

LANDSCAPE COMPLEXITY FOR THE EMPIRICAL RISK OF GENERALIZED LINEAR MODELS

Antoine Maillard, Gérard Ben Arous & Giulio Biroli

Mathematical and Scientific Machine Learning 2020



Universität Basel – October 20th 2021

GENERALIZED LINEAR MODELS

Goal: Recover $\mathbf{X}^* \in \mathbb{R}^n$ from $\{\Phi_\mu, Y_\mu\}_{\mu=1}^m$:

Observations $Y_\mu \in \mathbb{R}$

$$Y_\mu \sim P_{\text{out}}\left(\cdot \mid \frac{1}{\sqrt{n}} \sum_{i=1}^n \Phi_{\mu i} X_i^*\right) \quad \mu \in \{1, \dots, m\}$$

Channel: non-linearity
+ possible noise

Sensing matrix

Many examples: compressed sensing, perceptron learning, phase retrieval, ...

Goal: Fundamental limits of inference models with random input data in the typical case and in high dimension.

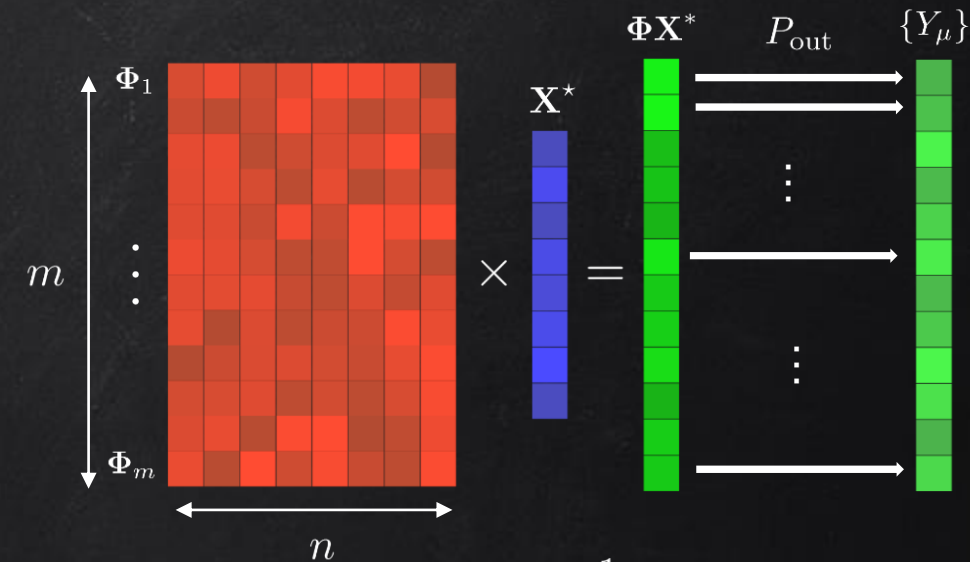


Different from “worst-case” injectivity studies. (e.g. [Bandeira&al ‘14])

“High-dimensional” limit

Number of parameters $n \rightarrow \infty$ + Number of data $m \rightarrow \infty$

In this presentation: $m/n \rightarrow \alpha > 0$ (sampling ratio).



$$Y_\mu = \frac{1}{n} |(\Phi \mathbf{X}^*)_\mu|^2$$

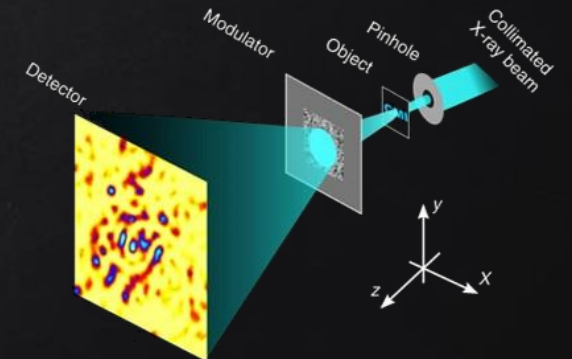


Image credits: [Zhang&al 16]

ESTIMATE X^* ?

