Asymmetric compressive learning guarantees with applications to quantized sketches



Vincent Schellekens CEA-List, Saclay



Laurent Jacques INMA, UCLouvain, Belgium

Curves and Surfaces 2022 Arcachon, France, June 20 - June 24





- (appetizer) Statistical learning intro
- (main course) and its fresh guarantees

Menu ·····•

(starter)

Symmetric Compressive Learning

Asymmetric Compressive Learning

(dessert)

Experiments

- (the bill)
- Conclusion



Minimize the <u>empirical risk</u> (ERM)

$$\widetilde{\boldsymbol{\theta}} \in \arg\min_{\boldsymbol{\theta}} \frac{1}{n} \sum_{i=1}^{n} \ell(\boldsymbol{\theta}, \boldsymbol{x}_i) \equiv \text{"Poole}$$



Sum of Squared Errors (SSE):

$$\sum_{i=1}^n \min_k \|oldsymbol{x}_i - oldsymbol{c}_k\|^2$$





oling of unhappiness"

 $\ell \equiv$ "How unhappy x_i is with the model θ "

Example: <u>GMM</u> fitting



Maximum likelihood: $\sum_{i} -\log\left(\sum_{k} w_{k} \mathcal{N}(\boldsymbol{x}_{i}; \boldsymbol{\mu}_{k}, \boldsymbol{\Gamma}_{k})\right)$

Minimize the <u>empirical risk</u> (ERM)

 $\widetilde{\boldsymbol{\theta}} \in \arg\min_{\boldsymbol{\theta}} \underbrace{\frac{1}{n} \sum_{i=1}^{n} \ell(\boldsymbol{\theta}, \boldsymbol{x}_{i})}_{\mathbb{E}_{\boldsymbol{x} \sim \widehat{\mathcal{P}}_{\boldsymbol{\mathcal{X}}}} \underbrace{\ell(\boldsymbol{\theta}, \boldsymbol{x})}_{\mathbb{E}_{\boldsymbol{x} \sim \widehat{\mathcal{P}}_{\boldsymbol{x}}}} \underbrace{\ell(\boldsymbol{\theta}, \boldsymbol{x})}_{\mathbb{E}_{\boldsymbol{x} \sim \widehat{\mathcal{P}}_{\boldsymbol{x}}}}} \underbrace{\ell(\boldsymbol{\theta}, \boldsymbol{x})}_{\mathbb{E}_{\boldsymbol{x} \sim \widehat{\mathcal{P}}_{\boldsymbol{x}}}}} \underbrace{\ell(\boldsymbol{\theta}, \boldsymbol{x})}_{\mathbb{E}_{\boldsymbol{x} \sim \widehat{\mathcal{P}}_{\boldsymbol{x}}}}} \underbrace{\ell(\boldsymbol{\theta}, \boldsymbol{x})}_{\mathbb{E}_{\boldsymbol{x} \sim \widehat{\mathcal{P}}_{\boldsymbol{x}}}} \underbrace{\ell(\boldsymbol{\theta}, \boldsymbol{x})}_{\mathbb{E}_{\boldsymbol{x} \sim \widehat{\mathcal{P}}_{\boldsymbol{x}}}}} \underbrace{\ell(\boldsymbol{\theta}, \boldsymbol{x})}_{\mathbb{E}_{\boldsymbol{x} \sim \widehat{\mathcal{P}}_{\boldsymbol{x}}}} \underbrace{\ell(\boldsymbol{\theta}, \boldsymbol{x})}_{\mathbb{E}_{\boldsymbol{x} \sim \widehat{\mathcal{P}}_{\boldsymbol{x}}}} \underbrace{\ell(\boldsymbol{\theta}, \boldsymbol{x})}} \underbrace{\ell(\boldsymbol{\theta}, \boldsymbol{x})}_{\mathbb{E}_$



Empirical distribution





Minimize the <u>empirical risk</u> (ERM)

