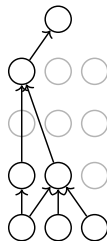
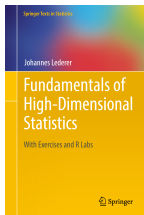


Sparse Deep Learning: From High-Dimensional Statistics to Neural Networks

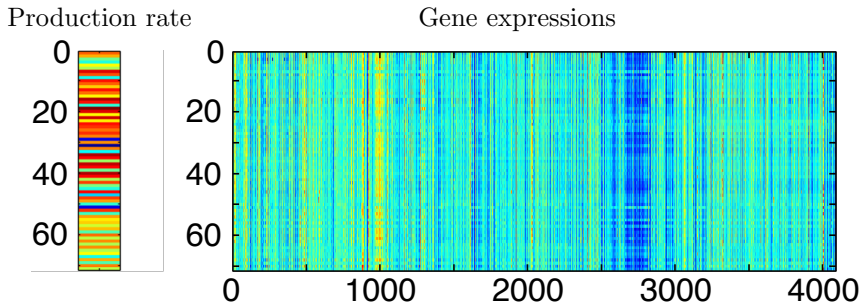


Extreme Value Theory
Optimization Online Learning
Econometrics Computer Science
Algorithms Tuning Parameters
Machine Learning Mathematics
Deep Learning
Genomics High Dimensions
Neuroscience

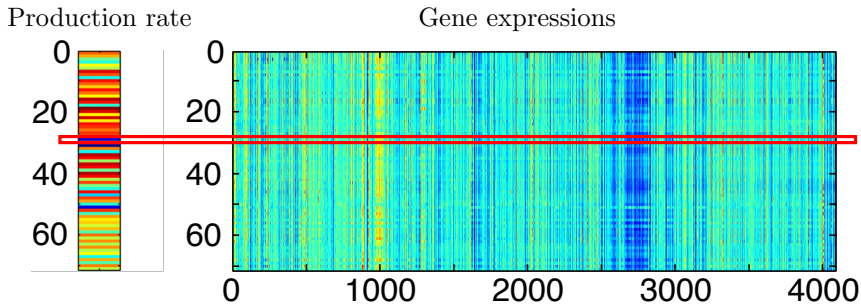
Johannes Lederer



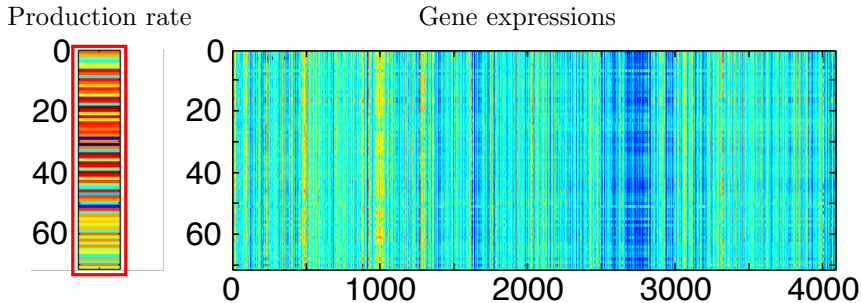
High-Dimensionality Is Common: Riboflavin Production in *B. subtilis*



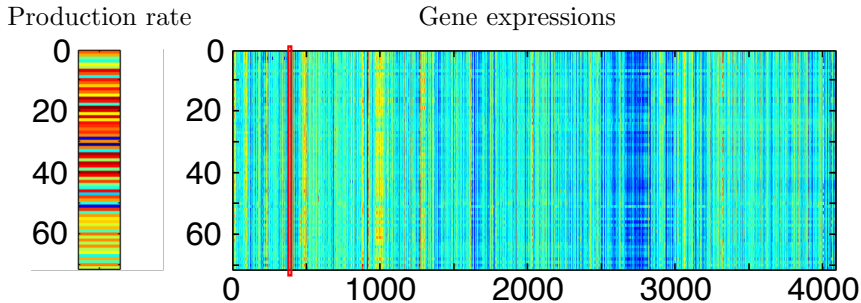
High-Dimensionality Is Common: Riboflavin Production in *B. subtilis*



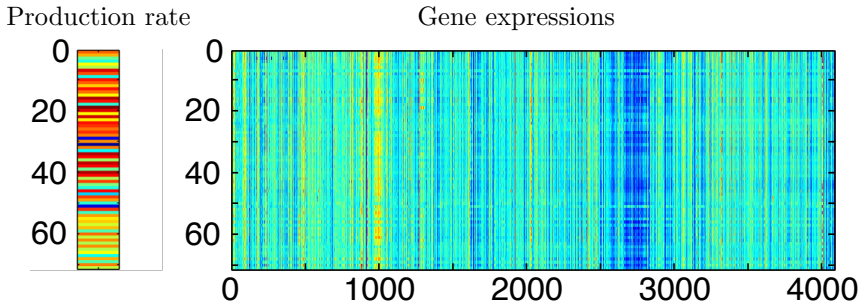
High-Dimensionality Is Common: Riboflavin Production in *B. subtilis*



