

Optimal kernel methods for large scale learning

Alessandro Rudi
INRIA - École Normale Supérieure, Paris

joint work with Luigi Carratino, Lorenzo Rosasco

6 Mar 2018 – École Polytechnique

Learning problem

The problem \mathcal{P}

Find

$$f_{\mathcal{H}} = \operatorname{argmin}_{f \in \mathcal{H}} \mathcal{E}(f), \quad \mathcal{E}(f) = \int d\rho(x, y) (y - f(x))^2$$

with ρ **unknown** but given $(x_i, y_i)_{i=1}^n$ i.i.d. samples.

Remarks:

- ▶ stochastic optimization problem
- ▶ \mathcal{H} is a space of candidate solutions.

Outline

Learning with kernels

Random projections

FALKON: Random projections and preconditioning

Kernel ridge regression

Let K p.d. kernel (e.g. $K(x, x') = e^{-\gamma\|x-x'\|^2}$) and

$$\mathcal{H} = \overline{\text{span}\{K(x, \cdot) | x \in X\}},$$

Problem $\hat{\mathcal{P}}_n$

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

KRR: Statistic (worst case)

Noise

$$\mathbb{E}[|Y|^p \mid X = x] \leq \frac{1}{2} p! \sigma^2 b^{p-2}, \quad \forall p \geq 2$$

Kernel boundness $\sup_x K(x, x) < \infty$.

Best model There exists $f_{\mathcal{H}}$ solving

$$\min_{f \in \mathcal{H}} \mathcal{E}(f).$$

KRR: Statistic (worst case)

Noise

$$\mathbb{E}[|Y|^p \mid X = x] \leq \frac{1}{2} p! \sigma^2 b^{p-2}, \quad \forall p \geq 2$$

Kernel boundness $\sup_x K(x, x) < \infty$.

Best model There exists $f_{\mathcal{H}}$ solving

$$\min_{f \in \mathcal{H}} \mathcal{E}(f).$$

Theorem[(Caponnetto, De Vito '05)] Under the assumptions above

$$\mathbb{E} \mathcal{E}(\hat{f}_{\lambda}) - \mathcal{E}(f_{\mathcal{H}}) \lesssim \frac{1}{\lambda n} + \lambda.$$

By selecting $\lambda_n = \frac{1}{\sqrt{n}}$

$$\mathbb{E} \mathcal{E}(\hat{f}_{\lambda_n}) - \mathcal{E}(f_{\mathcal{H}}) \lesssim \frac{1}{\sqrt{n}}$$

KRR: Statistics (refined case)

Let $Lf(x') = \mathbb{E}K(x', x)f(x)$ and $\mathcal{N}(\lambda) = \text{Trace}((L + \lambda I)^{-1}L)$

Capacity condition:

$$\mathcal{N}(\lambda) = O(\lambda^{-\gamma}), \quad \gamma \in [0, 1]$$

Source condition:

$$f_{\mathcal{H}} \in \text{Range}(L^r), \quad r \geq 1/2$$

KRR: Statistics (refined case)

Let $Lf(x') = \mathbb{E}K(x', x)f(x)$ and $\mathcal{N}(\lambda) = \text{Trace}((L + \lambda I)^{-1}L)$

Capacity condition:

$$\mathcal{N}(\lambda) = O(\lambda^{-\gamma}), \quad \gamma \in [0, 1]$$

Source condition:

$$f_{\mathcal{H}} \in \text{Range}(L^r), \quad r \geq 1/2$$

Theorem[(Caponnetto, De Vito '05)] Under (basic) and (refined)

$$\mathbb{E}\mathcal{E}(\hat{f}_{\lambda}) - \mathcal{E}(f_{\mathcal{H}}) \lesssim \frac{\mathcal{N}(\lambda)}{n} + \lambda^{2r}.$$

By selecting $\lambda_n = n^{-\frac{1}{2r+\gamma}}$

$$\mathbb{E}\mathcal{E}(\hat{f}_{\lambda_n}) - \mathcal{E}(f_{\mathcal{H}}) \lesssim n^{-\frac{2r}{2r+\gamma}}$$

KRR: Statistics (refined case)