- $\widetilde{\boldsymbol{\theta}} \in \arg\min_{\boldsymbol{\theta}} \frac{1}{n} \sum_{i=1}^{n} \ell(\boldsymbol{\theta}, \boldsymbol{x}_i) = \mathcal{R}(\boldsymbol{\theta}; \widehat{\mathcal{P}}_{\mathcal{X}})$ $\mathbb{E}_{oldsymbol{x}\sim\hat{\mathcal{P}}_{oldsymbol{\mathcal{X}}}}\ell(oldsymbol{ heta},oldsymbol{x})$
- ... as a proxy for the true risk:
 - $\boldsymbol{\theta}^* \in \arg\min_{\boldsymbol{\theta}} \mathcal{R}(\boldsymbol{\theta}; \mathcal{P}_0) = \mathbb{E}_{\boldsymbol{x} \sim \mathcal{P}_0} \ell(\boldsymbol{\theta}, \boldsymbol{x})$ if $x_i \sim_{\text{iid}} \mathcal{P}_0$ and $n \to +\infty$ $\mathcal{R}(\boldsymbol{\theta}, \widehat{\mathcal{P}}_{\mathcal{X}}) \approx \mathcal{R}(\boldsymbol{\theta}, \mathcal{P}_0)$
 - → Target statistical guarantee: Probably Approximately Correct PAC "excess risk"





 $\mathbb{P}\Big[\mathcal{R}\big(\hat{\boldsymbol{\theta}}; \mathcal{P}_0 \big) - \mathcal{R}\big(\boldsymbol{\theta}^*; \mathcal{P}_0 \big) \leq \eta \Big] \geq 1 - \delta \quad \text{(for small } \eta, \delta > 0 \text{)}$



Minimize the <u>empirical risk</u> (ERM) $\widetilde{\boldsymbol{\theta}} \in \arg\min_{\boldsymbol{\theta}} \frac{1}{n} \sum_{i=1}^{n} \ell(\boldsymbol{\theta}, \boldsymbol{x}_i) = \mathcal{R}(\boldsymbol{\theta}; \widehat{\mathcal{P}}_{\mathcal{X}})$

Problem:

 \rightarrow Require several passes over \mathcal{X} (*e.g.*, stochastic gradient descent);

 $\rightarrow n$ typically very large;



Computationally expensive on modern large-scale data











Gribonval, R., Blanchard, G., Keriven, N., & Traonmilin, Y. (2021). Gribonval, R., Chatalic, A., Keriven, N., Schellekens, V., Jacques, L., & Schniter, P. (2021)



Break the training down into two "cheaper" steps



- m coefficients
 - $d < m \ll n$
 - "Summary"

Sketching phase

Random features moments ("empirical average")

$$\boldsymbol{z} = \mathcal{A}_{\Phi}(\widehat{P}_{\mathcal{X}}) = \frac{1}{n} \sum_{i=1}^{n} \Phi(\boldsymbol{x}_i)$$

with random nonlinearity $\Phi : \mathbb{R}^d \to \mathbb{C}^m$







Sketching phase

Random features moments ("empirical average")

$$\boldsymbol{z} = \mathcal{A}_{\Phi}(\widehat{P}_{\mathcal{X}}) = \frac{1}{n} \sum_{i=1}^{n} \Phi(\boldsymbol{x}_i)$$

with random nonlinearity $\Phi : \mathbb{R}^d \to \mathbb{C}^m$

Sketching with RFF: $\{\mathcal{F}[\mathcal{P}_0](\boldsymbol{\omega}_j)\}_{j=1}^m$ $oldsymbol{z}_{\Phi_{\mathrm{RFF}}} = rac{1}{n} \sum \Phi_{\mathrm{RFF}}(oldsymbol{x}_i) \simeq [\mathbb{E}_{oldsymbol{x}\sim\mathcal{P}_0}oldsymbol{e}]$ i=1

Random sampling of the characteristic function (aka Fourier transform \mathcal{F}) of the data distribution \mathcal{P}_0



$$\rightarrow Example: random Fourier Features (RFF) \Phi_{\rm RFF}(\boldsymbol{x}) := \frac{1}{\sqrt{m}} \exp(i\boldsymbol{\Omega}^{\top}\boldsymbol{x}) \in \mathbb{C}^{m} \text{ with } \boldsymbol{\Omega} = (\boldsymbol{\omega}_{j})_{j=1}^{m}, \ \boldsymbol{\omega}_{j} \sim_{\rm i.i.d.} \Lambda(\boldsymbol{\omega}) = \underbrace{[\mathcal{F}\kappa^{\Delta}](\boldsymbol{\omega})}_{e.g., \text{ Gaussian k}} Embeds a "nonlinear geometry", e.g., \langle \Phi_{\rm RFF}(\boldsymbol{x}), \Phi_{\rm RFF}(\boldsymbol{x}') \rangle \simeq C \underbrace{\exp\left(-\frac{\|\boldsymbol{x}-\boldsymbol{x}'\|^{2}}{2\sigma^{2}}\right)}_{\kappa_{\Delta}} \\ \begin{bmatrix} \mathbb{R} \text{ Rhimi, Recket} \end{bmatrix}$$