Bayesian estimators: Leverage the posterior distribution $\mathbb{P}(\mathbf{x}|\mathbf{Y}, \Phi) = \frac{\mathbb{P}(\mathbf{x})\mathbb{P}(\mathbf{Y}|\mathbf{x}, \Phi)}{\mathbb{P}(\mathbf{Y}|\Phi)}$

➤ *Maximum A Posteriori* $\hat{\mathbf{X}}_{\text{MAP}} \equiv \arg \max_{\mathbf{x}} \mathbb{P}(\mathbf{x}|\mathbf{Y})$

➤ *Minimal Mean Squared Error* $\hat{\mathbf{X}}_{\text{MMSE}} \equiv \arg \min_{\mathbf{x}} \left\{ \mathbb{E}_{\mathbf{Y}} \int d\mathbf{x}' \mathbb{P}(\mathbf{x}'|\mathbf{Y}) \|\mathbf{x} - \mathbf{x}'\|^2 \right\} = \mathbb{E}_{\mathbf{Y}} \int d\mathbf{x} \mathbb{P}(\mathbf{x}|\mathbf{Y}) \mathbf{x}$

M-estimation:

$$\hat{\mathbf{X}} \equiv \arg \min_{\mathbf{x}} \sum_{\mu=1}^m \rho(\mathbf{x}, \Phi_{\mu}, y_{\mu})$$

$$y_{\mu} = \varphi \left[\frac{1}{\sqrt{n}} \sum_{i=1}^n \Phi_{\mu i} X_i^* \right]$$

(to simplify)

Example: Empirical Risk Minimizer, e.g. square loss

$$\hat{\mathbf{X}} \equiv \arg \min_{\mathbf{x}} \left\{ \frac{1}{m} \sum_{\mu=1}^m \left[y_{\mu} - \varphi \left(\frac{1}{\sqrt{n}} \Phi \cdot \mathbf{x} \right) \right]^2 \right\}$$

[This presentation](#)

TOPOLOGY OF HIGH-DIMENSIONAL LANDSCAPES

Empirical risk of a noiseless GLM

$$L_2(\mathbf{x}) = \frac{1}{2m} \sum_{\mu=1}^m \left[\varphi(\Phi_{\mu} \cdot \mathbf{x}) - \varphi(\Phi_{\mu} \cdot \mathbf{X}^*) \right]^2 \quad \begin{array}{l} \mathbf{X}^* \in \mathbb{R}^n, \|\mathbf{X}^*\|^2 = 1 \\ \mathbf{x} \in \mathbb{R}^n, \|\mathbf{x}\|^2 = 1 \end{array}$$

i.i.d. Gaussian input data

In many nonconvex problems: one can provably find regimes in which the optimization landscape is 'easy' (matrix decomposition, tensor factorization, neural nets...) [Soudry&al '16; Ge&al '16; Ge&Ma '17; ...]

In practice, local optimization algorithms work far beyond these regimes !

WHY ?

- The bounds on the simplicity of the landscape are not tight enough ?
- Optimization algorithms work in the "hard" regime (i.e. many spurious minima) ?
- Analyze the topological transition, and characterize the 'hard' regime ?



“COMPLEX” LANDSCAPES

Characterize complexity of a landscape

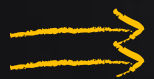
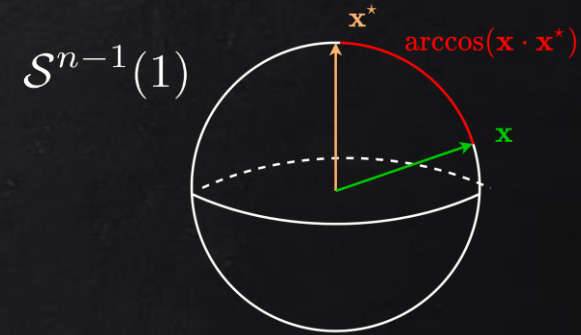


Count critical points with fixed “**energy**” $L_2(\mathbf{x})$ and **overlap** $q = \mathbf{X}^* \cdot \mathbf{x}$?

$$\text{Crit}_*(B, Q) \equiv \sum_{\mathbf{x}: \text{grad } L_2(\mathbf{x})=0} \mathbb{1}\{L_2(\mathbf{x}) \in B, \mathbf{x} \cdot \mathbf{X}^* \in Q\}$$

Random variable
(randomness of the data)

- Typically of size $e^{\Theta(n)}$.
- Strongly fluctuating!