High-Dimensionality Is Common: Riboflavin Production in *B. subtilis*



High-Dimensionality Is Common: Riboflavin Production in *B. subtilis*



$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta}^* + \mathbf{u}$$

Sparsity-Inducing Prior Terms Are Standard Tools in Statistics

(Fundamentals of High-Dimensional Statistics—With Exercises and R Labs, 2021)

$$\mathbf{y} = X\boldsymbol{\beta}^* + \mathbf{u}$$

$$\hat{\boldsymbol{\beta}}_{\text{lasso}} \in \underset{\boldsymbol{\beta} \in \mathbb{R}^p}{\operatorname{argmin}} \{ \|\mathbf{y} - X\boldsymbol{\beta}\|_2^2 + r\|\boldsymbol{\beta}\|_1 \}$$

$$\hat{\boldsymbol{\beta}}_{\text{grplasso}} \in \underset{\boldsymbol{\beta} \in \mathbb{R}^p}{\operatorname{argmin}} \{ \|\mathbf{y} - X\boldsymbol{\beta}\|_2^2 + r\|\boldsymbol{\beta}\|_{2,1} \}$$

$$\text{where } \|\boldsymbol{\beta}\|_1 := \sum_{j=1}^p |\beta_j| \quad \text{and} \quad \|\boldsymbol{\beta}\|_{2,1} := \sum_{k=1}^g \|\boldsymbol{\beta}_{\mathcal{G}_k}\|_2$$

Theory Exists: Oracle Inequalities

(Oracle inequalities for high-dimensional prediction, 2019)

Theorem (Power-One Bound): It holds with probability at least $1 - 1/n$ that

$$\frac{\|X\boldsymbol{\beta}^* - X\hat{\boldsymbol{\beta}}_{\text{lasso}}\|_2^2}{n} \leq 6\sigma\|\boldsymbol{\beta}^*\|_1\|\mathbf{x}\|_n \sqrt{\frac{\log[np]}{n}}.$$

$\sigma \leftrightarrow$ noise

$\|\mathbf{x}\|_n = 1$ (for normalized inputs)

$n =$ # samples

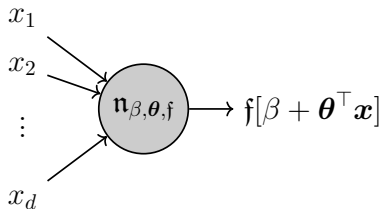
$p =$ # parameters

Neural Networks Consist of Neurons

(Activation functions in artificial neural networks: a systematic overview, 2021)

$$\mathbf{n}_{\beta, \boldsymbol{\theta}, \mathfrak{f}} : \mathbb{R}^d \rightarrow \mathbb{R}$$

$$\mathbf{x} \mapsto \mathfrak{f}[\beta + \boldsymbol{\theta}^\top \mathbf{x}] = \mathfrak{f}\left[\beta + \sum_{j=1}^d \theta_j x_j\right]$$



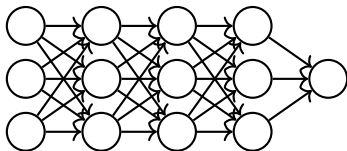
Neurons Can Be Readily Combined

(Activation functions in artificial neural networks: a systematic overview, 2021)

$$\mathbf{g}_{\Theta} : \mathbb{R}^d \rightarrow \mathbb{R}$$

$$\mathbf{x} \mapsto \mathbf{g}_{\Theta}[\mathbf{x}] := W^L \mathbf{f}^L [\dots W^1 \mathbf{f}^1 [W^0 \mathbf{x}]]$$

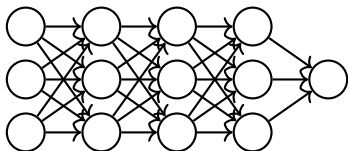
$$\mathcal{A} := \{ \Theta = (W^L, \dots, W^0) : W^l \in \mathbb{R}^{p_{l+1} \times p_l} \}$$



Neural Networks Are High-Dimensional

(Statistical guarantees for regularized neural networks, 2021)

$$\mathcal{A} := \left\{ \Theta = (W^L, \dots, W^0) : W^l \in \mathbb{R}^{p_{l+1} \times p_l} \right\}$$



Number of parameters in the toy network: $P = 30$

Number of parameters in general: $P = \sum_{l=0}^L p_{l+1} p_l$

We Invoke Sparsity I: Motivation

(Layer sparsity in neural networks, 2020)

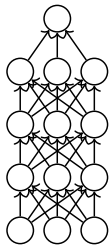
Sparsity is well established in statistics and machine learning:

- Avoid overfitting
- Save computations and memory
- Improve interpretability

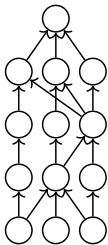
We Invoke Sparsity II: Concepts

(Layer sparsity in neural networks, 2020)