Let $Lf(x') = \mathbb{E}K(x', x)f(x)$ and $\mathcal{N}(\lambda) = \text{Trace}((L + \lambda I)^{-1}L)$

Capacity condition:

$$\mathcal{N}(\lambda) = O(\lambda^{-\gamma}), \quad \gamma \in [0, 1]$$

Source condition:

$$f_{\mathcal{H}} \in \text{Range}(L^r), \quad r \geq 1/2$$

Theorem[(Caponnetto, De Vito '05)] Under (basic) and (refined)

$$\mathbb{E}\mathcal{E}(\hat{f}_{\lambda}) - \mathcal{E}(f_{\mathcal{H}}) \lesssim \frac{\mathcal{N}(\lambda)}{n} + \lambda^{2r}.$$

By selecting $\lambda_n = n^{-\frac{1}{2r+\gamma}}$

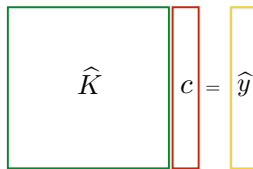
$$\mathbb{E}\mathcal{E}(\hat{f}_{\lambda_n}) - \mathcal{E}(f_{\mathcal{H}}) \lesssim n^{-\frac{2r}{2r+\gamma}}$$

KRR: Optimization

$$\hat{f}_\lambda(x) = \sum_{i=1}^n K(x, x_i) c_i$$

$$(\hat{K} + \lambda n I) c = \hat{y}$$

Linear System



The diagram illustrates a linear system $\hat{K}c = \hat{y}$. It features a large green square matrix labeled \hat{K} , a narrow red vertical rectangle labeled c , and a narrow yellow vertical rectangle labeled \hat{y} . The matrix and vector are separated by an equals sign.

Computations

Space $O(n^2)$

Kernel eval. $O(n^2)$

Time $O(n^3)$

Large scale ML:

Running out of time and space ... can we fix it?

Computations for optimal statistical accuracy

Model: $O(n)$

Space: $O(n^2)$

Kernel eval.: $O(n^2)$

Time: $O(n^3)$

Outline

Learning with kernels

Random projections

FALKON: Random projections and preconditioning

Random projections

Solve $\hat{\mathcal{P}}_n$ on $\mathcal{H}_M = \text{span}\{K(\tilde{x}_1, \cdot), \dots, K(\tilde{x}_M, \cdot)\}$

$$\hat{f}_{\lambda, M} = \underset{f \in \mathcal{H}_M}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2 + \lambda \|f\|_{\mathcal{H}}^2$$

Random projections

Solve $\hat{\mathcal{P}}_n$ on $\mathcal{H}_M = \text{span}\{K(\tilde{x}_1, \cdot), \dots, K(\tilde{x}_M, \cdot)\}$

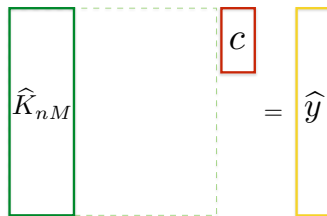
$$\hat{f}_{\lambda, M} = \underset{f \in \mathcal{H}_M}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2 + \lambda \|f\|_{\mathcal{H}}^2$$

► ... that is, pick M columns at random

$$\hat{f}_{\lambda, M}(x) = \sum_{i=1}^M K(x, \tilde{x}_i) c_i$$

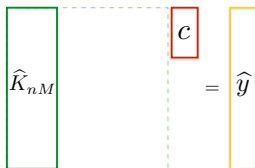
$$(\hat{K}_{nM}^\top \hat{K}_{nM} + \lambda n \hat{K}_{MM}) c = \hat{K}_{nM}^\top \hat{y}$$

Linear System



- **Nyström methods** (Smola, Scholköpfung '00)
- Gaussian processes: inducing inputs (Quionero-Candela et al '05)

Nystrom KRR: Computations



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

Computations (train)

- ▶ **Space** $\cancel{O(n^2)} \rightarrow O(M_n^2)$
- ▶ **Kernel eval.** $\cancel{O(n^2)} \rightarrow O(nM_n)$
- ▶ **Time** $\cancel{O(n^3)} \rightarrow O(nM_n^2)$