- Mean value: **annealed complexity** $\Sigma_*^{(\text{an.})}(B, Q) \equiv \lim_{n \rightarrow \infty} \frac{1}{n} \ln \mathbb{E} \text{Crit}_*(B, Q)$
- Typical value: **quenched complexity** $\Sigma_*^{(\text{qu.})}(B, Q) \equiv \lim_{n \rightarrow \infty} \frac{1}{n} \mathbb{E} \ln \text{Crit}_*(B, Q)$

⚠ Quenched \neq Annealed!
(Few exceptions [Subag '17])

Complexity at a fixed **index** $\text{Crit}_k(B, Q) \equiv \sum_{\mathbf{x}: \text{grad } L_2(\mathbf{x})=0} \mathbb{1}\{\text{i[Hess } L_2(\mathbf{x})] = k, L_2(\mathbf{x}) \in B, \mathbf{x} \cdot \mathbf{X}^* \in Q\}$

COUNTING CRITICAL POINTS

$$L_2(\mathbf{x}) = \frac{1}{2m} \sum_{\mu=1}^m \left[\varphi(\Phi_{\mu} \cdot \mathbf{x}) - \varphi(\Phi_{\mu} \cdot \mathbf{X}^*) \right]^2$$

Let $f : \mathcal{S}^{n-1}(1) \rightarrow \mathbb{R}$ a smooth function.

Count $N_f(\mathbf{u}) \equiv \#\{\mathbf{x} \in \mathcal{S}^{n-1}(1) : \text{grad } f(\mathbf{x}) = \mathbf{u}\} ?$

?

$$N_f(\mathbf{u}) = \int_{\text{grad } f(\mathcal{S}^{n-1}(1))} d\mathbf{v} \delta(\mathbf{v} - \mathbf{u}) \longrightarrow = \int_{\mathcal{S}^{n-1}(1)} d\mathbf{x} \delta(\text{grad } f(\mathbf{x}) - \mathbf{u}) |\det \text{Hess } f(\mathbf{x})|$$

Weak sense

Area formula

Test function

$$\int d\mathbf{u} g(\mathbf{u}) N_f(\mathbf{u}) = \int_{\mathcal{S}^{n-1}(1)} d\mathbf{x} g(\text{grad } f(\mathbf{x})) |\det \text{Hess } f(\mathbf{x})| \quad [\text{Federer '59}]$$

⚠ $L_2(\mathbf{x})$ is also a random function !

f is a.s. Morse (all critical points are non-degenerate)

+ technical regularity properties (hard for non-Gaussian functions)...

Strong sense at $\mathbf{u} = 0$

$$\mathbb{E} \text{Crit}_k(f) = \int_{\mathcal{S}^{n-1}(1)} \sigma(d\mathbf{x}) \varphi_{\text{grad } f(\mathbf{x})}(0) \times \mathbb{E}[|\det \text{Hess } f(\mathbf{x})| | \text{grad } f(\mathbf{x}) = 0; i(\text{Hess } f(\mathbf{x})) = k]$$

Density of the (random) gradient taken in 0

THE KAC-RICE FORMULA

Theorem (Kac-Rice)

e.g. [Adler&Taylor '09]

$f : \mathcal{S}^{n-1}(1) \rightarrow \mathbb{R}$ is a smooth random function that is a.s. Morse. Then:

$$\mathbb{E} \text{Crit}_k(f) = \int_{\mathcal{S}^{n-1}(1)} \sigma(d\mathbf{x}) \varphi_{\text{grad } f(\mathbf{x})}(0) \times \mathbb{E} \left[\underbrace{|\det \text{Hess } f(\mathbf{x})| \mid \text{grad } f(\mathbf{x}) = 0; i(\text{Hess } f(\mathbf{x})) = k} \right]$$

- Random differential geometry \longrightarrow **Random matrix theory**.
- Many possible refinements: fix the value of $f(\mathbf{x})$, higher-order moments, ...

- Takeaways:
- **Distribution of $\{\text{Hess } f(\mathbf{x}) \mid \text{grad } f(\mathbf{x}) = 0\}$** : intractable for “generic” functions !
 - To compute $\mathbb{E} \text{Crit}_k(f)$: we need the **large deviations of the k-th largest eigenvalue of the Hessian**.

\longrightarrow Applications limited to Gaussian random functions

Pure p-spin and variants

$$f(\mathbf{x}) = \sum_{i_1, \dots, i_p} \underbrace{J_{i_1, \dots, i_p}}_{\mathcal{N}(0,1)} x_{i_1} \cdots x_{i_p}$$

[Bray&Moore '80, Crisanti&al '95, Fyodorov&al '07, Auffinger&al '13, Ros&al '19, Belius&al '21,]

ANNEALED COMPLEXITY

$$\alpha = m/n = \Theta(1)$$

$$L_2(\mathbf{x}) = \frac{1}{2m} \sum_{\mu=1}^m \left[\varphi(\Phi_{\mu} \cdot \mathbf{x}) - \varphi(\Phi_{\mu} \cdot \mathbf{X}^*) \right]^2 \xrightarrow{\text{simplification}} L_1(\mathbf{x}) = \frac{1}{m} \sum_{\mu=1}^m \varphi(\Phi_{\mu} \cdot \mathbf{x})$$

Theorem & proof transfer to $L_2(\mathbf{x})$.

