Less is More: **Computational Regularization by Subsampling**

Lorenzo Rosasco University of Genova - Istituto Italiano di Tecnologia Massachusetts Institute of Technology lcsl_mit_edu

joint work with Alessandro Rudi, Raffaello Camoriano

Jan 13th. 2016 UCL. London



A Starting Point

Classically:

Statistics and optimization distinct steps in algorithm design

A Starting Point

Classically: Statistics and optimization distinct steps in algorithm design

Large Scale: Consider interplay between statistics and optimization! (Bottou, Bousquet '08)

A Starting Point

Classically: Statistics and optimization distinct steps in algorithm design

Large Scale: Consider interplay between statistics and optimization! (Bottou, Bousquet '08)

Computational Regularization: Computation "tricks" = regularization

Supervised Learning

Problem: Estimate f^*



Supervised Learning

Problem: Estimate f^* given $S_n = \{(x_1, y_1), \dots, (x_n, y_n)\}$



Supervised Learning

Problem: Estimate f^* given $S_n = \{(x_1, y_1), \dots, (x_n, y_n)\}$





$$y_i = f^*(x_i) + \varepsilon_i \qquad i \in \{1, \dots, n\}$$

• $\varepsilon_i \in \mathbb{R}, x_i \in \mathbb{R}^d$ random (with unknown distribution)

▶ f* unknown

Outline

Learning with dictionaries and kernels

Data Dependent Subsampling

Data Independent Subsampling

$$\widehat{f}(x) = \sum_{i=1}^{M} c_i q(x, w_i)$$

$$\widehat{f}(x) = \sum_{i=1}^{M} c_i q(x, w_i)$$

► q non linear function

$$\widehat{f}(x) = \sum_{i=1}^{M} c_i q(x, w_i)$$

- ► q non linear function • $w_i \in \mathbb{R}^d$ centers

$$\widehat{f}(x) = \sum_{i=1}^{M} c_i q(x, w_i)$$

- \blacktriangleright q non linear function
- $w_i \in \mathbb{R}^d$ centers
- $c_i \in \mathbb{R}$ coefficients

$$\widehat{f}(x) = \sum_{i=1}^{M} c_i q(x, w_i)$$

- ► q non linear function
- $w_i \in \mathbb{R}^d$ centers
- $c_i \in \mathbb{R}$ coefficients
- $M = M_n$ could/should grow with n

$$\widehat{f}(x) = \sum_{i=1}^{M} c_i q(x, w_i)$$

- ► q non linear function
- $w_i \in \mathbb{R}^d$ centers
- $c_i \in \mathbb{R}$ coefficients
- $M = M_n$ could/should grow with n

Question: How to choose w_i , c_i and M given S_n ?

Learning with Positive Definite Kernels

There is an *elegant* answer if:

- ► q is symmetric
- ▶ all the matrices $\hat{Q}_{ij} = q(x_i, x_j)$ are positive semi-definite¹

¹They have non-negative eigenvalues

Learning with Positive Definite Kernels

There is an *elegant* answer if:

- q is symmetric
- ▶ all the matrices $\widehat{Q}_{ij} = q(x_i, x_j)$ are **positive semi-definite**¹

Representer Theorem (Kimeldorf, Wahba '70; Schölkopf et al. '01)

- ▶ *M* = *n*,
- $w_i = x_i$,
- c_i by convex optimization!

¹They have non-negative eigenvalues

Kernel Ridge Regression (KRR)

a.k.a. Penalized Least Squares

$$\widehat{f}_{\lambda} = \operatorname*{argmin}_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^{n} (y_i - f(x_i))^2 + \lambda \|f\|^2$$

Kernel Ridge Regression (KRR)

a.k.a. Penalized Least Squares

$$\widehat{f}_{\lambda} = \operatorname*{argmin}_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^{n} (y_i - f(x_i))^2 + \lambda \|f\|^2$$

where

$$\mathcal{H} = \{ f \mid f(x) = \sum_{i=1}^{M} c_i q(x, w_i), \ c_i \in \mathbb{R}, \underbrace{w_i \in \mathbb{R}^d}_{\text{any center!}}, \ \underbrace{M \in \mathbb{N}}_{\text{any length!}} \}$$

Kernel Ridge Regression (KRR)

a.k.a. Penalized Least Squares

$$\widehat{f}_{\lambda} = \operatorname*{argmin}_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^{n} (y_i - f(x_i))^2 + \lambda \|f\|^2$$

where

$$\mathcal{H} = \{ f \mid f(x) = \sum_{i=1}^{M} c_i q(x, w_i), \ c_i \in \mathbb{R}, \underbrace{w_i \in \mathbb{R}^d}_{\text{any center!}}, \ \underbrace{M \in \mathbb{N}}_{\text{any length!}} \}$$

