

Extending Sparse Dictionary Learning Methods for Adversarial Robustness

36th OTDK

Mahmoud Aslan Balázs Mészáros

Wednesday 5th July, 2023



Outline

1 Theory

- Sparse Coding
- Thresholding & Iterative Thresholding Pursuit
- Layered Basis Pursuit
- Deep Pursuit
- Group Pursuit

2 Experiments

- Synthetic database
- MNIST
 - Basis Pursuit and Feedforward networks
 - Layered Basis Pursuit and Deep Pursuit

3 Conclusions

Outline

1 Theory

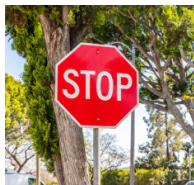
- Sparse Coding
- Thresholding & Iterative Thresholding Pursuit
- Layered Basis Pursuit
- Deep Pursuit
- Group Pursuit

2 Experiments

- Synthetic database
- MNIST
 - Basis Pursuit and Feedforward networks
 - Layered Basis Pursuit and Deep Pursuit

3 Conclusions

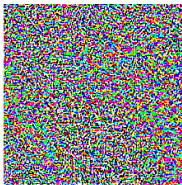
Adversarial Examples



x

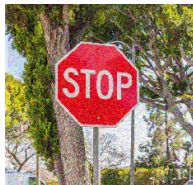
Prediction: "Stop"

$+ .007 \times$



$\text{sign}(\nabla_x J(\theta, x, y))$

$=$



$x + \epsilon \text{sign}(\nabla_x J(\theta, x, y))$

Prediction: "Turn Right" *

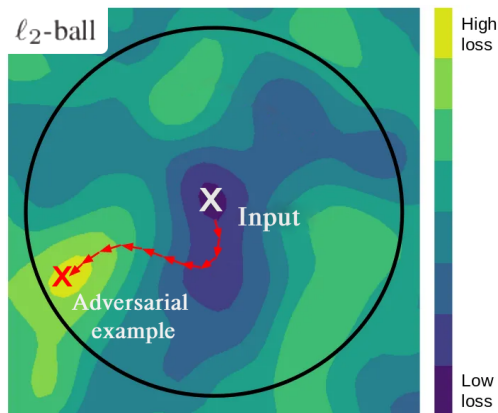
* illustrative example

Goodfellow, I.J., Shlens, J. and Szegedy, C., 2014. Explaining and harnessing adversarial examples. arXiv preprint arXiv:1412.6572.

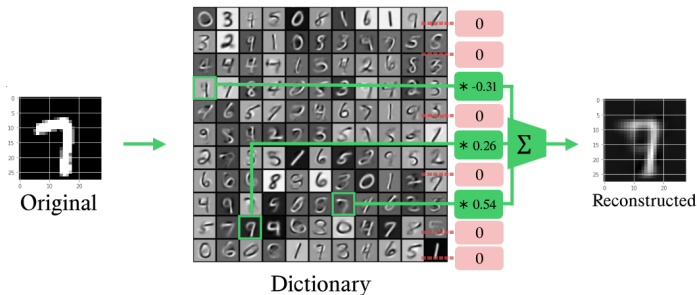
Kurakin, A., Goodfellow, I.J. and Bengio, S., 2018. Adversarial examples in the physical world. In Artificial intelligence safety and security (pp. 99-112). Chapman and Hall/CRC.

Adversarial Attacks

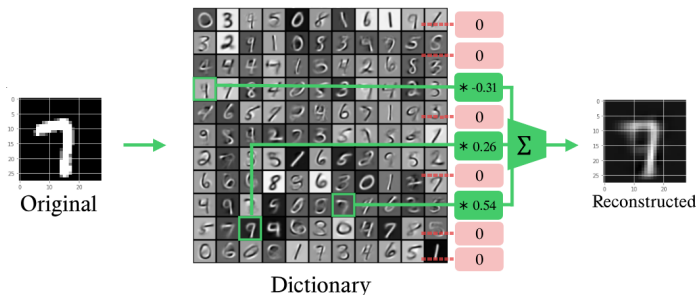
Figure: ℓ_2 bounded IFGSM attack



Sparse Coding



Sparse Coding - ℓ_1 Relaxation¹



$$\arg \min_{\hat{\mathbf{r}}} \frac{1}{2} \|\mathbf{D}\hat{\mathbf{r}} - \mathbf{x}\|_2^2 + \gamma \|\hat{\mathbf{r}}\|_1 \quad (\text{BP})$$