First exact high-dimensional result obtained with Kac-Rice for non-Gaussian functions!

Theorem

$$\Sigma_{\star}^{(\text{an})}(B) = \frac{1 + \ln \alpha}{2} + \sup_{\substack{\nu \in \mathcal{M}_1^+(\mathbb{R}) \\ \int \nu(dt) \varphi(t) \in B}} \left[-\frac{1}{2} \ln \left\{ \int \nu(dt) \varphi'(t)^2 \right\} - \underbrace{\alpha H(\nu | \mathcal{N}(0, 1))}_{\text{Relative entropy}} + \kappa_{\alpha}(\nu) \right]$$

Involved function: related to the logarithmic potential of the asymptotic spectral measure of \mathbf{zDz}^T/m if $\mathbf{z} \in \mathbb{R}^{n \times m}$ is a Gaussian i.i.d. matrix and $D_{\mu} = \varphi''(y_{\mu})$ with $y_{\mu} \stackrel{\text{i.i.d.}}{\sim} \nu$.

- Term $\kappa_{\alpha}(\nu) \implies$ Analytically very hard variational problem!
- Generic result, applies to mixture of Gaussians, binary classification, ...

SKETCH OF PROOF

$$L_1(\mathbf{x}) = \frac{1}{m} \sum_{\mu=1}^m \varphi(\Phi_{\mu} \cdot \mathbf{x})$$

Kac-Rice formula \implies Hessian conditioned by zero gradient.

Main idea: Condition everything by the i.i.d. Gaussian random variables $y_{\mu} \equiv \Phi_{\mu} \cdot \mathbf{x}$

$\mathbb{E}_{\Phi}[\cdots] = \mathbb{E}_{\mathbf{y}}\mathbb{E}[\cdots | \mathbf{y}]$ Under this conditional distribution, and under the gradient being zero:

$$\text{Hess } L_1(\mathbf{x}) \stackrel{\text{d}}{=} \frac{1}{m} \sum_{\mu=1}^m \varphi''(y_{\mu}) \mathbf{z}_{\mu} \mathbf{z}_{\mu}^{\top} + t(\mathbf{y}) \mathbb{1}_n (+ \text{finite rank term})$$

\mathbf{z}_{μ} : i.i.d. standard Gaussian vectors

“Generalized” version of a sample covariance matrix

- We prove fast enough concentration of $\ln |\det \text{Hess}|$ as a function of $\{y_{\mu}\}_{\mu=1}^m$: $\mathbb{E} |\det \text{Hess}| \simeq e^{\mathbb{E} \ln |\det \text{Hess}|}$
- The expectation only depends on the empirical distribution $\nu_{\mathbf{y}} \equiv (1/m) \sum_{\mu=1}^m \delta_{y_{\mu}}$: $\mathbb{E} \ln |\det \text{Hess}| \simeq m \kappa_{\alpha}(\nu_{\mathbf{y}})$
- We use **Sanov's theorem**: the law of $\nu_{\mathbf{y}}$ satisfies large deviations with rate function: $I(\nu) = \alpha H(\nu | \mathcal{N}(0, 1))$
- Kac-Rice formula and Varadhan's lemma $\implies \Sigma_{\star}^{(\text{an})} = \sup_{\nu \in \mathcal{M}_1^{+}(\mathbb{R})} [\kappa_{\alpha}(\nu) + \underline{G(\nu)} - \alpha H(\nu | \mathcal{N}(0, 1))]$ \blacksquare

(Gradient density term in Kac-Rice)

EXTENSION 1: THE QUENCHED COMPLEXITY

Replica trick

$$\mathbb{E}[\ln \text{Crit}(f)] = \lim_{r \downarrow 0} \frac{\mathbb{E} \text{Crit}(f)^r - 1}{r}$$

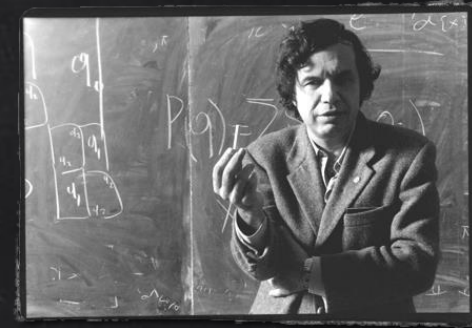


- Compute $\mathbb{E}[\text{Crit}(f)^r]$ for $r \in \mathbb{N}$.
- Perform an analytical continuation to consider $r \downarrow 0$.

