

LEARNING DISCRETE MEASUREMENT MATRICES AND ALGORITHM UNFOLDING FOR SIGNAL RECOVERY

SIYAKHULA - GERMAN RESEARCH DAYS AT AIMS' 20-YEARS
CELEBRATION, 2024

Peter Jung, Technical University Berlin

peter.jung@tu-berlin.de

March 22, 2024

Inverse problems: given observation \mathbf{y} and partially **known** forward operator A :

$$\text{find } \mathbf{x} \in \mathcal{X} \text{ s.t. } A(\mathbf{x}) \approx \mathbf{y}$$

enforce structure for \mathbf{x} via priors/regularization/models...

☺ Sparsity, low-rankness, atomic sets - **well-established theory**

Inverse problems: given observation \mathbf{y} and partially **known** forward operator A :

$$\text{find } \mathbf{x} \in \mathcal{X} \text{ s.t. } A(\mathbf{x}) \approx \mathbf{y}$$

enforce structure for \mathbf{x} via priors/regularization/models...

- 😊 Sparsity, low-rankness, atomic sets - **well-established theory**
- 😞 often **too generic** for real-world data
- 😞 algorithms Q are too **slow** and need handcrafted **tuning** for many real-world/time applications

Inverse problems: given observation \mathbf{y} and partially **known** forward operator \mathbf{A} :

$$\text{find } \mathbf{x} \in \mathcal{X} \quad \text{s.t.} \quad \mathbf{A}(\mathbf{x}) \approx \mathbf{y}$$

enforce structure for \mathbf{x} via priors/regularization/models...

- 😊 Sparsity, low-rankness, atomic sets - **well-established theory**
- 😞 often **too generic** for real-world data
- 😞 algorithms \mathbf{Q} are too **slow** and need handcrafted **tuning** for many real-world/time applications
- 😞 optimize \mathbf{A} under **practical constraints** is often difficult

Inverse problems: given observation \mathbf{y} and partially **known** forward operator \mathbf{A} :

$$\text{find } \mathbf{x} \in \mathcal{X} \quad \text{s.t.} \quad \mathbf{A}(\mathbf{x}) \approx \mathbf{y}$$

enforce structure for \mathbf{x} via priors/regularization/models...

- ☺ Sparsity, low-rankness, atomic sets - **well-established theory**
- ☹ often **too generic** for real-world data
- ☹ algorithms \mathbf{Q} are too **slow** and need handcrafted **tuning** for many real-world/time applications
- ☹ optimize \mathbf{A} under **practical constraints** is often difficult

“**Learn to sense and to recover**”

$$\text{find } (\mathbf{A}, \mathbf{Q}) \quad \text{s.t.} \quad \mathbb{E}_{\mathbf{x} \sim \mathcal{X}} \text{loss}(\mathbf{x}, \mathbf{Q}(\mathbf{A}(\mathbf{x}))) \quad \text{”is small”}$$

- ① Problem Statement
- ② Gradient-based Learning of Discrete Measurement Matrices
- ③ Application 1: Learning Masks for Single-Pixel Imaging
- ④ Application 2: Learning Pooling Matrices for Group Testing

Gradient-Based Learning of Discrete Structured Measurement Operators for Signal Recovery

Jonathan Sauder¹, Martin Genzel¹, and Peter Jung¹

Abstract—Countless signal processing applications include the reconstruction of signals from few indirect linear measurements. The design of effective measurement operators is typically constrained by the underlying hardware and physics, posing a challenging and often even discrete optimization task. While the potential of gradient-based learning via the unrolling of iterative recovery algorithms has been demonstrated, it has remained unclear how to leverage this technique when the set of admissible measurement operators is structured and discrete. We tackle this problem by combining unrolled optimization with Gumbel reparametrizations, which enable the computation of

discrete subset—to improve the performance of downstream tasks poses great computational challenges. While it is often easy to create a suitable random mask, it is not obvious how to optimize the measurement matrix in a way that is both efficient and respects the feasibility constraints. Classical approaches commonly use discrete optimization to find such sets, as no gradients can be directly computed.

On the other hand, gradient-based optimization via back-propagation through massive nonlinear computational graphs has shown impressive results in the field of machine learning

PROBLEM STATEMENT

- Recover $\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^n$ from (noisy) linear measurements $\mathbf{y} = \mathbf{A}(\mathbf{x}) = \Phi\mathbf{x} + \mathbf{e} \in \mathbb{R}^m$

- Recover $\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^n$ from (noisy) linear measurements $\mathbf{y} = \mathbf{A}(\mathbf{x}) = \Phi\mathbf{x} + \mathbf{e} \in \mathbb{R}^m$
- minimal acquisition time&storage $m \simeq \text{DoF}(\mathcal{X}) \ll n \Rightarrow$ “**sensing & compression**”

- Recover $\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^n$ from (noisy) linear measurements $\mathbf{y} = \mathbf{A}(\mathbf{x}) = \Phi \mathbf{x} + \mathbf{e} \in \mathbb{R}^m$
- minimal acquisition time&storage $m \simeq \text{DoF}(\mathcal{X}) \ll n \Rightarrow$ “**sensing & compression**”
- Φ lies in **finite set of admissible matrices** Φ (determined by physics / detectors)
- \mathbf{x} is recovered from \mathbf{y} using parametrized recovery algorithm $f_{\theta, \Phi} : \mathbb{R}^m \rightarrow \mathbb{R}^n$
- $\mathbf{x} \rightarrow f_{\theta, \Phi}(\Phi \mathbf{x})$ can be understood as autoencoder

- Recover $\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^n$ from (noisy) linear measurements $\mathbf{y} = \mathbf{A}(\mathbf{x}) = \Phi \mathbf{x} + \mathbf{e} \in \mathbb{R}^m$
- minimal acquisition time&storage $m \simeq \text{DoF}(\mathcal{X}) \ll n \Rightarrow$ “**sensing & compression**”
- Φ lies in **finite set of admissible matrices** Φ (determined by physics / detectors)
- \mathbf{x} is recovered from \mathbf{y} using parametrized recovery algorithm $f_{\theta, \Phi} : \mathbb{R}^m \rightarrow \mathbb{R}^n$
- $\mathbf{x} \rightarrow f_{\theta, \Phi}(\Phi \mathbf{x})$ can be understood as autoencoder
- find good and **admissible** $\Phi \in \Phi$ and algorithm parameters θ

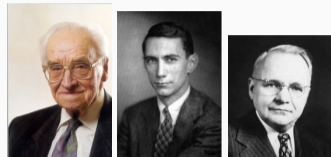
$$\min_{\Phi \in \Phi, \theta} \mathbb{E}_{\mathbf{x} \sim \mathcal{X}} \left[\text{loss}(f_{\theta, \Phi}(\Phi \mathbf{x}), \mathbf{x}) \right]$$

data-driven via “training”

prototypical playground is “**compressed sensing**” \Leftrightarrow recovering sparse \mathbf{x} from $\mathbf{y} = \Phi \mathbf{x} + \mathbf{e}$

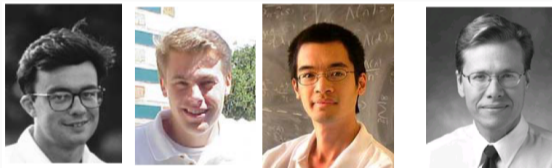
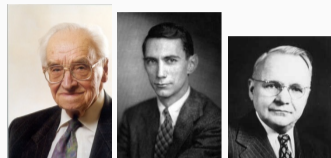
classical sampling/scanning “...if you sample densely enough, you can perfectly reconstruct the original analog data...”

