

**ICTP, Fall 2018**  
**Winter School on Learning and AI**

**Class 03: Learning with (Random) Projections**

Lorenzo Rosasco

## Learning algorithm design so far

► ERM + Optimization

$$\hat{w}_\lambda = \arg \min_{w \in \mathbb{R}^d} \underbrace{\frac{1}{n} \sum_{i=1}^n (y_i - w^\top x_i)^2 + \lambda \|w\|^2}_{\hat{L}^\lambda(w)}, \quad w_{t+1} = w_t - \gamma_t \nabla \hat{L}^\lambda(w_t).$$

► Learning by optimization (GD/SGD)

$$\hat{w}_{t+1} = \hat{w}_t - \gamma_t \nabla \hat{L}(\hat{w}_t), \quad \underbrace{\frac{1}{n} \sum_{i=1}^n (y_i - w^\top x_i)^2}_{\hat{L}(w)}.$$

Non linear extensions via features/kernels.

## Statistics and computations

- ▶ Regularization by penalization separates statistics and computations.
- ▶ Implicit regularization: training time controls statistics and computations.

*What about memory?*

## Large scale learning

In many modern applications, space is the real constraint,

$$\underbrace{\widehat{X}}_{n \times d},$$

$$\underbrace{\widehat{X}^T \widehat{X}}_{d \times d},$$

$$\underbrace{\widehat{X} \widehat{X}^T \text{ or } \widehat{K}}_{n \times n}.$$

Think  $n \sim d$  large!

## Projections and dimensionality reduction

Let  $S$  be a  $d \times M$  matrix and

$$\hat{X}_M = \hat{X}S.$$

Equivalently

$$x \in \mathbb{R}^d \mapsto x_M = (s_j^\top x)_{j=1}^M \in \mathbb{R}^m,$$

with  $s_1, \dots, s_M$  columns of  $S$ .

## Learning with projected data

$$\min_{w \in \mathbb{R}^M} \frac{1}{n} \sum_{i=1}^n \ell(y_i, w^\top (x_M)_i) + \lambda \|w\|^2, \quad \lambda \geq 0.$$

We will focus on ERM based learning and least squares in particular.

## Principal component analysis (PCA)

The SVD of  $\hat{X}$  is

$$\hat{X} = U\Sigma V^T.$$

Consider  $V_M$  the matrix  $d \times M$  of the first  $M$  columns of  $V$ .

A corresponding projection is given by

$$\hat{X}_M = \hat{X}S, \quad S = V_M.$$

## Representer theorem for PCA

Note that

$$\hat{X} = U\Sigma V^T \quad \Leftrightarrow \quad \hat{X}^T = V\Sigma U^T \quad \Leftrightarrow \quad V = \hat{X}^T U\Sigma^{-1}$$

and  $V_M = \hat{X}^T U_M \Sigma_M^{-1}$ .

Then

$$\hat{X}_M = \hat{X} V_M = \underbrace{\hat{X} \hat{X}^T}_{\hat{K}} U_M \Sigma_M^{-1} = U_M \Sigma_M$$

and for any  $x$

$$x^T v_j = \sum_{i=1}^n \underbrace{x^T x_i}_{k(x, x_i)} \frac{u_j^i}{\sigma_j},$$

with  $(u_j, \sigma_j^2)_j$  eigenvectors/eigenvalues of  $\hat{K}$ .



## Kernel PCA

If  $\Phi$  is a feature map, then the SVD in feature space is

$$\hat{\Phi} = U\Sigma V^T$$

and if  $V_M$  is the matrix  $d \times M$  of the first  $M$  columns of  $V$ ,

$$\hat{\Phi}_M = \hat{\Phi} V_M.$$

Equivalently using kernels

$$\hat{\Phi}_M = \hat{K} U_M \Sigma_M^{-1} = U_M \Sigma_M,$$

and for any  $x$

$$\Phi(x)^\top v_j = \sum_{i=1}^n k(x, x_i) \frac{u_j^i}{\sigma_j}.$$

## PCA+ERM for least squares

Consider (no penalization)

$$\min_{w \in \mathbb{R}^M} \frac{1}{n} \left\| \hat{X}_M w - \hat{Y} \right\|^2.$$

The solution is<sup>1</sup>

$$\hat{w}_M = (\hat{X}_M^\top \hat{X}_M)^{-1} \hat{X}_M^\top \hat{Y}.$$

---

<sup>1</sup>Assuming invertibility for simplicity. In general replace with pseudoinverse.