Solution

$$\widehat{f}_{\lambda} = \sum_{i=1}^{n} c_i q(x, \boldsymbol{x}_i) \text{ with } c = (\widehat{Q} + \lambda nI)^{-1} \widehat{y}$$

Well understood statistical properties:

Classical Theorem If $f^* \in \mathcal{H}$, then

$$\lambda_* = \frac{1}{\sqrt{n}} \qquad \mathbb{E}\left(\widehat{f}_{\lambda_*}(x) - f^*(x)\right)^2 \lesssim \frac{1}{\sqrt{n}}$$

Well understood statistical properties:

Classical Theorem If $f^* \in \mathcal{H}$, then

$$\lambda_* = \frac{1}{\sqrt{n}} \qquad \mathbb{E}\left(\widehat{f}_{\lambda_*}(x) - f^*(x)\right)^2 \lesssim \frac{1}{\sqrt{n}}$$

Remarks

Well understood statistical properties:

Classical Theorem If $f^* \in \mathcal{H}$, then

$$\lambda_* = \frac{1}{\sqrt{n}} \qquad \mathbb{E}\left(\widehat{f}_{\lambda_*}(x) - f^*(x)\right)^2 \lesssim \frac{1}{\sqrt{n}}$$

Remarks

1. Optimal nonparametric bound

Well understood statistical properties:

Classical Theorem If $f^* \in \mathcal{H}$, then

$$\lambda_* = \frac{1}{\sqrt{n}} \qquad \mathbb{E}\left(\widehat{f}_{\lambda_*}(x) - f^*(x)\right)^2 \lesssim \frac{1}{\sqrt{n}}$$

Remarks

- 1. Optimal nonparametric bound
- 2. Results for general kernels (e.g. splines/Sobolev etc.)

$$\lambda_* = n^{-\frac{1}{2s+1}}, \qquad \mathbb{E}\left(\widehat{f}_{\lambda_*}(x) - f^*(x)\right)^2 \lesssim n^{-\frac{2s}{2s+1}}$$

Well understood statistical properties:

Classical Theorem If $f^* \in \mathcal{H}$, then

$$\lambda_* = \frac{1}{\sqrt{n}} \qquad \mathbb{E}\left(\widehat{f}_{\lambda_*}(x) - f^*(x)\right)^2 \lesssim \frac{1}{\sqrt{n}}$$

Remarks

- 1. Optimal nonparametric bound
- 2. Results for general kernels (e.g. splines/Sobolev etc.)

$$\lambda_* = n^{-\frac{1}{2s+1}}, \qquad \mathbb{E}\left(\widehat{f}_{\lambda_*}(x) - f^*(x)\right)^2 \lesssim n^{-\frac{2s}{2s+1}}$$

- 3. Adaptive tuning via cross validation
- Proofs: analysis/linear algebra+ random matrices (Smale and Zhou + Caponnetto, De Vito, R.+ Steinwart)

KRR: Optimization

$$\widehat{f}_{\lambda} = \sum_{i=1}^{n} c_i q(x, x_i) \text{ with } c = (\widehat{Q} + \lambda nI)^{-1} \widehat{y}$$

Linear System

Complexity

- Space $O(n^2)$
- Time $O(n^3)$



KRR: Optimization

$$\widehat{f}_{\lambda} = \sum_{i=1}^{n} c_i q(x, x_i) \text{ with } c = (\widehat{Q} + \lambda nI)^{-1} \widehat{y}$$

Linear System

 \widehat{Q}



• Space
$$O(n^2)$$



 $C = \widehat{y}$

Can this be fixed?

 $(\hat{Q}+\lambda I)^{-1}$ approximation of \hat{Q}^{\dagger} controlled by λ

 $(\hat{Q}+\lambda I)^{-1}$ approximation of \hat{Q}^{\dagger} controlled by λ

Can we approximate \hat{Q}^{\dagger} by saving computations?

 $(\hat{Q}+\lambda I)^{-1}$ approximation of \hat{Q}^{\dagger} controlled by λ

Can we approximate \hat{Q}^{\dagger} by saving computations?

Yes!

 $(\hat{Q}+\lambda I)^{-1}$ approximation of \hat{Q}^{\dagger} controlled by λ

Can we approximate \hat{Q}^{\dagger} by saving computations? **Yes!**

Spectral filtering (Engl '96- inverse problems, Rosasco et al. 05- ML)

 $g_{\lambda}(\hat{Q}) \sim \hat{Q}^{\dagger}$

The filter function g_{λ} defines the form of the approximation

Spectral filtering

Examples

- Tikhonov- ridge regression
- Truncated SVD- principal component regression
- ▶ Landweber iteration- GD/ L₂-boosting
- nu-method- accelerated GD/Chebyshev method

▶ ...