- ☺ good, if signals “fill up” a subspace (linear structure)
- ☹ wasteful for compressible signals (non-linear structure)



classical sampling/scanning "...if you sample densely enough, you can perfectly reconstruct the original analog data..."

- ☺ good, if signals "fill up" a subspace (linear structure)
- ☹ wasteful for compressible signals (non-linear structure)

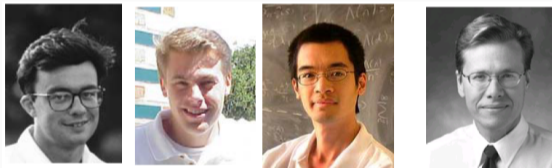
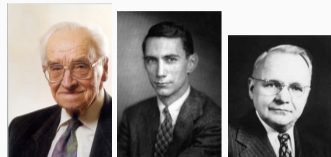


"... Can we not just directly measure the part that will not end up being thrown away?" [Donoho, 2006]

$$y = \Phi x$$

classical sampling/scanning “...if you sample densely enough, you can perfectly reconstruct the original analog data...”

- ☺ good, if signals “fill up” a subspace (linear structure)
- ☹ wasteful for compressible signals (non-linear structure)



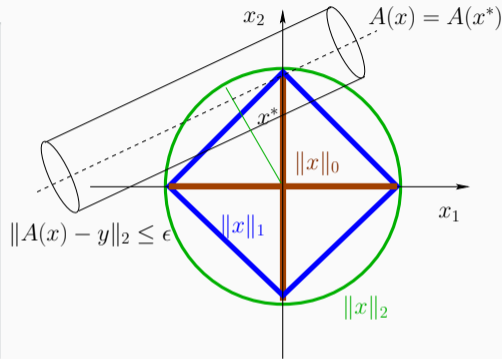
“... Can we not just directly measure the part that will not end up being thrown away?” [Donoho, 2006]

- ☺ Well-investigated, theoretical bounds, algorithms

$$y = \Phi x$$

Convex recovery approach with guarantees

$$\mathbf{x}^\sharp = \arg \min \|\mathbf{z}\|_1 \quad \text{s.t.} \quad \|\Phi \mathbf{z} - \mathbf{y}\| \leq \epsilon$$



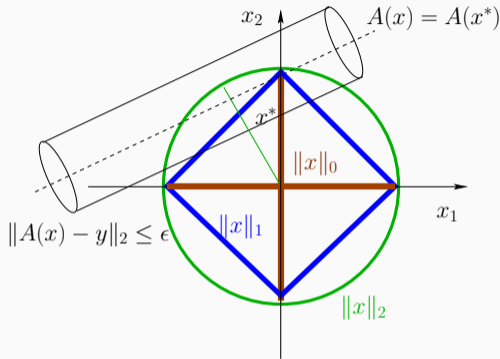
Convex recovery approach with guarantees

$$\mathbf{x}^\sharp = \arg \min \|z\|_1 \quad \text{s.t.} \quad \|\Phi z - \mathbf{y}\| \leq \epsilon$$

if sparse vectors are “well-separated” from nullspace of Φ
 (NSP \Leftarrow RIP \Leftarrow coherence):

$$\|\mathbf{x}^\sharp - \mathbf{x}\| \lesssim \epsilon$$

holds for all $\mathbf{y} = \Phi \mathbf{x} + \mathbf{e}$ with $\|\mathbf{e}\| \leq \epsilon$.



- depends on **hyper-parameter** ϵ (in this case the noise-level/SNR)
- ⊗ Problematic, if noise depends on \mathbf{x} (Poisson/multiplicative noise, covariance matching...)

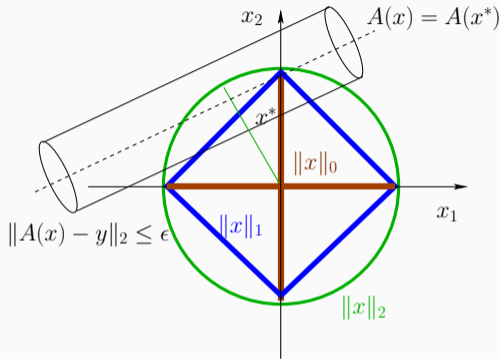
Convex recovery approach with guarantees

$$\mathbf{x}^\# = \arg \min \|\mathbf{z}\|_1 \quad \text{s.t.} \quad \|\Phi \mathbf{z} - \mathbf{y}\| \leq \epsilon$$

if sparse vectors are “well-separated” from nullspace of Φ
 (NSP \Leftarrow RIP \Leftarrow coherence):

$$\|\mathbf{x}^\# - \mathbf{x}\| \lesssim \epsilon$$

holds for all $\mathbf{y} = \Phi \mathbf{x} + \mathbf{e}$ with $\|\mathbf{e}\| \leq \epsilon$.



- depends on hyper-parameter ϵ (in this case the noise-level/SNR)
- ☹ Problematic, if noise depends on \mathbf{x} (Poisson/multiplicative noise, covariance matching...)
- ☺ Non-negativity & “biased” measurements are helpful, example later

task 1: find $\Phi \in \mathbb{R}^{m \times n}$ with minimal m for given sparsity s under practical constraints

task 1: find $\Phi \in \mathbb{R}^{m \times n}$ with minimal m for given sparsity s under practical constraints

- deterministic constructions known for $m = O(s^2)$ (based on [coherence](#))
- “sufficiently random” whp for $m = O(s \text{ polylog}(n))$ ([RIP/NSP](#) etc.) [[Candes & Tao, 2005](#)]
- ...many works on other random models ...

task 1: find $\Phi \in \mathbb{R}^{m \times n}$ with minimal m for given sparsity s under practical constraints

- deterministic constructions known for $m = O(s^2)$ (based on [coherence](#))
 - “sufficiently random” whp for $m = O(s \text{ polylog}(n))$ ([RIP/NSP](#) etc.) [[Candes & Tao, 2005](#)]
 - ...many works on other random models ...
 - random binary matrices [[Kueng & PJ, 2018](#)], designs and orthogonal arrays [[PJ, Kueng, Mixton, 2019](#)]
- ⇒ “structured derandomization”, find “concentrated measure” p on Φ

task 1: find $\Phi \in \mathbb{R}^{m \times n}$ with minimal m for given sparsity s under practical constraints