## PCA+ERM for least squares

It is easy to see that that , for all  $x$

$$f_M(x) = x_M^\top \hat{w}_M = \sum_{j=1}^M \frac{1}{\sigma_j} u_j^\top \hat{Y} v_j^\top x$$

where  $x_M = V_M x$ .

Essentially due to the fact that

$$\hat{X}_M^\top \hat{X}_M = V_M^\top \hat{X}^\top \hat{X} V_M$$

is the covariance matrix projected on its first  $M$  eigenvectors.

## PCR, TSVD, Filtering

$$f_M(x) = \sum_{j=1}^M \frac{1}{\sigma_j} u_j^\top \hat{Y} v_j^\top x$$

- ▶ PCA+ERM is called Principal component regression in statistics
- ▶ ... and truncated singular value decomposition in linear algebra.
- ▶ It corresponds to the spectral filter

$$F(\sigma_j) = \begin{cases} \frac{1}{\sigma_j}, & j \leq M \\ 0, & \text{oth.} \end{cases}$$

Compare to Tikhonov and Landweber,

$$F_{\text{Tik.}}(\sigma_j) = \sigma_j / (1 + \lambda \sigma_j) \quad F_{\text{Land.}}(\sigma_j) = (1 - (1 - \gamma \sigma_j)^t) \sigma_j^{-1}.$$

## Projection and complexity

Then,

- ▶ PCA + ERM = regularization.
- ▶ In principle, down stream learning is computationally cheaper. . .

. . . however SVD requires time

$$O(nD^2 \vee d^3)$$

or with kernel matrices

$$O(n^2 C_K \vee n^3).$$

## Sketching

Let  $S$  be a  $d \times M$  matrix s.t.  $S_{ij} \sim \mathcal{N}(0, 1)$  and

$$\hat{X}_M = \hat{X}S.$$

Computing  $\hat{X}_M$  is time  $O(ndM)$  and memory  $O(nd)$ .

## Dimensionality reduction with sketching

Note that if  $x_M = S^\top x$  and  $x'_M = S^\top x'$ , then

$$\frac{1}{M} \mathbb{E}[x_M^\top x'_M] = \frac{1}{M} \mathbb{E}[x^\top S S^\top x'] = x^\top \mathbb{E}[S S^\top] x' = \frac{1}{M} x^\top \sum_{j=1}^M \underbrace{\mathbb{E}[s_j s_j^\top]}_{\text{Identity}} x' = x^\top x'.$$

- ▶ Inner products, norms distances preserved in expectation..
- ▶ ... and with high probability for given  $M$  (Johnson-Linderstrauss Lemma).

## Least squares with sketching

Consider

$$\min_{w \in \mathbb{R}^M} \frac{1}{n} \left\| \widehat{X}_M w - \widehat{Y} \right\|^2 + \lambda \|w\|^2, \quad \lambda > 0.$$

Regularization is needed. For sketching

$$\widehat{X}_M^T \widehat{X}_M = S^T \widehat{X}^T \widehat{X} S,$$

is **not** the covariance matrix projected on its first  $M$  eigenvectors, but

$$\mathbb{E}[\widehat{X}_M \widehat{X}_M^T] = \mathbb{E}[\widehat{X} S S^T \widehat{X}^T] = \widehat{X} \widehat{X}^T.$$

There is extra variability.