Computations (test) $\cancel{O(n)} \rightarrow O(M_n)$

Nyström KRR: Statistics (worst case)

Theorem[Rudi, Camoriano, Rosasco '15] Under the basic assumptions

$$\mathbb{E}\mathcal{E}(\hat{f}_{\lambda,M}) - \mathcal{E}(f_{\mathcal{H}}) \lesssim \frac{1}{\lambda n} + \lambda + \frac{1}{M}.$$

Nyström KRR: Statistics (worst case)

Theorem[Rudi, Camoriano, Rosasco '15] Under the basic assumptions

$$\mathbb{E}\mathcal{E}(\hat{f}_{\lambda,M}) - \mathcal{E}(f_{\mathcal{H}}) \lesssim \frac{1}{\lambda n} + \lambda + \frac{1}{M}.$$

By selecting $\lambda_n = \frac{1}{\sqrt{n}}$, $M_n = \frac{1}{\lambda_n}$

$$\mathbb{E}\mathcal{E}(\hat{f}_{\lambda_n,M_n}) - \mathcal{E}(f_{\mathcal{H}}) \lesssim \frac{1}{\sqrt{n}}$$

Remarks

$$M = O(\sqrt{n}) \text{ suffices for optimal generalization}$$

- ▶ Previous works: only for fixed design
(Bach '13, Alaoui, Mahoney, '15, Yang et al. '15, Musco, Musco '16)
- ▶ Matches statistical minimax lower bounds [Caponnetto, De Vito '05].
- ▶ Special case: Sobolev spaces with $s = d/2$, e.g. exponential kernel and Fourier features.
- ▶ Same statistical bound of (kernel) ridge regression [Caponnetto, De Vito '05].

Nyström KRR: Statistics (refined)

Theorem[Rudi, Camoriano, Rosasco '15] Under (basic) and (refined)

$$\mathbb{E}\mathcal{E}(\hat{f}_{\lambda,M}) - \mathcal{E}(f_{\mathcal{H}}) \lesssim \frac{\mathcal{N}(\lambda)}{n} + \lambda^{2r} + \frac{1}{M}.$$

Nyström KRR: Statistics (refined)

Theorem[Rudi, Camoriano, Rosasco '15] Under (basic) and (refined)

$$\mathbb{E}\mathcal{E}(\hat{f}_{\lambda,M}) - \mathcal{E}(f_{\mathcal{H}}) \lesssim \frac{\mathcal{N}(\lambda)}{n} + \lambda^{2r} + \frac{1}{M}.$$

By selecting $\lambda_n = n^{-\frac{1}{2r+\gamma}}$, $M_n = \frac{1}{\lambda_n}$

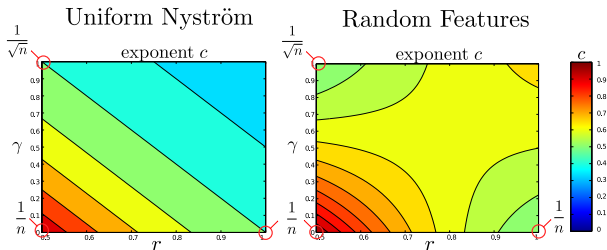
$$\mathbb{E}\mathcal{E}(\hat{f}_{\lambda_n,M_n}) - \mathcal{E}(f_{\mathcal{H}}) \lesssim n^{-\frac{2r}{2r+\gamma}}$$

Remarks

- ▶ The obtained rate is minmax optimal [Caponnetto, De Vito '05].
- ▶ Reduces to worst case for $\gamma = 1, r = 1/2$.

Comparison with Random Features: [Rudi, Rosasco '17]

$$M = n^c$$



Computations required for $1/\sqrt{n}$ rate

Model: $O(\sqrt{n})$

Space: $O(n)$

Kernel eval.: $O(n\sqrt{n})$

Time: $O(n^2)$

Possible improvements:

- ▶ adaptive sampling
- ▶ **optimization**

Outline

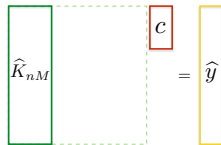
Learning with kernels

Random projections

FALKON: Random projections and preconditioning

Beyond $O(n^2)$ time?

