


A guided tour of machine learning (theory)

Lorenzo Rosasco



MaLGA, Università degli Studi di Genova, MIT, IIT

AI everywhere (literally...)

AI




iCub >1000 sensors





ERC_Interview_SLING8.pdf (page 3 of 28)

Data Science



Hep 30Pb/Y



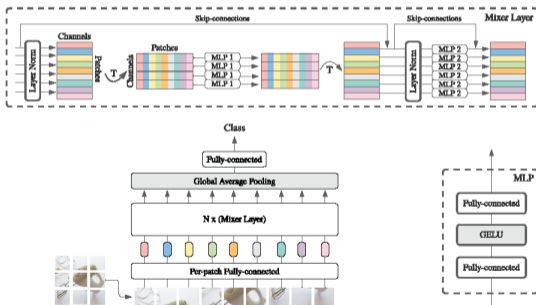
predictive maintenance, connected cars, precision agriculture, personalized fitness and wearables, smart housing, cities, healthcare, etc.

50 billion microcontrollers in 2019
>100 billions in service



Year	Market (\$M)	Units (M)	ASP (\$)
2012	~150	~180	~1.11
2013	~160	~190	~1.11
2014	~170	~19.18	~115.98
2015F	~180	~16.68	~117.78
2016F	~190	~11.98	~158.48
2017F	~200	~11.78	~170.58
2018F	~210	~11.78	~178.68
2019F	~220	~11.78	~187.78

The state of affairs



Rethinking machine learning:

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

Outline

The paradigm of learning from examples

Statistical learning theory (and optimization)

A theory crisis?

The basic picture

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

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

$$\mathbf{w} \in \mathbb{R}^p \mapsto f_{\mathbf{w}}$$

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

$$f_{\mathbf{w}}(\mathbf{x}) = \sum_{j=1}^p \beta^j \sigma(\mathbf{a}_j^{\top} \mathbf{x} + \alpha_j),$$

Model fitting

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

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

Model fitting

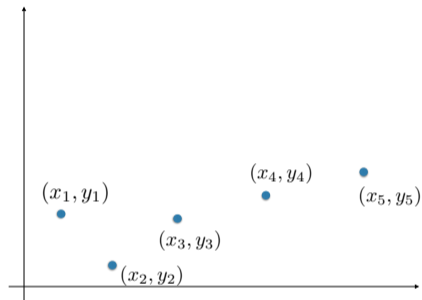
$$\min_w \frac{1}{n} \sum_{i=1}^n (f_w(x_i) - y_i)^2$$

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

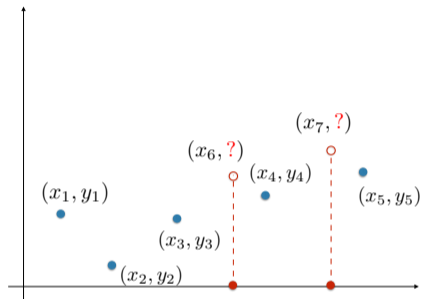
"With four parameters I can fit an elephant, and with five I can make him wiggle his trunk"

von Neumann:

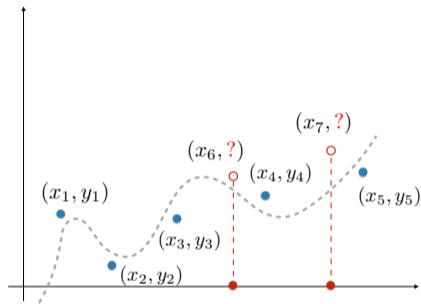
Learning is not (just) fitting, but prediction



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)

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

Learning pipeline

Model fitting (regularized)

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

Model tuning

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

Learning pipeline

Model fitting (regularized)

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

Model tuning

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

Model assessment

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

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

The paradigm of learning from examples

Statistical learning theory (and 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

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

ERM and its excess risk

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

Excess risk

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

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

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

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

Approximation error

Assume

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

Lemma

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

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

Approximation error

Assume

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

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.

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

□

Universality

A model is universal if for all f_*

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

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

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

where $\alpha(\theta, s)$ decays with 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, Györfi, Devroye, Lugosi, '96]

Capacity measures

Empirical process

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

Capacity measures

Empirical process

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

Lemma (Rademacher complexities)

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

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

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

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

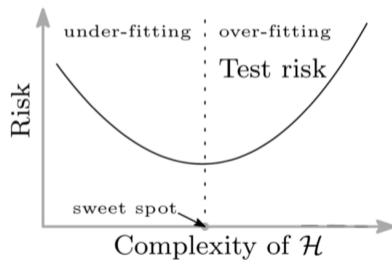
Results for nonlinear models can be similarly derived.