- deterministic constructions known for $m = O(s^2)$ (based on [coherence](#))
 - “sufficiently random” whp for $m = O(s \text{ polylog}(n))$ ([RIP/NSP](#) etc.) [[Candes & Tao, 2005](#)]
 - ...many works on other random models ...
 - random binary matrices [[Kueng & PJ, 2018](#)], designs and orthogonal arrays [[PJ, Kueng, Mixton, 2019](#)]
- ⇒ “structured derandomization”, find “concentrated measure” p on Φ

task 2: develop fast algorithms and hyper-parameter tuning

- most of algorithms are too complex and not suited for real-time applications

task 1: find $\Phi \in \mathbb{R}^{m \times n}$ with minimal m for given sparsity s under practical constraints

- deterministic constructions known for $m = O(s^2)$ (based on [coherence](#))
 - “sufficiently random” whp for $m = O(s \text{ polylog}(n))$ ([RIP/NSP](#) etc.) [[Candes & Tao, 2005](#)]
 - ...many works on other random models ...
 - random binary matrices [[Kueng & P, 2018](#)], designs and orthogonal arrays [[P, Kueng, Mixton, 2019](#)]
- ⇒ “structured derandomization”, find “concentrated measure” p on Φ

task 2: develop fast algorithms and hyper-parameter tuning

- most of algorithms are too complex and not suited for real-time applications

in “**machine learning**” terminology:

$$\min_{p, \theta} \mathbb{E}_{x, \Phi \sim p} [\mathbf{loss}(f_{\theta, \Phi}(\Phi x), x)]$$

two techniques: **Gumbel reparametrizations** and **algorithm unfolding**

GRADIENT-BASED LEARNING OF DISCRETE MEASUREMENT MATRICES

- learning **unstructured** $\Phi \Rightarrow$ use ∇_{Φ} [Adler et al 2016, Wu et al. 2019]

- learning **unstructured** $\Phi \Rightarrow$ use ∇_{Φ} [Adler et al 2016, Wu et al. 2019]
- here Φ is **discrete**, i.e., optimize instead **a discrete distribution p on Φ**

$$\min_{p, \theta} \mathbb{E}_{\mathbf{x}, \Phi \sim p} [\mathbf{loss}(f_{\theta, \Phi}(\Phi \mathbf{x}), \mathbf{x})]$$

aligns with derandomization principle

- learning **unstructured** $\Phi \Rightarrow$ use ∇_{Φ} [Adler et al 2016, Wu et al. 2019]
- here Φ is **discrete**, i.e., optimize instead **a discrete distribution p on Φ**

$$\min_{p, \theta} \mathbb{E}_{\mathbf{x}, \Phi \sim p} [\text{loss}(f_{\theta, \Phi}(\Phi \mathbf{x}), \mathbf{x})]$$

aligns with derandomization principle

- ① Gradient descent on p but estimate $\nabla_p \mathbb{E}_{\Phi \sim p} f(\Phi)$ via **sampling**
- ② perform sampling in a way that allows backprop - **Gumbel reparametrization**

- Consider computational graph with a *random node* v taking the values in $[a] := \{1, \dots, a\}$.
- inputs to node v are unnormalized log-probabilities $\varphi = (\varphi_1 \dots, \varphi_a) \in \mathbb{R}^a$

$$\mathbb{P}(v = i) = \frac{\exp(\varphi_i)}{\sum_{j=1}^a \exp(\varphi_j)} = \text{softmax}(\varphi)_i, \quad i = 1, \dots, a$$

- realization of v is then passed through a differentiable function f .

Estimate $\nabla_{\varphi} \mathbb{E}_v [f(v)]$ by sampling

- Consider computational graph with a *random node* v taking the values in $[a] := \{1, \dots, a\}$.
- inputs to node v are unnormalized log-probabilities $\varphi = (\varphi_1 \dots, \varphi_a) \in \mathbb{R}^a$

$$\mathbb{P}(v = i) = \frac{\exp(\varphi_i)}{\sum_{j=1}^a \exp(\varphi_j)} = \text{softmax}(\varphi)_i, \quad i = 1, \dots, a$$

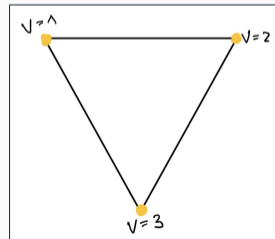
- realization of v is then passed through a differentiable function f .

Estimate $\nabla_{\varphi} \mathbb{E}_v [f(v)]$ by sampling

- efficiently sample from v by using [\[Gumbel 1954\]](#)

$$v = \text{one_hot}(\arg \max_i [\varphi_i + g_i])$$

where $g_i = -\log(-\log(u_i))$ with $u_1, \dots, u_a \sim_{\text{iid}} \text{unif}([0, 1])$



...so, we can sample v via

$$v = \text{one_hot}(\arg \max_i [\varphi_i + g_i]) \quad g_i \sim G(0, 1)$$

☹ But backprop doesn't work because of $\arg \max$.

...so, we can sample v via

$$v = \text{one_hot}(\arg \max_i [\varphi_i + g_i]) \quad g_i \sim G(0, 1)$$

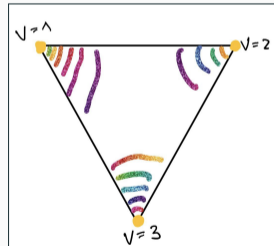
☹️ But backprop doesn't work because of $\arg \max$.

😊 **Gumbel-softmax** [Jang et al. 2016; Maddison et al. 2016]

replace $\arg \max$ in backward pass with softmax with temperature τ :

$$v_i = \frac{\exp((g_i + \varphi_i)/\tau)}{\sum_{j=1}^a \exp((g_j + \varphi_j)/\tau)} = \text{softmax}((\varphi + \mathbf{g})/\tau)_i$$

- As $\tau \rightarrow 0$, softmax approaches $\arg \max$



- from reservoir sampling: Gumbel softmax extends naturally to Gumbel top-k! [\[Vieira 2014\]](#)
- ☺ learnable component for "selecting k out of a set"

The Method (for structured binary $m \times n$ -matrix)

- from reservoir sampling: Gumbel softmax extends naturally to Gumbel top-k! [\[Vieira 2014\]](#)
- ☺ learnable component for "selecting k out of a set"

The Method (for structured binary $m \times n$ -matrix)

- ① let $\{I_i\}_{i \in [l]}$ a partition of $\mathcal{I} := [m] \times [n]$

- from reservoir sampling: Gumbel softmax extends naturally to Gumbel top-k! [\[Vieira 2014\]](#)
- ☺ learnable component for "selecting k out of a set"

The Method (for structured binary $m \times n$ -matrix)

- ① let $\{I_i\}_{i \in [l]}$ a partition of $\mathcal{I} := [m] \times [n]$
- ② for each $i \in [l]$ use Gumbel top-k to select k_i elements from I_i .

