

(Un)conventional regularization for efficient machine learning

Lorenzo Rosasco
Laboratory for Computational Statistical Learning (LCSL)
University of Genova
Massachusetts Institute of Technology
Istituto Italiano di Tecnologia
lcs1.mit.edu

Nov. 23rd, 2018 – Horizon Maths 2018: Intelligence Artificielle

joint work with L. Carratino, R. Camoriano (LCSL), J. Lin (EPFL), A. Rudi (INRIA)



A general motivation: efficient ML

Beyond the statistics vs optimization dichotomy:

Numerical resources budgeted to data **quality** (not just size)

[Bottou, Bousquet '08]

Outline

Classical regularization

Regularization by optimization

Regularization & Projections

Preconditioning

Statistical learning

Find

$$w_* = \operatorname{argmin}_{w \in \mathbb{R}^p} \mathbb{E}[(y - \langle w, \Phi(x) \rangle)^2]$$

given only i.i.d. samples $(x_i, y_i)_{i=1}^n$.

We assume throughout that:

- ▶ Φ a given high/infinite dimensional representation.
- ▶ If $p = \infty$, $\langle \Phi(x), \Phi(x') \rangle$ can be computed in $O(1)$ (kernel methods/Gaussian processes).

Empirical risk minimization (ERM)

$$\hat{w}_\lambda = \operatorname{argmin}_{w \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n (y_i - \langle w, \Phi(x_i) \rangle)^2 + \lambda \|w\|^2$$

Empirical risk minimization (ERM)

$$\hat{w}_\lambda = \operatorname{argmin}_{w \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n (y_i - \langle w, \Phi(x_i) \rangle)^2 + \lambda \|w\|^2$$

Theorem (Smale, Zhou '05, Caponetto De Vito '05)

If $p = \infty$, $\|\Phi(x)\|, |y| \leq 1$ a.s. and $\lambda = 1/\sqrt{n}$

$$\mathbb{E}[(\langle \Phi(x), \hat{w}_\lambda \rangle - \langle \Phi(x), w_* \rangle)^2] \lesssim \frac{1}{\sqrt{n}}$$

Empirical risk minimization (ERM)

$$\hat{w}_\lambda = \operatorname{argmin}_{w \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n (y_i - \langle w, \Phi(x_i) \rangle)^2 + \lambda \|w\|^2$$

Theorem (Smale, Zhou '05, Caponetto De Vito '05)

If $p = \infty$, $\|\Phi(x)\|, |y| \leq 1$ a.s. and $\lambda = 1/\sqrt{n}$

$$\mathbb{E}[(\langle \Phi(x), \hat{w}_\lambda \rangle - \langle \Phi(x), w_* \rangle)^2] \lesssim \frac{1}{\sqrt{n}}$$

Proof

$$\forall \lambda > 0, \quad \mathbb{E}[(\langle \Phi(x), \hat{w}_\lambda \rangle - \langle \Phi(x), w_* \rangle)^2] \lesssim \frac{1}{\lambda n} + \lambda$$

Remark: Optimal bound (can be improved under further assumptions)

ERM computations

$$\hat{w}_\lambda = \operatorname{argmin}_{w \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n (y_i - \langle w, \Phi(x_i) \rangle)^2 + \lambda \|w\|^2$$

Nonparametrics $n < p = \infty$

$$\hat{c}_\lambda = (\hat{\Phi} \hat{\Phi}^\top + \lambda n I)^{-1} \hat{y}$$

- ▶ $\hat{\Phi}$ is the n by p data/feature matrix
- ▶ $\hat{w}_\lambda = \hat{\Phi}^\top \hat{c}_\lambda$
- ▶ $\langle \hat{w}_\lambda, \Phi(x) \rangle = \sum_{i=1}^n \langle \Phi(x_i), \Phi(x) \rangle \hat{c}_\lambda^i$

$$\underbrace{O(n^3)}_{\text{time}} + \underbrace{O(n^2)}_{\text{space}}$$

time $O(n^3)$ + space $O(n^2)$ for optimal $O(1/\sqrt{n})$ learning bound

The rate $1/\sqrt{n}$ is optimal, can we improve time/space requirements?