Spectral filtering

Examples

- Tikhonov- ridge regression
- Truncated SVD- principal component regression
- ▶ Landweber iteration- GD/ L₂-boosting
- nu-method- accelerated GD/Chebyshev method

▶ ...

Landweber iteration (truncated power series)...

$$c_t = g_t(\hat{Q}) = \gamma \sum_{r=0}^{t-1} (I - \gamma \hat{Q})^r \hat{y}$$

Spectral filtering

Examples

- Tikhonov- ridge regression
- Truncated SVD- principal component regression
- ▶ Landweber iteration- GD/ L₂-boosting
- nu-method- accelerated GD/Chebyshev method

▶ ...

Landweber iteration (truncated power series)...

$$c_t = g_t(\hat{Q}) = \gamma \sum_{r=0}^{t-1} (I - \gamma \hat{Q})^r \hat{y}$$

... it's GD for ERM!!

$$r = 1 \dots t$$
 $c_r = c_{r-1} - \gamma (\hat{Q}c_{r-1} - \hat{y}), \quad c_0 = 0$

Semiconvergence



Early Stopping at Work


Early Stopping at Work



Early Stopping at Work



Early Stopping at Work



Statistics and computations with spectral filtering

The different filters achieve essentially the same optimal statistical error!

Statistics and computations with spectral filtering

The different filters achieve essentially the same optimal statistical error!

Difference is in computations

Filter	Time	Space
Tikhonov	n^3	n^2
GD	$n^2 \lambda_*^{-1}$	n^2
Accelerated GD	$n^2 \lambda_*^{-1/2}$	n^2
Truncated SVD	$n^2 \lambda_*^{-\gamma}$	n^2

Computational regularization: iterations control statistics and time complexity

Computational regularization: iterations control statistics and time complexity

Built-in regularization path

 Computational regularization: iterations control statistics and time complexity

Built-in regularization path

Is there an advantage going for on-line learning?

 Computational regularization: iterations control statistics and time complexity

Built-in regularization path

▶ Is there an advantage going for **on-line learning**? Not much, maybe $n^2 \log n$? (Bach Dieluevet '15 – R. Villa '15)

 Computational regularization: iterations control statistics and time complexity

Built-in regularization path

▶ Is there an advantage going for **on-line learning**? Not much, maybe $n^2 \log n$? (Bach Dieluevet '15 – R. Villa '15)

 Computational regularization: principles to control statistics, time and space complexity

Outline

Learning with dictionaries and kernels

Data Dependent Subsampling

Data Independent Subsampling

1. pick w_i at random...

1. pick w_i at random... from training set (Smola, Scholköpf '00)

$$\tilde{w}_1,\ldots,\tilde{w}_M\subset x_1,\ldots x_n$$
 $M\ll n$



1. pick w_i at random... from training set (Smola, Scholköpf '00)

$$\tilde{w}_1,\ldots,\tilde{w}_M\subset x_1,\ldots x_n$$
 $M\ll n$



2. perform KRR on

$$\mathcal{H}_M = \{ f \mid f(x) = \sum_{i=1}^{M} c_i q(x, \tilde{w}_i), \ c_i \in \mathbb{R}, \ w_i \in \mathbb{R}^d, \ \mathcal{M} \in \mathbb{N} \}.$$

1. pick w_i at random... from training set (Smola, Scholköpf '00)

$$\tilde{w}_1,\ldots,\tilde{w}_M\subset x_1,\ldots x_n$$
 $M\ll n$



2. perform KRR on

$$\mathcal{H}_M = \{ f \mid f(x) = \sum_{i=1}^{M} c_i q(x, \tilde{w}_i), \ c_i \in \mathbb{R}, \ w_i \in \mathbb{R}^d, \ \mathcal{M} \in \mathbb{N} \}.$$

Linear System



Complexity

- ▶ Space $O(n^2) \rightarrow O(nM)$ ▶ Time $O(n^3) \rightarrow O(nM^2)$

1. pick w_i at random... from training set (Smola, Scholköpf '00)

$$\tilde{w}_1,\ldots,\tilde{w}_M\subset x_1,\ldots x_n$$
 $M\ll n$



2. perform KRR on

$$\mathcal{H}_M = \{ f \mid f(x) = \sum_{i=1}^{M} c_i q(x, \tilde{\boldsymbol{w}}_i), \ c_i \in \mathbb{R}, \ \boldsymbol{w}_i \in \mathbb{R}^d, \ \mathcal{M} \in \mathbb{N} \}.$$

Linear System



Complexity

- ▶ Space $O(n^2) \to O(nM)$ ▶ Time $Q(n^3) \to O(nM^2)$

What about statistics? What's the price for efficient computations?