$$e^{\mathbf{i}\boldsymbol{\omega}_{j}^{\top}\boldsymbol{x}}]_{j=1}^{m} =: \mathcal{A}_{\Phi_{\mathrm{RFF}}}(\mathcal{P}_{0})$$

"Compressive sensing of prob. distributions"













Theoretical guarantees:

<u>Assumption</u>: The sketch encodes the risk, *i.e.*,

 \mathcal{A}_{Φ} respects the Lower Restricted Isometry Property (LRIP):

 $\forall \boldsymbol{\theta}, \boldsymbol{\theta}' \in \Theta, \| \mathcal{P}_{\boldsymbol{\theta}} - \mathcal{P}_{\boldsymbol{\theta}'} \|_{\mathcal{R}} \leq \gamma \| \mathcal{A}_{\Phi}(\mathcal{P}_{\boldsymbol{\theta}}) - \mathcal{A}_{\Phi}(\mathcal{P}_{\boldsymbol{\theta}'}) \|_{2}$



(given $\|\mathcal{P} - \mathcal{P}'\|_{\mathcal{R}} := \sup_{\alpha \in \Theta} |\mathcal{R}(\alpha, \mathcal{P}) - \mathcal{R}(\alpha, \mathcal{P}')|)$

= distance between $\mathcal{P} \otimes \mathcal{P}'$ wrt task-specific risk \mathcal{R}



Theoretical guarantees:

<u>Assumption</u>: The sketch encodes the risk, *i.e.*,

 \mathcal{A}_{Φ} respects the Lower Restricted Isometry Property (LRIP):

$$\forall \boldsymbol{\theta}, \boldsymbol{\theta}' \in \Theta, \ \|\mathcal{P}_{\boldsymbol{\theta}} - \mathcal{P}_{\boldsymbol{\theta}'}\|_{\mathcal{R}} \leq \gamma \|\mathcal{A}_{\Phi}(\mathcal{P}_{\boldsymbol{\theta}}) - \mathcal{A}_{\Phi}(\mathcal{P}_{\boldsymbol{\theta}'})\|_{2}$$

e estimate $\hat{\boldsymbol{\theta}} \in \arg\min_{\boldsymbol{\theta} \in \Theta} \|\widehat{\mathcal{A}}_{\Phi}(\widehat{\mathcal{P}}_{\mathcal{X}}) - \mathcal{A}_{\Phi}(\mathcal{P}_{\boldsymbol{\theta}})\|$, (assuming we can exactly solve it)

Then, given the

$$\mathcal{R}(\widehat{\boldsymbol{\theta}}; \mathcal{P}_0) - \mathcal{R}(\boldsymbol{\theta}^{\star}; \mathcal{P}_0) \leq \inf_{\boldsymbol{\theta} \in \Theta} d(\mathcal{P}_0, \mathcal{P}_{\boldsymbol{\theta}}) + 4\gamma \|\mathcal{A}_{\Phi}(\widehat{\mathcal{P}}_{\mathcal{X}}) - \mathcal{A}_{\Phi}(\mathcal{P}_0)\|_2$$

Excess risk

Task (in Satisfying the LRIP? (w.h.p) k-m

Gribonval, R., Blanchard, G., Keriven, N., & Traonmilin, Y. (2021).



(given
$$\|\mathcal{P} - \mathcal{P}'\|_{\mathcal{R}} := \sup_{\alpha \in \Theta} |\mathcal{R}(\alpha, \mathcal{P}) - \mathcal{R}(\alpha, \mathcal{P})|$$

= distance between $\mathcal{P} \otimes \mathcal{P}'$ wrt task-specific risk \mathcal{R}

(for a certain distance d)

Modeling error

Sampling error

in \mathbb{R}^d)	Sample complexity (up to log.)
means	$m\gtrsim k^2 d$
GMM	$m \gtrsim k^2 d \text{ (or } m \gtrsim ke^d \text{)}$





Asymmetric compressive learning



Possible advantages:

- avoid complex exponential at sketching;

 - reduced transmission cost;
- still <u>differentiable</u> learning cost & relaxed conditions.



opt for hardware friendly, <u>quantized</u> sketching procedure;

Asymmetric CL: sketching phase?