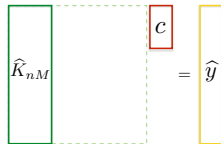
$$(\hat{K}_{nM}^\top \hat{K}_{nM} + \lambda n \hat{K}_{MM}) c = \hat{K}_{nM}^\top \hat{y}.$$



Bottleneck: compute $\hat{K}_{nM}^\top \hat{K}_{nM}$ requires $O(nM^2)$ time.

Optimization to rescue

$$\underbrace{\hat{K}_{nM}^\top \hat{K}_{nM} + \lambda n \hat{K}_{MM}}_H c = \underbrace{\hat{K}_{nM}^\top \hat{y}}_b.$$



Idea: First order methods

$$c_t = c_{t-1} - \frac{\tau}{n} \left[\hat{K}_{nM}^\top (\hat{K}_{nM} c_{t-1} - y_n) + \lambda n \hat{K}_{MM} c_{t-1} \right]$$

Pros: requires $O(nMt)$

Cons: $t \propto \kappa(H)$ arbitrarily large- $\kappa(H) = \sigma_{\max}(H)/\sigma_{\min}(H)$ condition number.

Preconditioning

Idea: solve an equivalent linear system with better condition number

Preconditioning

$$Hc = b \quad \mapsto \quad \textcolor{red}{P}^\top H \textcolor{red}{P} \beta = \textcolor{red}{P}^\top b, \quad c = \textcolor{red}{P} \beta.$$

Ideally $PP^\top = H^{-1}$, so that

$$t = O(\kappa(H)) \quad \mapsto \quad t = O(1)!$$

Computing a good preconditioning can be hard!

Remarks

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

$$H = K + \lambda nI$$

Can we precondition Nystrom-KRR?

Preconditioning Nystom-KRR

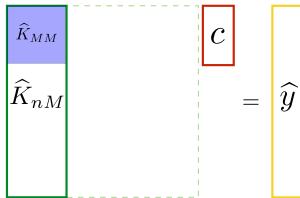
Consider $H := \hat{K}_{nM}^\top \hat{K}_{nM} + \lambda n \hat{K}_{MM}$

Proposed Preconditioning

$$PP^\top = \left(\frac{n}{M} \hat{K}_{MM}^2 + \lambda n \hat{K}_{MM} \right)^{-1}$$

Compare to naive preconditioning

$$PP^\top = \left(\hat{K}_{nM}^\top \hat{K}_{nM} + \lambda n \hat{K}_{MM} \right)^{-1}.$$



Baby FALKON

Proposed Preconditioning

$$PP^\top = \left(\frac{n}{M} \hat{K}_{MM}^2 + \lambda n \hat{K}_{MM} \right)^{-1},$$

Gradient descent

$$\hat{f}_{\lambda, M, t}(x) = \sum_{i=1}^M K(x, \tilde{x}_i) c_{t, i}, \quad c_t = \textcolor{red}{P} \beta_t$$

$$\beta_t = \beta_{t-1} - \frac{\tau}{n} \textcolor{red}{P}^\top \left[\hat{K}_{nM}^\top (\hat{K}_{nM} \textcolor{red}{P} \beta_{t-1} - y_n) + \lambda n \hat{K}_{MM} \textcolor{red}{P} \beta_{t-1} \right]$$

FALKON

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

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

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



Falkon statistics (worst case)

Theorem

Under (basic), when $M > \frac{\log n}{\lambda}$,

$$\mathbb{E}\mathcal{E}(\hat{f}_{\lambda_n, M_n, t_n}) - \mathcal{E}(f_{\mathcal{H}}) \lesssim \frac{1}{\lambda n} + \lambda + \frac{1}{M} + \exp \left[-t \left(1 - \frac{\log n}{\lambda M} \right)^{1/2} \right]$$

Falkon statistics (worst case)

Theorem

Under (basic), when $M > \frac{\log n}{\lambda}$,

$$\mathbb{E}\mathcal{E}(\hat{f}_{\lambda_n, M_n, t_n}) - \mathcal{E}(f_{\mathcal{H}}) \lesssim \frac{1}{\lambda n} + \lambda + \frac{1}{M} + \exp \left[-t \left(1 - \frac{\log n}{\lambda M} \right)^{1/2} \right]$$

