq \theta} |\widehat{L}(f_w) - L(f_w)|$$

Capacity measures

Empirical process

$$\sup_{\|w\| \leq \theta} |\hat{L}(f_w) - L(f_w)|$$

Lemma (Rademacher complexities)

If $\sigma_i \in \{\pm 1\}$, $P(1) = P(-1) = 1/2$, $i = 1, \dots, n$ (Rademacher random variables), then

$$\mathbb{E} \left[\sup_{\|w\| \leq \theta} |\hat{L}(f_w) - L(f_w)| \right] \leq 2C_\ell \underbrace{\mathbb{E} \left[\frac{1}{n} \sup_{\|w\| \leq \theta} \sum_{i=1}^n \sigma_i f_w(x_i) \right]}_{\text{Rademacher complexity}},$$

Capacity measures for linear models

$$f_w = \sum_{j=1}^{\infty} w^j \phi_j$$

Capacity measures for linear models

$$f_w = \sum_{j=1}^{\infty} w^j \phi_j$$

Lemma

If

$$\sup_x \left| \sum_{j=1}^{\infty} \phi_j(x) \right|^2 \leq \kappa^2$$

then

$$\mathbb{E} \left[\frac{1}{n} \sup_{\|w\| \leq \theta} \sum_{i=1}^n \sigma_i f_w(x_i) \right] \leq \frac{\theta \kappa}{\sqrt{n}}$$

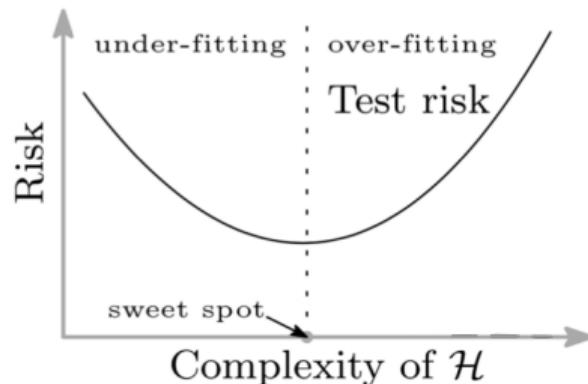
Results for nonlinear models can be similarly derived.

The bias-variance trade-off

$$L(\hat{f}_\theta) - \min L(f) \lesssim \frac{\theta}{\sqrt{n}} + a(\theta, s)$$

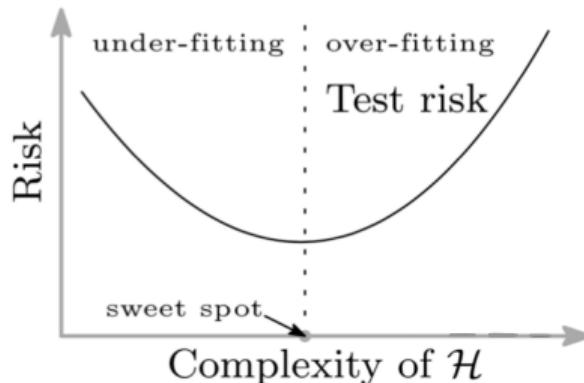
The bias-variance trade-off

$$L(\hat{f}_\theta) - \min L(f) \lesssim \frac{\theta}{\sqrt{n}} + a(\theta, s)$$



The bias-variance trade-off

$$L(\hat{f}_\theta) - \min L(f) \lesssim \frac{\theta}{\sqrt{n}} + a(\theta, s)$$



$$\theta_* = \theta(s, n) \implies L(\hat{f}_{\theta_*}) - \min L(f) \lesssim \epsilon(n, s)$$

where $\epsilon(\theta, a)$ decays with n increasing and rate depending on s , e.g. $n^{-\frac{2s}{2s+d}}$

Outline

Data driven modeling paradigm

Statistical learning, bias and variance

From statistics to optimization

A theory crisis?