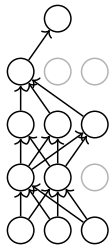
no sparsity



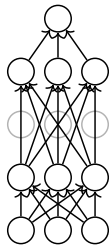
connections



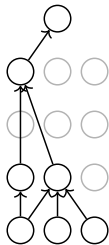
nodes



layers



combined



We Invoke Sparsity III: Implementations

(Layer sparsity in neural networks, 2020)

$$\mathfrak{r}^{\text{con}}[\Theta] := \sum_{l=0}^L \lVert\lVert W^l \rVert_1 \rVert := \sum_{l=0}^L \sum_{v=1}^{p_{l+1}} \sum_{w=1}^{p_l} |(W^l)_{vw}|$$

$$\mathfrak{r}^{\text{node}}[\Theta] := \sum_{l=0}^L \lVert\lVert W^l \rVert_{2,1} \rVert := \sum_{l=0}^L \sum_{v=1}^{p_{l+1}} \sqrt{\sum_{w=1}^{p_l} |(W^l)_{vw}|^2}$$

$$\mathfrak{r}^{\text{lay}}[\Theta] := \sum_{l=1}^L \lVert\lVert W^l \rVert_{2,-} \rVert := \sum_{l=1}^L \sqrt{\sum_{v=1}^{p_{l+1}} \sum_{w=1}^{p_l} (\text{neg}[(W^l)_{vw}])^2}$$

Sparsity Leads to Accurate Deep Learning

(Statistical guarantees for regularized neural networks, 2021)

Theorem: It holds with probability at least $1 - 1/N$ that

$$\text{err}^2[\hat{\kappa}_{\mathfrak{h}} \mathfrak{g}_{\hat{\Omega}_{\mathfrak{h}}}] \leq a\sigma\kappa_* \left(\frac{2a_{\text{Lip}}}{L}\right)^L \|\mathbf{x}\|_N \sqrt{L \log(2P)} \frac{\log(2N)}{\sqrt{N}}.$$

a = constant

$\sigma \leftrightarrow$ noise

$\kappa_* \leftrightarrow \|\beta^*\|_1$

$a_{\text{Lip}} = 1$ (for relu)

$\|\mathbf{x}\|_N = 1$ (for normalized inputs)

$N = \#$ samples

$P = \#$ parameters

$L = \#$ hidden layers

We Use Standard Measures for the Accuracy

(Statistical guarantees for regularized neural networks, 2021)

$$\text{err}[\hat{\kappa}_{\mathfrak{h}} \mathfrak{g}_{\hat{\Omega}_{\mathfrak{h}}}] := \sqrt{\frac{1}{N} \sum_{i=1}^N (\hat{\kappa}_{\mathfrak{h}} \mathfrak{g}_{\hat{\Omega}_{\mathfrak{h}}}[\mathbf{x}_i] - \kappa_* \mathfrak{g}_{\Omega_*}[\mathbf{x}_i])^2}$$

$$\text{risk}[\hat{\kappa}_{\mathfrak{h}} \mathfrak{g}_{\hat{\Omega}_{\mathfrak{h}}}] := \mathbb{E}_{(\mathbf{x}, y)} \left[(\hat{\kappa}_{\mathfrak{h}} \mathfrak{g}_{\hat{\Omega}_{\mathfrak{h}}}[\mathbf{x}] - y)^2 \right]$$

Using a Parametric Model for Illustration

(Statistical guarantees for regularized neural networks, 2021)

$$y_i = \kappa_* \mathfrak{g}_{\Omega_*}[\mathbf{x}_i] + u_i$$

$$u_1, \dots, u_N \text{ i.i.d. } \mathcal{N}[0, \sigma^2]$$

Using a Parametric Model for Illustration

(Statistical guarantees for regularized neural networks, 2021)

$$y_i = \kappa_* \mathbf{g}_{\Omega_*}[\mathbf{x}_i] + u_i$$

Proposition: Assume relu activation and define for a norm \mathfrak{h}

$$\mathcal{A}_{\mathfrak{h}} := \{ \Theta \in \mathcal{A} : \mathfrak{h}[\Theta] \leq 1 \}.$$

Then, for every $\Theta \in \mathcal{A}$, there exists a pair of $\kappa \in [0, \infty)$ and $\Omega \in \mathcal{A}_{\mathfrak{h}}$ such that

$$\mathbf{g}_{\Theta}[\mathbf{x}] = \kappa \mathbf{g}_{\Omega}[\mathbf{x}] \quad \text{for all } \mathbf{x} \in \mathbb{R}^d;$$

and vice versa, for every pair of $\kappa \in [0, \infty)$ and $\Omega \in \mathcal{A}_{\mathfrak{h}}$, there exists a $\Theta \in \mathcal{A}$ such that the above equality holds.

Scale Regularization Reduces Ambiguity

(Statistical guarantees for regularized neural networks, 2021)

$$(\hat{\kappa}_{\mathfrak{h}}, \hat{\Omega}_{\mathfrak{h}}) \in \underset{\substack{\kappa \in [0, \infty) \\ \Omega \in \mathcal{A}_{\mathfrak{h}}}}{\operatorname{argmin}} \left\{ \frac{1}{N} \sum_{i=1}^N (y_i - \kappa \mathfrak{g}_{\Omega}[\mathbf{x}_i])^2 + r\kappa \right\}$$

r is an “appropriate” tuning parameter

Our Proofs Connect Different Fields

(Statistical guarantees for regularized neural networks, 2021)