Putting our Result in Context

Many different subsampling schemes

(Smola, Scholkopf '00; Williams, Seeger '01; ... 20+)

Putting our Result in Context

 Many different subsampling schemes (Smola, Scholkopf '00; Williams, Seeger '01; ... 20+)

Theoretical guarantees mainly on matrix approximation (Mahoney and Drineas '09; Cortes et al '10, Kumar et al.'12...10+)

$$\|\widehat{Q} - \widehat{Q}_M\| \lesssim \frac{1}{\sqrt{M}}$$

Putting our Result in Context

 Many different subsampling schemes (Smola, Scholkopf '00; Williams, Seeger '01; ... 20+)

 Theoretical guarantees mainly on matrix approximation (Mahoney and Drineas '09; Cortes et al '10, Kumar et al.'12 ... 10+)

$$\|\widehat{Q} - \widehat{Q}_M\| \lesssim \frac{1}{\sqrt{M}}$$

Few prediction guarantees either suboptimal or in restricted setting (Cortes et al. '10; Jin et al. '11, Bach '13, Alaoui, Mahoney '14)

(Rudi, Camoriano, Rosasco, '15)

Theorem If $f^* \in \mathcal{H}$, then

$$\lambda_* = \frac{1}{\sqrt{n}} \quad , M_* = \frac{1}{\lambda_*}, \quad \mathbb{E}\left(\widehat{f}_{\lambda_*, M_*}(x) - f^*(x)\right)^2 \lesssim \frac{1}{\sqrt{n}}$$

(Rudi, Camoriano, Rosasco, '15)

Theorem If $f^* \in \mathcal{H}$, then

$$\lambda_* = \frac{1}{\sqrt{n}} \quad , M_* = \frac{1}{\lambda_*}, \quad \mathbb{E} \left(\widehat{f}_{\lambda_*, M_*}(x) - f^*(x) \right)^2 \lesssim \frac{1}{\sqrt{n}}$$

Remarks

(Rudi, Camoriano, Rosasco, '15)

Theorem If $f^* \in \mathcal{H}$, then

$$\lambda_* = \frac{1}{\sqrt{n}} \quad , M_* = \frac{1}{\lambda_*}, \quad \mathbb{E}\left(\widehat{f}_{\lambda_*, M_*}(x) - f^*(x)\right)^2 \lesssim \frac{1}{\sqrt{n}}$$

Remarks

1. Subsampling achives **optimal** bound...

(Rudi, Camoriano, Rosasco, '15)

Theorem If $f^* \in \mathcal{H}$, then

$$\lambda_* = \frac{1}{\sqrt{n}} \quad , M_* = \frac{1}{\lambda_*}, \quad \mathbb{E}\left(\widehat{f}_{\lambda_*,M_*}(x) - f^*(x)\right)^2 \lesssim \frac{1}{\sqrt{n}}$$

Remarks

- 1. Subsampling achives **optimal** bound...
- 2. ... with $M_* \sim \sqrt{n}$!!

(Rudi, Camoriano, Rosasco, '15)

Theorem If $f^* \in \mathcal{H}$, then

$$\lambda_* = \frac{1}{\sqrt{n}} \quad , M_* = \frac{1}{\lambda_*}, \quad \mathbb{E}\left(\widehat{f}_{\lambda_*,M_*}(x) - f^*(x)\right)^2 \lesssim \frac{1}{\sqrt{n}}$$

Remarks

- 1. Subsampling achives **optimal** bound...
- 2. ... with $M_* \sim \sqrt{n}$!!
- 3. More generally,

$$\lambda_* = n^{-\frac{1}{2s+1}}, \quad M_* = \frac{1}{\lambda_*}, \quad \mathbb{E}_x \left(\widehat{f}_{\lambda_*, M_*}(x) - f^*(x) \right)^2 \lesssim n^{-\frac{2s}{2s+1}}$$

(Rudi, Camoriano, Rosasco, '15)

Theorem If $f^* \in \mathcal{H}$, then

$$\lambda_* = \frac{1}{\sqrt{n}} \quad , M_* = \frac{1}{\lambda_*}, \quad \mathbb{E} \left(\widehat{f}_{\lambda_*, M_*}(x) - f^*(x) \right)^2 \lesssim \frac{1}{\sqrt{n}}$$

Remarks

- 1. Subsampling achives optimal bound...
- 2. ... with $M_* \sim \sqrt{n}$!!
- 3. More generally,

$$\lambda_* = n^{-\frac{1}{2s+1}}, \quad M_* = \frac{1}{\lambda_*}, \quad \mathbb{E}_x \left(\widehat{f}_{\lambda_*, M_*}(x) - f^*(x) \right)^2 \lesssim n^{-\frac{2s}{2s+1}}$$

Note: An interesting insight is obtained rewriting the result...