$\Psi \equiv$ random periodic Features (RPF)



Nonlinear periodic "activation" f:









Dithering "smoothes" quantization/discontinuities (see later)

random projections



Random modulo feature

Asymmetric CL: learning phase?

Plug the " Ψ " sketch into the " Φ " cost function! Solve "as before", without changing anything!

$$\widehat{oldsymbol{ heta}}' \in rg\min_{oldsymbol{ heta}\in\Theta} \mathcal{C}_{\Phi}(oldsymbol{ heta};oldsymbol{z}_{\Psi}) := \|oldsymbol{z}_{\Psi} - oldsymbol{ heta}_{ heta} + oldsymbol{ heta}_{$$



(Similar philosophy in non-linear compressive sensing) Plan, Y., & Vershynin, R. (2016)

 $-\mathcal{A}_{\Phi}(\mathcal{P}_{\theta})\|$ e.g., full-precision RFF 7





Asymmetric CL: learning phase?

Plug the " Ψ " sketch into the " Φ " cost function! Solve "as before", without changing anything!

$$\widehat{\boldsymbol{\theta}}' \in \arg\min_{\boldsymbol{\theta}\in\Theta} \mathcal{C}_{\Phi}(\boldsymbol{\theta}; \boldsymbol{z}_{\Psi}) := \|\boldsymbol{z}_{\Psi} - \mathcal{A}_{\Phi}(\mathcal{P}_{\theta})\|$$

e.g., binary features \mathbb{P} e.g., full-precision RFF $\widetilde{\mathcal{N}}$

Intuitive explanation: (the dither magic trick)

$$\mathcal{C}_{\Phi}(\boldsymbol{\theta}; \boldsymbol{z}_{\Phi}) \approx \| \widehat{\boldsymbol{z}_{\Phi}} - \mathcal{A}_{\Phi}(\mathcal{P}_{\boldsymbol{\theta}}) \| + \sum_{j=1}^{m} \sum_{k \neq 1}^{k} C_{k,j} \exp(\mathrm{i}(k-1)\xi_{j}) + \sum_{j=1}^{m} \mathrm{i}(k-1)\xi_{j}) + \mathrm{i}(k-1)\xi_{j} + \mathrm{i}(k-1)\xi_{j} + \mathrm{i}(k-1)\xi_{j}) + \mathrm{i}(k-1)\xi_{j} + \mathrm{i}(k-1)\xi_{j} + \mathrm{i}(k-1)\xi_{j} + \mathrm{i}(k-1)\xi_{j}) + \mathrm{i}(k-1)\xi_{j} + \mathrm{$$

Good of CL Cost (what we want)





(Similar philosophy in non-linear compressive sensing) Plan, Y., & Vershynin, R. (2016)

Fourier Series ______ decomposition

いいい WWW + **MMMM** +

. . .

Fundamental frequency

Higher frequencies

Artefacts of quantization

 $\rightarrow \mathbf{0} \text{ in expectation on } \xi \text{ since}$ $\mathbb{E}_{\xi} \exp(\mathrm{i}(k-1)\xi) = 0 \quad \forall k \neq 1$











Asymmetric CL: theoretical guarantees?

<u>Assumption 1</u>: The (*learning*) sketch Φ encodes the risk, *i.e.*,

 \mathcal{A}_{Φ} respects the Lower Restricted Isometry Property (LRIP): $\exists \gamma > 0$ s.t. $\forall \boldsymbol{\theta}, \boldsymbol{\theta}' \in \Theta, \| \mathcal{P}_{\boldsymbol{\theta}} - \mathcal{P}_{\boldsymbol{\theta}'} \|_{\mathcal{R}} \leq \gamma \| \mathcal{A}_{\Phi}(\mathcal{P}_{\boldsymbol{\theta}}) - \mathcal{A}_{\Phi}(\mathcal{P}_{\boldsymbol{\theta}'}) \|_{2}$

