

LANDSCAPE COMPLEXITY FOR THE EMPIRICAL RISK OF GENERALIZED LINEAR MODELS

A.M, Gérard Ben Arous, Giulio Biroli

arXiv:1912.02143



ENS

PSL



NYU

1

LANDSCAPE COMPLEXITY

Main goal: Understand the landscape of the empirical risk in statistical estimation

To analyze local optimization algorithms
(e.g. gradient descent and stochastic variants)

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

In practice, algorithms work far beyond these regimes ! Why ?

→ The bounds on the simplicity of the landscape are not tight enough ?

→ Optimization algorithms can work in a 'hard' regime (i.e. many spurious local minima) ?

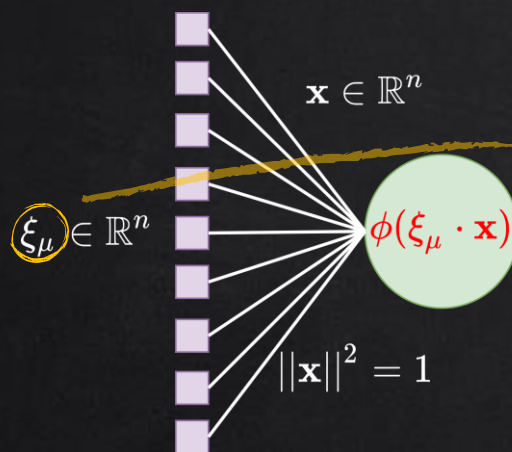
→ Can we analyze the topological transition in the landscape, and characterize the 'hard' regime ?

2

GENERALIZED LINEAR MODELS

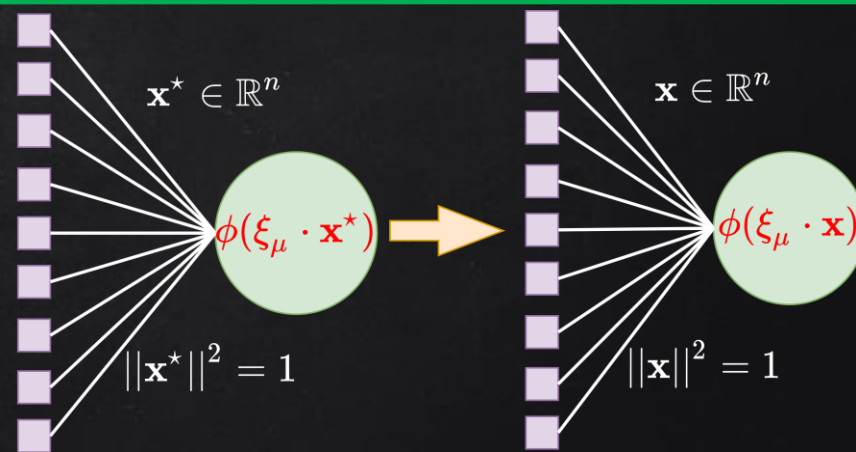
“Perceptron” energy

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

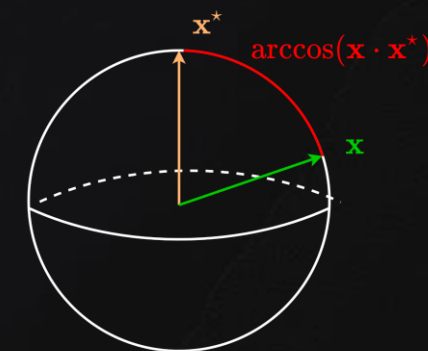


Squared loss of a generalized linear model

$$L_2(\mathbf{x}) = \frac{1}{2m} \sum_{\mu=1}^m [\phi(\xi_{\mu} \cdot \mathbf{x}) - \phi(\xi_{\mu} \cdot \mathbf{x}^*)]^2$$



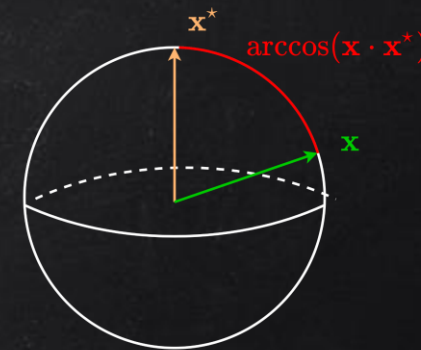
“Teacher-student” setup



- High dimensional limit $m, n \rightarrow \infty$ with $\alpha = m/n = \Theta(1)$.
- Can we count the number of critical points of these functions with a definite “energy” $L(\mathbf{x})$ and overlap $q = \mathbf{x} \cdot \mathbf{x}^*$?

1,

LANDSCAPE COMPLEXITY



Number of critical points with “energy” $L(\mathbf{x})$ and overlap $q = \mathbf{x} \cdot \mathbf{x}^*$?