Replica theory is a prolific field of statistical physics



Physics
2021



Giorgio Parisi

Notably for describing the possible breaking of the symmetry between replicas in disordered systems

Refined Kac–Rice for $\mathbb{E}[\text{Crit}(f)^r]$



Replica trick



$$\Sigma_{\star}^{(\text{qu.})}(B, Q) = \sup_{\nu \in \mathcal{M}_1(\mathbb{R})} [\dots]$$

Closed formula for the quenched complexity

- Heuristic result, under a replica-symmetric ansatz.
- Drawback: Hard to solve, and even to interpret all terms so far...

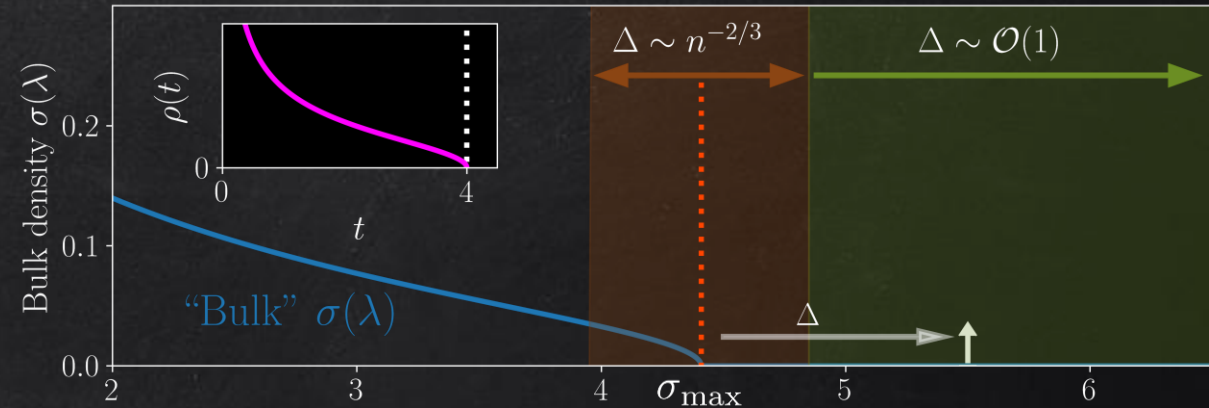
EXTENSION 2: COUNTING MINIMA?

Goal: count only **local minima**, not all critical points !

Kac-Rice + restriction to local minima



Large deviations: $\frac{1}{n} \ln \mathbb{P}[\lambda_{\min}(\text{Hess } f) \simeq x] \quad ??$



LDP for smallest/largest eigenvalue of “generalized covariance matrix” $\mathbf{M} = \frac{1}{m} \sum_{\mu=1}^m \rho_{\mu} \mathbf{z}_{\mu} \mathbf{z}_{\mu}^{\dagger}$?

Solved in [\[A.M. 21\]](#)

Real/complex

$$x \geq \sigma_{\max} : \lim_{n \rightarrow \infty} \left\{ \frac{1}{n} \ln \mathbb{P}(\lambda_{\max}(\mathbf{M}) \simeq x) \right\} = -\frac{\beta}{2} \int_{\sigma_{\max}}^x [\overline{G}_{\sigma}(u) - G_{\sigma}(u)] du$$

$(G_{\sigma}(x), \overline{G}_{\sigma}(x))$ solutions to

$$x = \frac{1}{G} + \alpha \int dt \rho(t) \frac{t}{\alpha - tG}$$

Dyson/Marchenko–Pastur equation

Using a technique based on a tilting of the measure.

[Biroli&Guionnet’20, Belinschi&al ’20, Guionnet&al ’20, Husson ’20, Augeri&al ’21]...

SUMMARY & OUTLOOK

- ❖ First exact high-dimensional result obtained with Kac–Rice for non-Gaussian functions !
- ❖ Both annealed and quenched computations for the total complexity of critical points.
- ❖ Generalizes to other models: mixture of two Gaussians, binary linear classification... → Neural networks ?

Physical discussion is still hard to reach. Many problems ahead:

- Numerically solve the variational problem? Hints in [A.M.&al '20]...

Sign of the complexity given α, φ ? ...

- Count local minima ? We need a LDP for $\lambda_{\min}(\text{Hess } L(\mathbf{x}))$...

→ Obtained in [A.M., EPL 2021] ! To be continued...

- Discrete systems ?

- ❖ TAP approach ?

- ❖ Recent algorithmic progress on F-RSB spin glasses [Subag '21, Montanari '21, El Alaoui & al '20].