Enter optimization

$$\min_{\|w\| \leq \theta} \widehat{L}(f_w)$$

Enter optimization

$$\min_{\|w\| \leq \theta} \widehat{L}(f_w)$$

Gradient methods

$$\widehat{w}_{\theta,t+1} = P_\theta \left(\widehat{w}_{\theta,t} - \gamma_t \nabla \widehat{L}(f_{\widehat{w}_{\theta,t}}) \right)$$

Hard to prove convergence unless the model is linear...

Linear models

$$f_w = \sum_{j=1}^{\infty} w^j \phi_j + \ell \text{ convex}$$

$$\implies \min_{\|w\| \leq \theta} \widehat{L}(f_w) \text{ convex problem!}$$

Linear models

$$f_w = \sum_{j=1}^{\infty} w^j \phi_j + \ell \text{ convex}$$

$$\implies \min_{\|w\| \leq \theta} \widehat{L}(f_w) \text{ convex problem!}$$

Optimization results ensures that $\widehat{f}_{\theta,t} = f_{\widehat{w}_{\theta,t}}$

$$\widehat{L}(\widehat{f}_{\theta,t}) - \min_{\|w\| \leq \theta} \widehat{L}(f) \leq \delta_t$$

with δ_t decaying in t .

[Nesterov, '03]

Estimation + optimization error

Lemma

$$L(\hat{f}_{\theta,t}) - L(f_\theta) \leq C_\ell \sup_{\|w\| \leq \theta} |\hat{L}(f_w) - L(f_w)| + \delta_t$$

Estimation + optimization error

Lemma

$$L(\widehat{f}_{\theta,t}) - L(f_\theta) \leq C_\ell \sup_{\|w\| \leq \theta} |\widehat{L}(f_w) - L(f_w)| + \delta_t$$

Proof.

$$L(\widehat{f}_{\theta,t}) - L(f_\theta) = L(\widehat{f}_{\theta,t}) - \widehat{L}(\widehat{f}_{\theta,t}) + \underbrace{\widehat{L}(\widehat{f}_{\theta,t}) - \widehat{L}(f_\theta)}_{\leq \delta_t} + \widehat{L}(f_\theta) - L(f_\theta)$$

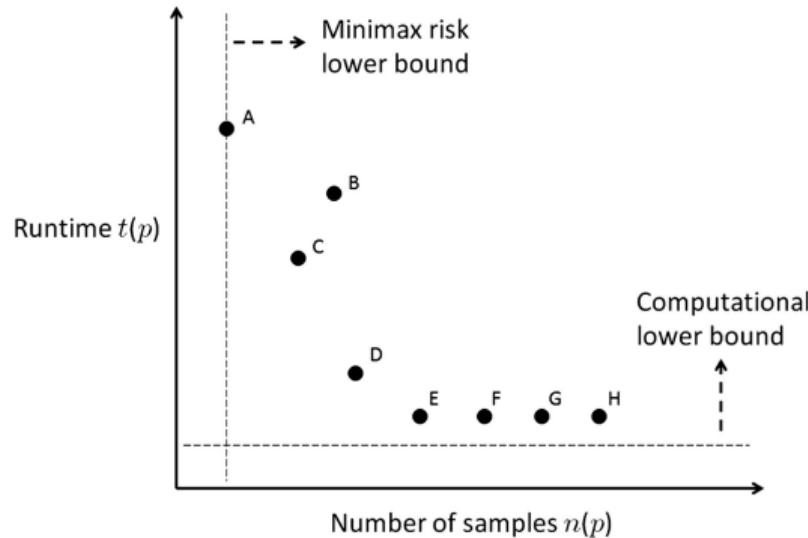
□

The statistical-optimization trade-off

$$L(\widehat{f}_{\theta,t}) - \min L(f) \lesssim \frac{\theta}{\sqrt{n}} + \alpha(\theta, s) + \delta_t$$

The statistical-optimization trade-off

$$L(\hat{f}_{\theta,t}) - \min L(f) \lesssim \frac{\theta}{\sqrt{n}} + \alpha(\theta, s) + \delta_t$$



