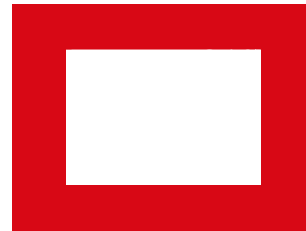


On the loss function landscape in the simplest constrained least-square optimization

Yan V Fyodorov

Department of Mathematics



Project supported by the EPSRC grant EP/N009436/1

"XIV Brunel-Bielefeld RMT Workshop " London, 14th of December 2018

¹Based on: [YVF](#) & [Rashel Tublin](#) , under preparation.

Background:

The simplest optimization problem of the **least-square** type on the sphere $x \in \mathbb{R}^N$; $x^2 = \text{const}$ arises in the **Multiple Factor Data Analysis** and is known as the **Oblique Procrustes Problem** :

For a given pair of $M \times N$ matrices A and B and such $N \times N$ matrix X that the equality $B = AX$ holds as close as possible and columns $x_i \in \mathbb{R}^N$; $i = 1; \dots; N$ are of unit length.

For $M > N$ this system of linear equations is overcomplete and a solution can be found separately for each column x by minimizing the **loss/cost function**

$$H(x) = \frac{1}{2} \|Ax - b\|^2 := \frac{1}{2} \sum_{k=1}^M \left(\sum_{j=1}^N A_{kj} x_j - b_k \right)^2 ; \quad x^2 = \text{const}$$

The problem was first analysed in that setting by **M. W. BROWNE** in 1967, and then independently by numerical mathematicians (e.g. **W. GANDER** 1981) who used the **Lagrange multiplier** to take care of the spherical constraint. Introducing the Lagrangian $L(x) = H(x) - \frac{\lambda}{2}(x^2 - \text{const})$, with real λ being the Lagrange multiplier, the stationary conditions $\nabla L(x) = 0$ yields linear system:

$$A^T [Ax - b] = \lambda x ; \quad x = (A^T A - \lambda I_N)^{-1} A^T b$$

Setting of the problem:

The spherical constraint $x^2 = N$ yields the equation for λ in the form:

$$b^T A \frac{1}{(A^T A - I_N)^2} A^T b = N$$

which is equivalent to a polynomial equation of degree $2N$ in λ . Each **real** solution for the **Lagrange multiplier** λ_i corresponds to a **stationary point** x_i of the loss function $H(x) = \frac{1}{2} \|Ax - b\|^2$ on the sphere $x^2 = N$ and one can show that the order $\lambda_1 < \lambda_2 < \dots < \lambda_N$ implies $H(x_1) < H(x_2) < \dots < H(x_N)$. Thus the **minimal loss** is given by $E_{\min} = H(x_1)$.

Our goal: To count the **stationary points** via the Lagrange multipliers

$$\lambda_i; i = 1, \dots, N \in \mathbb{R}$$

and eventually find the **minimal loss** E_{\min} after assuming the entries A_{kj} of $M \times N$; $M > N$ matrix A to be i.i.d. normal real variables such that $A^T A = W$ is $N \times N$

Wishart with the probability density

$$P_{N;M}(W) = C_{N;M} e^{-\frac{N}{2} \text{Tr} W} (dW)^{\frac{M-N-1}{2}}$$

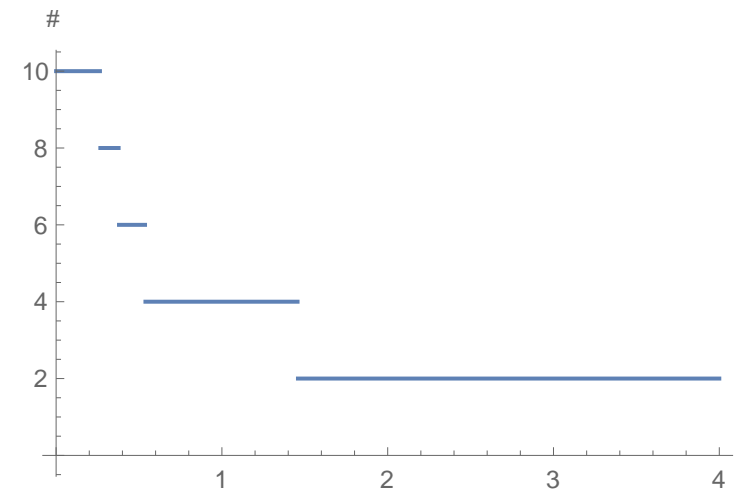
We will also assume for convenience that the vector b is normally distributed: $b =$