The bias-variance trade-off

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

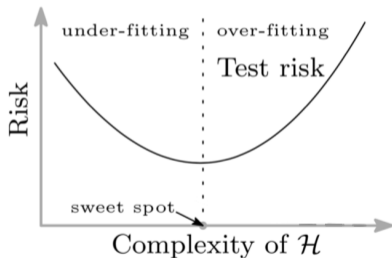
The bias-variance trade-off

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



The bias-variance trade-off

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



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

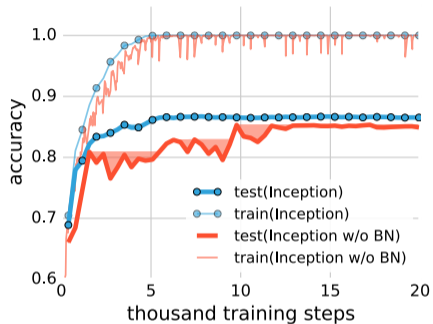
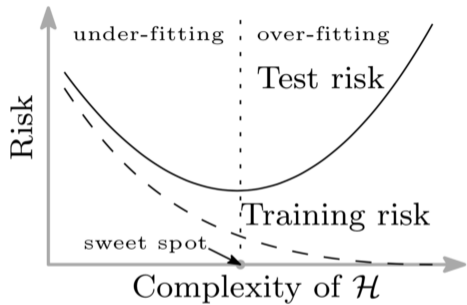
where $\epsilon(\theta, \alpha)$ decays with n increasing and rate depending on s , e.g. $n^{-\frac{2s}{2s+d}}$

Outline

The paradigm of learning from examples

Statistical learning theory (and optimization)

A theory crisis?



Explicit regularization

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

Explicit regularization

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

Implicit regularization

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

Explicit regularization

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

Implicit regularization

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

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

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

Inexact optimization with linear models

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

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

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

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 idea¹.





"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 \mapsto FALKON

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

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

Function $\text{LinOp}(\beta)$:

$v \leftarrow A^{-1}\beta;$

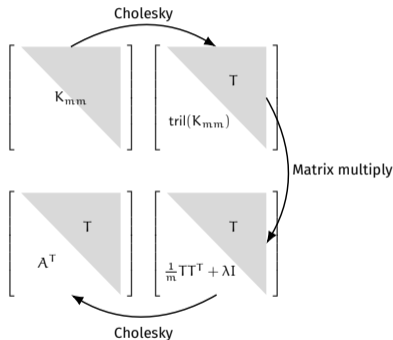
$c \leftarrow k(X_m, X)k(X, X_m)T^{-1}v;$

return $A^{-T}T^{-T}c + \lambda nv;$

$\text{rhs} \leftarrow A^{-T}T^{-T}k(X, X_m)y;$

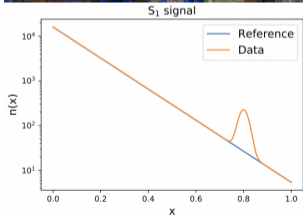
$\beta \leftarrow \text{ConjugateGradient}(\text{LinOp}, \text{rhs}, t);$

return $T^{-1}A^{-1}\beta;$

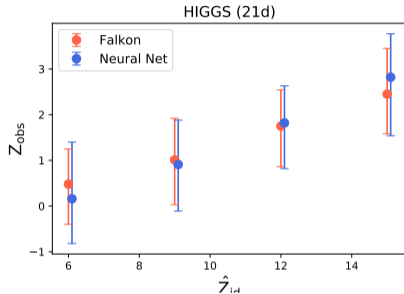


[Meanti, Carratino, R., Rudi '20, Meanti, Carratino, De Vito, R. '21]

Efficient linear models in practice: HEP



[Wulzer, D'Agnolo '18]



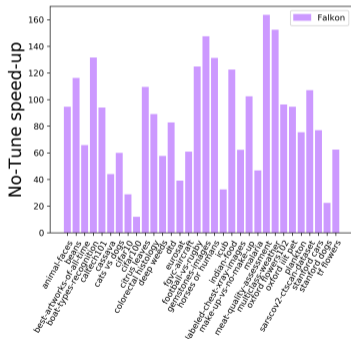
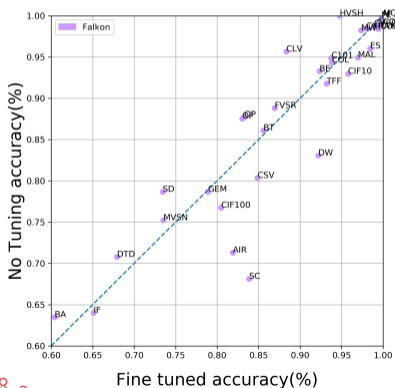
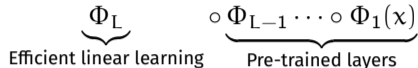
Model	DIMUON	SUSY	HIGGS
Falcon	(53.8 ± 1.9) s	(44.8 ± 1.5) s	(88.7 ± 2.2) s
Neural Net	(4.23 ± 0.73) h	(73.1 ± 10) h	(112 ± 9) h

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

[Letizia et al. '21]

Efficient linear models in practice: vision

$$f(x) = \langle w, \Phi(x) \rangle, \quad x \mapsto$$



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

Wrapping up

- ▶ A guided tour of statistical learning theory
- ▶ Statistics and optimization under the lens of linear models
- ▶ Modern gist to classic ideas (hopefuyly!)

What's next?

- ▶ Data driven + mechanistic modeling
- ▶ Efficient implementation for other loss functions.
- ▶ Random projections+ multiscale approaches [Chen, Avron, Sindawhani '16].



PhD/Postdoc positions available!