- from reservoir sampling: Gumbel softmax extends naturally to Gumbel top-k! [\[Vieira 2014\]](#)
- ☺ learnable component for "selecting k out of a set"

The Method (for structured binary $m \times n$ -matrix)

- 1 let $\{I_i\}_{i \in [l]}$ a partition of $\mathcal{I} := [m] \times [n]$
- 2 for each $i \in [l]$ use Gumbel top-k to select k_i elements from I_i .
- 3 forward pass: gives a random $m \times n$ binary matrix Φ via top- k_i and ...

- from reservoir sampling: Gumbel softmax extends naturally to Gumbel top-k! [\[Vieira 2014\]](#)
- ☺ learnable component for "selecting k out of a set"

The Method (for structured binary $m \times n$ -matrix)

- 1 let $\{I_i\}_{i \in [l]}$ a partition of $\mathcal{I} := [m] \times [n]$
- 2 for each $i \in [l]$ use Gumbel top-k to select k_i elements from I_i .
- 3 forward pass: gives a random $m \times n$ binary matrix Φ via top- k_i and ...
- 4 backward pass: well-defined gradient prop. via softmax instead of hard top- k_i [\[Bengio et.al, 2013\]](#)
- 5 construction can be freely used in conjunction with automatic differentiation

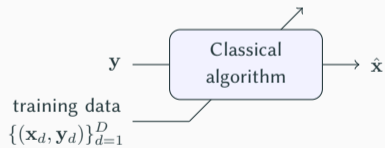
- from reservoir sampling: Gumbel softmax extends naturally to Gumbel top-k! [\[Vieira 2014\]](#)
- ☺ learnable component for "selecting k out of a set"

The Method (for structured binary $m \times n$ -matrix)

- ① let $\{I_i\}_{i \in [l]}$ a partition of $\mathcal{I} := [m] \times [n]$
- ② for each $i \in [l]$ use Gumbel top-k to select k_i elements from I_i .
- ③ forward pass: gives a random $m \times n$ binary matrix Φ via top- k_i and ...
- ④ backward pass: well-defined gradient prop. via softmax instead of hard top- k_i [\[Bengio et.al, 2013\]](#)
- ⑤ construction can be freely used in conjunction with automatic differentiation

ok, now we also need a "trainable" recovery algorithm ...

ALGORITHM UNFOLDING



Computational graph of many iterative algorithms can be viewed as neural networks...

- consider for example:

$$\min_x \|\mathbf{y} - \Phi \mathbf{x}\|_2^2 + \lambda g(\mathbf{x})$$

and use proximal (Landweber) iterations:

$$\hat{\mathbf{x}}^{(t+1)} = \text{prox}_{\lambda g} \left(\hat{\mathbf{x}}^{(t)} + \gamma \nabla (\|\mathbf{y} - \Phi \hat{\mathbf{x}}^{(t)}\|_2^2) \right),$$

Computational graph of many iterative algorithms can be viewed as neural networks...

- consider for example:

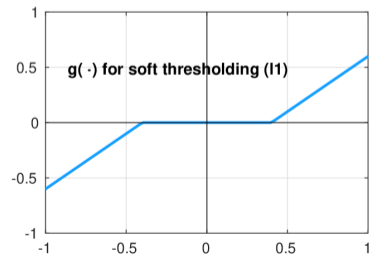
$$\min_{\mathbf{x}} \|\mathbf{y} - \Phi \mathbf{x}\|_2^2 + \lambda g(\mathbf{x})$$

and use proximal (Landweber) iterations:

$$\hat{\mathbf{x}}^{(t+1)} = \text{prox}_{\lambda g} \left(\hat{\mathbf{x}}^{(t)} + \gamma \nabla (\|\mathbf{y} - \Phi \hat{\mathbf{x}}^{(t)}\|_2^2) \right),$$

- $g(\mathbf{x}) = \|\mathbf{x}\|_1 \Rightarrow$ element-wise soft-thresholding, known as *Iterative Shrinkage and Thresholding Algorithm* [Daubechies et al., 2004]

$$\text{(ISTA)} \quad \hat{\mathbf{x}}^{(t+1)} = \eta_{\frac{\lambda}{L}} \left(\hat{\mathbf{x}}^{(t)} + \gamma \Phi^* (\Phi \hat{\mathbf{x}}^{(t)} - \mathbf{y}) \right)$$



Unfolding iterations into network $f_{\theta, \Phi}$, several proposals [Gregor&LeCun 2010, Liu 2019, Chen 2018,...]:

$$\text{(LISTA-CP)} \quad \mathbf{x}^{(k+1)} = \eta_{\alpha^{(k)}} \left(\mathbf{x}^{(k)} - \mathbf{B}^{(k)*} \left(\Phi \mathbf{x}^{(k)} - \mathbf{y} \right) \right) \quad \{\alpha^{(k)}, \mathbf{B}^{(k)}\}$$

Unfolding iterations into network $f_{\theta, \Phi}$, several proposals [Gregor&LeCun 2010, Liu 2019, Chen 2018,...]:

$$\text{(LISTA-CP)} \quad \mathbf{x}^{(k+1)} = \eta_{\alpha^{(k)}} \left(\mathbf{x}^{(k)} - \mathbf{B}^{(k)*} \left(\Phi \mathbf{x}^{(k)} - \mathbf{y} \right) \right) \quad \{\alpha^{(k)}, \mathbf{B}^{(k)}\}$$

$$\text{(ALISTA)} \quad \mathbf{x}^{(k+1)} = \eta_{\alpha^{(k)}} \left(\mathbf{x}^{(k)} - \gamma^{(k)} \mathbf{B}^* \left(\Phi \mathbf{x}^{(k)} - \mathbf{y} \right) \right) \quad \{\alpha^{(k)}, \gamma^{(k)}\}$$

- supervised training on \mathcal{X} with $\mathbb{E}_{\mathbf{x} \sim \mathcal{X}, e} \mathbf{loss}(f_{\theta, \Phi}(\Phi \mathbf{x} + e), \mathbf{x})$, unsupervised is possible as well

Unfolding iterations into network $f_{\theta, \Phi}$, several proposals [Gregor&LeCun 2010, Liu 2019, Chen 2018,...]:

$$\text{(LISTA-CP)} \quad \mathbf{x}^{(k+1)} = \eta_{\alpha^{(k)}} \left(\mathbf{x}^{(k)} - \mathbf{B}^{(k)*} \left(\Phi \mathbf{x}^{(k)} - \mathbf{y} \right) \right) \quad \{\alpha^{(k)}, \mathbf{B}^{(k)}\}$$

$$\text{(ALISTA)} \quad \mathbf{x}^{(k+1)} = \eta_{\alpha^{(k)}} \left(\mathbf{x}^{(k)} - \gamma^{(k)} \mathbf{B}^* \left(\Phi \mathbf{x}^{(k)} - \mathbf{y} \right) \right) \quad \{\alpha^{(k)}, \gamma^{(k)}\}$$