¹Chen, S.S., Donoho, D.L. and Saunders, M.A., 2001. Atomic decomposition by basis pursuit. SIAM review, 43(1), pp.129-159.

Thresholding Pursuit

$$\phi_{\gamma}(\mathbf{D}^T \mathbf{X}) \tag{1}$$

Thresholding Pursuit

$$\phi_{\gamma}(\mathbf{D}^T \mathbf{X}) \quad (2)$$

ϕ_{γ} can be the soft, hard, or non-neg soft thresholding.

Thresholding Operators

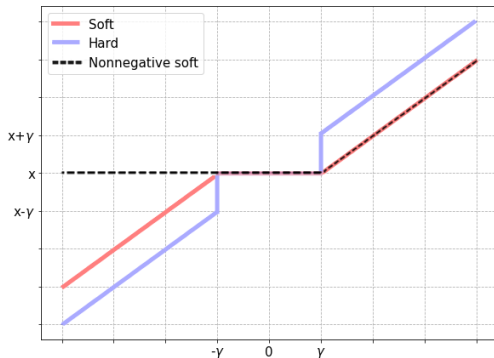


Figure: Soft vs hard vs nonnegative soft thresholding

²Papayan, V., Romano, Y. and Elad, M., 2017. Convolutional neural networks analyzed via convolutional sparse coding. The Journal of Machine Learning Research, 18(1), pp.2887-2938.

Thresholding Operators

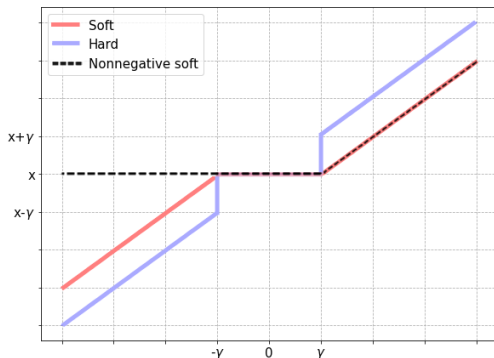


Figure: Soft vs hard vs nonnegative soft thresholding

$$S_{\gamma}^{+}(\mathbf{x}) = \max(\mathbf{x} - \gamma, 0) = \text{ReLU}(\mathbf{x} - \gamma) \quad (3^2)$$

²Papayan, V., Romano, Y. and Elad, M., 2017. Convolutional neural networks analyzed via convolutional sparse coding. The Journal of Machine Learning Research, 18(1), pp.2887-2938.

Thresholding Pursuit: Simplicity vs. Recovery

Problem: Thresholding pursuit doesn't recover exact support.³

³Donoho, D.L. and Elad, M., 2003. Optimally sparse representation in general (nonorthogonal) dictionaries via ℓ_1 minimization. Proceedings of the National Academy of Sciences, 100(5), pp.2197-2202.

FISTA - Fast Iterative Shrinkage & Thresholding Algorithm⁴

Initialize:

$$\mathbf{r}^0 := \phi_{\gamma}(\mathbf{D}^T \mathbf{X})$$

⁴Beck, A. and Teboulle, M., 2009. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. SIAM journal on imaging sciences, 2(1), pp.183-202.

FISTA - Fast Iterative Shrinkage & Thresholding Algorithm⁴

Initialize:

$$\mathbf{\Gamma}^0 := \phi_{\gamma}(\mathbf{D}^T \mathbf{X})$$

iterate:

$$\mathbf{\Gamma}^t := \phi_{\gamma}(\mathbf{\Gamma}^{t-1} - \alpha \mathbf{D}^T (\mathbf{D} \mathbf{\Gamma}^{t-1} - \mathbf{X}))$$

⁴Beck, A. and Teboulle, M., 2009. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. SIAM journal on imaging sciences, 2(1), pp.183-202.