with $\sigma > 0$ and the components of $b = (\lambda_1, \dots, \lambda_M)^T$ are mean zero standard normals.

Qualitative considerations:

The equation for the Lagrange multiplier can be conveniently written in terms of N nonzero eigenvalues $s_1; \dots; s_N$ of $M \times M$ matrix $W^{(a)} = AA^T$ and the associated eigenvectors v_i :

$$\sum_{i=1}^N \frac{s_i}{(s_i)^2} (v_i^T v_i)^2 = \frac{N}{2}$$



Case $N = 5$

Counting Lagrange multipliers via the Kac-Rice formula:

The number $N_{st}[a; b]$ of real solutions of the equation $A^T [Ax - b] = 0$ on the sphere $x^2 = N$ such that $2 \in [a; b]$ can be counted by employing the **Kac-Rice** type formula

$$N_{st}[a; b] = \int_a^b \int_{\mathbb{R}^N} \det \begin{pmatrix} A^T A & I_N \\ x^T & 0 \end{pmatrix} dx$$

Using Gaussianity of both the matrix entries $A_{ij} \sim N(0; 1)$ and the vector components $b \sim N_M(0; I_M)$ and introducing the parameter $\beta = \frac{1}{2} \ln(1 + \frac{2}{N})$ one can eventually find the mean number of solutions as

$$E[N_{st}[a; b]] = \int_a^b \dots$$

Counting Lagrange multipliers via the Kac-Rice formula :

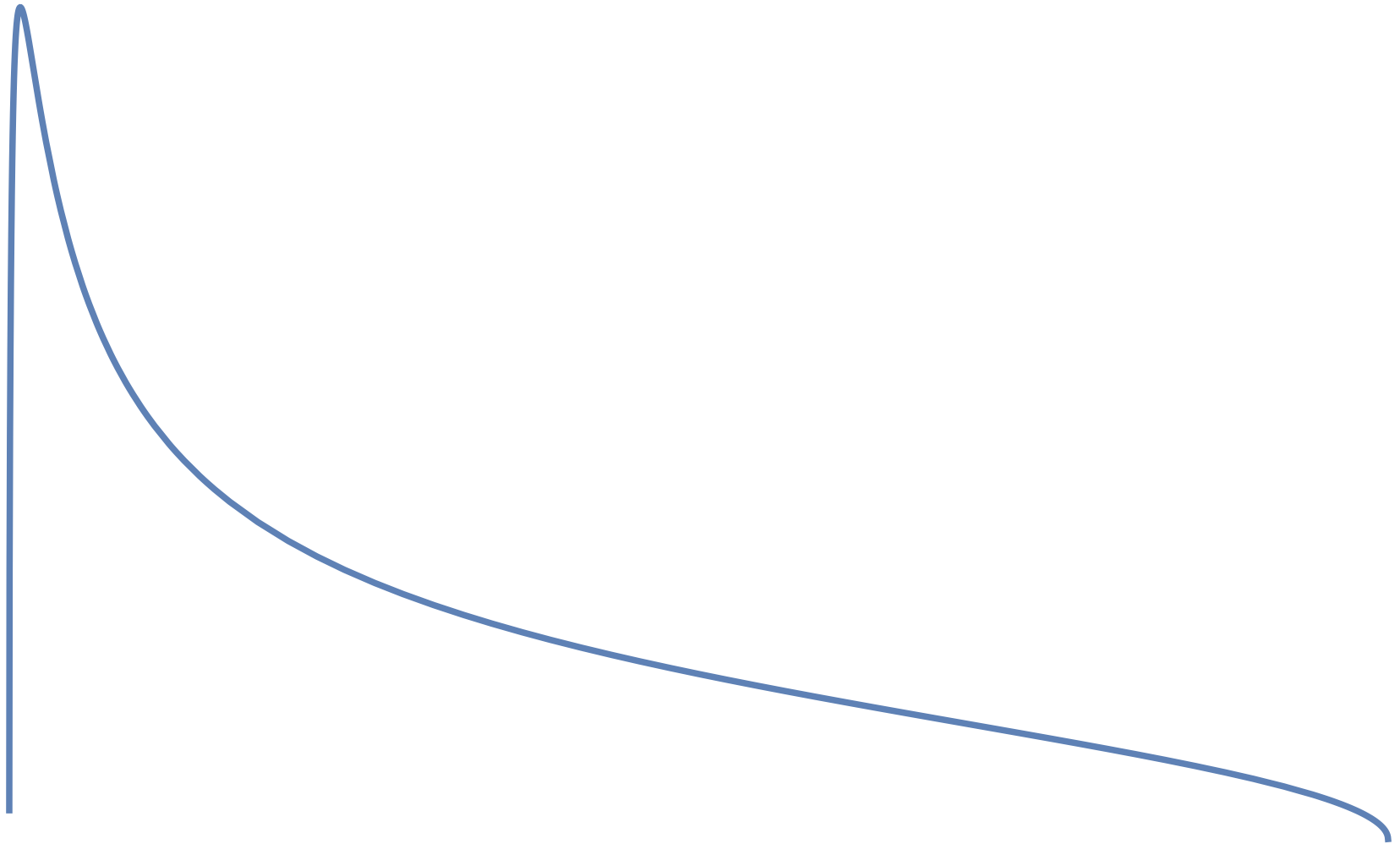
For negative values of the Lagrange multiplier we have instead:

