

# Fundamentals of ESs' Statistical Learning

Ofer M. Shir (Tel-Hai College & Migal Institute)

Joint work with Amir Yehudayoff, Mathematics/Technion

ofersh@telhai.ac.il



Dagstuhl Seminar 17191

Theory of Randomized Optimization Heuristics

May 11 2017, Schloss Dagstuhl, DE

# Outline

- 1 Introduction  
Model
- 2 The Analytical Covariance Matrix  
A Single Winner:  $(1, \lambda)$ -Selection  
 $(\mu, \lambda)$ -Truncation Selection
- 3 The Inverse Relation  
Proving  $\lim_{\lambda \rightarrow \infty} \alpha \mathcal{C}\mathcal{H} = \mathbf{I}$
- 4 Discussion
- 5 Backup Slides

# ESs' statistical landscape learning

the classical hypothesis  $\mathcal{C} \rightarrow \mathcal{H}^{-1}$ 

- Open question since the early development of ESs; widely discussed [Rudolph1992]
- Sheer amount of empirical evidence for this relation, in addition to extensive branding “C=inv(H)” made this hypothesis a practical *postulate* throughout the years

# the classical hypothesis $\mathcal{C} \rightarrow \mathcal{H}^{-1}$

- Open question since the early development of ESs; widely discussed [Rudolph1992]
- Sheer amount of empirical evidence for this relation, in addition to extensive branding “ $\mathcal{C}=\text{inv}(\mathcal{H})$ ” made this hypothesis a practical *postulate* throughout the years
- Recent proofs published, yet limited to Derandomization (or Natural Gradient); they exercise IGO [Akimoto2012, Beyer2014]
- We seek the fundamentals of this learning capability and consider a theoretical model – which is not likely to reflect an everyday’s heuristic – e.g.,  $\mathcal{C}$ ’s eigenvalues are  $\Omega(1/\lambda^2)$  and  $\lambda$  tends to infinity

# so why bother?

- IGO is an elephant gun – can we target an equivalent result in more fundamental ways? (David's stone+sling to knock Goliath down)

## so why bother?

- IGO is an elephant gun – can we target an equivalent result in more fundamental ways? (David's stone+sling to knock Goliath down)
- “Going back to basics” using first principles of *probability theory* and *calculus* on a basic ES model

## so why bother?

- IGO is an elephant gun – can we target an equivalent result in more fundamental ways? (David's stone+sling to knock Goliath down)
- “Going back to basics” using first principles of *probability theory* and *calculus* on a basic ES model
- Mathematically beautiful, but may also serve as a tool elsewhere in the future
- This work concerns the absolutely continuous case, but any theory should find it interesting :-)



# statistical sampling by $(1, \lambda)$ -selection

```

1  $t \leftarrow 0$ 
2  $\mathcal{S} \leftarrow \emptyset$ 
3 repeat
4   for  $k \leftarrow 1$  to  $\lambda$  do
5      $\vec{x}_k^{(t+1)} \leftarrow \vec{x}_0 + \vec{z}_k, \quad \vec{z}_k \sim \mathcal{N}(\vec{0}, \mathbf{I})$ 
6      $J_k^{(t+1)} \leftarrow \text{evaluate} \left( \vec{x}_k^{(t+1)} \right)$ 
7   end
8    $m_{t+1} \leftarrow \arg \min \left( \left\{ J_i^{(t+1)} \right\}_{i=1}^{\lambda} \right)$ 
9    $\mathcal{S} \leftarrow \mathcal{S} \cup \left\{ \vec{x}_{m_{t+1}}^{(t+1)} \right\}$ 
10   $t \leftarrow t + 1$ 
11 until  $t \geq N_{iter}$ 
output:  $\mathcal{C}^{\text{stat}} = \text{statCovariance}(\mathcal{S})$ 

```

first results [FOGA'17]