By selecting

$$\lambda_n = 1/\sqrt{n}, \quad M_n = \frac{2 \log n}{\lambda}, \quad t_n = \log n,$$

then

$$\mathbb{E}\mathcal{E}(\hat{f}_{\lambda_n, M_n, t_n}) - \mathcal{E}(f_{\mathcal{H}}) \lesssim \frac{1}{\sqrt{n}}$$

Falkon statistics (refined results)

Theorem

Under (basic) and (refined), when $M > \frac{\log n}{\lambda}$,

$$\mathbb{E}\mathcal{E}(\hat{f}_{\lambda_n, M_n, t_n}) - \mathcal{E}(f_{\mathcal{H}}) \lesssim \frac{\mathcal{N}(\lambda)}{n} + \lambda^{2r} + \frac{1}{M} + \exp \left[-t \left(1 - \frac{\log n}{\lambda M} \right)^{1/2} \right]$$

Falkon statistics (refined results)

Theorem

Under (basic) and (refined), when $M > \frac{\log n}{\lambda}$,

$$\mathbb{E}\mathcal{E}(\hat{f}_{\lambda_n, M_n, t_n}) - \mathcal{E}(f_{\mathcal{H}}) \lesssim \frac{\mathcal{N}(\lambda)}{n} + \lambda^{2r} + \frac{1}{M} + \exp \left[-t \left(1 - \frac{\log n}{\lambda M} \right)^{1/2} \right]$$

By selecting

$$\lambda_n = n^{-\frac{1}{2r+\gamma}}, \quad M_n = \frac{2 \log n}{\lambda}, \quad t_n = \log n,$$

then

$$\mathbb{E}\mathcal{E}(\hat{f}_{\lambda_n, M_n, t_n}) - \mathcal{E}(f_{\mathcal{H}}) \lesssim n^{-\frac{2r}{2r+\gamma}}$$

Remarks

Relevant works

- ▶ SGD
 - ▶ RF-KRR (Rahimi, Recht '07; Bach '15; Rudi, Rosasco '17)
 - ▶ Divide and conquer (Zhang et al. '13)
 - ▶ NYTRO (Angles et al '16)
 - ▶ Nyström SGD (Lin, Rosasco '16)
-
- ▶ Same statistical properties/memory requirements
 - ▶ Much smaller time complexity

Proof: bridging statistics and optimization

Lemma

Let $\delta > 0$, $\kappa_P := \kappa(P^\top H P)$, $c_\delta = c_0 \log \frac{1}{\delta}$. When $\lambda \geq \frac{1}{n}$

$$\mathcal{E}(\hat{f}_{\lambda, M, t}) - \mathcal{E}(f_{\mathcal{H}}) \leq \mathcal{E}(\hat{f}_{\lambda, M}) - \mathcal{E}(f_{\mathcal{H}}) + c_\delta \exp(-t/\sqrt{\kappa_P}).$$

with probability $1 - \delta$.

Proof: bridging statistics and optimization

Lemma

Let $\delta > 0$, $\kappa_P := \kappa(P^\top H P)$, $c_\delta = c_0 \log \frac{1}{\delta}$. When $\lambda \geq \frac{1}{n}$

$$\mathcal{E}(\hat{f}_{\lambda, M, t}) - \mathcal{E}(f_{\mathcal{H}}) \leq \mathcal{E}(\hat{f}_{\lambda, M}) - \mathcal{E}(f_{\mathcal{H}}) + c_\delta \exp(-t/\sqrt{\kappa_P}).$$

with probability $1 - \delta$.

Lemma

Let $\delta \in (0, 1]$, $\lambda > 0$. When

$$M = \frac{2 \log \frac{1}{\delta}}{\lambda},$$

then

$$\kappa(P^\top H P) \leq \left(1 - \frac{\log \frac{1}{\delta}}{\lambda M}\right)^{-1} < 4$$

with probability $1 - \delta$.