Outline

Classical regularization

Regularization by optimization

Regularization & Projections

Preconditioning

Gradient descent (GD) learning aka L2 boosting, Landweber iteration

GD for

$$\frac{1}{n} \sum_{i=1}^n (y_i - \langle w, \Phi(x_i) \rangle)^2 = \frac{1}{n} \|\hat{\Phi}w - \hat{y}\|^2$$

Gradient descent (GD) learning aka L2 boosting, Landweber iteration

GD for

$$\frac{1}{n} \sum_{i=1}^n (y_i - \langle w, \Phi(x_i) \rangle)^2 = \frac{1}{n} \|\widehat{\Phi}w - \widehat{y}\|^2$$

Nonparametrics $n \leq p = \infty$

$$\widehat{c}_{t+1} = \widehat{c}_t - \frac{\gamma}{n} (\widehat{\Phi} \widehat{\Phi}^\top \widehat{c}_t - \widehat{y})$$

- ▶ $\widehat{w}_\lambda = \widehat{\Phi}^\top \widehat{c}_t$
- ▶ $\langle \widehat{w}_t, \Phi(x) \rangle = \sum_{i=1}^n \langle \Phi(x_i), \Phi(x) \rangle \widehat{c}_t^i$

$$\underbrace{O(n^2 t)}_{\text{time}} + \underbrace{O(n^2)}_{\text{space}}$$

Gradient descent (GD) learning aka L2 boosting, Landweber iteration

GD for

$$\frac{1}{n} \sum_{i=1}^n (y_i - \langle w, \Phi(x_i) \rangle)^2 = \frac{1}{n} \|\widehat{\Phi}w - \widehat{y}\|^2$$

Nonparametrics $n \leq p = \infty$

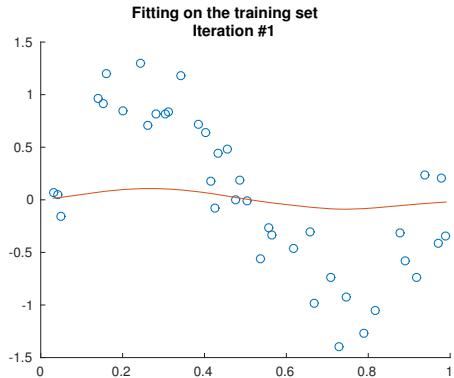
$$\widehat{c}_{t+1} = \widehat{c}_t - \frac{\gamma}{n} (\widehat{\Phi} \widehat{\Phi}^\top \widehat{c}_t - \widehat{y})$$

- ▶ $\widehat{w}_\lambda = \widehat{\Phi}^\top \widehat{c}_t$
- ▶ $\langle \widehat{w}_t, \Phi(x) \rangle = \sum_{i=1}^n \langle \Phi(x_i), \Phi(x) \rangle \widehat{c}_t^i$

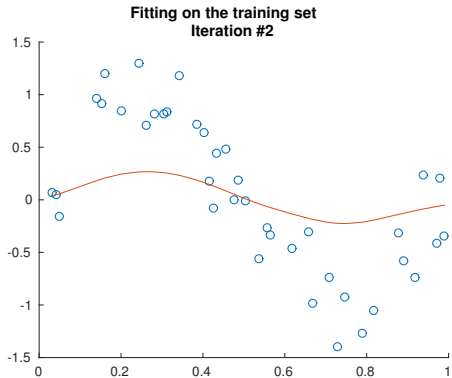
$$\underbrace{O(n^2 t)}_{\text{time}} + \underbrace{O(n^2)}_{\text{space}}$$

Why should this work??

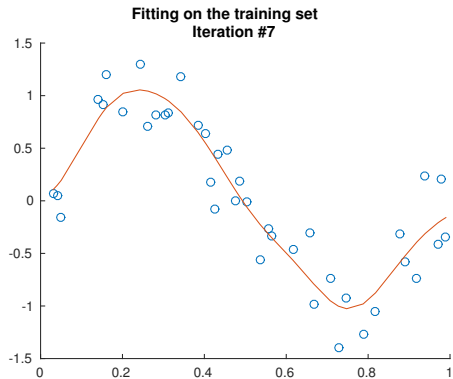
An intuition



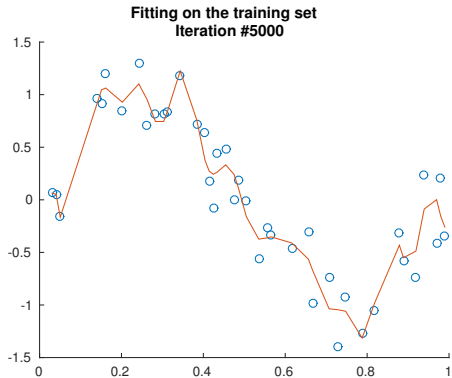
An intuition



An intuition



An intuition



Regularization and stability

$$\hat{w}_t = \frac{\gamma}{n} \sum_{j=0}^{t-1} \left(I - \frac{\gamma}{n} \hat{\Phi}^\top \hat{\Phi}\right)^j \Phi^\top \hat{y}$$