## Least squares with sketching (cont.)

Consider

$$\min_{w \in \mathbb{R}^M} \frac{1}{n} \left\| \widehat{X}_M w - \widehat{Y} \right\|^2 + \lambda \|w\|^2, \quad \lambda > 0.$$

The solution is

$$\widehat{w}_{\lambda, M} = (\widehat{X}_M^T \widehat{X}_M + \lambda n I)^{-1} \widehat{X}_M^T \widehat{Y}.$$

Computing  $\widehat{w}_{\lambda, M}$  is time  $O(nM^2 + ndM)$  and memory  $O(nM)$ .

## Beyond linear sketching

Let  $S$  be a  $d \times M$  random matrix and

$$\hat{X}_M = \sigma(\hat{X}S)$$

where  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  is a given nonlinearity.

Then consider functions of the form,

$$f_M(x) = x_M^\top w = \sum_{j=1}^M w^j \sigma(s_j^\top x).$$

## Learning with random weights networks

$$f_M(x) = x_M^\top w = \sum_{j=1}^M w^j \sigma(s_j^\top x).$$

Here,  $w^1, \dots, w^M$  can be computed solving a convex problem

$$\min_{w \in \mathbb{R}^M} \frac{1}{n} \sum_{i=1}^n (y_i - f_M(x_i))^2 + \lambda \|w\|^2, \quad \lambda > 0,$$

in time  $O(nM^2 + ndM)$  and memory  $O(nM)$ .

## Neural networks, random features and kernels

$$f_M(x) = \sum_{j=1}^M w^j \sigma(s_j^\top x)$$

- ▶ It is a *one hidden layer neural network with random weights*.
- ▶ It is defined by a random feature map  $\Phi_M(x) = \sigma(S^\top x)$ .
- ▶ There are a number of cases in which

$$\mathbb{E}[\Phi_M(x)^\top \Phi_M(x')] = k(x, x')$$

with  $k$  a suitable pos. def. kernel  $k$ .

## Random Fourier features

Let  $X = \mathbb{R}$ ,  $s \sim \mathcal{N}(0, 1)$  and

$$\Phi_M^j(x) = \frac{1}{\sqrt{M}} \underbrace{e^{is_j x}}_{\text{complex exp.}} .$$

For  $k(x, x') = e^{-|x-x'|^2 \gamma}$  it holds

$$\mathbb{E}[\Phi_M(x)^\top \Phi_M(x')] = k(x, x').$$

Proof: from basic properties of the Fourier transform

$$e^{-|x-x'|^2 \gamma} = \text{const.} \int ds \underbrace{e^{isx}}_{\text{Inv. transf.}} \underbrace{e^{-isx'}}_{\text{Transl.}} \underbrace{e^{-\frac{s^2}{\gamma}}}_{\text{Tranf. of Gaussian}} .$$

## Random Fourier features (cont.)

- ▶ The above reasoning immediately extends to  $X = \mathbb{R}^d$ .
- ▶ Using symmetry one can show the same result holds for

$$\Phi_M^j(x) = \frac{1}{\sqrt{M}} \cos(s_j^\top x + b_j)$$

with  $b_j$  uniformly distributed.

## Other random features

The relation

$$\mathbb{E}[\Phi_M(x)^\top \Phi_M(x')] = k(x, x').$$

is satisfied by a number of nonlinearities and corresponding kernels:

- ▶ ReLU  $\sigma(a) = |a|_+ \dots$
- ▶ Sigmoidal  $\sigma(a)$ ,  $\dots$
- ▶  $\dots$

As for all feature map the relation with kernels is not one to one.