 \mathcal{A}_{Ψ} and \mathcal{A}_{Φ} respect the Limited Projected Distortion (LPD): $\exists \epsilon > 0$ s.t. $\forall \mathcal{X}, \forall \boldsymbol{\theta} \in \Theta, \ |\langle \mathcal{A}_{\Psi}(\widehat{\mathcal{P}}_{\mathcal{X}}) - \mathcal{A}_{\Phi}(\widehat{\mathcal{P}}_{\mathcal{X}}), \mathcal{A}_{\Phi}(\mathcal{P}_{\boldsymbol{\theta}}) \rangle| \leq \epsilon$



- We do not generalize de LRIP (too hard), we rather complete it! (easier)

 - (as for SCL)

<u>Assumption 2</u>: The (*encoding*) sketch not too much "distorted" along specific directions, *i.e.*, (new condition)





Asymmetric CL: theoretical guarantees?

- If both LRIP (with $\gamma > 0$) & LPD (with $\epsilon > 0$) hold then,
 - given the estimate $\hat{\theta}' \in \arg\min_{\theta \in \Theta} \| \hat{\theta} \|$

 $\mathcal{R}(\widehat{\boldsymbol{\theta}}'; \mathcal{P}_0) - \mathcal{R}(\boldsymbol{\theta}^*; \mathcal{P}_0) \leq \inf_{\boldsymbol{\theta} \in \Theta} d(\mathcal{P}_0,$

Excess risk

Modeling e (for a certain distance d)



We do not generalize de LRIP (too hard), we rather complete it! (easier)

$$\mathcal{A}_{\Psi}(\widehat{\mathcal{P}}_{\mathcal{X}}) - \mathcal{A}_{\Phi}(\mathcal{P}_{m{ heta}}) \|$$
, (assuming we can exactly solve it) z_{Ψ}

(as before) (new term)

$$(\mathcal{P}_{\theta}) + 4\gamma \|\mathcal{A}_{\Phi}(\widehat{\mathcal{P}}_{\mathcal{X}}) - \mathcal{A}_{\Phi}(\mathcal{P}_{0})\|_{2} + 4\gamma \sqrt{\epsilon}$$
(new term)

18

(Sketch distortion)

Proving the LPD

What we want: \mathcal{A}_{Ψ} and \mathcal{A}_{Φ} respect the Limited Projected Distortion (LPD): $\exists \epsilon > 0$ s.t. $\forall \mathcal{X}, \forall \boldsymbol{\theta} \in \Theta, \ |\langle \mathcal{A}_{\Psi}(\widehat{\mathcal{P}}_{\mathcal{X}}) - \mathcal{A}_{\Phi}(\widehat{\mathcal{P}}_{\mathcal{X}}), \mathcal{A}_{\Phi}(\mathcal{P}_{\boldsymbol{\theta}}) \rangle| \leq \epsilon$



Proving the LPD

Assuming $\mathbb{P}[x \in \Sigma] \geq 1 - \zeta$ if $x \sim \widehat{\mathcal{P}}_{\mathcal{X}}$ or \mathcal{P}_{θ} , we first note that $\left| \langle \mathcal{A}_{\Psi}(\widehat{\mathcal{P}}_{\mathcal{X}}) - \mathcal{A}_{\Phi}(\widehat{\mathcal{P}}_{\mathcal{X}}), \mathcal{A}_{\Phi}(\mathcal{P}_{\theta}) \rangle \right| \leq \left| \int \left| \langle \Psi(\boldsymbol{x}) - \Phi(\boldsymbol{x}), \Phi(\boldsymbol{x}') \rangle \right| d\widehat{\mathcal{P}}_{\mathcal{X}}(\boldsymbol{x}) d\mathcal{P}_{\theta}(\boldsymbol{x}') \leq C \epsilon' + D \boldsymbol{\zeta},$



Data domain



What we want : \mathcal{A}_{Ψ} and \mathcal{A}_{Φ} respect the Limited Projected Distortion (LPD): $\exists \epsilon > 0$ s.t. $\forall \mathcal{X}, \forall \boldsymbol{\theta} \in \Theta, \ |\langle \mathcal{A}_{\Psi}(\widehat{\mathcal{P}}_{\mathcal{X}}) - \mathcal{A}_{\Phi}(\widehat{\mathcal{P}}_{\mathcal{X}}), \mathcal{A}_{\Phi}(\mathcal{P}_{\boldsymbol{\theta}}) \rangle| \leq \epsilon$

 $|\langle \Psi(\boldsymbol{x}) - \Phi(\boldsymbol{x}), \Phi(\boldsymbol{x}') \rangle| \leq \epsilon', \ \forall \boldsymbol{x}, \boldsymbol{x}' \in \Sigma.$ signal-LPD (sLPD) so, slPD \Rightarrow LPD

Ok, how to prove the sLPD now?