- supervised training on \mathcal{X} with $\mathbb{E}_{\mathbf{x} \sim \mathcal{X}, e} \text{loss}(f_{\theta, \Phi}(\Phi \mathbf{x} + e), \mathbf{x})$, unsupervised is possible as well
- \mathcal{X} = bounded s -sparse vectors and cross-coherence $\mu := \mu(\mathbf{B}, \Phi)$

Theorem: [Liu et al., 2019 ; Hauffen, PJ, Mücke 2022]

For any $\mathbf{x} \in \mathcal{X}$ with $s < (\mu^{-1} + 1)/2$, $\|\mathbf{e}\|_2 \leq \epsilon$ and $\alpha^{(k)} \gtrsim \mu + \epsilon$ **ALISTA** yields:

$$\|\mathbf{x}^{(k)} - \mathbf{x}\|_2 \lesssim c^k s + (1 + kc^k) \|\mathbf{e}\|_2$$

for $c = c(\{\alpha_{k'}\}_{k'=1}^k) < 1$ and constants depend on \mathcal{X} and $\{\alpha_{k'}\}_{k'=1}^k$.

- adaptive parameters depend on $\|\mathbf{x}^{(k)} - \mathbf{x}\|_1$

- adaptive parameters depend on $\|\mathbf{x}^{(k)} - \mathbf{x}\|_1$
- use $r^{(k)} = \|\Phi \mathbf{x}^{(k)} - \mathbf{y}\|_1$ and $u^{(k)} = \|\mathbf{B}^* (\Phi \mathbf{x}^{(k)} - \mathbf{y})\|_1$
- *neurally augmented ALISTA* (NA-ALISTA) [Behrens, Sauder, PJ 2020]:

$$\mathbf{x}^{(k+1)} = \eta_{\alpha(r^{(k)}, u^{(k)})} \left(\mathbf{x}^{(k)} - \gamma(r^{(k)}, u^{(k)}) \left[\mathbf{B}^* (\Phi \mathbf{x}^{(k)} - \mathbf{y}) \right] \right)$$

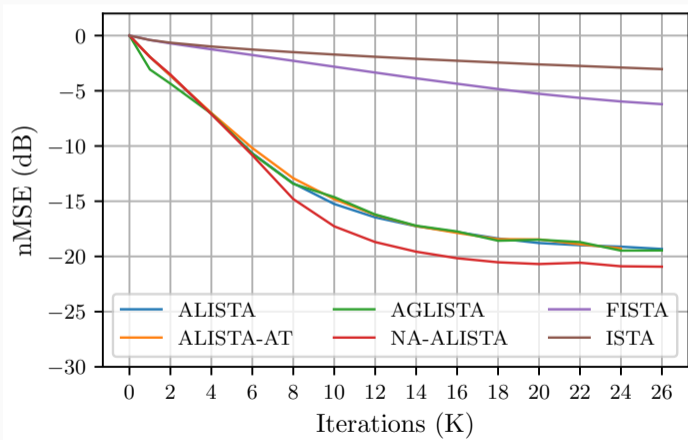
train LSTMs to represent $\alpha(\cdot, \cdot)$ and $\gamma(\cdot, \cdot)$

- adaptive parameters depend on $\|\mathbf{x}^{(k)} - \mathbf{x}\|_1$
- use $r^{(k)} = \|\Phi \mathbf{x}^{(k)} - \mathbf{y}\|_1$ and $u^{(k)} = \|\mathbf{B}^*(\Phi \mathbf{x}^{(k)} - \mathbf{y})\|_1$
- *neurally augmented ALISTA* (NA-ALISTA) [Behrens, Sauder, PJ 2020]:

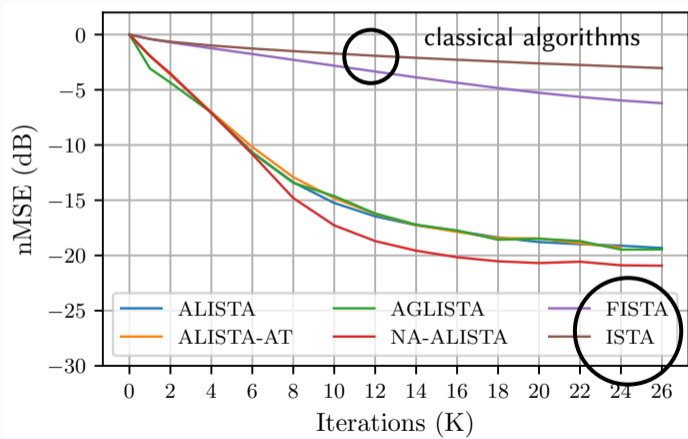
$$\mathbf{x}^{(k+1)} = \eta_{\alpha(r^{(k)}, u^{(k)})} \left(\mathbf{x}^{(k)} - \gamma(r^{(k)}, u^{(k)}) \left[\mathbf{B}^* (\Phi \mathbf{x}^{(k)} - \mathbf{y}) \right] \right)$$

train LSTMs to represent $\alpha(\cdot, \cdot)$ and $\gamma(\cdot, \cdot)$

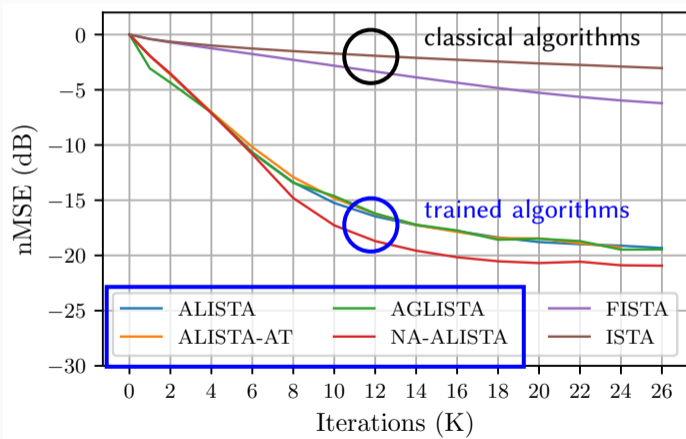
Significantly reduces number of iterations required! Enabling real-time applications



Significantly reduces number of iterations required! Enabling real-time applications

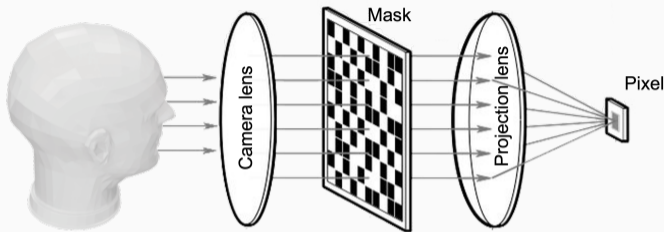


Significantly reduces number of iterations required! Enabling real-time applications



APPLICATION 1: LEARNING MASKS FOR SINGLE-PIXEL IMAGING

APPLICATION 1: LEARNING MASKS FOR SINGLE-PIXEL IMAGING



[adapted from [\[Bacca et al. 2019\]](#)]