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

Random variable
(randomness of the data)

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



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



In general, quenched and annealed are very different ! (Few exceptions [Subag '17])

3

THE KAC-RICE FORMULA

For a 1D function $f(x)$, one would like to write: $\#\{x \text{ s.t. } f(x) = 0\} \stackrel{?}{=} \int \delta(f(x)) |f'(x)| dx$

The **Kac-Rice formula** makes this intuition precise \implies mean number of critical points of random $f : \mathbb{R}^n \rightarrow \mathbb{R}$:

$$\mathbb{E} \text{Crit}(f) = \int_{\|\mathbf{x}\|^2=1} d\sigma(\mathbf{x}) \underbrace{\varphi_{\text{grad } f(\mathbf{x})}(0)}_{\text{Density of the gradient taken in 0}} \mathbb{E} [\underbrace{|\det \text{Hess } f(\mathbf{x})|}_{\text{random matrix theory}} | \text{grad } f(\mathbf{x}) = 0]$$

Density of the gradient taken in 0

- Turns a random differential geometry problem into a **random matrix theory** problem.
- One can also fix the value of $f(\mathbf{x})$, get higher-order moments...

We need the distribution of the Hessian, conditioned by the gradient being zero : intractable for “generic” functions !

\longrightarrow Applications limited so far to Gaussian random functions

Pure p-spin and variants

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

- Large literature on the applications to Gaussian models, in physics and mathematics [Bray&Moore '80, Crisanti&al '95, Fyodorov&al '07, Auffinger&al '13, Ros&al '19,]. See [Adler&Taylor '07, Azais&Wschebor '09] for a mathematical introduction to Kac-Rice.

4

MAIN RESULT

A first high-dimensional exact result with the Kac-Rice formula for a non-Gaussian random function !

Theorem (for L_1) $\Sigma_{\star}^{(\text{an.})}(B) = \frac{1 + \ln \alpha}{2} + \sup_{\substack{\nu \in \mathcal{M}_1^+(\mathbb{R}) \\ \int \nu(dt) \phi(t) \in B}} \left[-\frac{1}{2} \ln \left\{ \int \nu(dt) \phi'(t)^2 \right\} - \underbrace{\alpha H(\nu | \mathcal{N}(0, 1))}_{\text{Relative entropy}} + \kappa_{\alpha}(\nu) \right].$

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

There is a similar theorem for L_2 : $\Sigma_{\star}^{(\text{an.})}(B, Q) = \sup_{q \in Q} \sup_{\nu \in \mathcal{M}_1^+(\mathbb{R}^2)} [\dots]$

- Because of $\kappa_{\alpha}(\nu)$: very hard to solve in general ! (apart from trivial activation functions)
- We derive a heuristic closed formula of $\kappa_{\alpha}(\nu)$, based on [Marchenko&Pastur '67, Silverstein&Bai '95], which leads to simpler scalar fixed point equations (for L_1 and L_2) \longrightarrow Numerical solutions ? (Ongoing work)

5

OVERVIEW OF THE PROOF

We focus on L_1 (same method for L_2)

Kac-Rice formula \Rightarrow We need to study the Hessian, conditioned by the gradient being zero.

Main idea: Condition (gradient, Hessian) by the i.i.d. Gaussian random variables $y_\mu \equiv \xi_\mu \cdot \mathbf{x}$

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

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

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

“Generalized” version of a sample covariance matrix
($\phi''(y_\mu)$ can be negative).

- 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** $\nu_y \equiv (1/m) \sum_{\mu=1}^m \delta_{y_\mu}$: $\mathbb{E} \ln |\det \text{Hess}| = nF(\nu_y)$
- We use **Sanov's theorem**: ν_y satisfies large deviations with rate function: $I(\nu) = \alpha H(\nu | \mathcal{N}(0, 1))$
- **Kac-Rice formula and Varadhan's lemma**: $\Sigma^{(\text{an.})} = \sup_{\nu} [F(\nu) + \underline{G(\nu)} - \alpha H(\nu | \mathcal{N}(0, 1))]$

(Gradient density term in Kac-Rice)

6

CONCLUSION & PERSPECTIVES

Additional results

- We derive a similar **closed formula for the quenched complexity** $\Sigma_{\star}^{(\text{qu.})}(B, Q)$. The derivation is based on the Kac-Rice formula for higher order moments and the non-rigorous replica method [Parisi&al '87, Ros&al '19].
- We generalize our results to other models : mixture of two Gaussians, binary linear classification, a simple unsupervised learning problem \longrightarrow Can we generalize to neural networks ? (Open question)

Some future directions

- We present a first theoretical step : can we solve the variational problem and obtain numerical curves ? (Ongoing work, cf the "Main Result" slide)
- We only derived formulas for the **total number of critical points**. Can we count **only the local minima** ? To do so, we need the large deviations of the lowest eigenvalue of the Hessian.
 \longrightarrow We believe we can ! (Ongoing work [A.M., to appear in 2020])

THANK YOU !