Large t

$$\hat{w}_t = \frac{\gamma}{n} \sum_{j=0}^{t-1} \left(I - \frac{\gamma}{n} \hat{\Phi}^\top \hat{\Phi}\right)^j \Phi^\top \hat{y} \approx \frac{\gamma}{n} \sum_{j=0}^{\infty} \left(I - \frac{\gamma}{n} \hat{\Phi}^\top \hat{\Phi}\right)^j \Phi^\top \hat{y} = (\hat{\Phi}^\top \hat{\Phi})^{-1} \hat{\Phi}^\top \hat{y}$$

Small t

$$\hat{w}_t = \left(I - \left(I - \frac{\gamma}{n} \hat{\Phi}^\top \hat{\Phi}\right)^t\right) (\hat{\Phi}^\top \hat{\Phi})^{-1} \Phi^\top \hat{y} \stackrel{t=1}{\propto} \Phi^\top \hat{y}$$

compare to

$$\hat{w}_\lambda = (\hat{\Phi}^\top \hat{\Phi} + \lambda n I)^{-1} \Phi^\top \hat{y} \stackrel{\lambda \text{ large}}{\propto} \Phi^\top \hat{y}$$

Statistics

Theorem (Bauer, Pereverzev, R. '07)

If $p = \infty$, $\|\Phi(x)\|, |y| \leq 1$ a.s. and $t = \sqrt{n}$

$$\mathbb{E}[(\langle \Phi(x), \hat{w}_t \rangle - \langle \Phi(x), w_* \rangle)^2] \lesssim \frac{1}{\sqrt{n}}$$

Statistics

Theorem (Bauer, Pereverzev, R. '07)

If $p = \infty$, $\|\Phi(x)\|, |y| \leq 1$ a.s. and $t = \sqrt{n}$

$$\mathbb{E}[(\langle \Phi(x), \hat{w}_t \rangle - \langle \Phi(x), w_* \rangle)^2] \lesssim \frac{1}{\sqrt{n}}$$

Proof

$$\forall t > 1, \quad \mathbb{E}[(\langle \Phi(x), \hat{w}_t \rangle - \langle \Phi(x), w_* \rangle)^2] \lesssim \frac{t}{n} + \frac{1}{t}$$

Remarks:

- ▶ Same bound as Tikhonov regularization.
- ▶ Known as iterative regularization (1955), early stopping (1985), implicit regularization (2018).

Computational regularization

time $O(n^2\sqrt{n})$ + space $O(n^2)$ for optimal $O(1/\sqrt{n})$ learning bound

Regularization by stochastic optimization

control statistics and time at once

What about other optimization methods?

Beyond gradient descent

What about other optimization methods?

- ▶ Accelerated methods (Conjugate Gradient, Nesterov, Heavyball).
- ▶ **Stochastic methods (SGD).**

SGD with minibatches

$$\hat{w}_{t+1} = \hat{w}_t - \gamma_t \frac{1}{b} \sum_{i=b(t-1)+1}^{bt} (\langle \hat{w}_t, \Phi(x_{j_i}) \rangle - y_{j_i}) \Phi(x_{j_i})$$

- ▶ b minibatch size ($b = 1$ SGD, $b = n$ GD)
- ▶ $t = \lceil \frac{n}{b} \rceil$ one pass

Statistics+Optimization

Theorem (Lin R. '16)

If $p = \infty$, $\|\Phi(x)\|, |y| \leq 1$ a.s. if

1. $b = 1$, $\gamma_t \simeq \frac{1}{\sqrt{n}}$, and $t = n$ iterations (1 pass over the data);
2. $b = \sqrt{n}$, $\gamma_t \simeq 1$, and $t = \sqrt{n}$ iterations (1 pass over the data);
3. $b = n$, $\gamma_t \simeq 1$, and $t = \sqrt{n}$ iterations (\sqrt{n} passes over the data);

then,

$$\mathbb{E}[(\langle \Phi(x), \hat{w}_t \rangle - \langle \Phi(x), w_* \rangle)^2] \lesssim \frac{1}{\sqrt{n}}.$$