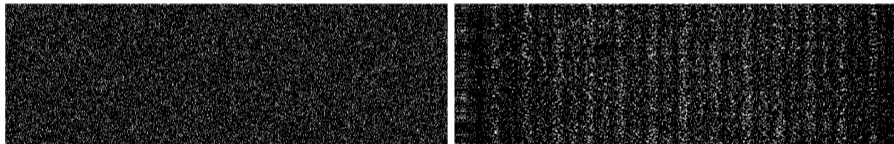
- $\mathcal{X} = \text{MNIST}$ ($n = 28^2 = 784$)
- Φ is the set of **binary matrices** with $k_i = 32$ ones per row i (“on”-pixels per mask)
- $f_{\theta, \Phi}$ is NA-ALISTA with $T = 20$ layers/iterations
- $\mathcal{L}(\hat{x}, x) = \|\hat{x} - x\|_2^2$

APPLICATION 1: LEARNING MASKS FOR SINGLE-PIXEL IMAGING



(a) Random Masks

(b) Learned Masks



(c) Random Φ

(d) Learned Φ



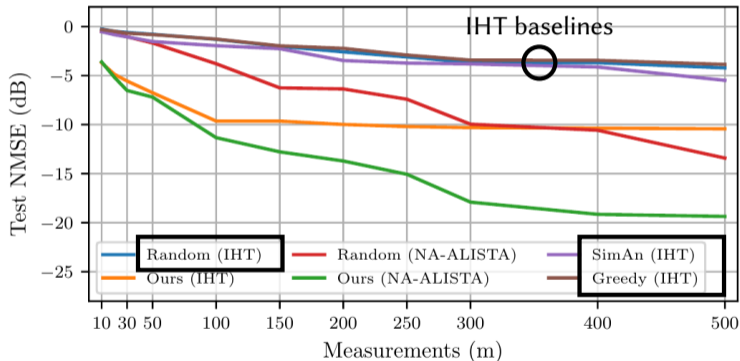
(e) Reconstructed (random Φ)

(f) Reconstructed (learned Φ)

APPLICATION 1: LEARNING MASKS FOR SINGLE-PIXEL IMAGING

baseline: Iterative Hard-thresholding (IHT) $\mathbf{x}^{(k+1)} = H_s (\mathbf{x}^{(k)} - \Phi^* (\Phi \mathbf{x}^{(k)} - \mathbf{y}))$

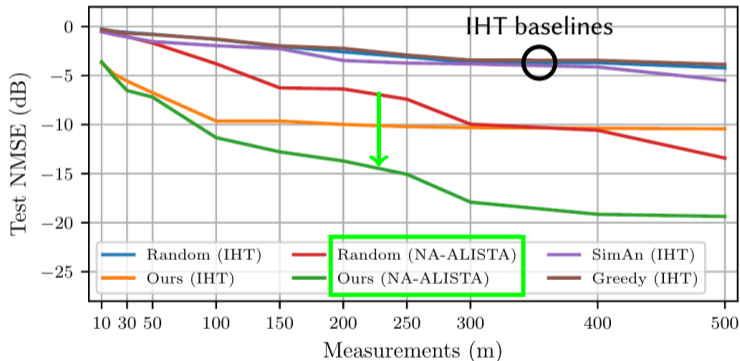
- $n = 28^2 = 784$ (\mathcal{X} =MNIST)
- random $I \subset \mathcal{I} = [m] \times [n]$
- swap a random $1 \Leftrightarrow 0$
- accept if loss L decreases
- **Greedy:** not accept if L increase
- **SimAn:** accept if $\exp[(L - L')/\tau] \leq u \sim U(0, 1)$ and $\tau \rightarrow 0$ with training



APPLICATION 1: LEARNING MASKS FOR SINGLE-PIXEL IMAGING

baseline: Iterative Hard-thresholding (IHT) $\mathbf{x}^{(k+1)} = H_s (\mathbf{x}^{(k)} - \Phi^* (\Phi \mathbf{x}^{(k)} - \mathbf{y}))$

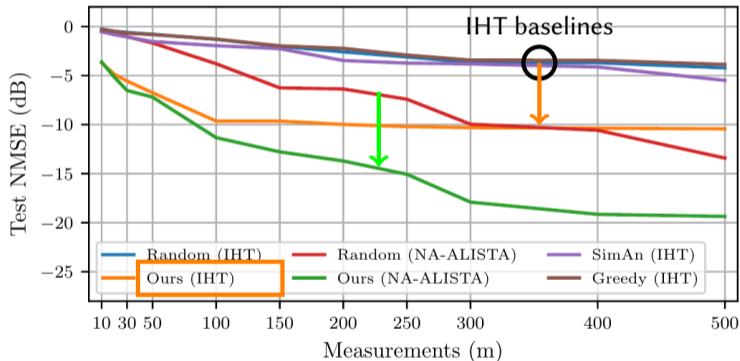
- $n = 28^2 = 784$ (\mathcal{X} =MNIST)
- random $I \subset \mathcal{I} = [m] \times [n]$
- swap a random $1 \Leftrightarrow 0$
- accept if loss L decreases
- **Greedy:** not accept if L increase
- **SimAn:** accept if $\exp[(L - L')/\tau] \leq u \sim U(0, 1)$ and $\tau \rightarrow 0$ with training



APPLICATION 1: LEARNING MASKS FOR SINGLE-PIXEL IMAGING

baseline: Iterative Hard-thresholding (IHT) $\mathbf{x}^{(k+1)} = H_s (\mathbf{x}^{(k)} - \Phi^* (\Phi \mathbf{x}^{(k)} - \mathbf{y}))$

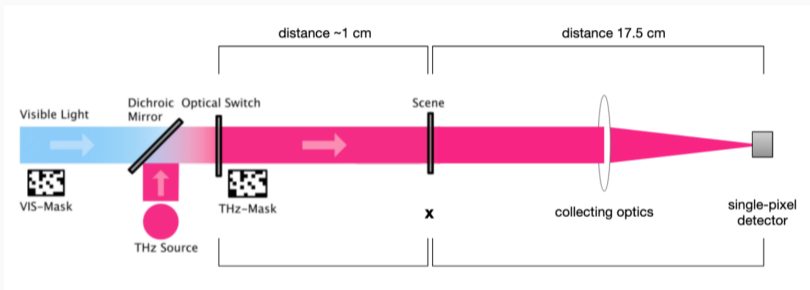
- $n = 28^2 = 784$ (\mathcal{X} =MNIST)
- random $I \subset \mathcal{I} = [m] \times [n]$
- swap a random $1 \Leftrightarrow 0$
- accept if loss L decreases
- **Greedy:** not accept if L increase
- **SimAn:** accept if $\exp[(L - L')/\tau] \leq u \sim U(0, 1)$ and $\tau \rightarrow 0$ with training