FISTA - Fast Iterative Shrinkage & Thresholding Algorithm⁵

Initialize:

$$\mathbf{\Gamma}^0 := \phi_{\gamma}(\mathbf{D}^T \mathbf{X})$$

iterate:

$$\mathbf{\Gamma}^t := \phi_{\gamma}(\mathbf{\Gamma}^{t-1} - \alpha \mathbf{D}^T (\mathbf{D} \mathbf{\Gamma}^{t-1} - \mathbf{X}))$$

⁵Beck, A. and Teboulle, M., 2009. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. SIAM journal on imaging sciences, 2(1), pp.183-202.

Layered Basis Pursuit⁶

$$\arg \min_{\mathbf{r}_i} \frac{1}{2} \|\mathbf{D}_i \mathbf{r}_i - \hat{\mathbf{r}}_{i-1}\|_2^2 + \gamma_i \|\mathbf{r}_i\|_1 \quad (\text{LBP})$$

⁶Papayan, V., Romano, Y. and Elad, M., 2017. Convolutional neural networks analyzed via convolutional sparse coding. The Journal of Machine Learning Research, 18(1), pp.2887-2938.

Layered Basis Pursuit⁶

$$\arg \min_{\mathbf{r}_i} \frac{1}{2} \|\mathbf{D}_i \mathbf{r}_i - \hat{\mathbf{r}}_{i-1}\|_2^2 + \gamma_i \|\mathbf{r}_i\|_1 \quad (\text{LBP})$$

LBP suffers from **error accumulation** as we go deeper, and doesn't offer support for **skip connections**.

⁶Papayan, V., Romano, Y. and Elad, M., 2017. Convolutional neural networks analyzed via convolutional sparse coding. The Journal of Machine Learning Research, 18(1), pp.2887-2938.

Deep Pursuit⁷

$$\arg \min_{\mathbf{r}_j, j \in \{1, \dots, l\}} \frac{1}{2} \sum_{j=1}^l \|\mathbf{r}_{j-1} - \mathbf{D}_j \mathbf{r}_j\|_2^2 + \gamma_j \|\mathbf{r}_j\|_1 \quad (\text{DP})$$

⁷Cazenavette, G., Murdock, C. and Lucey, S., 2021. Architectural adversarial robustness: The case for deep pursuit. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition: (pp. 7150-7158).

Deep Pursuit⁷

$$\arg \min_{\mathbf{r}_j, j \in \{1, \dots, l\}} \frac{1}{2} \sum_{j=1}^l \|\mathbf{r}_{j-1} - \mathbf{D}_j \mathbf{r}_j\|_2^2 + \gamma_j \|\mathbf{r}_j\|_1 \quad (\text{DP})$$

or in a matrix form:

$$\arg \min_{\mathbf{r}_j, j \in \{1, \dots, l\}} \frac{1}{2} \left\| \begin{bmatrix} \mathbf{x} \\ \mathbf{0} \\ \vdots \\ \mathbf{0} \end{bmatrix} - \begin{bmatrix} \mathbf{D}_1 & \dots & \mathbf{0} \\ -\mathbf{I} & \mathbf{D}_2 & \vdots \\ \vdots & \ddots & \ddots \\ \mathbf{0} & -\mathbf{I} & \mathbf{D}_l \end{bmatrix} \begin{bmatrix} \mathbf{r}_1 \\ \mathbf{r}_2 \\ \vdots \\ \mathbf{r}_l \end{bmatrix} \right\|_2^2 + \sum_{j=1}^l \gamma_j \|\mathbf{r}_j\|_1$$

⁷Cazenavette, G., Murdock, C. and Lucey, S., 2021. Architectural adversarial robustness: The case for deep pursuit. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (pp. 7150-7158).