Computational implications

Cost for optimal generalization with FALKON

$$M_n = O(\sqrt{n}), \quad t_n = \log n$$

\Downarrow

Model: $O(\sqrt{n})$

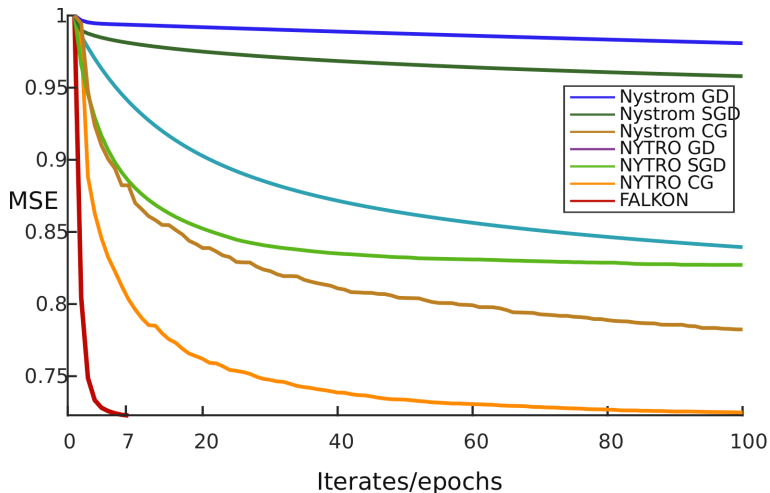
Space: $O(n)$

Kernel eval.: $O(n\sqrt{n})$

Time: $\cancel{O(n^2)} \rightarrow O(n\sqrt{n})$

In practice

Higgs dataset: $n = 10,000,000$, $M = 50,000$



Some experiments

	MillionSongs ($n \sim 10^6$)			YELP ($n \sim 10^6$)		TIMIT ($n \sim 10^6$)	
	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 ($n \sim 10^6$)			HIGGS ($n \sim 10^7$)		IMAGENET ($n \sim 10^6$)	
	c-err	AUC	Time(m)	AUC	Time(h)	c-err	Time(h)
FALKON	19.6%	0.877	4	0.833	3	20.7%	4
EigenPro	19.8%	-	6^{\dagger}	-	-	-	-
Hierarchical	20.1%	-	40^{\dagger}	-	-	-	-
Boosted Decision Tree	-	0.863	-	0.810	-	-	-
Neural Network	-	0.875	-	0.816	-	-	-
Deep Neural Network	-	0.879	4680^{\ddagger}	0.885	78^{\ddagger}	-	-
Inception-V4	-	-	-	-	-	20.0%	-

Table: Architectures: \dagger cluster with IBM POWER8 12-core cpu, 512 GB RAM, \dagger single machine with two Intel Xeon E5-2620, one Nvidia GTX Titan X GPU, 128GB RAM, \ddagger single machine.

Contributions

- ▶ Best computations so far for optimal statistics

Space $O(n)$	Time $O(n\sqrt{n})$
---------------------	----------------------------

- ▶ Random projections+iterative solvers+preconditioning
- ▶ ... fast rates
- ▶ ... adaptive sampling

Next?

- ▶ Distributed architectures...
- ▶ Open the **kernel blackbox!**

Proving $\kappa(P^\top HP) \approx 1$

Let $K_x = K(x, \cdot) \in \mathcal{H}$,

$$C = \int K_x \otimes K_x d\rho_X(x), \quad \hat{C}_n = \frac{1}{n} \sum_{i=1}^n K_{x_i} \otimes K_{x_i}, \quad \hat{C}_M = \frac{1}{M} \sum_{j=1}^M K_{\tilde{x}_j} \otimes K_{\tilde{x}_j}.$$

Recall that $P = \frac{1}{\sqrt{n}} T^{-1} A^{-1}$, $T = \text{chol}(K_{MM})$, $A = \text{chol}(\frac{1}{M} T T^\top + \lambda I)$.

Steps

1. $P^\top HP = A^{-\top} V^* (\hat{C}_n + \lambda I) V A^{-1}$

Proving $\kappa(P^\top HP) \approx 1$

Let $K_x = K(x, \cdot) \in \mathcal{H}$,

$$C = \int K_x \otimes K_x d\rho_X(x), \quad \hat{C}_n = \frac{1}{n} \sum_{i=1}^n K_{x_i} \otimes K_{x_i}, \quad \hat{C}_M = \frac{1}{M} \sum_{j=1}^M K_{\tilde{x}_j} \otimes K_{\tilde{x}_j}.$$

Recall that $P = \frac{1}{\sqrt{n}} T^{-1} A^{-1}$, $T = \text{chol}(K_{MM})$, $A = \text{chol}(\frac{1}{M} TT^\top + \lambda I)$.

