# Data Science: Statistics Meets Optimization







ComputerScience StatisticalTheory Optimization TuningParameters Mathematics MachineLearning DeepLearning Remotistics Augorithms All HighDimensions DataScience DiffusionModels Astrophysics

# Johannes Lederer

Universität Hamburg der forschung | der lehre | der Bildung

Funded by



GEFÖRDERT VOM Bundesministerium für Bildung und Forschung





Lambro, one year ago

# Extreme-Value Theory Is Not Extrapolation (Extremes in high dimensions: methods and scalable algorithms, 2023)



# Extreme-Value Theory Is Not Extrapolation (Extremes in high dimensions: methods and scalable algorithms, 2023)



# Extreme-Value Theory Is Not Extrapolation (Extremes in high dimensions: methods and scalable algorithms, 2023)



 $\lim_{t \to \infty} \mathbb{P}\left\{ X/t \in \mathcal{A} \mid \|X\| > t \right\} \sim \mathfrak{m}[\mathcal{A}]$ 

## Networks Are High Dimensional

(Extremes in high dimensions: methods and scalable algorithms, 2023)



## Networks Are High Dimensional

(Extremes in high dimensions: methods and scalable algorithms, 2023)



Key challenge:  $p \approx d^2$  large, n usually small

$$\frac{1}{c_{\Gamma}} \frac{1}{x_m} \left( \prod_{j=1}^d \frac{1}{x_j} \right) \mathfrak{n}_{d-1} \left[ \log[\boldsymbol{x}_{-m}/x_m]; -\boldsymbol{\Gamma}_{-m,m}/2, \boldsymbol{\Sigma} \right]$$

• 
$$x \in \{a \in (0,\infty)^d : \|a\| > 1\}$$

- $\mathfrak{n}_{d-1}$ : Gauss density
- $\mathbf{\Gamma} \in \mathbb{R}^{d \times d}$ : variogram matrix
- $\Sigma \equiv \Sigma[\Gamma, m] \in \mathbb{R}^{(d-1) \times (d-1)}$ : covariance matrix
- $c_{\Gamma}$ : normalizing constant

$$\frac{1}{c_{\Gamma}} \frac{1}{x_m} \left( \prod_{j=1}^d \frac{1}{x_j} \right) \mathfrak{n}_{d-1} \left[ \log[\boldsymbol{x}_{-m}/x_m]; -\boldsymbol{\Gamma}_{-m,m}/2, \boldsymbol{\Sigma} \right]$$

C Almost like Gaussian graphical models



$$\frac{1}{c_{\Gamma}} \frac{1}{x_m} \left( \prod_{j=1}^d \frac{1}{x_j} \right) \mathfrak{n}_{d-1} \left[ \log[\boldsymbol{x}_{-m}/x_m]; \underbrace{-\boldsymbol{\Gamma}_{-m,m}/2}_{\sim \boldsymbol{\mu}}, \underbrace{\boldsymbol{\Sigma}}_{\sim \boldsymbol{\Lambda}^{-1}} \right]$$

- C Almost like Gaussian graphical models
- $\bigcirc$  Not really like Gaussian graphical models  $\longrightarrow$  we establish generalized "densities"  $\bigcirc$

$$\frac{1}{c_{\Gamma}} \frac{1}{x_m} \left( \prod_{j=1}^d \frac{1}{x_j} \right) \mathfrak{n}_{d-1} \left[ \log[\boldsymbol{x}_{-m}/x_m]; \underbrace{-\boldsymbol{\Gamma}_{-m,m}/2}_{\sim \boldsymbol{\mu}}, \underbrace{\boldsymbol{\Sigma}}_{\sim \boldsymbol{\Lambda}^{-1}} \right]$$

- C Almost like Gaussian graphical models
- $\bigcirc$  Not really like Gaussian graphical models  $\longrightarrow$  we establish generalized "densities"  $\bigcirc$