Proving the sLPD (for RPF & RFF)

Considering Random Periodic Features (RPF), with: $\Psi_f(\boldsymbol{x}) := rac{1}{F_1\sqrt{m}} f(\boldsymbol{\Omega}^{ op} \boldsymbol{x} + \boldsymbol{\xi}) \in \mathbb{C}^{ op}$

Key property of *f*: *mean Lipschitz smoothness*

$$\mathbb{E}_{t\sim\mathcal{U}[0,2\pi]} \sup_{r\in[-\delta,+\delta]} |f(t-t)| \leq C$$

Sampling condition

$$m \gtrsim \epsilon^{-2} \mathcal{H}_{c\epsilon}(\Sigma) \Rightarrow |\langle \Psi(\boldsymbol{x}) - \boldsymbol{\mu} \rangle|$$

"Dimension" of Σ ,
e.g., $k \log(k/\epsilon)$ if dim $(\Sigma) = k$

Schellekens, V., & Jacques, L. (2022a).





$$\Sigma^m, \ \Phi_{\mathrm{RFF}} = \Psi_{\exp(i\cdot)}, \ f(t) = \sum_k F_k e^{ikt}$$

(more permissive than Lipschitz smoothness) $|+r) - f(t)| \leq L_f^{\mu} \cdot \delta$ (e.g., for $q = \Box \Box$, $L_q^{\mu} = \frac{8}{\pi}$)

nooth on average"

signal-LPD (sLPD)

 $-\Phi(\boldsymbol{x}), \Phi(\boldsymbol{x}')\rangle \leq \epsilon, \ \forall \boldsymbol{x}, \boldsymbol{x}' \in \Sigma, \text{ w.h.p.}$

Asymmetric CL: theoretical guarantees with RPF

Let us consider <u>k-means</u> or <u>GMM fitting</u>, and the estimate

$$\widehat{\boldsymbol{ heta}}' \in rgmin_{\boldsymbol{ heta}\in\Theta} \mathcal{C}_{\Phi}(\boldsymbol{ heta};$$

<u>Assumptions:</u> + bounded data domain (in ℓ_{∞}) + mode variance $\leq S$ (for GMM) + LRIP satisfied for the task (with $\gamma > 0$)

If, for k-means, $m \gtrsim \epsilon^{-2} d \log(\frac{c\sqrt{dr}}{\epsilon})$, or, for

<u>Then</u>, with proba $\geq 1 - C \exp(-cm\epsilon^2)$,

SCL part (as before)

$$\mathcal{R}(\widehat{\boldsymbol{\theta}}';\mathcal{P}_0) - \mathcal{R}(\boldsymbol{\theta}^*;\mathcal{P}_0) \leq \inf_{\boldsymbol{\theta}\in\Theta} d(\mathcal{P}_0,\mathcal{P}_{\boldsymbol{\theta}}) + 4\gamma \|\mathcal{A}_{\Phi}(\widehat{\mathcal{P}}_{\mathcal{X}}) - \mathcal{A}_{\Phi}(\mathcal{P}_0)\|_2 +$$

Excess risk

Modeling error (for a certain distance d)



Schellekens, V., & Jacques, L. (2022b).

- $oldsymbol{z}_{oldsymbol{\Psi}}) := \|oldsymbol{z}_{oldsymbol{\Psi}} \mathcal{A}_{\Phi}(\mathcal{P}_{oldsymbol{ heta}})\|$







$$\frac{\text{GMM fitting,}}{\epsilon} \quad m \gtrsim \epsilon^{-2} d \log(\frac{c \sqrt{d(2\rho S + r)}}{\epsilon})$$
 (for some $\rho > d$)

Distortion error (Sketch distortion)

Sampling error

Smoother f (in mean Lipschitz & Fourier) \Rightarrow better performances







Some numerical results* (synthetic data)

How large should the sketch be to reach desired learning performance?