☺ **Improvements for IHT while trained with NA-ALISTA \Rightarrow due to Φ**

APPLICATION 1: LEARNING MASKS FOR SINGLE-PIXEL IMAGING

SUPER-PIXEL MASKS FOR SPI

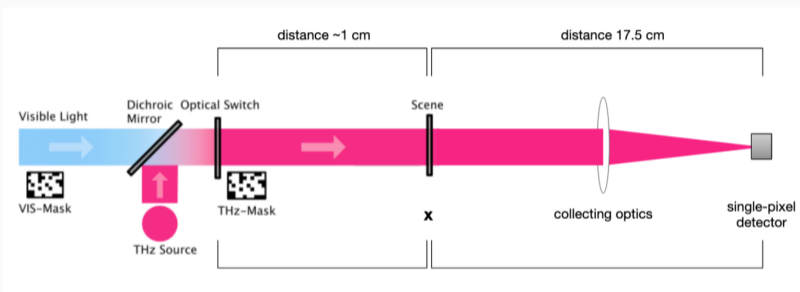


Terahertz SPI [Augustin, PJ, Frohmann, Hübers, 2019] - image: [Reiche, PJ 2020]

- wavelength is in order of pixel size, diffraction degrades image quality \Rightarrow use super-pixels!

APPLICATION 1: LEARNING MASKS FOR SINGLE-PIXEL IMAGING

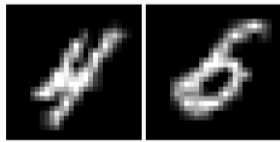
SUPER-PIXEL MASKS FOR SPI



(b) Learned Masks



(d) Learned Φ



Terahertz SPI [Augustin, PJ, Frohmann, Hübers, 2019] - image: [Reiche, PJ 2020]

- wavelength is in order of pixel size, diffraction degrades image quality \Rightarrow use super-pixels!
- sample binary masks with d ones forming super-pixels of size Δ

APPLICATION 2: LEARNING POOLING MATRICES FOR GROUP TESTING

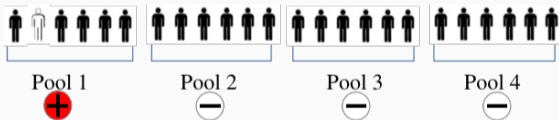
APPLICATION 2: LEARNING POOLING MATRICES FOR GROUP TESTING

- **Non-adaptive viral testing of n individuals with minimal amount of m qPCR-tests**
- $s \ll n$ individuals are positive, i.e., vector of viral loads $\mathbf{x} = (x_1, \dots, x_n)$ **is sparse**
thus, viral-load recovery is a **compressed sensing problem** [Bah, Petersen & PJ, 2024]

APPLICATION 2: LEARNING POOLING MATRICES FOR GROUP TESTING

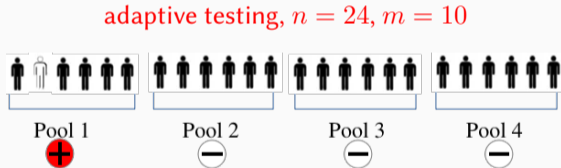
- Non-adaptive viral testing of n individuals with minimal amount of m qPCR-tests
- $s \ll n$ individuals are positive, i.e., vector of viral loads $\mathbf{x} = (x_1, \dots, x_n)$ is sparse
thus, viral-load recovery is a **compressed sensing problem** [Bah, Petersen & PJ, 2024]

adaptive testing, $n = 24$, $m = 10$



APPLICATION 2: LEARNING POOLING MATRICES FOR GROUP TESTING

- **Non-adaptive viral testing of n individuals with minimal amount of m qPCR-tests**
- $s \ll n$ individuals are positive, i.e., vector of viral loads $\mathbf{x} = (x_1, \dots, x_n)$ **is sparse**
thus, viral-load recovery is a **compressed sensing problem** [Bah, Petersen & PJ, 2024]

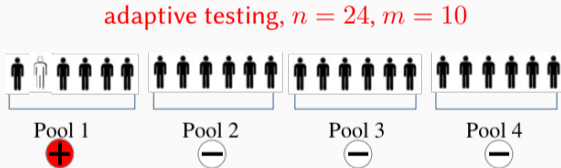


	+	-	-	-	-	Test result
4	●	●	●	●	●	-
3	●	●	●	●	●	-
2	●	●	●	●	●	-
1	○	●	●	●	●	+
0	●	●	●	●	●	-
Grid points	0	1	2	3	4	
		Pools				

non-adaptive testing
 $n = 25, m = 10$

APPLICATION 2: LEARNING POOLING MATRICES FOR GROUP TESTING

- Non-adaptive viral testing of n individuals with minimal amount of m qPCR-tests
- $s \ll n$ individuals are positive, i.e., vector of viral loads $\mathbf{x} = (x_1, \dots, x_n)$ is sparse
thus, viral-load recovery is a **compressed sensing problem** [Bah, Petersen & PJ, 2024]



	+	-	-	-	-	Test result
4	●	●	●	●	●	-
3	●	●	●	●	●	-
2	●	●	●	●	●	-
1	○	●	●	●	●	+
0	●	●	●	●	●	-
Grid points	0	1	2	3	4	
	Pools					

non-adaptive testing
 $n = 25, m = 10$

- binary pooling matrix $\Phi \in \{0, 1\}^{m \times n}$ has k non-zeros per row (pool size)
- pooled measurements of viral loads are $\mathbf{y} = \Phi \mathbf{x} + \mathbf{e}$ with sparse $\mathbf{x} \geq 0$

APPLICATION 2: LEARNING POOLING MATRICES FOR GROUP TESTING

- qPCR modeling is complicated and involves multiplicative noise, yielding heavy-tailed models
- **non-negative least absolute deviations** (NNLAD) is a compressed sensing algorithm

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x} \geq 0} \|\Phi \mathbf{x} - \mathbf{y}\|_1$$

if $\exists \mathbf{t}$ such that $\Phi^T \mathbf{t} > 0$ and $\|\hat{\mathbf{x}} - \mathbf{x}\|_1 \lesssim O(\|e\|_1)$ if Φ has **NSP** [Petersen, Bah, PJ, 2021]

- ℓ_1 -regularization is superflous \Rightarrow **no hyperparameter tuning !**

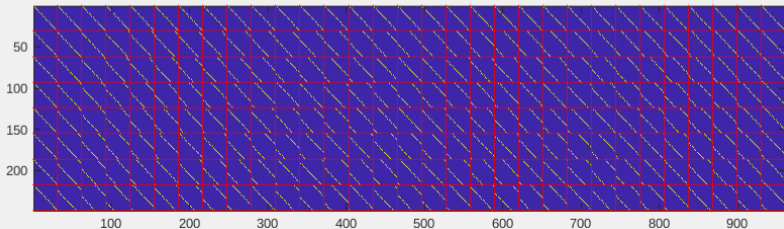
APPLICATION 2: LEARNING POOLING MATRICES FOR GROUP TESTING