[Bousquet, Bottou '07, Chendrasekaran Jordan, '12]

Outline

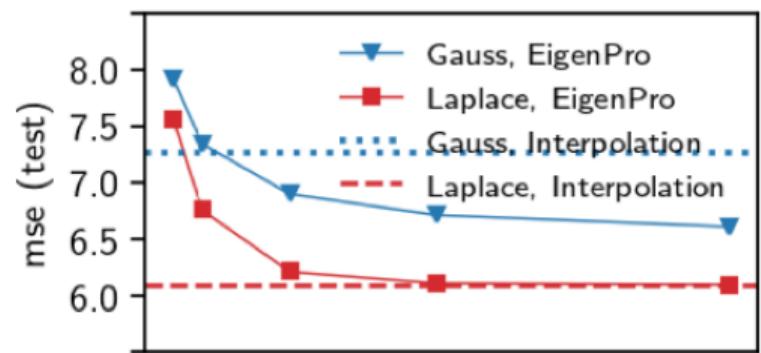
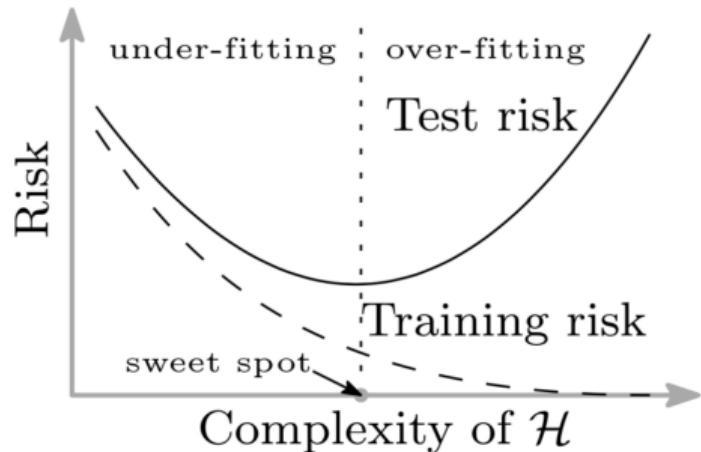
Data driven modeling paradigm

Statistical learning, bias and variance

From statistics to optimization

A theory crisis?

Was it all wrong?



Explicit regularization

$$\min_{\|w\| \leq \theta} \widehat{L}(f_w) \quad \widehat{w}_{\theta,t+1} = P_\theta \left(\widehat{w}_{\theta,t} - \gamma_t \nabla \widehat{L}(f_{\widehat{w}_{\theta,t}}) \right)$$

Explicit regularization

$$\min_{\|w\| \leq \theta} \widehat{L}(f_w) \quad \widehat{w}_{\theta,t+1} = P_\theta \left(\widehat{w}_{\theta,t} - \gamma_t \nabla \widehat{L}(f_{\widehat{w}_{\theta,t}}) \right)$$

Implicit regularization

$$\widehat{w}_{t+1} = \widehat{w}_t - \gamma_t \nabla \widehat{L}(f_{\widehat{w}_t})$$

Explicit regularization

$$\min_{\|w\| \leq \theta} \widehat{L}(f_w) \quad \widehat{w}_{\theta,t+1} = P_\theta \left(\widehat{w}_{\theta,t} - \gamma_t \nabla \widehat{L}(f_{\widehat{w}_{\theta,t}}) \right)$$

Implicit regularization

$$\widehat{w}_{t+1} = \widehat{w}_t - \gamma_t \nabla \widehat{L}(f_{\widehat{w}_t})$$