Statistics+Optimization

Theorem (Lin R. '16)

If $p = \infty$, $\|\Phi(x)\|, |y| \leq 1$ a.s. if

1. $b = 1$, $\gamma_t \simeq \frac{1}{\sqrt{n}}$, and $t = n$ iterations (1 pass over the data);
2. $b = \sqrt{n}$, $\gamma_t \simeq 1$, and $t = \sqrt{n}$ iterations (1 pass over the data);
3. $b = n$, $\gamma_t \simeq 1$, and $t = \sqrt{n}$ iterations (\sqrt{n} passes over the data);

then,

$$\mathbb{E}[(\langle \Phi(x), \hat{w}_t \rangle - \langle \Phi(x), w_* \rangle)^2] \lesssim \frac{1}{\sqrt{n}}.$$

Proof

$$\forall \gamma > 0, t > 1, \quad \mathbb{E}[(\langle \Phi(x), \hat{w}_t \rangle - \langle \Phi(x), w_* \rangle)^2] \lesssim \lesssim \frac{1}{\gamma t} + \frac{1}{\sqrt{n}} \left(\frac{\gamma t}{\sqrt{n}} \right)^2 + \frac{\gamma}{b} \left(1 + \frac{\gamma t}{\sqrt{n}} \right)$$

Other Flavors of SGD

- ▶ Cyclic, reshuffle SGD (R. Villa, '15)
- ▶ SGD with averaging (Bach et al. 14-..., Pillaud, Rudi, Bach '18)- larger step-size.
- ▶ Averaging as regularization (Neu, R. '18)

Computational regularization

time $O(n^2)$ + space $O(n^2)$ for optimal $O(1/\sqrt{n})$ learning bound

Regularization by stochastic optimization

improved time while keeping optimal statistical error

What about memory costs?

Outline

Classical regularization

Regularization by optimization

Regularization & Projections

Preconditioning

Tackling memory with random projections

- ▶ Sketching & random features
- ▶ **Nyström methods & subsampling**

Subsampling aka Nyström methods

Consider

$$\bar{x}_1, \dots, \bar{x}_M \subset x_1, \dots, x_n$$

and

$$w_M = \sum_{j=1}^M \Phi(\bar{x}_j) c_j,$$

a random projection on a subspace.

Connections to

- ▶ Nyström methods
- ▶ Galerkin methods
- ▶ Column subsampling

Linear algebra perspective

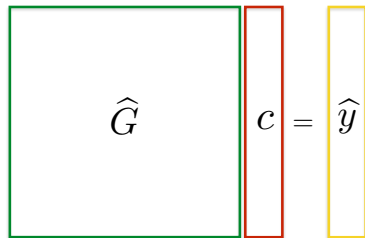
RandNLA=randomized numerical linear algebra (Halko et al. '11)

Back to ridge regression.

$$\langle \hat{w}_\lambda, \Phi(x) \rangle = \sum_{i=1}^n \langle \Phi(x_i), \Phi(x) \rangle c^i$$

$$c = \underbrace{(\hat{\Phi}\hat{\Phi}^\top)}_{\hat{G}} + \lambda n I)^{-1} \hat{y}$$

Linear System



The diagram illustrates a linear system $\hat{G}c = \hat{y}$. It features three vertical rectangular boxes. The first box on the left is a large square with a green border, containing the symbol \hat{G} . To its right is a tall, narrow vertical rectangle with a red border, containing the symbol c . To the right of the red box is an equals sign, followed by another tall, narrow vertical rectangle with a yellow border, containing the symbol \hat{y} .

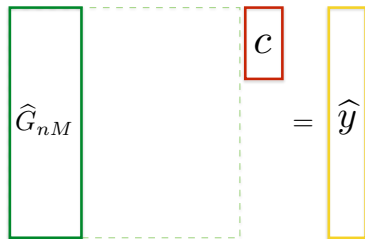
Nystrom/Column subsampling

Take $\bar{x}_1, \dots, \bar{x}_M \subset x_1, \dots, x_n$, $M < n$.

$$\langle \hat{w}_{\lambda, M}, \Phi(x) \rangle = \sum_{i=1}^M \langle \Phi(\bar{x}_i), \Phi(x) \rangle c^i$$

$$(\hat{G}_{nM}^\top \hat{G}_{nM} + \lambda n \hat{G}_{MM}) c = \hat{G}_{nM}^\top \hat{y}$$

Linear System



Statistical guarantees

Theorem (Rudi, Camoriano, R. '16)

If $p = \infty$, $\|\Phi(x)\|, |y| \leq 1$ a.s. if

$$\lambda = 1/\sqrt{n}, \quad M = \sqrt{n}$$

then,

$$\mathbb{E}[(\langle \Phi(x), \hat{w}_{\lambda, M} \rangle - \langle \Phi(x), w_* \rangle)^2] \lesssim \frac{1}{\sqrt{n}}$$

Remarks:

- ▶ Same bound again...
- ▶ Improve previous bounds (Bach et al. '12, Alaoui, Mahoney '14)
- ▶ Regularization by projection!