(Rudi, Camoriano, Rosasco, '15)

A simple idea: "swap" the role of λ and M...

(Rudi, Camoriano, Rosasco, '15)

A simple idea: "swap" the role of λ and M...

Theorem If $f^* \in \mathcal{H}$, then

$$M_* = n^{\frac{1}{2s+1}}, \quad \lambda_* = \frac{1}{M_*}, \quad \mathbb{E}_x \left(\widehat{f}_{\lambda_*, M_*}(x) - f^*(x) \right)^2 \lesssim n^{-\frac{2s}{2s+1}}$$

(Rudi, Camoriano, Rosasco, '15)

A simple idea: "swap" the role of λ and M...

Theorem If $f^* \in \mathcal{H}$, then

$$M_* = n^{\frac{1}{2s+1}}, \quad \lambda_* = \frac{1}{M_*}, \quad \mathbb{E}_x \left(\widehat{f}_{\lambda_*, M_*}(x) - f^*(x) \right)^2 \lesssim n^{-\frac{2s}{2s+1}}$$

• λ and M play the same role. .

... new interpretation: subsampling regularizes!

(Rudi, Camoriano, Rosasco, '15)

A simple idea: "swap" the role of λ and M...

Theorem If $f^* \in \mathcal{H}$, then

$$M_* = n^{\frac{1}{2s+1}}, \quad \lambda_* = \frac{1}{M_*}, \quad \mathbb{E}_x \left(\widehat{f}_{\lambda_*, M_*}(x) - f^*(x) \right)^2 \lesssim n^{-\frac{2s}{2s+1}}$$

λ and M play the same role... ...new interpretation: subsampling regularizes!

New natural incremental algorithm...

Algorithm

(Rudi, Camoriano, Rosasco, '15)

A simple idea: "swap" the role of λ and M...

Theorem If $f^* \in \mathcal{H}$, then

$$M_* = n^{\frac{1}{2s+1}}, \quad \lambda_* = \frac{1}{M_*}, \quad \mathbb{E}_x \left(\widehat{f}_{\lambda_*, M_*}(x) - f^*(x) \right)^2 \lesssim n^{-\frac{2s}{2s+1}}$$

- λ and M play the same role... ...new interpretation: subsampling regularizes!
- New natural incremental algorithm...

Algorithm

1. Pick a center + compute solution

(Rudi, Camoriano, Rosasco, '15)

A simple idea: "swap" the role of λ and M...

Theorem If $f^* \in \mathcal{H}$, then

$$M_* = n^{\frac{1}{2s+1}}, \quad \lambda_* = \frac{1}{M_*}, \quad \mathbb{E}_x \left(\widehat{f}_{\lambda_*, M_*}(x) - f^*(x) \right)^2 \lesssim n^{-\frac{2s}{2s+1}}$$

- λ and M play the same role... ...new interpretation: subsampling regularizes!
- New natural incremental algorithm...

Algorithm

- 1. Pick a center + compute solution
- 2. Pick another center + rank one update

(Rudi, Camoriano, Rosasco, '15)

A simple idea: "swap" the role of λ and M...

Theorem If $f^* \in \mathcal{H}$, then

$$M_* = n^{\frac{1}{2s+1}}, \quad \lambda_* = \frac{1}{M_*}, \quad \mathbb{E}_x \left(\widehat{f}_{\lambda_*, M_*}(x) - f^*(x) \right)^2 \lesssim n^{-\frac{2s}{2s+1}}$$

- λ and M play the same role...
 ... new interpretation: subsampling regularizes!
- New natural incremental algorithm...

Algorithm

- 1. Pick a center + compute solution
- 2. Pick another center + rank one update
- 3. Pick another center . . .

CoRe Illustrated





Computation controls stability!