<u>K-means</u>

d = 5 dims, K = 10 clusters, $n = 10^5$ samples.

*: Real data, audio classification task in extra slides

<u>GMM fitting</u>

d = 5 dims, K = 10 modes, $n = 10^5$ samples.



Some numerical results* (synthetic data)

How large should the sketch be to reach desired learning performance?



*: Real data, audio classification task in extra slides



Conclusion

Take away messages:

- We can make CL asymmetric & keep theoretical guarantees (with LPD) Relaxed conditions on sketching vs learning Hardware friendly sketching (e.g., quantized)

- New sketching are possible (such as modulo)

Limitation & open questions:

- Hidden (large) constants for LRIP and LPD (\equiv the "compressive sensing crime") Sample complexity for LPD in function of model set dimension (and not Σ)? Gap between the theoretical guarantees and empirical performances (SCL & ACL)



Thank you for your attention

Related references:

- scale learning: Keeping only what you need." IEEE Signal Processing Magazine, 38(5), 12-36.
- processing systems, 20.
- clustering and compressive mixture modeling." Mathematical Statistics and Learning, 3(2), 165-257.
- theory, 62(3), 1528-1537.
- conference on Multimedia, pages 1015–1018, 2015.
- sketches." IEEE Transactions on Signal Processing, 70, 1348-1360.
- machines." Information and Inference: A Journal of the IMA, 11(1), 385-421.

Gribonval, R., Chatalic, A., Keriven, N., Schellekens, V., Jacques, L., & Schniter, P. (2021). "Sketching data sets for large-

Rahimi, A., & Recht, B. (2007). "Random features for large-scale kernel machines." Advances in neural information

Gribonval, R., Blanchard, G., Keriven, N., & Traonmilin, Y. (2021). "Statistical learning guarantees for compressive

Plan, Y., & Vershynin, R. (2016). "The generalized lasso with non-linear observations." IEEE Transactions on information

Karol J Piczak. "ESC: Dataset for environmental sound classification." In Proceedings of the 23rd ACM international

Anurag Kumar and Bhiksha Raj. "Features and kernels for audio event recognition." arXiv preprint arXiv:1607.05765, 2016. Schellekens, V., & Jacques, L. (2022b). "Asymmetric compressive learning guarantees with applications to quantized

Schellekens, V., & Jacques, L. (2022a). "Breaking the waves: asymmetric random periodic features for low-bitrate kernel







Extra slides

Real data, audio classification task

ESC-50 dataset: 🖬 Karol J Piczak, 2015.

- J = 2500 audio clips lasting 5s; C = 50 classes (e.g., animals, water sounds, urban noises)

Scenario:

- We compute $z_f = \frac{1}{n} \sum_{i=1}^n \Psi_f(x_i)$

for SCL (with $f(\cdot) = \exp(i \cdot)$) and ACL (with f = q)

- Minimizing transmission cost
- Compare (64bits) SCL, and (1bit) ACL
 - \rightarrow Objective \neq competitive classification

Method:

- GMM-fit with K = 32 modes extracted from SCL & ACL.
- Classification using "alpha features" (on a test set)

A. Kumar, B. Raj, 2016.

 $\mathscr{X} = \{x_i\}_{i=1}^n \subset \mathbb{R}^{d=10} \equiv \text{Mel Frequency Cepstral Coefficients (MFCC) of 15ms } (n = 334J = 835\,000)$

Schellekens, V., & Jacques, L. (2022b).



Real data, audio classification task

ESC-50 dataset: 🖬 Karol J Piczak, 2015.