**quadratic approximation; optimum's vicinity**

$$J(\vec{x}) = J(\vec{x} - \vec{x}^*) = \vec{x}^T \cdot \mathcal{H} \cdot \vec{x} \quad (1)$$

## first results [FOGA'17]

**quadratic approximation; optimum's vicinity**

$$J(\vec{x}) = J(\vec{x} - \vec{x}^*) = \vec{x}^T \cdot \mathcal{H} \cdot \vec{x} \quad (1)$$

**main results**

- i. Rigorous formulation of the covariance matrix  $\mathcal{C}$  over ESs' selected individuals (“winners”)
- ii. The covariance matrix and the Hessian commute and are simultaneously diagonalizable for any  $\lambda$  for  $(\mu, \lambda)$ -selection
- iii. For every invertible  $\mathcal{H}$  and  $\lambda \in \mathbb{N}$ , there exists a constant  $\alpha = \alpha(\mathcal{H}, \lambda) > 0$  such that

$$\lim_{\lambda \rightarrow \infty} \alpha \mathcal{C} \mathcal{H} = \mathbf{I}.$$

## current model

**general quadratic approximation**

$$J(\vec{x}) = (\vec{x} - \vec{x}^*)^T \cdot \mathcal{H} \cdot (\vec{x} - \vec{x}^*) \quad (2)$$

## current model

## general quadratic approximation

$$J(\vec{x}) = (\vec{x} - \vec{x}^*)^T \cdot \mathcal{H} \cdot (\vec{x} - \vec{x}^*) \quad (2)$$

## truncation selection (“winners”)

$$\vec{y} = \arg \min \{J(\vec{x}_1), J(\vec{x}_2), \dots, J(\vec{x}_\lambda)\} \quad (3)$$

$$\omega = J(\vec{y}) = \min \{J(\vec{x}_1), J(\vec{x}_2), \dots, J(\vec{x}_\lambda)\} \quad (4)$$

# the covariance matrix

## the analytical form

The expectation vector of the winner is defined by its  $i^{th}$  element:

$$\mathcal{E}_i = \int x_i \text{PDF}_{\vec{y}}(\vec{x}) d\vec{x} , \quad (5)$$

$$\mathcal{C}_{ij} = \int (x_i - \mathcal{E}_i)(x_j - \mathcal{E}_j) \text{PDF}_{\vec{y}}(\vec{x}) d\vec{x} . \quad (6)$$

$\text{PDF}_{\vec{y}}(\vec{x})$  is an  $n$ -dimensional density function characterizing the *winning* decision variables about the optimum.

## the analytical form

The expectation vector of the winner is defined by its  $i^{\text{th}}$  element:

$$\mathcal{E}_i = \int x_i \text{PDF}_{\vec{y}}(\vec{x}) d\vec{x} , \quad (5)$$

$$\mathcal{C}_{ij} = \int (x_i - \mathcal{E}_i)(x_j - \mathcal{E}_j) \text{PDF}_{\vec{y}}(\vec{x}) d\vec{x} . \quad (6)$$

$\text{PDF}_{\vec{y}}(\vec{x})$  is an  $n$ -dimensional density function characterizing the *winning* decision variables about the optimum.

One of the primary goals is to fully understand this expression and utilize it.



winners' density in  $(1, \lambda)$ -selection**Proposition 0**

$$\text{PDF}_{\vec{y}}(\vec{x}) = \text{PDF}_{\omega}(J(\vec{x})) \cdot \frac{\text{PDF}_{\vec{z}}(\vec{x})}{\text{PDF}_{\psi}(J(\vec{x}))} \quad (7)$$

- $\text{PDF}_{\omega}$  : density of the *winning* value  $\omega$
- $\text{PDF}_{\vec{z}}$  : density for generating an individual by *mutation*
- $\text{PDF}_{\psi}$  : density of the objective function values (Eqs. 13 or 16)

winners' density in  $(1, \lambda)$ -selection**Proposition 0**

$$\text{PDF}_{\vec{y}}(\vec{x}) = \text{PDF}_{\omega}(J(\vec{x})) \cdot \frac{\text{PDF}_{\vec{z}}(\vec{x})}{\text{PDF}_{\psi}(J(\vec{x}))} \quad (7)$$