Computational regularization

time $O(n^2)$ + space $O(n\sqrt{n})$ for optimal $O(1/\sqrt{n})$ learning bound

Regularization by projection

control statistics, time and memory costs at once

Can we improve computational costs?

Outline

Classical regularization

Regularization by optimization

Regularization & Projections

Preconditioning

Preconditioning

Idea: define equivalent linear system with better condition number

Preconditioning

Idea: define equivalent linear system with better condition number

Preconditioning

$$(\widehat{G} + \lambda n I)c = \hat{y} \quad \mapsto \quad B^\top (\widehat{G} + \lambda n I) B \beta = B^\top \hat{y}, \quad c = B \beta.$$

Ideally $BB^\top = (\widehat{G} + \lambda n I)^{-1}$, so that

$$t = O(1/\lambda) \quad \mapsto \quad t = O(1)!$$

(Fasshauer et al '12, Avron et al '16, Cutaját '16, Ma, Belkin '17)

Baby FALKON

Recall Nyström $(\widehat{G}_{nM}^\top \widehat{G}_{nM} + \lambda n \widehat{G}_{MM})c = \widehat{G}_{nM}^\top \widehat{y}$

Preconditioning

$$BB^\top = \left(\frac{n}{M} \widehat{G}_{MM}^2 + \lambda n \widehat{G}_{MM} \right)^{-1},$$

Baby FALKON

$$\langle \widehat{w}_{\lambda, M, t}, \Phi(x) \rangle = \sum_{i=1}^M \langle \Phi(\tilde{x}_i), \Phi(x) \rangle c^i \quad c_t = B\beta_t$$

$$\beta_t = \beta_{t-1} - \frac{\tau}{n} B^\top \left[\widehat{G}_{nM}^\top (\widehat{G}_{nM} B\beta_{t-1} - y_n) + \lambda n \widehat{G}_{MM} B\beta_{t-1} \right]$$

FALKON

- ▶ Gradient descent \mapsto conjugate gradient
- ▶ Computing B

$$B = \frac{1}{\sqrt{n}}T^{-1}A^{-1}, \quad T = \text{chol}(G_{MM}), \quad A = \text{chol}\left(\frac{1}{M}TT^{\top} + \lambda I\right),$$

where $\text{chol}(\cdot)$ is the Cholesky decomposition.



Some Theory

Theorem (Rudi, Carratino, R. '17)

If $p = \infty$, $\|\Phi(x)\|, |y| \leq 1$ a.s. if

$$\lambda = 1/\sqrt{n}, \quad M = \sqrt{n}, \quad t = \log n$$

then

$$\mathbb{E}[(\langle \Phi(x), \hat{w}_{\lambda, M, t} \rangle - \langle \Phi(x), w_* \rangle)^2] \lesssim \frac{1}{\sqrt{n}}$$

Remarks:

- ▶ Same bound again... improved time cost!
- ▶ Improved results by considering adaptive sampling.

Computational regularization

time $O(n\sqrt{n})$ + space $O(n\sqrt{n})$ for optimal $O(1/\sqrt{n})$ learning bound

Maybe optimal?

Some experiments

	MillionSongs			YELP		TIMIT	
	MSE	Relative error	Time(s)	RMSE	Time(m)	c-err	Time(h)
FALKON	80.30	4.51×10^{-3}	55	0.833	20	32.3%	1.5
Prec. KRR	-	4.58×10^{-3}	289 [†]	-	-	-	-
Hierarchical	-	4.56×10^{-3}	293 [*]	-	-	-	-
D&C	80.35	-	737 [*]	-	-	-	-
Rand. Feat.	80.93	-	772 [*]	-	-	-	-
Nyström	80.38	-	876 [*]	-	-	-	-
ADMM R. F.	-	5.01×10^{-3}	958 [†]	-	-	-	-
BCD R. F.	-	-	-	0.949	42 [‡]	34.0%	1.7 [‡]
BCD Nyström	-	-	-	0.861	60 [‡]	33.7%	1.7 [‡]
KRR	-	4.55×10^{-3}	-	0.854	500 [‡]	33.5%	8.3 [‡]
EigenPro	-	-	-	-	-	32.6%	3.9 [‡]
Deep NN	-	-	-	-	-	32.4%	-
Sparse Kernels	-	-	-	-	-	30.9%	-
Ensemble	-	-	-	-	-	33.5%	-