Score Matching Circumvents Normalizations (Extremes in high dimensions: methods and scalable algorithms, 2023)

Our densities are of the form

$$\mathfrak{h}[\boldsymbol{x}; \boldsymbol{\mu}, \boldsymbol{\Lambda}] = e^{\mathfrak{g}[\boldsymbol{x}; \boldsymbol{\mu}, \boldsymbol{\Lambda}]} / c_{\boldsymbol{\mu}, \boldsymbol{\Lambda}}$$

with explicit  $\mathfrak{g}[\boldsymbol{x};\boldsymbol{\mu},\boldsymbol{\Lambda}]$  and integral normalization  $c_{\boldsymbol{\mu},\boldsymbol{\Lambda}}$ .



Score Matching Circumvents Normalizations (Extremes in high dimensions: methods and scalable algorithms, 2023)

Our densities are of the form

$$\mathfrak{h}[\boldsymbol{x}; \boldsymbol{\mu}, \boldsymbol{\Lambda}] = e^{\mathfrak{g}[\boldsymbol{x}; \boldsymbol{\mu}, \boldsymbol{\Lambda}]} / c_{\boldsymbol{\mu}, \boldsymbol{\Lambda}}$$

with explicit  $\mathfrak{g}[\boldsymbol{x};\boldsymbol{\mu},\boldsymbol{\Lambda}]$  and integral normalization  $c_{\boldsymbol{\mu},\boldsymbol{\Lambda}}$ .

Maximum likelihood

$$\begin{aligned} \operatorname*{argmax}_{\boldsymbol{\mu},\boldsymbol{\Lambda}} \mathbb{E} \Big[ \log \mathfrak{h}[\boldsymbol{x};\boldsymbol{\mu},\boldsymbol{\Lambda}] \Big] \\ = \ \operatorname*{argmax}_{\boldsymbol{\mu},\boldsymbol{\Lambda}} \mathbb{E} \Big[ \mathfrak{g}[\boldsymbol{x};\boldsymbol{\mu},\boldsymbol{\Lambda}] - \log c_{\boldsymbol{\mu},\boldsymbol{\Lambda}} \Big] \end{aligned}$$

is computationally infeasible.

Score Matching Circumvents Normalizations (Extremes in high dimensions: methods and scalable algorithms, 2023)

Our densities are of the form

$$\mathfrak{h}[\boldsymbol{x}; \boldsymbol{\mu}, \boldsymbol{\Lambda}] = e^{\mathfrak{g}[\boldsymbol{x}; \boldsymbol{\mu}, \boldsymbol{\Lambda}]} / c_{\boldsymbol{\mu}, \boldsymbol{\Lambda}}$$

with explicit  $\mathfrak{g}[\boldsymbol{x};\boldsymbol{\mu},\boldsymbol{\Lambda}]$  and integral normalization  $c_{\boldsymbol{\mu},\boldsymbol{\Lambda}}$ .

Score matching

$$\underset{\boldsymbol{\mu}, \boldsymbol{\Lambda}}{\operatorname{argmin}} \mathbb{E} \Big[ \Big\| \underbrace{\nabla_{\boldsymbol{x}} \log \mathfrak{h}[\boldsymbol{x}; \boldsymbol{\mu}, \boldsymbol{\Lambda}]}_{\text{"score"}} - \nabla_{\boldsymbol{x}} \log \mathfrak{h}^{*}[\boldsymbol{x}] \Big\|_{2}^{2} \Big] \\ = \operatorname{argmin}_{\boldsymbol{\mu}, \boldsymbol{\Lambda}} \mathbb{E} \Big[ \Big\| \nabla_{\boldsymbol{x}} \mathfrak{g}[\boldsymbol{x}; \boldsymbol{\mu}, \boldsymbol{\Lambda}] - \nabla_{\boldsymbol{x}} \log \mathfrak{h}^{*}[\boldsymbol{x}] \Big\|_{2}^{2} \Big]$$

looks much more promising.