$$p(\lambda < 0) = \frac{N! N^{(M+N-3)/2}}{2^{(M+N-3)/2} \binom{N}{2} \binom{M}{2}} \frac{1}{p} \frac{e^{-(M+N-1)/2}}{\sinh} e^{-\frac{1}{2} N \lambda^2}$$

"Bulk" Scaling Regime: extensive number of stationary points:

As N & $M \rightarrow 1$ in such a way that $1 < \frac{M}{N} < 1$ the number of stationary points in the loss function landscapes shows **three different regimes** depending on the magnitude of the parameter $\beta = \frac{1}{2} \ln(1 + \frac{M}{N})$.

"Bulk" Scaling Regime: for small enough β $\frac{M}{N} \rightarrow 1$ so that



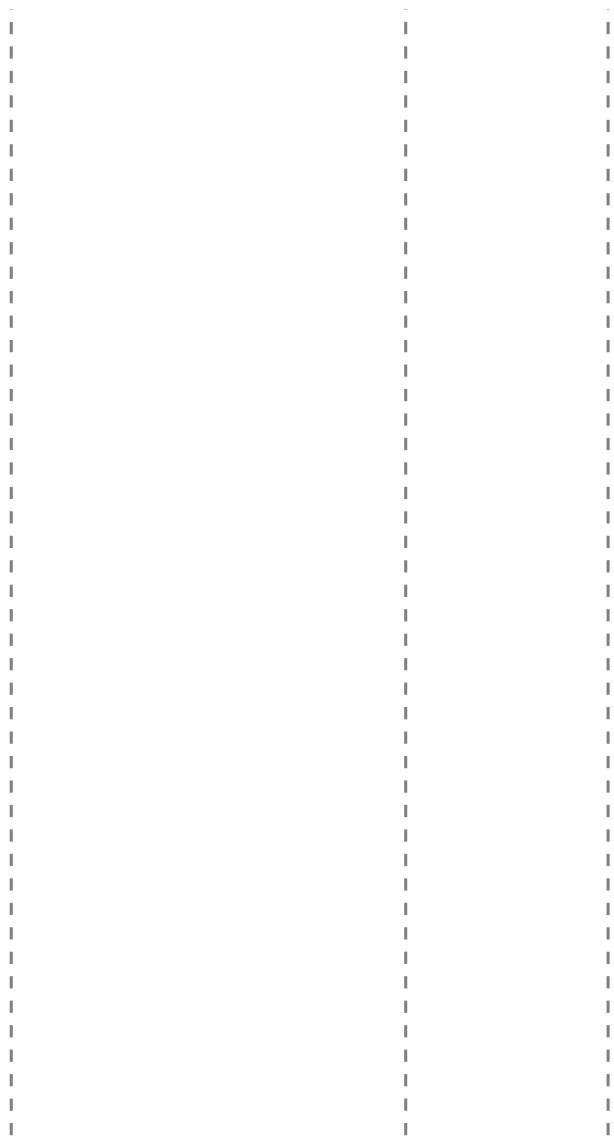
Evolution of the density ρ_B () in the 'bulk scaling' regime.