Deep Pursuit's Gradient

$$\hat{\mathbf{g}}_j^t := \begin{cases} \mathbf{D}_j^T (\mathbf{D}_j \hat{\mathbf{r}}_j^{t-1} - \mathbf{r}_{j-1}^t) + (\hat{\mathbf{r}}_j^{t-1} - \mathbf{D}_{j+1} \mathbf{r}_{j+1}^{t-1}) & j < l \\ \mathbf{D}_j^T (\mathbf{D}_j \hat{\mathbf{r}}_j^{t-1} - \mathbf{r}_{j-1}^t) & j = l \end{cases}$$

Deep Pursuit with Skip Connections

$$\arg \min_{\mathbf{r}_j, j \in \{1, \dots, l\}} \frac{1}{2} \left\| \begin{bmatrix} \mathbf{X} \\ \mathbf{0} \\ \vdots \\ \mathbf{0} \end{bmatrix} - \begin{bmatrix} \mathbf{D}_1 & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{B}_{21} & \mathbf{D}_2 & \ddots & \vdots \\ \vdots & \ddots & \ddots & \mathbf{0} \\ \mathbf{B}_{l1} & \dots & \mathbf{B}_{l(l-1)} & \mathbf{D}_l \end{bmatrix} \begin{bmatrix} \mathbf{r}_1 \\ \mathbf{r}_2 \\ \vdots \\ \mathbf{r}_l \end{bmatrix} \right\|_2^2 + \sum_{j=1}^l \gamma_j \|\mathbf{r}_j\|_1$$

Group Pursuit⁸

$$\arg \min_{\hat{\mathbf{f}}} \frac{1}{2} \|\mathbf{X} - \mathbf{D}\hat{\mathbf{f}}\|_2^2 + \langle \gamma, l(\hat{\mathbf{f}}) \rangle \quad (\text{GBP})$$

where $\langle \gamma, l(\hat{\mathbf{f}}) \rangle$ is a generalized regularizer to include ℓ_1 , ℓ_2 , $\ell_{1,2}$, and $\ell_{\beta,1,2}$ on groups of $\hat{\mathbf{f}}$.

⁸Szeghy, D., Aslan, M., Fóthi, Á., Mészáros, B., Milacski, Z.Á. and Lőrincz, A., 2022. Structural Extensions of Basis Pursuit: Guarantees on Adversarial Robustness. arXiv preprint arXiv:2205.08955.

Group Pursuit Generalization

$$\arg \min_{\hat{\mathbf{r}}_j, j \in \{1, \dots, l\}} \frac{1}{2} \left\| \begin{bmatrix} \mathbf{x} \\ \mathbf{0} \\ \vdots \\ \mathbf{0} \end{bmatrix} - \begin{bmatrix} \mathbf{D}_1 & \mathbf{F}_{12} & \dots & \mathbf{F}_{1l} \\ \mathbf{B}_{21} & \mathbf{D}_2 & \ddots & \vdots \\ \vdots & \ddots & \ddots & \mathbf{F}_{(l-1)l} \\ \mathbf{B}_{l1} & \dots & \mathbf{B}_{l(l-1)} & \mathbf{D}_l \end{bmatrix} \begin{bmatrix} \hat{\mathbf{r}}_1 \\ \hat{\mathbf{r}}_2 \\ \vdots \\ \hat{\mathbf{r}}_l \end{bmatrix} \right\|_2^2 + \langle \gamma, l(\hat{\mathbf{r}}) \rangle$$

Outline

1 Theory

- Sparse Coding
- Thresholding & Iterative Thresholding Pursuit
- Layered Basis Pursuit
- Deep Pursuit
- Group Pursuit

2 Experiments

- Synthetic database
- MNIST
 - Basis Pursuit and Feedforward networks
 - Layered Basis Pursuit and Deep Pursuit