Can we characterize $\widehat{f}_t = f_{\widehat{w}_t}$

$$L(\widehat{f}_t) - L(f_*)$$

Inexact optimization with linear models

Linear models: $f_w = \sum_{j=1}^{\infty} w^j \phi_j$, ℓ convex and

$$w_{t+1} = w_t - \gamma_t \nabla L(f_{w_t}),$$

then for $f_t = f_{w_t}$

$$L(f_t) - L(f_*) \leq \delta_t.$$

Inexact optimization with linear models

Linear models: $f_w = \sum_{j=1}^{\infty} w^j \phi_j$, ℓ convex and

$$w_{t+1} = w_t - \gamma_t \nabla L(f_{w_t}),$$

then for $f_t = f_{w_t}$

$$L(f_t) - L(f_*) \leq \delta_t.$$

Idea: consider

$$\hat{w}_{t+1} = \hat{w}_t - \gamma_t (\nabla L(f_{\hat{w}_t}) + e_t)$$

with

$$e_t = \nabla \hat{L}(f_{\hat{w}_{\theta,t}}) - \nabla L(f_{\hat{w}_{\theta,t}}).$$

[Rockafellar, '76, Salzo, Villa '11, Schmidt, Le Roux, Bach '11]

Excess risk control with inexact gradient

Lemma

$$L(\hat{f}_t) - L(f_*) \leq \delta_t + \sum_{j=1}^t \langle e_t, \hat{f}_t - f_* \rangle.$$

Excess risk control with inexact gradient

Lemma

$$L(\hat{f}_t) - L(f_*) \leq \delta_t + \sum_{j=1}^t \langle e_t, \hat{f}_t - f_* \rangle.$$

Need to control:

- ▶ gradient error e_t ,
- ▶ path $(\hat{f}_j)_j$ around f_* .

Gradient concentration

$$\mathbb{E} \left[\sup_{\|w\| \leq \theta} \|\nabla \widehat{L}(f_w) - \nabla L(f_w)\| \right] \lesssim \frac{\theta}{\sqrt{n}}$$

Gradient concentration

$$\mathbb{E} \left[\sup_{\|w\| \leq \theta} \|\nabla \widehat{L}(f_w) - \nabla L(f_w)\| \right] \lesssim \frac{\theta}{\sqrt{n}}$$

Path control

For $j \lesssim \sqrt{n}$

$$\|\widehat{f}_t - f_*\| \lesssim \|f_*\|.$$

[Stankewitz, Mücke, R. '21, see also Lin R. '17]

Excess risk control with inexact gradient

Theorem (Stankewitz, Mücke, R. '21)

] For $t \lesssim \sqrt{n}$,

$$\mathbb{E} \left[L(\hat{f}_t) - L(f_*) \right] \lesssim \frac{1}{\sqrt{n}}$$

Same as explicit regularization: implicit regularization a *new algorithmic idea*¹.





"Looking for the lost keys under the lamp, because that's where the light is.", Yann Lecun

- ▶ Can we explain the lack of variance? Learning & interpolation?
- ▶ **Are linear model of any practical use?**
- ▶ Can linear model explain deep learning?

ML meets large scale computing

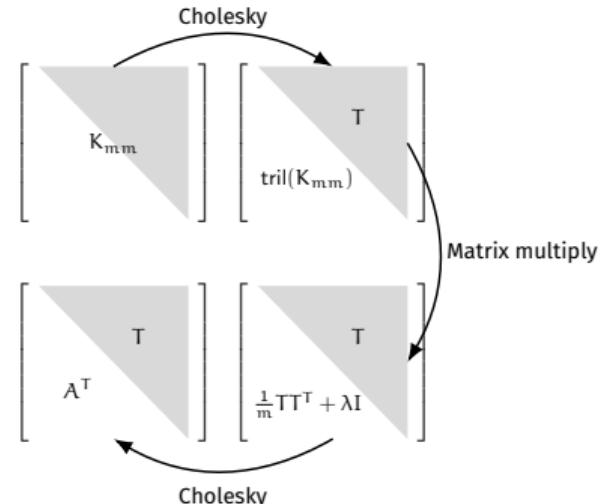
Scalable implementations needed ↫ FALKON