- J = 2500 audio clips lasting 5s; C = 50 classes (e.g., animals, water sounds, urban noises)

Scenario:

- We compute $z_f = \frac{1}{n} \sum_{i=1}^n \Psi_f(x_i)$

for SCL (with $f(\cdot) = \exp(i \cdot)$) and ACL (with f = q)

- Minimizing transmission cost
- Compare (64bits) SCL, and (1bit) ACL
 - \rightarrow Objective \neq competitive classification

Method:

- GMM-fit with K = 32 modes extracted from SCL & ACL.
- Classification using "alpha features" (on a test set)

A. Kumar, B. Raj, 2016.

 $\mathscr{X} = \{x_i\}_{i=1}^n \subset \mathbb{R}^{d=10} \equiv \text{Mel Frequency Cepstral Coefficients (MFCC) of 15ms } (n = 334J = 835\,000)$







Basic proof idea

First, show that the LPD and a bit of algebra imply

$$\|\mathcal{A}_{\Phi}(\widehat{\mathcal{P}}_{\mathcal{X}}) - \mathcal{A}_{\Phi}(\mathcal{P}_{\widehat{\theta}'})\| \leq \|\mathcal{A}_{\Phi}(\mathcal{P}_{\widehat{\theta}'})\| \leq \|\mathcal{A}_{\Phi}(\mathcal{P}_{\widehat{\theta}'})\| \leq \|\mathcal{A}_{\Phi}(\mathcal{P}_{\widehat{\theta}'})\| \leq \|\mathcal{P}_{\Phi}(\mathcal{P}_{\widehat{\theta}'})\| \leq \|\mathcal{P}_{\Phi}(\mathcal{P}_{\widehat{\theta$$

Asymmetric CL solution

$$\widehat{\boldsymbol{\theta}}' \in \arg\min_{\boldsymbol{\theta}\in\Theta} \|\mathcal{A}_{\Psi}(\widehat{\mathcal{P}}_{\mathcal{X}}) - \mathcal{A}_{\Phi}(\mathcal{P}_{\boldsymbol{\theta}})\|$$

Then, plug this into the usual LRIP \Rightarrow performance guarantee proof







Some numerical results. (synthetic data)

Influence of the dataset size n?









Theoretical guarantees:

<u>Assumption</u>: The sketch encodes the risk, *i.e.*,

 \mathcal{A}_{Φ} respects the Lower Restricted Isometry Property (LRIP):

$$\forall \boldsymbol{\theta}, \boldsymbol{\theta}' \in \Theta, \ \|\mathcal{P}_{\boldsymbol{\theta}} - \mathcal{P}_{\boldsymbol{\theta}'}\|_{\mathcal{R}} \leq \gamma \|\mathcal{A}_{\Phi}(\mathcal{P}_{\boldsymbol{\theta}}) - \mathcal{A}_{\Phi}(\mathcal{P}_{\boldsymbol{\theta}'})\|_{2}$$

e estimate $\hat{\boldsymbol{\theta}} \in \arg\min_{\boldsymbol{\theta} \in \Theta} \|\widehat{\mathcal{A}}_{\Phi}(\widehat{\mathcal{P}}_{\mathcal{X}}) - \mathcal{A}_{\Phi}(\mathcal{P}_{\boldsymbol{\theta}})\|$, (assuming we can exactly solve it)

Then, given the

$$\mathcal{R}(\widehat{\boldsymbol{\theta}}; \mathcal{P}_0) - \mathcal{R}(\boldsymbol{\theta}^{\star}; \mathcal{P}_0) \leq \inf_{\boldsymbol{\theta} \in \Theta} d(\mathcal{P}_0, \mathcal{P}_{\boldsymbol{\theta}}) + 4\gamma \|\mathcal{A}_{\Phi}(\widehat{\mathcal{P}}_{\mathcal{X}}) - \mathcal{A}_{\Phi}(\mathcal{P}_0)\|_2$$

Excess risk

Task (in \mathbb{R}^d) Satisfying the LRIP? (w.h.p) *k*-means GMM



(given
$$\|\mathcal{P} - \mathcal{P}'\|_{\mathcal{R}} := \sup_{\alpha \in \Theta} |\mathcal{R}(\alpha, \mathcal{P}) - \mathcal{R}(\alpha, \mathcal{P})|$$

= distance between $\mathcal{P} \otimes \mathcal{P}'$ wrt task-specific risk \mathcal{R}

(for a certain distance d)

Modeling error

Sampling error

Sample complexity	Condition
$m\gtrsim k^2d\lograc{kdR}{arepsilon}$	$\min_{j \neq k} \ \boldsymbol{c}_j - \boldsymbol{c}_k \ \geq \varepsilon$
$m \gtrsim k^2 d \ (\text{or} \ m \gtrsim k e^d)$	sufficiently far modes