- $\text{PDF}_{\omega}$  : density of the *winning* value  $\omega$
- $\text{PDF}_{\vec{z}}$  : density for generating an individual by *mutation*
- $\text{PDF}_{\psi}$  : density of the objective function values (Eqs. 13 or 16)

**sketch:** consider the distribution of  $[\vec{y}; \omega]$  on  $\mathbb{R}^{n+1}$

- sample  $\{J_1, \dots, J_{\lambda}\}$  according to  $\text{PDF}_{\psi}$  independently
- sample  $\{\vec{x}_1, \dots, \vec{x}_{\lambda}\}$  conditioned on  $J_1, \dots, J_{\lambda}$  independently
- $\omega$  is set to the minimum  $J_{\ell}$ , and  $\vec{y}$  is set to  $\vec{x}_{\ell}$

$(\mu, \lambda)$ -selection

- $J_{1:\lambda} \leq J_{2:\lambda} \leq \dots \leq J_{\lambda:\lambda}$  are the order statistics obtained by sorting the objective function values.
- $\omega_{1:\lambda}, \dots, \omega_{\mu:\lambda}$  are the first  $\mu$  values from this list.
- $\vec{y}_{1:\lambda}, \dots, \vec{y}_{\mu:\lambda}$  are their corresponding vectors.

$(\mu, \lambda)$ -selection

- $J_{1:\lambda} \leq J_{2:\lambda} \leq \dots \leq J_{\lambda:\lambda}$  are the order statistics obtained by sorting the objective function values.
- $\omega_{1:\lambda}, \dots, \omega_{\mu:\lambda}$  are the first  $\mu$  values from this list.
- $\vec{y}_{1:\lambda}, \dots, \vec{y}_{\mu:\lambda}$  are their corresponding vectors.

To study the covariance in this case, we consider the pairwise density of the  $k^{th}$ -degree and  $\ell^{th}$ -degree winners ( $\ell > k$ ):

$$\text{PDF}_{\vec{y}_{k:\lambda}, \vec{y}_{\ell:\lambda}}(\vec{x}_k, \vec{x}_\ell) = \text{PDF}_{\omega_{k:\lambda}, \omega_{\ell:\lambda}}(J(\vec{x}_k), J(\vec{x}_\ell)) \times \left( \frac{\text{PDF}_{\vec{z}}(\vec{x}_k)}{\text{PDF}_{\psi}(J(\vec{x}_k))} \right) \cdot \left( \frac{\text{PDF}_{\vec{z}}(\vec{x}_\ell)}{\text{PDF}_{\psi}(J(\vec{x}_\ell))} \right) \cdot \quad (8)$$

# the inverse relation

## winning values' density &amp; proposition 1

For simplicity, we consider  $(1, \lambda)$ -selection.

## winning values' density &amp; proposition 1

For simplicity, we consider  $(1, \lambda)$ -selection.

$$\text{CDF}_\omega(v) = 1 - (1 - \text{CDF}_\psi(v))^\lambda \quad (9)$$

$$\text{PDF}_\omega(v) = \lambda \cdot (1 - \text{CDF}_\psi(v))^{\lambda-1} \cdot \text{PDF}_\psi(v) \quad (10)$$

## winning values' density &amp; proposition 1

For simplicity, we consider  $(1, \lambda)$ -selection.

$$\text{CDF}_\omega(v) = 1 - (1 - \text{CDF}_\psi(v))^\lambda \quad (9)$$

$$\text{PDF}_\omega(v) = \lambda \cdot (1 - \text{CDF}_\psi(v))^{\lambda-1} \cdot \text{PDF}_\psi(v) \quad (10)$$

**Proposition 1:**

For every invertible  $\mathcal{H}$  and  $\lambda \in \mathbb{N}$ , there exists a constant  $\alpha = \alpha(\mathcal{H}, \lambda) > 0$  such that

$$\lim_{\lambda \rightarrow \infty} \alpha \mathcal{C}\mathcal{H} = \mathbf{I}.$$



# intuition for proving proposition 1

We first target a diagonal  $\mathcal{H}$

For a large  $\lambda$ , the winner  $\vec{y}$  is close to the optimum, which in turn implies that  $(\mathcal{CH})_{ii}$  does not actually depend on  $i$ .