Function Falkon($X \in \mathbb{R}^{n \times d}$, $y \in \mathbb{R}^n$, λ , m , t):

```
 $X_m \leftarrow \text{RandomSubsample}(X, m);$ 
 $T, A \leftarrow \text{Preconditioner}(X_m, \lambda);$ 
```

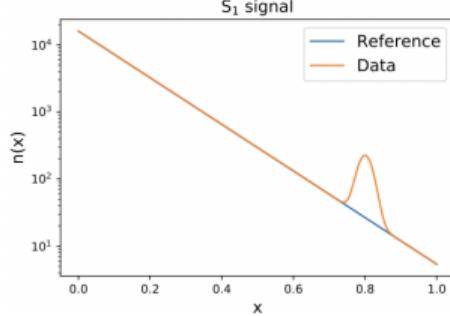
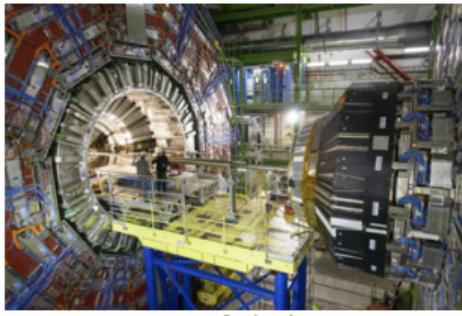
Function LinOp(β):

```
 $v \leftarrow A^{-1}\beta;$ 
 $c \leftarrow k(X_m, X)k(X, X_m)T^{-1}v;$ 
return  $A^{-T}T^{-T}c + \lambda nv;$ 
 $\text{rhs} \leftarrow A^{-T}T^{-T}k(X, X_m)y;$ 
 $\beta \leftarrow \text{ConjugateGradient}(\text{LinOp}, \text{rhs}, t);$ 
return  $T^{-1}A^{-1}\beta;$ 
```

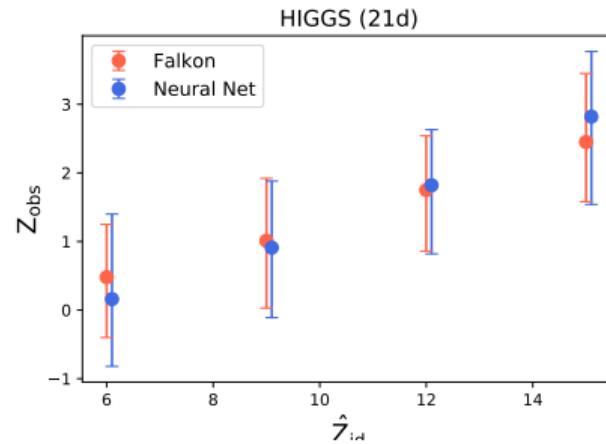


[Meanti, Carratino, R., Rudi '20, Meanti, Carratino, De Vito, R. '21]

Efficient linear models in practice: HEP



[Wulzer, D'Agnolo '18]