High-dimensional statistics

- Oracle inequalities and effective noise

Analysis

- “Scale-free” NNs are bounded and Lipschitz

Empirical-process theory

- Metric entropy and concentration inequalities

Oracle Inequality Due to Scale Trick

(Statistical guarantees for regularized neural networks, 2021)

$$\frac{1}{N} \sum_{i=1}^N (y_i - \hat{\kappa}_{\mathfrak{h}} \mathfrak{g}_{\hat{\Omega}_{\mathfrak{h}}}[\mathbf{x}_i])^2 + r \hat{\kappa}_{\mathfrak{h}} \leq \frac{1}{N} \sum_{i=1}^N (y_i - \kappa_{\Omega} \mathfrak{g}_{\Omega}[\mathbf{x}_i])^2 + r \kappa$$

Oracle Inequality Due to Scale Trick

(Statistical guarantees for regularized neural networks, 2021)

$$\frac{1}{N} \sum_{i=1}^N (y_i - \hat{\kappa}_{\mathfrak{h}} \mathfrak{g}_{\hat{\Omega}_{\mathfrak{h}}}[\mathbf{x}_i])^2 + r \hat{\kappa}_{\mathfrak{h}} \leq \frac{1}{N} \sum_{i=1}^N (y_i - \kappa \mathfrak{g}_{\Omega}[\mathbf{x}_i])^2 + r \kappa$$

$$\frac{1}{N} \sum_{i=1}^N (\kappa_* \mathfrak{g}_{\Omega_*}[\mathbf{x}_i] + u_i - \hat{\kappa}_{\mathfrak{h}} \mathfrak{g}_{\hat{\Omega}_{\mathfrak{h}}}[\mathbf{x}_i])^2 + r \hat{\kappa}_{\mathfrak{h}} \leq \frac{1}{N} \sum_{i=1}^N (\kappa_* \mathfrak{g}_{\Omega_*}[\mathbf{x}_i] + u_i - \kappa \mathfrak{g}_{\Omega}[\mathbf{x}_i])^2 + r \kappa$$

Oracle Inequality Due to Scale Trick

(Statistical guarantees for regularized neural networks, 2021)

$$\frac{1}{N} \sum_{i=1}^N (y_i - \hat{\kappa}_{\mathfrak{h}} \mathfrak{g}_{\hat{\Omega}_{\mathfrak{h}}}[\mathbf{x}_i])^2 + r \hat{\kappa}_{\mathfrak{h}} \leq \frac{1}{N} \sum_{i=1}^N (y_i - \kappa \mathfrak{g}_{\Omega}[\mathbf{x}_i])^2 + r \kappa$$

$$\frac{1}{N} \sum_{i=1}^N (\kappa_* \mathfrak{g}_{\Omega_*}[\mathbf{x}_i] + u_i - \hat{\kappa}_{\mathfrak{h}} \mathfrak{g}_{\hat{\Omega}_{\mathfrak{h}}}[\mathbf{x}_i])^2 + r \hat{\kappa}_{\mathfrak{h}} \leq \frac{1}{N} \sum_{i=1}^N (\kappa_* \mathfrak{g}_{\Omega_*}[\mathbf{x}_i] + u_i - \kappa \mathfrak{g}_{\Omega}[\mathbf{x}_i])^2 + r \kappa$$

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N (\hat{\kappa}_{\mathfrak{h}} \mathfrak{g}_{\hat{\Omega}_{\mathfrak{h}}}[\mathbf{x}_i] - \kappa_* \mathfrak{g}_{\Omega_*}[\mathbf{x}_i])^2 &\leq \frac{1}{N} \sum_{i=1}^N (\kappa \mathfrak{g}_{\Omega}[\mathbf{x}_i] - \kappa_* \mathfrak{g}_{\Omega_*}[\mathbf{x}_i])^2 \\ &+ \frac{2}{N} \sum_{i=1}^N \hat{\kappa}_{\mathfrak{h}} \mathfrak{g}_{\hat{\Omega}_{\mathfrak{h}}}[\mathbf{x}_i] u_i - \frac{2}{N} \sum_{i=1}^N \kappa \mathfrak{g}_{\Omega}[\mathbf{x}_i] u_i + r \kappa - r \hat{\kappa}_{\mathfrak{h}} \end{aligned}$$

Oracle Inequality Due to Scale Trick

(Statistical guarantees for regularized neural networks, 2021)

$$\frac{1}{N} \sum_{i=1}^N (y_i - \hat{\kappa}_{\mathfrak{h}} \mathfrak{g}_{\hat{\Omega}_{\mathfrak{h}}}[\mathbf{x}_i])^2 + r\hat{\kappa}_{\mathfrak{h}} \leq \frac{1}{N} \sum_{i=1}^N (y_i - \kappa \mathfrak{g}_{\Omega}[\mathbf{x}_i])^2 + r\kappa$$

$$\frac{1}{N} \sum_{i=1}^N (\kappa_* \mathfrak{g}_{\Omega_*}[\mathbf{x}_i] + u_i - \hat{\kappa}_{\mathfrak{h}} \mathfrak{g}_{\hat{\Omega}_{\mathfrak{h}}}[\mathbf{x}_i])^2 + r\hat{\kappa}_{\mathfrak{h}} \leq \frac{1}{N} \sum_{i=1}^N (\kappa_* \mathfrak{g}_{\Omega_*}[\mathbf{x}_i] + u_i - \kappa \mathfrak{g}_{\Omega}[\mathbf{x}_i])^2 + r\kappa$$