Table: MillionSongs, YELP and TIMIT Datasets. Times obtained on: ‡ = cluster of 128 EC2 r3.2xlarge machines, † = cluster of 8 EC2 r3.8xlarge machines, † = single machine with two Intel Xeon E5-2620, one Nvidia GTX Titan X GPU and 128GB of RAM, * = cluster with 512 GB of RAM and IBM POWER8 12-core processor, * = unknown platform.

Some more experiments

	SUSY			HIGGS	
	c-err	AUC	Time(<i>m</i>)	AUC	Time(<i>h</i>)
FALKON	19.6%	0.877	4	0.825	3
EigenPro	19.8%	-	6 [‡]	-	-
SVM	26.4%	-	9*	-	-
Hierarchical	20.1%	-	40 [†]	-	-
Boosted Decision Tree	-	0.863	-	0.810	-
Neural Network	-	0.875	-	0.816	-
Deep Neural Network	-	0.879	4680 [‡]	0.885	78 [†]

Table: SUSY and HIGGS Datasets. Time obtained working on : † = cluster with 512 GB of RAM and IBM POWER8 12-core processor, ‡ = single machine with two Intel Xeon E5-2620, one Nvidia GTX Titan X GPU and 128GB of RAM, ‡ = single machine, * = 14 workers.

Image classification

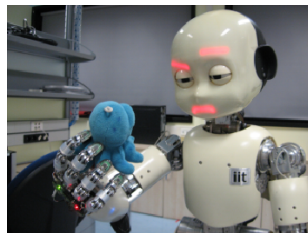
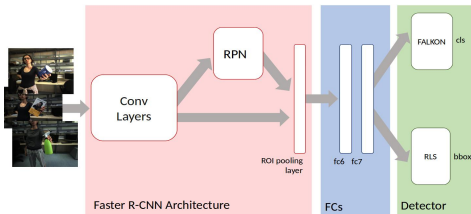
$$f(x) = \langle w, \Phi(x) \rangle, \quad x \mapsto \underbrace{\Phi_L}_{\text{Kernel representation}} \circ \underbrace{\Phi_{L-1} \cdots \circ \Phi_1(x)}_{\text{Convolutional}}$$

Imagenet

	Top-1 class error
FALKON + I-v3 feat.	22.1%
Inception-v3	21.2%
Inception-v2	23.4%
BN-Inception	25.2%
BN-GoogLeNet	26.8%
GoogLeNet	29.0%

Table: Single crop experimental results on the validation set of ILSVRC 2012.

Object detection



Method	mAP [%]	Train Time
Faster R-CNN	51,9	~25 min
FALKON + Full Bootstrap (~ 1K×1000)	51,5	~8 min
FALKON + Random BKG (0 × 7000)	47,7	~25 sec



Method	Train Time	mAP [%]	soda bottle	mug	pencil case	ring binder	wallet	flower	book	body lotion	hair clip	sprayer
Faster R-CNN Fine-tuning	~40 min	49,7	63,2	68,4	23,3	29,6	49,9	66,1	35,3	56,2	60,2	45,8
FALKON + Random BKG (0×6000)	~25 sec	40,5	57,7	67,9	17,5	23,1	23,8	59,5	26,8	39,6	48,5	40,5
FALKON + Mini Bootstrap (4×2500)	~40 sec	48,1	63,1	67,2	18,4	25,7	47,4	70,3	36,1	52,3	58,8	41,5
FALKON + Mini Bootstrap (10×1500)	~50 sec	51,3	64,7	71,2	27,2	31,7	56,9	69,4	39,6	54,0	60,7	37,1

Summing up

- ▶ Computational regularization for efficient learning.
- ▶ Faster GP/Kernel solver ever.

Looking ahead

- ▶ Other loss functions, norms, learning problems
- ▶ Parallelization.
- ▶ Non convex problems.
- ▶ Optimal complexity.

check papers on arxiv.org