Steps

$$2. \quad P^\top HP = A^{-\top} V^* (\hat{C}_M + \lambda I) V A^{-1} + A^{-\top} V^* (\hat{C}_n - \hat{C}_M) V A^{-1}$$

Proving $\kappa(P^\top HP) \approx 1$

Let $K_x = K(x, \cdot) \in \mathcal{H}$,

$$C = \int K_x \otimes K_x d\rho_X(x), \quad \hat{C}_n = \frac{1}{n} \sum_{i=1}^n K_{x_i} \otimes K_{x_i}, \quad \hat{C}_M = \frac{1}{M} \sum_{j=1}^M K_{\tilde{x}_j} \otimes K_{\tilde{x}_j}.$$

Recall that $P = \frac{1}{\sqrt{n}} T^{-1} A^{-1}$, $T = \text{chol}(K_{MM})$, $A = \text{chol}(\frac{1}{M} T T^\top + \lambda I)$.

Steps

$$3. \quad P^\top HP = I + A^{-\top} V^* (\hat{C}_n - \hat{C}_M) V A^{-1}$$

Proving $\kappa(P^\top HP) \approx 1$

Let $K_x = K(x, \cdot) \in \mathcal{H}$,

$$C = \int K_x \otimes K_x d\rho_X(x), \quad \hat{C}_n = \frac{1}{n} \sum_{i=1}^n K_{x_i} \otimes K_{x_i}, \quad \hat{C}_M = \frac{1}{M} \sum_{j=1}^M K_{\tilde{x}_j} \otimes K_{\tilde{x}_j}.$$

Recall that $P = \frac{1}{\sqrt{n}} T^{-1} A^{-1}$, $T = \text{chol}(K_{MM})$, $A = \text{chol}(\frac{1}{M} T T^\top + \lambda I)$.

Steps

$$3. \quad P^\top HP = I + E \quad \text{with} \quad E = A^{-\top} V^* (\hat{C}_n - \hat{C}_M) V A^{-1}$$

Proving $\kappa(P^\top HP) \approx 1$

Let $K_x = K(x, \cdot) \in \mathcal{H}$,

$$C = \int K_x \otimes K_x d\rho_X(x), \quad \hat{C}_n = \frac{1}{n} \sum_{i=1}^n K_{x_i} \otimes K_{x_i}, \quad \hat{C}_M = \frac{1}{M} \sum_{j=1}^M K_{\tilde{x}_j} \otimes K_{\tilde{x}_j}.$$

Recall that $P = \frac{1}{\sqrt{n}} T^{-1} A^{-1}$, $T = \text{chol}(K_{MM})$, $A = \text{chol}(\frac{1}{M} TT^\top + \lambda I)$.

Steps

$$4. \quad \kappa(P^\top HP) = \kappa(I + E) \leq \frac{1 + \|E\|}{1 - \|E\|}, \quad \text{when } \|E\| < 1,$$

Proving $\kappa(P^\top HP) \approx 1$

Let $K_x = K(x, \cdot) \in \mathcal{H}$,

$$C = \int K_x \otimes K_x d\rho_X(x), \quad \hat{C}_n = \frac{1}{n} \sum_{i=1}^n K_{x_i} \otimes K_{x_i}, \quad \hat{C}_M = \frac{1}{M} \sum_{j=1}^M K_{\tilde{x}_j} \otimes K_{\tilde{x}_j}.$$

Recall that $P = \frac{1}{\sqrt{n}} T^{-1} A^{-1}$, $T = \text{chol}(K_{MM})$, $A = \text{chol}(\frac{1}{M} T T^\top + \lambda I)$.

Steps

$$5. E = A^{-\top} V^* (\hat{C}_n - \hat{C}_M) V A^{-1} \leq 1/2 \text{ w.h.p. when } M \geq \frac{c_0 \log \frac{1}{\delta}}{\lambda}$$