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N (\hat{\kappa}_{\mathfrak{h}} \mathfrak{g}_{\hat{\Omega}_{\mathfrak{h}}}[\mathbf{x}_i] - \kappa_* \mathfrak{g}_{\Omega_*}[\mathbf{x}_i])^2 &\leq \frac{1}{N} \sum_{i=1}^N (\kappa \mathfrak{g}_{\Omega}[\mathbf{x}_i] - \kappa_* \mathfrak{g}_{\Omega_*}[\mathbf{x}_i])^2 \\ &\quad + \frac{2}{N} \sum_{i=1}^N \hat{\kappa}_{\mathfrak{h}} \mathfrak{g}_{\hat{\Omega}_{\mathfrak{h}}}[\mathbf{x}_i] u_i - \frac{2}{N} \sum_{i=1}^N \kappa \mathfrak{g}_{\Omega}[\mathbf{x}_i] u_i + r\kappa - r\hat{\kappa}_{\mathfrak{h}} \end{aligned}$$

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N (\hat{\kappa}_{\mathfrak{h}} \mathfrak{g}_{\hat{\Omega}_{\mathfrak{h}}}[\mathbf{x}_i] - \kappa_* \mathfrak{g}_{\Omega_*}[\mathbf{x}_i])^2 &\leq \frac{1}{N} \sum_{i=1}^N (\kappa \mathfrak{g}_{\Omega}[\mathbf{x}_i] - \kappa_* \mathfrak{g}_{\Omega_*}[\mathbf{x}_i])^2 \\ &\quad + \hat{\kappa}_{\mathfrak{h}} \sup_{\Omega \in \mathcal{A}_{\mathfrak{h}}} \left| \frac{2}{N} \sum_{i=1}^N \mathfrak{g}_{\Omega}[\mathbf{x}_i] u_i \right| + \kappa \sup_{\Omega \in \mathcal{A}_{\mathfrak{h}}} \left| \frac{2}{N} \sum_{i=1}^N \mathfrak{g}_{\Omega}[\mathbf{x}_i] u_i \right| + r\kappa - r\hat{\kappa}_{\mathfrak{h}}. \end{aligned}$$

Oracle Inequality Due to Scale Trick

(Statistical guarantees for regularized neural networks, 2021)

$$\frac{1}{N} \sum_{i=1}^N (y_i - \hat{\kappa}_{\mathfrak{h}} \mathfrak{g}_{\hat{\Omega}_{\mathfrak{h}}}[\mathbf{x}_i])^2 + r\hat{\kappa}_{\mathfrak{h}} \leq \frac{1}{N} \sum_{i=1}^N (y_i - \kappa \mathfrak{g}_{\Omega}[\mathbf{x}_i])^2 + r\kappa$$

$$\frac{1}{N} \sum_{i=1}^N (\kappa_* \mathfrak{g}_{\Omega_*}[\mathbf{x}_i] + u_i - \hat{\kappa}_{\mathfrak{h}} \mathfrak{g}_{\hat{\Omega}_{\mathfrak{h}}}[\mathbf{x}_i])^2 + r\hat{\kappa}_{\mathfrak{h}} \leq \frac{1}{N} \sum_{i=1}^N (\kappa_* \mathfrak{g}_{\Omega_*}[\mathbf{x}_i] + u_i - \kappa \mathfrak{g}_{\Omega}[\mathbf{x}_i])^2 + r\kappa$$

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N (\hat{\kappa}_{\mathfrak{h}} \mathfrak{g}_{\hat{\Omega}_{\mathfrak{h}}}[\mathbf{x}_i] - \kappa_* \mathfrak{g}_{\Omega_*}[\mathbf{x}_i])^2 &\leq \frac{1}{N} \sum_{i=1}^N (\kappa \mathfrak{g}_{\Omega}[\mathbf{x}_i] - \kappa_* \mathfrak{g}_{\Omega_*}[\mathbf{x}_i])^2 \\ &\quad + \frac{2}{N} \sum_{i=1}^N \hat{\kappa}_{\mathfrak{h}} \mathfrak{g}_{\hat{\Omega}_{\mathfrak{h}}}[\mathbf{x}_i] u_i - \frac{2}{N} \sum_{i=1}^N \kappa \mathfrak{g}_{\Omega}[\mathbf{x}_i] u_i + r\kappa - r\hat{\kappa}_{\mathfrak{h}} \end{aligned}$$