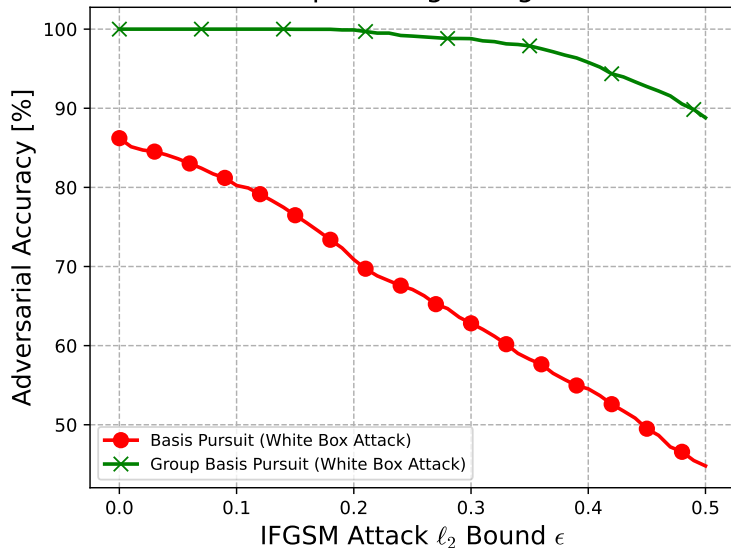
3 Conclusions

Shallow Experiments

- Databases:
 - ▶ Synthetic database
 - ▶ MNIST
- Architectures:
 - ▶ Basis Pursuit
 - ▶ Group Basis Pursuit
 - ▶ Pooled Group Basis Pursuit, Transformer, Shallow dense, Deep dense