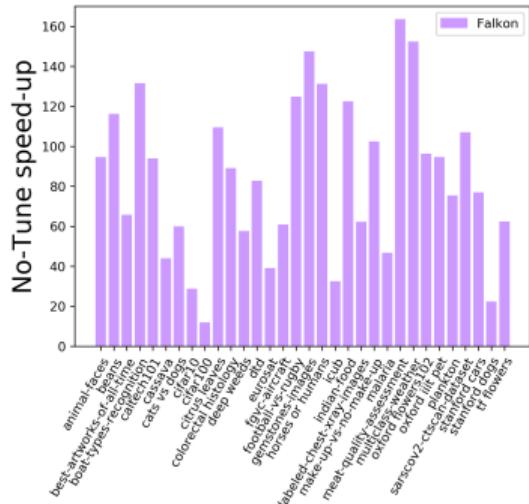
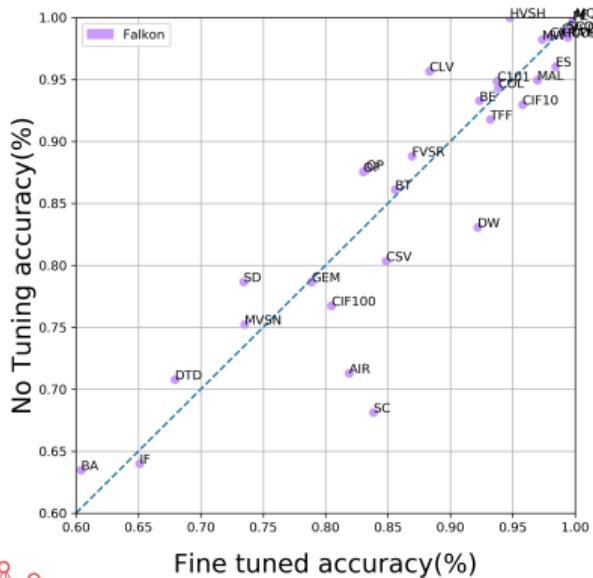
Model	DIMUON	SUSY	HIGGS
Falkon	(53.8 ± 1.9) s	(44.8 ± 1.5) s	(88.7 ± 2.2) s
Neural Net	(4.23 ± 0.73) h	(73.1 ± 10) h	(112 ± 9) h

Table 4: Average training times per single run with standard deviations.

[Letizia et al. '21]

Efficient linear models in practice: vision

$$f(x) = \langle w, \Phi(x) \rangle, \quad x \mapsto \underbrace{\Phi_L}_{\text{Efficient linear learning}} \circ \underbrace{\Phi_{L-1} \cdots \circ \Phi_1(x)}_{\text{Pre-trained layers}}$$



[Alfano, Pastore, Odone, R. '21]

- ▶ Can we explain the lack of variance? Learning & interpolation.
- ▶ Are linear model of any practical use?
- ▶ **Can linear model explain deep learning?**

Neural nets as a linear model?

$$f(x) = \sum_{j=1}^p \beta^j \sigma(a_j^\top x + \alpha_j)$$

Neural nets as a linear model?

$$f(x) = \sum_{j=1}^p \beta^j \sigma(a_j^\top x + \alpha_j)$$

Key idea: Let $\rho(x, \underbrace{(a, \alpha)}_\xi) = \sigma(a_j^\top x + \alpha_j)$.

$$f_\mu(x) = \int \rho(x, \xi) d\mu(\xi), \quad \mu = \sum_{j=1}^q \delta_{\xi_j} \beta_j$$

[Bach '14, Montanari et al. 18, ...]

Optimality of neural nets

$$f_\mu(x) = \int \rho(x, \xi) d\mu(\xi) = \langle \rho(x, \cdot), \mu \rangle, \quad \|f\|_{\mathcal{B}} = \|\mu\|_{TV}$$

Optimality of neural nets

$$f_\mu(x) = \int \rho(x, \xi) d\mu(\xi) = \langle \rho(x, \cdot), \mu \rangle, \quad \|f\|_{\mathcal{B}} = \|\mu\|_{TV}$$

Theorem (Bartolucci, De Vito, R., Vigogna, '21)

The problem

$$\min_{f \in \mathcal{B}} \widehat{L}(f) + \|f\|_{\mathcal{B}}$$

has a solution of the form

$$f(x) = \sum_{i=1}^u \beta^i \rho(x, \lambda_i), \quad u \leq n.$$

see also [Bredies and Carioni '20, Unser et al. '17, Unser '20, Parhi and Nowak '20].

Wrapping up

- ▶ From bias-variance...
- ▶ ...to statistics-optimization trade-off.
- ▶ Statistical learning 2.0?

What's next?

- ▶ Physics (equation?) informed ML.
- ▶ ML for Inverse problems.
- ▶ BIGGER linear models: ML meets scientific computing/HPC.



Multiple post-docs/PhD positions **@MaLGa!**



malga.unige.it