For the general case, we ought to show that both  $\mathcal{H}$  and  $\mathcal{C}$  are diagonalizable in the same base under the same limit conditions (not shown here).

## proof sketch for proposition 1

i. firstly, assume  $\mathcal{H}$  is diagonal and apply change of variables

$$r_i = \sqrt{\Delta_i} \cdot (x_i - x_i^*)$$

$$\text{ii. } \mathcal{E}_i - x_i^* = \frac{c_{\mathcal{H}}}{\sqrt{\Delta_i}} \int r_i \lambda (1 - \text{CDF}_{\psi}(\|\vec{r}^*\|^2))^{\lambda-1} \exp\left(-\hat{J}(\vec{r})\right) d\vec{r}$$

$$\text{iii. show that } |\mathcal{E}_i - x_i^*| \leq \epsilon_1 \sqrt{C_{ii}}$$

$$\text{iv. bound the off-diagonal terms } C_{ij} \leq \epsilon_2 \sqrt{C_{ii}C_{jj}}$$

v. show that  $\alpha \Delta_i C_{ii} \geq 1 - \epsilon_3$  and  $\alpha \Delta_i C_{ii} \leq 1 + \epsilon_4$  ( $\epsilon_3$  and  $\epsilon_4$  tend to zero as  $\lambda$  tends to infinity)

vi. secondly, for a non-diagonal  $\mathcal{H}$ ,

$$\lim_{\lambda \rightarrow \infty} \alpha \mathcal{C}\mathcal{H} - \mathbf{I} = \lim_{\lambda \rightarrow \infty} \mathbf{U} (\alpha \mathcal{T}\mathcal{D} - \mathbf{I}) \mathbf{U}^{-1} = 0$$

## discussion

i.  $\mathcal{C}$  and  $\mathcal{H}$  commute (for any  $\lambda$  near the optimum, for  $\lambda \rightarrow \infty$  elsewhere).

this learning capability stems only from two components:

(1) isotropic Gaussian mutations, and (2) rank-based selection.

\* learning the landscape is an inherent property of classical ESs.

\*\* it does not require Derandomization (adaptation) nor IGO (proofs)

## discussion

i.  $\mathcal{C}$  and  $\mathcal{H}$  commute (for any  $\lambda$  near the optimum, for  $\lambda \rightarrow \infty$  elsewhere).

this learning capability stems only from two components:

(1) isotropic Gaussian mutations, and (2) rank-based selection.

\* learning the landscape is an inherent property of classical ESs.

\*\* it does not require Derandomization (adaptation) nor IGO (proofs)

ii.  $\lim_{\lambda \rightarrow \infty} \alpha \mathcal{C} \mathcal{H} = \mathbf{I}$  ; this approximation has two parts:

(1) guaranteeing that  $\mathcal{C}^{\text{stat}}$  is pointwise  $\epsilon$ -close to  $\mathcal{C}$  with confidence  $1 - \delta$ . the eigenvalues of  $\mathcal{C}$  are at least  $\Omega(1/\lambda^2)$ ; for  $\mathcal{C}^{\text{stat}}$  to meaningfully approach  $\mathcal{C}$  it requires  $\epsilon \ll 1/\lambda^2$ .

$\implies$  number of samples required for this part is polynomial in  $\lambda, 1/\epsilon, \ln(n)$  and  $\ln(1/\delta)$ .