## Infinite networks and large scale kernel methods

- ▶ One hidden layer network with infinite random weights = kernels.
- ▶ Random features are an approach to scaling kernel methods:

from

$$\text{time } O(n^2 C_k \vee n^3) \quad \text{memory } O(n^2)$$

to

$$\text{time } O(ndM \vee nM^2) \quad \text{memory } O(nM).$$



## Subsampling aka Nyström method

Through the representer theorem, the ERM solution has the form,

$$w = \sum_{i=1}^n x_i c_i = \hat{X}^\top c.$$

For  $M < n$ , choose a set of *centers*  $\{\tilde{x}_1, \dots, \tilde{x}_M\} \subset \{x_1, \dots, x_n\}$  and let

$$w_M = \sum_{i=1}^M x_i (c_M)_i = \tilde{X}_M^\top c_M.$$

## Least squares with Nyström centers

Consider

$$\min_{w_M \in \mathbb{R}^d} \frac{1}{n} \left\| \widehat{X} w_M - \widehat{Y} \right\|^2 + \lambda \|w_M\|^2, \quad \lambda > 0.$$

Equivalently

$$\min_{c \in \mathbb{R}^M} \frac{1}{n} \left\| \underbrace{\widehat{X} \widetilde{X}_M^T}_{\widehat{K}_{nM}} c_M - \widehat{Y} \right\|^2 + \lambda c_M^T \underbrace{\widetilde{X}_M \widetilde{X}_M^T}_{\widehat{K}_M} c_M, \quad \lambda > 0.$$

## Least squares with Nyström centers

$$\min_{c \in \mathbb{R}^M} \frac{1}{n} \left\| \underbrace{\widehat{X} \widetilde{X}_M^T}_{\widehat{K}_{nM}} c_M - \widehat{Y} \right\|^2 + \lambda c_M^T \underbrace{\widetilde{X}_M \widetilde{X}_M^T}_{\widehat{K}_M} c_M, \quad \lambda > 0.$$

The solutions is

$$\widehat{c}_{\lambda, M} = (\widehat{K}_{nM}^T \widehat{K}_M + n\lambda \widehat{K}_M)^{-1} \widehat{K}_{nM}^T \widehat{Y}$$

requiring

time  $O(ndM \vee nM^2)$

memory  $O(nM)$ .

## Nyström centers and sketching

Note that Nyström corresponds to sketching

$$\hat{X}_M = \hat{X}S,$$

with

$$S = \tilde{X}_M.$$

## Regularization with sketching and Nyström centers

Considering regularization as we did for sketching leads to

$$\min_{c \in \mathbb{R}^M} \frac{1}{n} \|\widehat{X} \widetilde{X}_M^\top c_M - \widehat{Y}\|^2 + \lambda c_M^\top c_M, \quad \lambda > 0.$$

In the Nyström derivation we ended up with Equivalently

$$\min_{c \in \mathbb{R}^M} \frac{1}{n} \|\widehat{X} \widetilde{X}_M^\top c_M - \widehat{Y}\|^2 + \lambda c_M^\top \widetilde{X}_M \widetilde{X}_M^\top c_M, \quad \lambda > 0.$$

Different regularizers are considered.

## Nyström approximation

A classical discrete approximation to integral equations.

For all  $x$

$$\int k(x, x')c(x')dx' = y(x) \quad \mapsto \quad \sum_{j=1}^M k(x, \tilde{x}_j)c(\tilde{x}_j) = y(\tilde{x}_j).$$

Related to to quadrature methods.

From operators to matrices.

For all  $i = 1, \dots, n$

$$\sum_{j=1}^n k(x_i, x_j)c_j = y_j \quad \mapsto \quad \sum_{j=1}^M k(x_i, \tilde{x}_j)c_j = y_j.$$

## Nyström approximation and subsampling

For all  $i = 1, \dots, n$

$$\sum_{j=1}^n k(x_i, x_j) c_j = y_i \quad \mapsto \quad \sum_{j=1}^M k(x_i, \tilde{x}_j) c_j = y_i.$$

The above formulation highlights connection to columns subsampling

$$\widehat{K}c = \widehat{Y} \quad \mapsto \quad \widehat{K}_{nM}c_M = \widehat{Y}.$$

## In summary

- ▶ Projection (dim. reductions) regularizes.
- ▶ Reducing computations by sketching.
- ▶ Nyström approximation and columns subsampling.