#### Score Matching Provides Feasible Objectives (Extremes in high dimensions: methods and scalable algorithms, 2023)

Proposition (Score Matching): Under some assumptions  $\operatorname{argmin}_{\mu,\Lambda} \mathbb{E}\Big[ \|\nabla_{\boldsymbol{x}} \,\mathfrak{g}[\boldsymbol{x};\boldsymbol{\mu},\Lambda] - \nabla_{\boldsymbol{x}} \log \mathfrak{h}^*[\boldsymbol{x}] \|_2^2 \Big]$  $= \operatorname{argmin}_{\mu,\Lambda} \mathbb{E}\Big[ \sum_{j=1}^d \partial_{x_j} \mathfrak{g}[\boldsymbol{x};\boldsymbol{\mu},\Lambda] + 2 \sum_{j=1}^d \left( \partial_{x_j}^2 \mathfrak{g}[\boldsymbol{x};\boldsymbol{\mu},\Lambda] \right)^2 \Big].$ 

# The Final Estimator Is Sparse and Scalable (Extremes in high dimensions: methods and scalable algorithms, 2023)

$$\begin{split} \widehat{\boldsymbol{\mu}}, \widehat{\boldsymbol{\Lambda}} &\in \operatorname{argmin}_{\boldsymbol{\mu}, \boldsymbol{\Lambda}} \\ &\left\{ \sum_{i=1}^{n} \left\| \left( \boldsymbol{\mu} - \mathbf{1} - (\boldsymbol{\Lambda} + \boldsymbol{\Lambda}^{\top} - \operatorname{diag}[\boldsymbol{\Lambda} \mathbf{1} + \boldsymbol{\Lambda}^{\top} \mathbf{1}]) \log[\boldsymbol{x}_{i}] \right) \otimes \mathfrak{f}_{1}[\boldsymbol{x}_{i}] \right\|_{2}^{2} \\ &+ \sum_{i=1}^{n} \left( \boldsymbol{\mu} - \mathbf{1} - (\boldsymbol{\Lambda} + \boldsymbol{\Lambda}^{\top} - \operatorname{diag}[\boldsymbol{\Lambda} \mathbf{1} + \boldsymbol{\Lambda}^{\top} \mathbf{1}]) \log[\boldsymbol{x}_{i}] \right) \otimes \mathfrak{f}_{2}[\boldsymbol{x}_{i}] \\ &- \sum_{i=1}^{n} \operatorname{trace} \left[ (\boldsymbol{\Lambda} + \boldsymbol{\Lambda}^{\top} - \operatorname{diag}[\boldsymbol{\Lambda} \mathbf{1} + \boldsymbol{\Lambda}^{\top} \mathbf{1}]) \boldsymbol{F}[\boldsymbol{x}_{i}] \right] + r \operatorname{prior}[\boldsymbol{\mu}, \boldsymbol{\Lambda}] \right\} \end{split}$$

Can be computed within seconds  $(n \approx 1000, d \approx 20, p \approx 200)$  to minutes  $(n \approx 100\,000, d \approx 100, p \approx 5000)$  on a standard laptop.

# The Estimator Satisfies Finite-Sample Guarantees Even in High Dimensions

(Extremes in high dimensions: methods and scalable algorithms, 2023)

**Theorem (Extremes):** Under reasonable assumptions, such as r large enough, it holds that  $\|\widehat{\mu} - \mu^*\|_2^2 + \|\widehat{\Lambda} - \Lambda^*\|_{\mathrm{F}}^2 \lesssim \frac{(|\mathcal{S}_{\mu^*}| + |\mathcal{S}_{\Lambda^*}|)r^2}{n}$ with high probability

with high probability.