Time/space requirement tailored to generalization

Experiments

comparable/better w.r.t. the state of the art

Dataset	n_{tr}	d	Incremental CoRe	Standard KRLS	Standard Nyström	Random Features	Fastfood RF
Ins. Co.	5822	85	$0.23180 \pm 4 \times 10^{-5}$	0.231	0.232	0.266	0.264
CPU	6554	21	$\bf 2.8466 \pm 0.0497$	7.271	6.758	7.103	7.366
CT slices	42800	384	$\bf 7.1106 \pm 0.0772$	NA	60.683	49.491	43.858
Year Pred.	463715	90	$0.10470 \pm 5 imes 10^{-5}$	NA	0.113	0.123	0.115
Forest	522910	54	0.9638 ± 0.0186	NA	0.837	0.840	0.840

Random Features (Rahimi, Recht '07)

► Fastfood (Le et al. '13)

Summary so far

- Optimal learning with data dependent subsampling
- Computational regularization: subsampling regularizes!

Summary so far

- Optimal learning with data dependent subsampling
- Computational regularization: subsampling regularizes!

Few more questions:

- Can one do better than **uniform** sampling?
- What about data independent sampling?
(Approximate) Leverage scores

Leverage scores

$$l_i(t) = (\widehat{Q}(\widehat{Q} + tnI)^{-1})_{ii}$$

(Approximate) Leverage scores

Leverage scores

$$l_i(t) = (\widehat{Q}(\widehat{Q} + tnI)^{-1})_{ii}$$

ALS With probability at least $1 - \delta$,

$$\frac{1}{T}l_i(t) \le \tilde{l}_i(t) \le Tl_i(t)$$

(Approximate) Leverage scores

Leverage scores

$$l_i(t) = (\widehat{Q}(\widehat{Q} + tnI)^{-1})_{ii}$$

$$\label{eq:alsolution} \begin{split} & \mathsf{ALS} \\ & \mathsf{With \ probability \ at \ least \ } 1-\delta, \end{split}$$

$$\frac{1}{T}l_i(t) \le \tilde{l}_i(t) \le Tl_i(t)$$

ALS subsampling

Pick $\tilde{w}_1, \ldots, \tilde{w}_M \subset x_1, \ldots x_n$ with replacement, and probability

$$P_t(i) = \tilde{l}_i(t) / \sum_j \tilde{l}_j(t).$$

(Rudi, Camoriano, Rosasco, '15)

Theorem

$$\lambda_* = n^{-\frac{1}{1+\gamma}}, \quad M_* = \frac{1}{\lambda_*^{\gamma}}, \quad \mathbb{E} \, (\widehat{f}_{\lambda_*, M_*}(x) - f^*(x))^2 \lesssim n^{-\frac{1}{1+\gamma}}$$

(Rudi, Camoriano, Rosasco, '15)

Theorem

If $f_*\in\mathcal{H}$ and the integral operator associated to q has eigenvalue decay $i^{-\frac{1}{\gamma}},\,\gamma\in(0,1)$ then

$$\lambda_* = n^{-\frac{1}{1+\gamma}}, \quad M_* = \frac{1}{\lambda_*^{\gamma}}, \quad \mathbb{E} \, (\widehat{f}_{\lambda_*, M_*}(x) - f^*(x))^2 \lesssim n^{-\frac{1}{1+\gamma}}$$

Non uniform subsampling achieves optimal bound

(Rudi, Camoriano, Rosasco, '15)

Theorem

$$\lambda_* = n^{-\frac{1}{1+\gamma}}, \quad M_* = \frac{1}{\lambda_*^{\gamma}}, \quad \mathbb{E} \, (\widehat{f}_{\lambda_*, M_*}(x) - f^*(x))^2 \lesssim n^{-\frac{1}{1+\gamma}}$$

- Non uniform subsampling achieves optimal bound
- Most importantly: potentially much fewer samples needed

(Rudi, Camoriano, Rosasco, '15)

Theorem

$$\lambda_* = n^{-\frac{1}{1+\gamma}}, \quad M_* = \frac{1}{\lambda_*^{\gamma}}, \quad \mathbb{E} \, (\widehat{f}_{\lambda_*, M_*}(x) - f^*(x))^2 \lesssim n^{-\frac{1}{1+\gamma}}$$

- Non uniform subsampling achieves optimal bound
- Most importantly: potentially much fewer samples needed
- ▶ Need efficient ALS computation (...)