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N (\hat{\kappa}_{\mathfrak{h}} \mathfrak{g}_{\hat{\Omega}_{\mathfrak{h}}}[\mathbf{x}_i] - \kappa_* \mathfrak{g}_{\Omega_*}[\mathbf{x}_i])^2 &\leq \frac{1}{N} \sum_{i=1}^N (\kappa \mathfrak{g}_{\Omega}[\mathbf{x}_i] - \kappa_* \mathfrak{g}_{\Omega_*}[\mathbf{x}_i])^2 \\ &\quad + \hat{\kappa}_{\mathfrak{h}} \sup_{\Omega \in \mathcal{A}_{\mathfrak{h}}} \left| \frac{2}{N} \sum_{i=1}^N \mathfrak{g}_{\Omega}[\mathbf{x}_i] u_i \right| + \kappa \sup_{\Omega \in \mathcal{A}_{\mathfrak{h}}} \left| \frac{2}{N} \sum_{i=1}^N \mathfrak{g}_{\Omega}[\mathbf{x}_i] u_i \right| + r\kappa - r\hat{\kappa}_{\mathfrak{h}}. \end{aligned}$$

$$\frac{1}{N} \sum_{i=1}^N (\hat{\kappa}_{\mathfrak{h}} \mathfrak{g}_{\hat{\Omega}_{\mathfrak{h}}}[\mathbf{x}_i] - \kappa_* \mathfrak{g}_{\Omega_*}[\mathbf{x}_i])^2 \leq \frac{1}{N} \sum_{i=1}^N (\kappa \mathfrak{g}_{\Omega}[\mathbf{x}_i] - \kappa_* \mathfrak{g}_{\Omega_*}[\mathbf{x}_i])^2 + 2r\kappa$$

Connecting Statistics and Computations: Stationary Points Instead of Global Optima

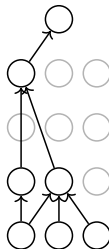
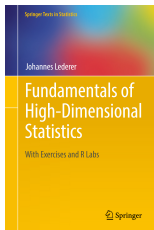
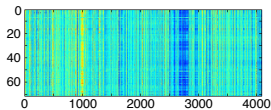
(Upcoming)

$$y_i = \underbrace{\alpha\beta}_{=:\gamma} x_i + u_i$$

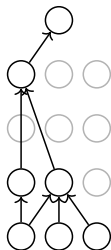
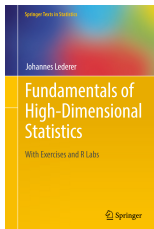
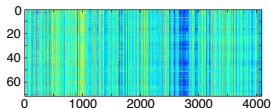
$$(\hat{\alpha}, \hat{\beta}) \in \operatorname{argmin}_{\alpha, \beta \in \mathbb{R}} \left\{ \sum_{i=1}^n (y_i - \alpha\beta x_i)^2 \right\}$$

$$\hat{\gamma} \in \operatorname{argmin}_{\gamma \in \mathbb{R}} \left\{ \sum_{i=1}^n (y_i - \gamma x_i)^2 \right\}$$

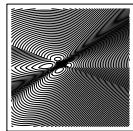
Sparsity Is a Promising Concept in Deep Learning



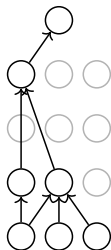
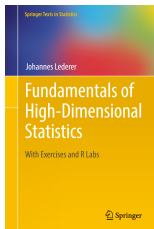
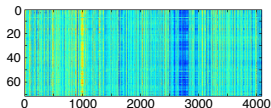
Sparsity Is a Promising Concept in Deep Learning



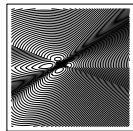
johanneslederer.com
github.com/LedererLab



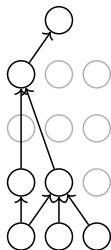
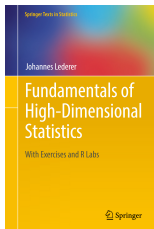
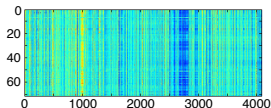
Sparsity Is a Promising Concept in Deep Learning



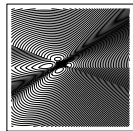
johanneslederer.com
github.com/LedererLab

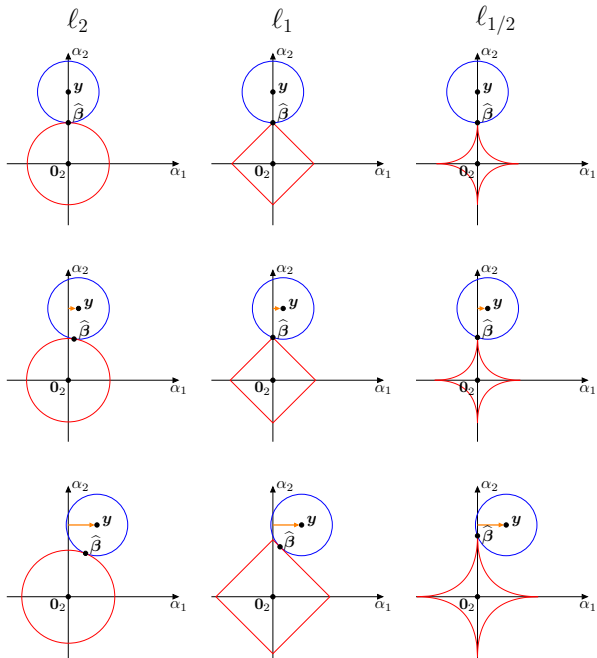


Sparsity Is a Promising Concept in Deep Learning



johanneslederer.com
github.com/LedererLab





We Use Standard Measures for the Accuracy

(Statistical guarantees for regularized neural networks, 2021)