 $\implies$  "sparsistency"

# The Statistical Proofs Are "Geometric"

(Extremes in high dimensions: methods and scalable algorithms, 2023)





#### Diffusion Models Destroy and Rebuild (Regularization can make diffusion models more efficient, 2025)



The Forward Process Injects Noise (Regularization can make diffusion models more efficient, 2025)

$$oldsymbol{x}_j \; \coloneqq \; \sqrt{1-eta_j} \, oldsymbol{x}_{j-1} + \sqrt{eta_j} \, oldsymbol{u}_{j-1}$$

with white noise  $\boldsymbol{u}_j \sim \mathcal{N}[\boldsymbol{0}_d, I_d]$  and noise schedule  $\beta_j \in (0, 1)$ 

The distribution of  $\boldsymbol{x}_j$  given an "original" sample  $\boldsymbol{x}_0$  is

$$\mathbb{Q}_{j}[oldsymbol{x}_{j} \mid oldsymbol{x}_{0}] \;=\; \mathcal{N}iggl[oldsymbol{x}_{j}; iggl(\prod_{l=1}^{j}\sqrt{1-eta_{l}}iggr)oldsymbol{x}_{0}, iggl(1-\prod_{l=1}^{j}(1-eta_{l})iggr)I_{d}iggr]$$

The Backward Process Recovers an Original (Regularization can make diffusion models more efficient, 2025)

$$oldsymbol{x}_{j-1} \;=\; rac{1}{\sqrt{1-eta_j}}oldsymbol{x}_j+eta_j \underbrace{
abla_{oldsymbol{x}_j}\log oldsymbol{\mathfrak{q}}_j[oldsymbol{x}_j]}_{ ext{"score"}}oldsymbol{)}+\sqrt{rac{eta_j}{1-eta_j}}\,oldsymbol{z}_j$$

with the marginal density  $q_j[\boldsymbol{x}_j] = \int q_j[\boldsymbol{x}_j \mid \boldsymbol{x}] p_0[\boldsymbol{x}] d\boldsymbol{x}$ and white noise  $\boldsymbol{z}_j \sim \mathcal{N}[\boldsymbol{0}_d, I_d]$ 

# Standard Estimation Suffers From the Curse of Dimensionality

(Regularization can make diffusion models more efficient, 2025)

$$\widehat{\boldsymbol{\Theta}} \in \operatorname{argmin}_{\boldsymbol{\Theta} \in \mathcal{A}} \left\{ \mathbb{E}_{\substack{j \sim \mathcal{U}_{[0,t]} \\ X_j \sim \mathbb{Q}_j}} \left[ \left\| \boldsymbol{s}_{\boldsymbol{\Theta}}[X_j, j] - \nabla_{X_j} \log \mathfrak{q}_j[X_j] \right\|_2^2 \right] \right\}$$

**Theorem ("Classical"):** It holds with high probability that  $\mathfrak{d}_{\mathrm{KL}}[\mathfrak{p}_0,\widehat{\mathfrak{p}}_0] \lesssim \frac{d^2}{t} \,.$ 

### Sparsity Breaks the Curse

(Regularization can make diffusion models more efficient, 2025)

$$(\widehat{\boldsymbol{\Theta}}, \widehat{\kappa}) \in \underset{\substack{\boldsymbol{\Theta} \in \mathcal{A}_1\\\kappa \in (0,\infty)}}{\operatorname{argmin}} \left\{ \underset{\substack{X_j \sim \mathcal{U}_{[0,t]}\\X_j \sim \mathbb{Q}_j}}{\mathbb{E}_j \sim \mathcal{U}_{[0,t]}} \left[ \left\| \kappa \boldsymbol{s}_{\boldsymbol{\Theta}}[X_j, j] - \nabla_{X_j} \log \mathfrak{q}_j[X_j] \right\|_2^2 \right] + r \kappa \right\}$$

**Theorem (Regularized):** It holds with high probability that  $\mathfrak{d}_{\mathrm{KL}}[\mathfrak{p}_0, \widehat{\mathfrak{p}}_0] \lesssim \frac{s^2}{t}.$ 