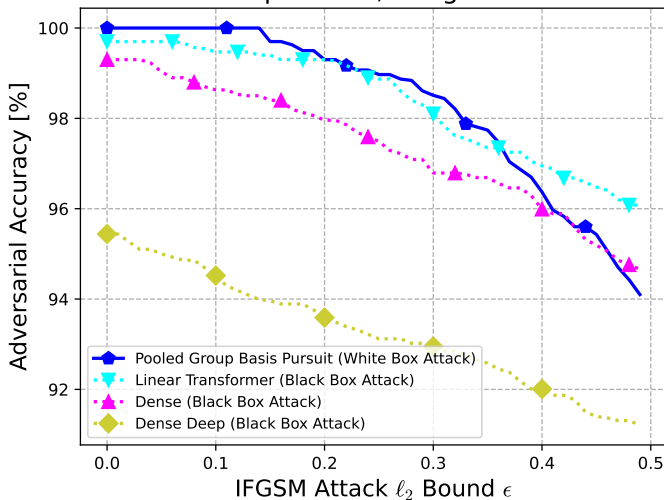
Synthetic database

Synthetic Robustness to IFGSM Attack No Group Pooling, Margin=0.1

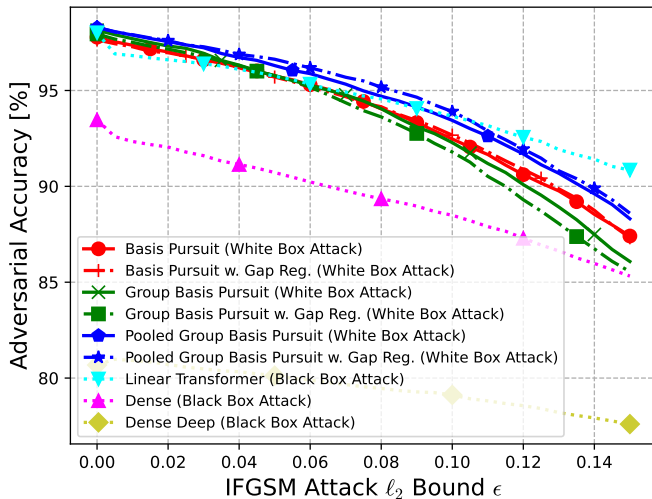


Synthetic database

Synthetic Robustness to IFGSM Attack
Group Pooled, Margin=0.1

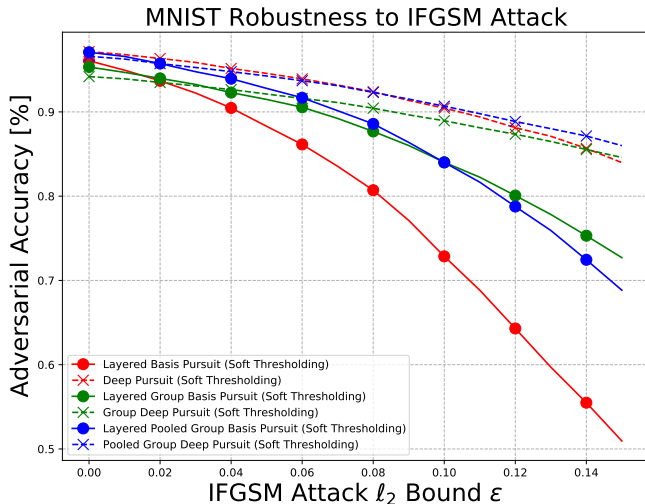


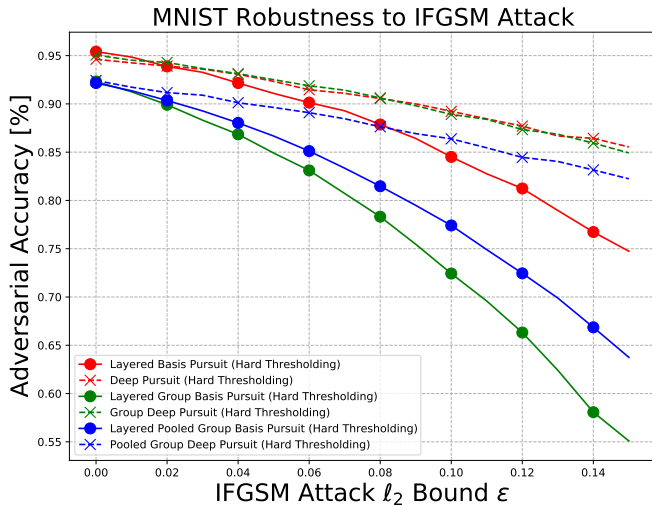
MNIST Robustness to IFGSM Attack



Deep Experiments

- Databases:
 - ▶ MNIST
- Architectures:
 - ▶ Layered Basis Pursuit
 - ▶ Layered Group Basis Pursuit
 - ▶ Layered Pooled Group Basis Pursuit
 - ▶ Deep Pursuit
 - ▶ Group Deep Pursuit
 - ▶ Pooled Group Deep Pursuit





Outline

1 Theory

- Sparse Coding
- Thresholding & Iterative Thresholding Pursuit
- Layered Basis Pursuit
- Deep Pursuit
- Group Pursuit

2 Experiments

- Synthetic database
- MNIST
 - Basis Pursuit and Feedforward networks
 - Layered Basis Pursuit and Deep Pursuit

3 Conclusions

Conclusions

- Our group and pooled group methods managed to overcome non-group pursuit methods in many cases
- As expected DP is more robust than LBP but not necessarily better than a single BP layer
- Feedforward estimations (especially the Transformer) were efficient
- Introduced MC and Gap terms in loss function
- Future directions:
 - ▶ More challenging data: e.g. ImageNet, CIFAR10, ...
 - ▶ More complex architectures: e.g. ResNet50, Transformers, ...
 - ▶ Other downstream tasks: e.g. 3D pose estimation, ...
- Szeghy, D.; Aslan, M.; Fóthi, Á.; Mészáros, B.; Milacski, Z. and Lőrincz, A. (2022). **Structural Extensions of Basis Pursuit: Guarantees on Adversarial Robustness.** In Proceedings of the 3rd International Conference on Deep Learning Theory and Applications - DeLTA.

Thank you for your attention, thank you András Lőrincz and Dávid Szeghy for supervising our research and thank you Zoltán Ádám Milacski, Áron Fóthi and Ellák Somfai for all of your contributions to our work!