"Edge" Scaling Regime: finite number of stationary points:

The density of Lagrange multipliers for $N^{1=3}$ is dominated by the vicinities of the spectral edges

$$j \quad s \quad j \quad N^{2=3} \quad \frac{4s^2}{s_+ \quad s} \quad 1=3$$

where the



Counting stationary points in the edge regime.

Large Deviations for the smallest Lagrange multiplier:

For large $N \rightarrow \infty$, $\text{with } 1 < \beta = M/N < \infty$ and $\text{with } \beta^2 > 0$ the probability density for the smallest Lagrange multiplier λ_{\min} has the **Large Deviation** form:

$$p(\lambda_{\min} < s) \sim e^{-\frac{N}{2} I(s)}; \quad I(s) = L_1(s) + L_2(s) + \frac{(\beta+1)}{2} \ln(1 + \beta^2),$$

where $s = \left(\beta - \frac{1}{\beta}\right)$ is the 'Marchenko-Pastur' left edge and for $\beta = \frac{(\beta-1)^2}{2}$

$$L_1(s) = \left(\beta - \frac{1}{\beta}\right) \frac{\beta - \frac{1}{\beta}}{2 + \beta^2} \ln \left(\frac{\beta - \frac{1}{\beta}}{2 + \beta^2} \right) + \frac{\beta - \frac{1}{\beta}}{2 + \beta^2} \frac{\beta - \frac{1}{\beta}}{(\beta - 1)}$$

$$L_2(s) = \frac{\beta - \frac{1}{\beta}}{(\beta - s)(\beta - s_+)} - 2 \ln \frac{(\beta + 1 + \frac{\beta - \frac{1}{\beta}}{2 + \beta^2})}{2\beta - \frac{\beta - \frac{1}{\beta}}{(\beta - s)(\beta - s_+)}}$$