#### Regularization Makes a Visible Difference (Regularization can make diffusion models more efficient, 2025)



t = 500

#### Regularization Makes a Visible Difference (Regularization can make diffusion models more efficient, 2025)



$$t = 50$$

# Regularization Makes a Visible Difference

(Regularization can make diffusion models more efficient, 2025)



t = 500

# Regularization Makes a Visible Difference

(Regularization can make diffusion models more efficient, 2025)

#### Vanilla Diffusion



t = 50



Example:  $\widehat{\boldsymbol{\theta}} \in \underset{\boldsymbol{\theta} \in \mathbb{R}^p}{\operatorname{argmin}} \{ \| \boldsymbol{y} - X \boldsymbol{\theta} \|_2^2 + r \| \boldsymbol{\theta} \|_1 \}$ 



Example: 
$$\widehat{\boldsymbol{\theta}} \in \operatorname*{argmin}_{\boldsymbol{\theta} \in \mathbb{R}^p} \{ \| \boldsymbol{y} - X \boldsymbol{\theta} \|_2^2 + r \| \boldsymbol{\theta} \|_1 \}; \ \boldsymbol{y} = X \boldsymbol{\theta}^* + \boldsymbol{u}$$

## Goal: Select t and r simultaneously

(Balancing statistical and computational precision: a general theory and applications to sparse regression, 2023)

$$\widehat{\mathcal{F}}_{\mathcal{R}} = \{\widehat{\boldsymbol{\theta}}^r : r \in \mathcal{R}\}$$
 theoretical estimators  
 $\mathfrak{d}[\widehat{\boldsymbol{\theta}}^r, \boldsymbol{\theta}^*] \leq \mathfrak{f}[r] \quad \forall r \geq r^*$ 

 $\widetilde{\mathcal{F}}_{\mathcal{R}} = \{ \widetilde{\boldsymbol{\theta}}_t^r : r \in \mathcal{R}, t \in \{1, 2, \dots\} \} \text{ practical estimators} \\ \mathfrak{d}[\widetilde{\boldsymbol{\theta}}_t^r, \widehat{\boldsymbol{\theta}}^r] \leq \mathfrak{g}[r, t] \quad \forall r \in \mathcal{R}$ 

# A Surprisingly Concise Algorithm

(Balancing statistical and computational precision: a general theory and applications to sparse regression, 2023)

# We Have Optimal Theoretical Guarantees

(Balancing statistical and computational precision: a general theory and applications to sparse regression, 2023)

Theorem (Optimality): It holds that

$$\mathfrak{D}[\widetilde{oldsymbol{ heta}}, oldsymbol{ heta}^*] \leq 6 \mathfrak{f}[r^*].$$

The Theorem Has a Short, Elementary Proof

Consider  $\hat{r}$  with  $\tilde{\boldsymbol{\theta}}^{\hat{r}} = \tilde{\boldsymbol{\theta}}$  (we omit  $\tilde{t}$  for ease of notation)

If  $\hat{r} > r^*$ , then  $\exists r', r'' \ge r^*$  such that

$$\mathfrak{d}[\widetilde{\boldsymbol{ heta}}^{r'},\widetilde{\boldsymbol{ heta}}^{r''}] > \mathfrak{f}[r'] + \mathfrak{g}[r'] + \mathfrak{f}[r''] + \mathfrak{g}[r''].$$

The Theorem Has a Short, Elementary Proof

Consider  $\hat{r}$  with  $\tilde{\boldsymbol{\theta}}^{\hat{r}} = \tilde{\boldsymbol{\theta}}$  (we omit  $\tilde{t}$  for ease of notation)