- qPCR modeling is complicated and involves multiplicative noise, yielding heavy-tailed models
- **non-negative least absolute deviations** (NNLAD) is a compressed sensing algorithm

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x} \geq 0} \|\Phi \mathbf{x} - \mathbf{y}\|_1$$

if $\exists \mathbf{t}$ such that $\Phi^T \mathbf{t} > 0$ and $\|\hat{\mathbf{x}} - \mathbf{x}\|_1 \lesssim O(\|e\|_1)$ if Φ has **NSP** [Petersen, Bah, PJ, 2021]

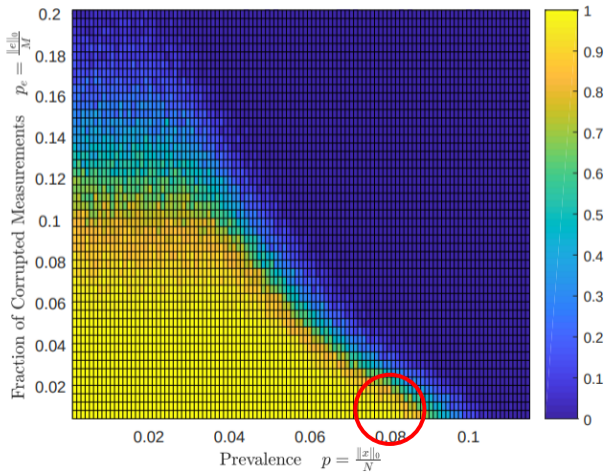
- ℓ_1 -regularization is superflous \Rightarrow **no hyperparameter tuning !**



Example: disjunct matrix $\Phi \in \{0, 1\}^{248 \times 961}$ with $k = 31$ and has NSP for $s \leq 7$ (LDPC/Array code) [Lofti&Vidyasagar, 2020]

APPLICATION 2: LEARNING POOLING MATRICES FOR GROUP TESTING

robust against corrupted tests (sparse noise) and **empirical success up to $s \lesssim 70$!!**



Learning pooling matrices (k ones per row) and tuning viral load recovery

- projected subgradient descent for>NNLAD can be unfolded
- slow convergence (not relevant for testing)
- supervised training with 200 iterations/layers
- inference with 1000 iterations

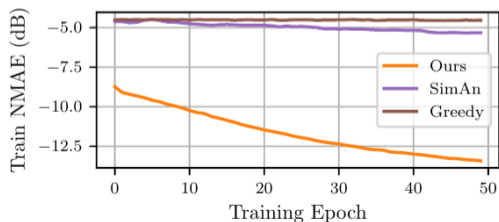
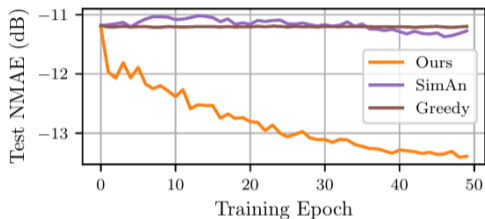
Learning pooling matrices (k ones per row) and tuning viral load recovery

- projected subgradient descent for>NNLAD can be unfolded
- slow convergence (not relevant for testing)
- supervised training with 200 iterations/layers
- inference with 1000 iterations
- $\Phi \in \{0, 1\}^{248 \times 961}$, $k = 31$
- $\mathcal{X} = \{s = 80 \text{ and } x_i^{\text{non-zero}} \sim \text{Beta iid}\}$
- $\text{NMAE} = \frac{\|x - \hat{x}\|_1}{\|x\|_1}$

APPLICATION 2: LEARNING POOLING MATRICES FOR GROUP TESTING

Learning pooling matrices (k ones per row) and tuning viral load recovery

- projected subgradient descent for>NNLAD can be unfolded
- slow convergence (not relevant for testing)
- supervised training with 200 iterations/layers
- inference with 1000 iterations
- $\Phi \in \{0, 1\}^{248 \times 961}$, $k = 31$
- $\mathcal{X} = \{s = 80 \text{ and } x_i^{\text{non-zero}} \sim \text{Beta iid}\}$
- $\text{NMAE} = \frac{\|x - \hat{x}\|_1}{\|x\|_1}$



- ⊖ optimal measurement matrices can not be constructed analytically, only random method
- ⊖ classical algorithms well-understood but too slow and generic for realtime applications

- ☹ optimal measurement matrices can not be constructed analytically, only random method
- ☹ classical algorithms well-understood but too slow and generic for realtime applications

algorithm unfolding allows backprop and data-driven optimization of Φ

- 😊 overcome handcrafted hyper-parameter tuning
- 😊 learn better and constrained random models for $\Phi \Rightarrow$ **derandomization**

“learn to sense and to recover”

- ☹ optimal measurement matrices can not be constructed analytically, only random method
- ☹ classical algorithms well-understood but too slow and generic for realtime applications

algorithm unfolding allows backprop and data-driven optimization of Φ

- 😊 overcome handcrafted hyper-parameter tuning
- 😊 learn better and constrained random models for $\Phi \Rightarrow$ **derandomization**

“learn to sense and to recover”

- 😊 examples in single-pixel/detector imaging and pooling

Thank you!

- [1] E. J. Gumbel. *Statistical theory of extreme values and some practical applications: a series of lectures*. Vol. 33. US Government Printing Office, 1954.
- [2] E. Jang, S. Gu, and B. Poole. “Categorical reparameterization with gumbel-softmax”. *arXiv preprint arXiv:1611.01144* (2016).
- [3] C. J. Maddison, A. Mnih, and Y. W. Teh. “The concrete distribution: A continuous relaxation of discrete random variables”. *arXiv preprint arXiv:1611.00712* (2016).
- [4] R. Kueng and P. Jung. “Robust Nonnegative Sparse Recovery and the Nullspace Property of 0/1 Measurements”. *IEEE Transactions on Information Theory* 64.2 (Feb. 2018), 689–703.
- [5] P. Jung, R. Kueng, and D. G. Mixon. “Derandomizing Compressed Sensing With Combinatorial Design”. *Frontiers in Applied Mathematics and Statistics* 5 (June 2019). eprint: 1812.08130.
- [6] H. B. Petersen, B. Bah, and P. Jung. “Practical high-throughput, non-adaptive and noise-robust SARS-CoV-2 testing”. *arXiv preprint arXiv:2007.09171* (2020).
- [7] F. Behrens, J. Sauder, and P. Jung. “Neurally Augmented ALISTA”. *International Conference on Learning Representations (ICLR)* (2021).
- [8] H. B. Petersen, B. Bah, and P. Jung. “Efficient Tuning-Free l1-Regression of Nonnegative Compressible Signals”. *Frontiers in Applied Mathematics and Statistics* 7 (June 2021).
- [9] B. Bah, H. B. Petersen, and P. Jung. “Compressed sensing-based SARS-CoV-2 pool testing”. *Notices of the American Mathematical Society* 2 (2024).