(Rudi, Camoriano, Rosasco, '15)

Theorem

$$\lambda_* = n^{-\frac{1}{1+\gamma}}, \quad M_* = \frac{1}{\lambda_*^{\gamma}}, \quad \mathbb{E} \, (\widehat{f}_{\lambda_*, M_*}(x) - f^*(x))^2 \lesssim n^{-\frac{1}{1+\gamma}}$$

- Non uniform subsampling achieves optimal bound
- Most importantly: potentially much fewer samples needed
- Need efficient ALS computation (...)
- Extensions to more general sampling schemes are possible

Outline

Learning with dictionaries and kernels

Data Dependent Subsampling

Data Independent Subsampling

$$\widehat{f}(x) = \sum_{i=1}^{M} c_i q(x, w_i)$$

 \blacktriangleright q general non linear function

$$\widehat{f}(x) = \sum_{i=1}^{M} c_i q(x, w_i)$$

- \blacktriangleright q general non linear function
- pick \tilde{w}_i at random according to a distribution μ

$$\tilde{w}_1,\ldots,\tilde{w}_M\sim\mu$$

$$\widehat{f}(x) = \sum_{i=1}^{M} c_i q(x, w_i)$$

 \blacktriangleright q general non linear function

• pick \tilde{w}_i at random according to a distribution μ

$$\tilde{w}_1,\ldots,\tilde{w}_M\sim\mu$$

▶ perform KRR on

$$\mathcal{H}_M = \{ f \mid f(x) = \sum_{i=1}^M c_i q(x, \tilde{w}_i), \ c_i \in \mathbb{R} \}.$$

(Rahimi, Recht '07)

Consider

$$q(x,w) = e^{iw^T x},$$

(Rahimi, Recht '07)

Consider

$$q(x,w) = e^{iw^T x}, \qquad \qquad w \sim \mu(w) = \mathcal{N}(0,I)$$

(Rahimi, Recht '07)

Consider

$$q(x,w) = e^{iw^T x}, \qquad \qquad w \sim \mu(w) = \mathcal{N}(0,I)$$

Then

$$\mathbb{E}_{w}[q(x,w)q(x',w)] = e^{-\|x-x'\|^{2}\gamma} = K(x,x')$$

(Rahimi, Recht '07)

Consider

$$q(x,w) = e^{iw^T x}, \qquad \qquad w \sim \mu(w) = \mathcal{N}(0,I)$$

Then

$$\mathbb{E}_{w}[q(x,w)q(x',w)] = e^{-\|x-x'\|^{2}\gamma} = K(x,x')$$

By sampling $\tilde{w}_1, \ldots, \tilde{w}_M$ we are considering the **approximating kernel**

$$\frac{1}{M}\sum_{i=1}^{M} \left[q(x,\tilde{w}_i)q(x',\tilde{w}_i)\right] = \widetilde{K}(x,x')$$

▶ translation invariant kernels K(x, x') = H(x - x'),

• translation invariant kernels K(x, x') = H(x - x'),

$$q(x,w) = e^{iw^T x},$$

• translation invariant kernels K(x, x') = H(x - x'),

$$q(x,w) = e^{iw^T x}, \qquad \qquad w \sim \mu = \mathcal{F}(H)$$

 $\blacktriangleright \ \ {\rm translation \ invariant \ kernels \ } K(x,x') = H(x-x'),$

$$q(x,w) = e^{iw^T x}, \qquad \qquad w \sim \mu = \mathcal{F}(H)$$

infinite neural nets kernels

• translation invariant kernels K(x, x') = H(x - x'),

$$q(x,w) = e^{iw^T x}, \qquad \qquad w \sim \mu = \mathcal{F}(H)$$

infinite neural nets kernels

$$q(x,w) = |w^T x + b|_+,$$

• translation invariant kernels K(x, x') = H(x - x'),

$$q(x,w) = e^{iw^T x}, \qquad \qquad w \sim \mu = \mathcal{F}(H)$$

infinite neural nets kernels

$$q(x,w) = |w^T x + b|_+,$$
 $(w,b) \sim \mu = U[\mathbb{S}^d]$

• translation invariant kernels K(x, x') = H(x - x'),

$$q(x,w) = e^{iw^T x}, \qquad w \sim \mu = \mathcal{F}(H)$$

infinite neural nets kernels

$$q(x,w) = |w^T x + b|_+,$$
 $(w,b) \sim \mu = U[\mathbb{S}^d]$

- infinite dot product kernels
- homogeneous additive kernels
- group invariant kernels

▶ ...

Note: Connections with hashing and sketching techniques.

Properties of Random Features

Properties of Random Features

Optimization

► Time: $O(n^3)$ $O(nM^2)$ ► Space: $O(n^2)$ O(nM)

Properties of Random Features

Optimization

▶ Time: $O(n^3)$ $O(nM^2)$ ▶ Space: $O(n^2)$ O(nM)

Statistics

As before: do we pay a price for efficient computations?