$$\text{err}[\hat{\kappa}_{\mathfrak{h}} \mathfrak{g}_{\hat{\Omega}_{\mathfrak{h}}}] := \sqrt{\frac{1}{N} \sum_{i=1}^N (\hat{\kappa}_{\mathfrak{h}} \mathfrak{g}_{\hat{\Omega}_{\mathfrak{h}}}[\mathbf{x}_i] - \kappa_* \mathfrak{g}_{\Omega_*}[\mathbf{x}_i])^2}$$

$$\text{risk}[\hat{\kappa}_{\mathfrak{h}} \mathfrak{g}_{\hat{\Omega}_{\mathfrak{h}}}] := \mathbb{E}_{(\mathbf{x}, y)} \left[(\hat{\kappa}_{\mathfrak{h}} \mathfrak{g}_{\hat{\Omega}_{\mathfrak{h}}}[\mathbf{x}] - y)^2 \right]$$

Using a Parametric Model for Illustration

(Statistical guarantees for regularized neural networks, 2021)

$$y_i = \kappa_* \mathfrak{g}_{\Omega_*}[\mathbf{x}_i] + u_i$$

$$u_1, \dots, u_N \text{ i.i.d. } \mathcal{N}[0, \sigma^2]$$

Using a Parametric Model for Illustration

(Statistical guarantees for regularized neural networks, 2021)

$$y_i = \kappa_* \mathbf{g}_{\Omega_*}[\mathbf{x}_i] + u_i$$

Proposition: Assume relu activation and define for a norm \mathfrak{h}

$$\mathcal{A}_{\mathfrak{h}} := \{ \Theta \in \mathcal{A} : \mathfrak{h}[\Theta] \leq 1 \}.$$

Then, for every $\Theta \in \mathcal{A}$, there exists a pair of $\kappa \in [0, \infty)$ and $\Omega \in \mathcal{A}_{\mathfrak{h}}$ such that

$$\mathbf{g}_{\Theta}[\mathbf{x}] = \kappa \mathbf{g}_{\Omega}[\mathbf{x}] \quad \text{for all } \mathbf{x} \in \mathbb{R}^d;$$

and vice versa, for every pair of $\kappa \in [0, \infty)$ and $\Omega \in \mathcal{A}_{\mathfrak{h}}$, there exists a $\Theta \in \mathcal{A}$ such that the above equality holds.

Scale Regularization Reduces Ambiguity

(Statistical guarantees for regularized neural networks, 2021)

$$(\hat{\kappa}_{\mathfrak{h}}, \hat{\Omega}_{\mathfrak{h}}) \in \underset{\substack{\kappa \in [0, \infty) \\ \Omega \in \mathcal{A}_{\mathfrak{h}}}}{\operatorname{argmin}} \left\{ \frac{1}{N} \sum_{i=1}^N (y_i - \kappa \mathfrak{g}_{\Omega}[\mathbf{x}_i])^2 + r\kappa \right\}$$

r is an “appropriate” tuning parameter

Neural Networks Have Lipschitz Properties

(Statistical guarantees for regularized neural networks, 2021)

Proposition: It holds for every $\mathbf{x} \in \mathbb{R}^d$ and $\Theta = (W^L, \dots, W^0), \Gamma = (V^L, \dots, V^0) \in \mathcal{A}$ that

$$|\mathbf{g}_\Theta[\mathbf{x}] - \mathbf{g}_\Gamma[\mathbf{x}]| \leq c_{\text{Lip}}[\mathbf{x}] \|\Theta - \Gamma\|_F$$

with

$$c_{\text{Lip}}[\mathbf{x}] := 2(a_{\text{Lip}})^L \sqrt{L} \|\mathbf{x}\|_2 \max_{l \in \{0, \dots, L\}} \prod_{j \in \{0, \dots, L\}, j \neq l} (\|W^j\|_2 \vee \|V^j\|_2).$$

Neural Networks Have Lipschitz Properties

(Statistical guarantees for regularized neural networks, 2021)

$$S_l \mathbf{g}_\Theta : \mathbb{R}^d \rightarrow \mathbb{R}^{p_l}$$

$$\mathbf{x} \mapsto S_l \mathbf{g}_\Theta[\mathbf{x}] := \mathbf{f}^l[\dots W^1 \mathbf{f}^1[W^0 \mathbf{x}]]$$

$$\|S_{l-1} \mathbf{g}_\Theta[\mathbf{x}]\|_2 \leq (a_{\text{Lip}})^{l-1} \|\mathbf{x}\|_2 \prod_{j=0}^{l-2} \|W^j\|_2$$

$$S^l \mathbf{g}_\Theta : \mathbb{R}^{p_l} \rightarrow \mathbb{R}$$

$$\mathbf{z} \mapsto S^l \mathbf{g}_\Theta[\mathbf{z}] := W^L \mathbf{f}^L[\dots W^l \mathbf{f}^l[W^{l-1} \mathbf{z}]]$$

$$|S^{l+1} \mathbf{g}_\Theta[\mathbf{z}_1] - S^{l+1} \mathbf{g}_\Theta[\mathbf{z}_2]| \leq (a_{\text{Lip}})^{L-l} \|\mathbf{z}_1 - \mathbf{z}_2\|_2 \prod_{j=l}^L \|W^j\|_2$$