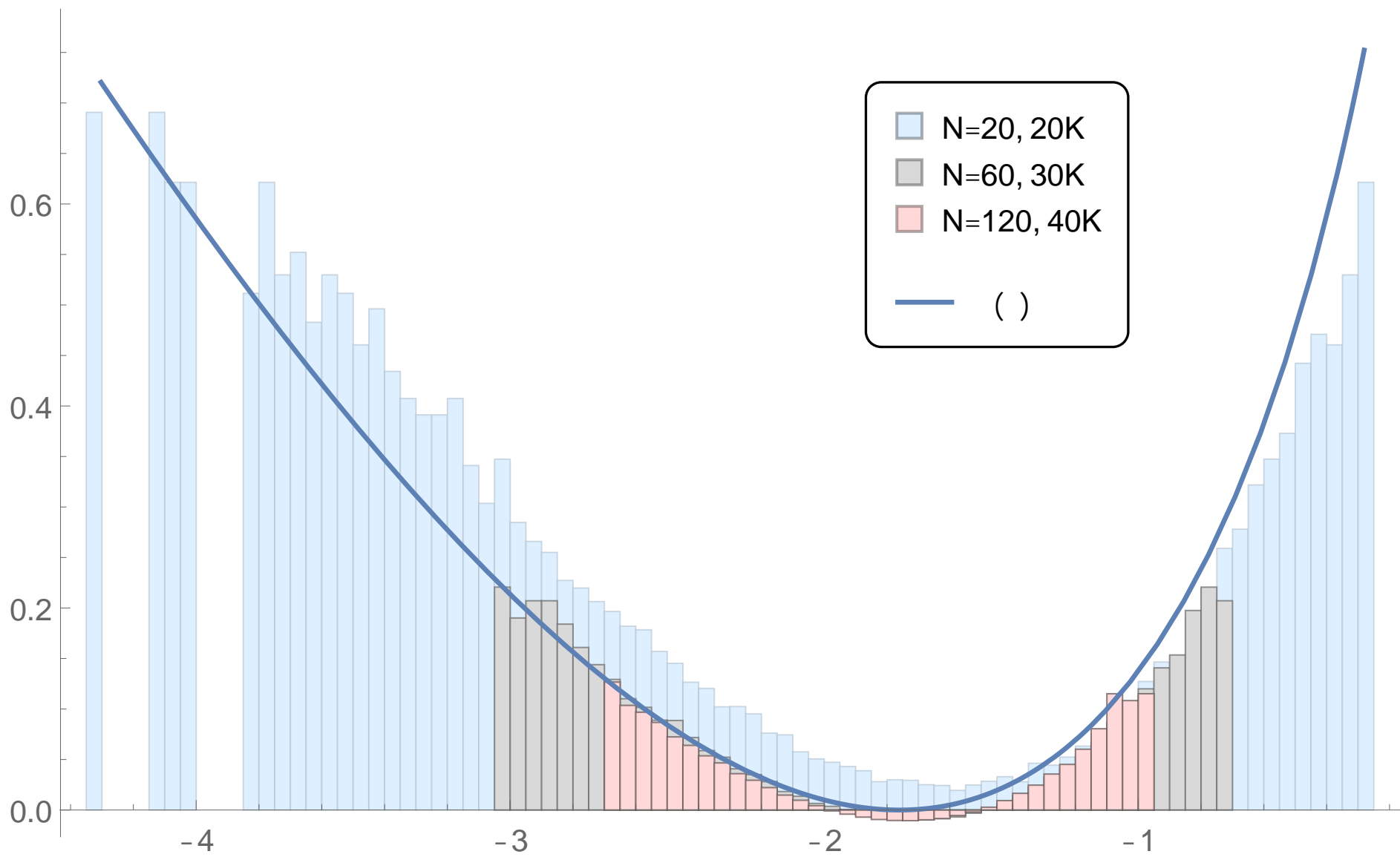
$$+ 2\left(\beta - \frac{1}{\beta}\right) \ln \frac{(\beta + 1 + \frac{\beta - \frac{1}{\beta}}{2 + \beta^2})}{2\beta - \frac{\beta - \frac{1}{\beta}}{(\beta - s)(\beta - s_+)}}$$

One finds that $I(s)$ is **minimized** for

$$s = \left(\beta - \frac{1}{\beta}\right) \frac{\beta - \frac{1}{\beta}}{1 + \beta^2}$$

which eventually implies the **most probable** value of the **minimal loss/error** :

$$\lim_{N \rightarrow \infty} \frac{E_{\min}}{N} = \frac{1}{2} \frac{\beta - \frac{1}{\beta}}{(1 + \beta^2)}$$



Conclusions:

We counted the mean number of **stationary points** of the simplest '**least-square**' optimization problem on a sphere via the Lagrange multipliers in various scaling regimes, and found the **typical** minimal loss E_{\min} .

Open questions:

- Fluctuations of the counting function,
- **large/small deviations** of the minimal loss E_{\min}
- Gradient search dynamics on the sphere
- Landscape for a **nonlinear** 'least-square' optimization, etc.

THANK YOU!