Previous works

Many different random features for different kernels (Rahimi, Recht '07, Vedaldi, Zisserman, ...10+)

Previous works

 Many different random features for different kernels (Rahimi, Recht '07, Vedaldi, Zisserman, ... 10+)

 Theoretical guarantees: mainly kernel approximation (Rahimi, Recht '07, ..., Sriperumbudur and Szabo '15)

$$|K(x,x') - \widetilde{K}(x,x')| \lesssim \frac{1}{\sqrt{M}}$$

Previous works

 Many different random features for different kernels (Rahimi, Recht '07, Vedaldi, Zisserman, ... 10+)

 Theoretical guarantees: mainly kernel approximation (Rahimi, Recht '07, ..., Sriperumbudur and Szabo '15)

$$|K(x,x') - \widetilde{K}(x,x')| \lesssim \frac{1}{\sqrt{M}}$$

 Statistical guarantees suboptimal or in restricted setting (Rahimi, Recht '09, Yang et al. '13 ..., Bach '15)

Let

$$q(x,w) = e^{iw^T x},$$

Let

$$q(x,w) = e^{iw^T x}, \qquad w \sim \mu(w) = c_d \left(\frac{1}{1+\|w\|^2}\right)^{\frac{d+1}{2}}$$

Let

$$q(x,w) = e^{iw^T x}, \qquad w \sim \mu(w) = c_d \left(\frac{1}{1+\|w\|^2}\right)^{\frac{d+1}{2}}$$

Theorem If $f_* \in \mathcal{H}_s$ Sobolev space, then

$$\lambda_* = n^{-\frac{1}{2s+1}}, \quad M_* = \frac{1}{\lambda_*^{2s}}, \quad \mathbb{E}\left(\widehat{f}_{\lambda_*,M_*}(x) - f^*(x)\right)^2 \lesssim n^{-\frac{2s}{2s+1}}$$

Let

$$q(x,w) = e^{iw^T x}, \qquad w \sim \mu(w) = c_d \left(\frac{1}{1+\|w\|^2}\right)^{\frac{d+1}{2}}$$

Theorem If $f_* \in \mathcal{H}_s$ Sobolev space, then

$$\lambda_* = n^{-\frac{1}{2s+1}}, \quad M_* = \frac{1}{\lambda_*^{2s}}, \quad \mathbb{E}\left(\widehat{f}_{\lambda_*,M_*}(x) - f^*(x)\right)^2 \lesssim n^{-\frac{2s}{2s+1}}$$

Random feature achieves optimal bound!

Let

$$q(x,w) = e^{iw^T x}, \qquad w \sim \mu(w) = c_d \left(\frac{1}{1+\|w\|^2}\right)^{\frac{d+1}{2}}$$

Theorem If $f_* \in \mathcal{H}_s$ Sobolev space, then

$$\lambda_* = n^{-\frac{1}{2s+1}}, \quad M_* = \frac{1}{\lambda_*^{2s}}, \quad \mathbb{E}\left(\widehat{f}_{\lambda_*,M_*}(x) - f^*(x)\right)^2 \lesssim n^{-\frac{2s}{2s+1}}$$

- Random feature achieves optimal bound!
- Efficient worst case subsampling $M_* \sim \sqrt{n}$, but cannot exploit smoothness.
Remarks & Extensions

Nÿstrom vs Random features

- Both achieve optimal rates
- Nÿstrom seems to need fewer samples (random centers)

Remarks & Extensions

Nÿstrom vs Random features

- Both achieve optimal rates
- Nÿstrom seems to need fewer samples (random centers)

How tight are the results?

Remarks & Extensions

Nÿstrom vs Random features

- Both achieve optimal rates
- Nÿstrom seems to need fewer samples (random centers)

How tight are the results?



Test Error

$\log M$



Contributions

- > Optimal bounds for data dependent/independent subsampling
- Subsampling: Nÿstrom vs Random features
- Beyond ridge regression: early stopping and multiple passes SGD (coming up in AISTATS!)

Contributions

- Optimal bounds for data dependent/independent subsampling
- Subsampling: Nÿstrom vs Random features
- Beyond ridge regression: early stopping and multiple passes SGD (coming up in AISTATS!)

Some questions:

- Quest for the **best** sampling
- ▶ Regularization by projection: inverse problems and preconditioning
- Beyond randomization: non convex optimization?

Contributions

- Optimal bounds for data dependent/independent subsampling
- Subsampling: Nÿstrom vs Random features
- Beyond ridge regression: early stopping and multiple passes SGD (coming up in AISTATS!)

Some questions:

- Quest for the **best** sampling
- Regularization by projection: inverse problems and preconditioning
- Beyond randomization: non convex optimization?

Some perspectives:

- Computational regularization: subsampling regularizes
- ► Algorithm design: Control stability with computations