The Relevant Dudley Integrals Are Bounded

(Statistical guarantees for regularized neural networks, 2021)

$$\mathcal{G}_1 := \{\mathbf{g}_\Omega : \Omega \in \mathcal{A}_{\|\cdot\|_1}\}$$

$$\|\mathbf{g}_\Theta - \mathbf{g}_\Gamma\|_N := \sqrt{\sum_{i=1}^N (\mathbf{g}_\Theta[\mathbf{x}_i] - \mathbf{g}_\Gamma[\mathbf{x}_i])^2 / N}$$

H usual metric entropy and J corresponding $\delta/(8\sigma)$ -Dudley integral

Proposition: It holds for every $v \in (0, \infty)$ and $\delta, \sigma \in (0, \infty)$ that satisfy $\delta \leq 8\sigma c_{\text{Lip1}}$ that

$$H[v, \mathcal{G}_1, \|\cdot\|_N] \leq \frac{6(c_{\text{Lip1}})^2}{v^2} \log \left[\frac{ePv^2}{(c_{\text{Lip1}})^2} \vee 2e \right]$$

and

$$J[v, \mathcal{G}_1, \|\cdot\|_N] \leq \frac{5c_{\text{Lip1}}}{2} \sqrt{\log[eP \vee 2e]} \log \left[\frac{8\sigma c_{\text{Lip1}}}{\delta} \right],$$

where $c_{\text{Lip1}} := 2(2a_{\text{Lip}}/L)^L \sqrt{L} \|\mathbf{x}\|_N$.

The Relevant Dudley Integrals Are Bounded

(Statistical guarantees for regularized neural networks, 2021)

Step 1: $H[v, \mathcal{G}_1, \|\cdot\|_N]$ to $H[v/c_{\text{Lip1}}, \mathcal{A}_1, \|\cdot\|_F]$ via Lipschitzness

Step 2: “Sparse covering”

Step 3: Approximations in the Dudley integral

Theory Exists: Oracle Inequalities

(Oracle inequalities for high-dimensional prediction, 2019)

Theorem (Power-One Bound): As long as $r \geq 2\|X^\top u\|_\infty$, it holds that

$$\frac{\|X\Theta_* - X\hat{\omega}_{\text{lasso}}\|_2^2}{n} \leq 2\|\Theta_*\|_1 \frac{r}{n}.$$

Theory Exists: Oracle Inequalities

(Oracle inequalities for high-dimensional prediction, 2019)

$$\|\mathbf{y} - X\hat{\omega}_{\text{lasso}}\|_2^2 + r\|\hat{\omega}_{\text{lasso}}\|_1 \leq \|\mathbf{y} - X\Theta_*\|_2^2 + r\|\Theta_*\|_1$$

$$\Rightarrow \|X\Theta_* - X\hat{\omega}_{\text{lasso}}\|_2^2 + 2\langle X^\top u, \Theta_* - \hat{\omega}_{\text{lasso}} \rangle + r\|\hat{\omega}_{\text{lasso}}\|_1 \leq r\|\Theta_*\|_1$$

$$\begin{aligned} \Rightarrow \|X\Theta_* - X\hat{\omega}_{\text{lasso}}\|_2^2 &\leq 2\|X^\top u\|_\infty \|\Theta_*\|_1 + 2\|X^\top u\|_\infty \|\hat{\omega}_{\text{lasso}}\|_1 \\ &\quad + r\|\Theta_*\|_1 - r\|\hat{\omega}_{\text{lasso}}\|_1 \end{aligned}$$

$$\Rightarrow \|X\Theta_* - X\hat{\omega}_{\text{lasso}}\|_2^2 \leq 2r\|\Theta_*\|_1$$

Theory Exists: Concentration Inequalities

(Fundamentals of High-Dimensional Statistics—With Exercises and R Labs, 2021)

Theorem (Effective Noise): Consider a fixed design matrix $X \in \mathbb{R}^{n \times p}$ with $\max_j (X^\top X)_{jj}/n = 1$, and consider Gauss-distributed noise $u \sim \mathcal{N}_n[\mathbf{0}_n, \sigma^2 \mathbf{I}_{n \times n}]$. Then, for all $t \in (0, 1)$, it holds that

$$\mathbb{P}\{2\|X^\top u\|_\infty/n \leq \sigma\sqrt{8\log[p/t]/n}\} \geq 1 - t.$$

Theory Exists: Concentration Inequalities

(Fundamentals of High-Dimensional Statistics—With Exercises and R Labs, 2021)

$$\mathbb{P}\{2\|X^\top u\|_\infty > r\} \leq \sum_{j=1}^p \mathbb{P}\{2|(X^\top u)_j| > r\}.$$

$$\mathbb{P}\{2\|X^\top u\|_\infty > r\} \leq p \max_{j \in \{1, \dots, p\}} \mathbb{P}\left\{ \frac{|(X^\top u)_j|}{\sigma \sqrt{(X^\top X)_{jj}}} > \frac{r}{2\sigma \sqrt{n}} \right\}$$

$$\mathbb{P}\{|z| \geq a\} \leq e^{-\frac{a^2}{2}} \quad \text{for all } a \in [0, \infty)$$