If  $\hat{r} > r^*$ , then  $\exists r', r'' \ge r^*$  such that

$$\mathfrak{d}[\widetilde{\boldsymbol{ heta}}^{r'},\widetilde{\boldsymbol{ heta}}^{r''}] > \mathfrak{f}[r'] + \mathfrak{g}[r'] + \mathfrak{f}[r''] + \mathfrak{g}[r''].$$

But



### The Theorem Has a Short, Elementary Proof

Recall that  $\hat{r} \leq r^*$ . Hence,



# Works Like a Charm for the Lasso

(Balancing statistical and computational precision: a general theory and applications to sparse regression, 2023)





Workshop August 2025: datascienceminds.com



johanneslederer.com C LedererLab in johanneslederer O thedata.scienceguy





Workshop August 2025: datascienceminds.com



johanneslederer.com C LedererLab in johanneslederer O thedata.scienceguy





Workshop August 2025: datascienceminds.com



johanneslederer.com C LedererLab in johanneslederer O thedata.scienceguy



#### Extreme-Value Theory Is Not Extrapolation (Extremes in high dimensions: methods and scalable algorithms, 2023)



 $\lim_{t \to \infty} \mathbb{P}\left\{ X/t \in \mathcal{A} \mid \|X\| > t \right\} \sim \mathfrak{m}[\mathcal{A}]$ 

New Densities Disentangle the Parameters (Extremes in high dimensions: methods and scalable algorithms, 2023)

**Theorem:** Every Hüsler-Reiss density is of the form

$$\frac{1}{c_{\boldsymbol{\mu},\boldsymbol{\Lambda}}} \left( \prod_{j=1}^{d} \frac{1}{x_j} \right) \exp \left[ \boldsymbol{\mu}^{\top} \log[\boldsymbol{x}] - \frac{1}{2} \log[\boldsymbol{x}] (\boldsymbol{\Lambda} + \boldsymbol{\Lambda}^{\top} - \operatorname{diag}(\boldsymbol{\Lambda} \mathbf{1} + \boldsymbol{\Lambda}^{\top} \mathbf{1})) \log[\boldsymbol{x}] \right].$$

- $\boldsymbol{\mu} \in \mathbb{R}^d$ : shifts
- $\mathbf{\Lambda} \in \mathbb{R}^{d \times d}$  (upper triangular): dependences
- $c_{\mu,\Lambda}$ : normalization

#### We Develop a Score-Matching Approach That Satisfies the Boundary Conditions (Extremes in high dimensions: methods and scalable algorithms, 2023)



### The Estimator Behaves as Intended and can Be Computed Super Fast

(Extremes in high dimensions: methods and scalable algorithms, 2023)

 $\mathbf{n} = 500, \ \mathbf{d} = 20, \ \mathbf{p} = 190$   $r \qquad 1000\sqrt{\frac{\log[d]}{n}} \qquad 100\sqrt{\frac{\log[d]}{n}} \qquad 10\sqrt{\frac{\log[d]}{n}} \qquad 1\sqrt{\frac{\log[d]}{n}} \qquad 0.1\sqrt{\frac{\log[d]}{n}} \qquad 0.1\sqrt{\frac{\log[d]}{n}} \qquad 0$   $\|\widehat{\mathbf{A}}\|_0/(d(d-1)/2) \qquad 33.9\% \pm 2.7\% \qquad 63.1\% \pm 3.7\% \qquad 96.0\% \pm 1.3\% \qquad 99.7\% \pm 0.5\% \qquad 100.0\% \pm 0.1\% \qquad 100.0\% \pm 0.1\% \qquad 100.0\% \pm 0.1\%$   $t_{\text{pre}} + t_{\text{opt}} \ (\text{s}) \qquad \qquad < 0.01 \ s \pm 0.01 \ s \qquad + \qquad 0.83 \ s \pm 0.11 \ s$ 

#### Plenary talk at ICORS, May 21, 2025

Image credits: River Lambro May 2024: LaPress/REX/Shutterstock realDonaldTrump on Truth Social