(2) guaranteeing that  $\mathcal{C}$  is pointwise  $\epsilon$ -close to  $\alpha \mathcal{H}^{-1}$  ,  $\alpha(\lambda, \mathcal{H}) > 0$ .

$\implies$  upper bound on the number of samples required for this part depends on  $\epsilon, \lambda$  and on the spectrum of  $\mathcal{H}$ .

# limit distributions of order statistics

In order to calculate  $\mathcal{C}_{ij}$  when  $\lambda$  tends to infinity, it is possible to approximate  $\text{PDF}_\omega(J(\vec{x}))$ , by considering  $\mathcal{L}_\lambda(v) = 1 - (1 - \text{CDF}_\psi(v))^\lambda$  at  $\lim_{\lambda \rightarrow \infty} \mathcal{L}_\lambda(v)$ :

# limit distributions of order statistics

In order to calculate  $\mathcal{C}_{ij}$  when  $\lambda$  tends to infinity, it is possible to approximate  $\text{PDF}_\omega(J(\vec{x}))$ , by considering  $\mathcal{L}_\lambda(v) = 1 - (1 - \text{CDF}_\psi(v))^\lambda$  at  $\lim_{\lambda \rightarrow \infty} \mathcal{L}_\lambda(v)$ :

**theorem [Fisher-Tippett]**

the generalized extreme value distributions (GEVD) are the only non-degenerate family of distributions satisfying this limit:

$$\mathcal{L}_\kappa(v; \kappa_1, \kappa_2, \kappa_3) = 1 - \exp \left\{ - \left[ 1 + \kappa_3 \left( \frac{v - \kappa_1}{\kappa_2} \right) \right]^{1/\kappa_3} \right\} \quad (11)$$

$\implies \text{CDF}_\psi$  belongs to Weibull

## limit distributions of order statistics

In order to calculate  $\mathcal{C}_{ij}$  when  $\lambda$  tends to infinity, it is possible to approximate  $\text{PDF}_\omega(J(\vec{x}))$ , by considering  $\mathcal{L}_\lambda(v) = 1 - (1 - \text{CDF}_\psi(v))^\lambda$  at  $\lim_{\lambda \rightarrow \infty} \mathcal{L}_\lambda(v)$ :

**theorem [Fisher-Tippett]**

the generalized extreme value distributions (GEVD) are the only non-degenerate family of distributions satisfying this limit:

$$\mathcal{L}_\kappa(v; \kappa_1, \kappa_2, \kappa_3) = 1 - \exp \left\{ - \left[ 1 + \kappa_3 \left( \frac{v - \kappa_1}{\kappa_2} \right) \right]^{1/\kappa_3} \right\} \quad (11)$$

$\implies \text{CDF}_\psi$  belongs to Weibull

This tool is hardly ever exercised amongst our scholars; Rudolph utilized it in his book [Rudolph1997].

Acknowledgements to Jonathan Roslund.

**danke**



## probability functions

**isotropic case:**  $\mathcal{H} = \mathbf{I}$

$\psi = J(\vec{z})$  is a random variable obeying the  $\chi^2$ -distribution:

$$F_{\chi^2}(\psi) = \frac{1}{2^{n/2}\Gamma(n/2)} \int_0^\psi t^{\frac{n}{2}-1} \exp\left(-\frac{t}{2}\right) dt \quad (12)$$

$$f_{\chi^2}(\psi) = \frac{1}{2^{n/2}\Gamma(n/2)} \psi^{n/2-1} \exp\left(-\frac{\psi}{2}\right) \quad (13)$$

## probability functions

**isotropic case:**  $\mathcal{H} = \mathbf{I}$

$\psi = J(\vec{z})$  is a random variable obeying the  $\chi^2$ -distribution:

$$F_{\chi^2}(\psi) = \frac{1}{2^{n/2}\Gamma(n/2)} \int_0^\psi t^{\frac{n}{2}-1} \exp\left(-\frac{t}{2}\right) dt \quad (12)$$

$$f_{\chi^2}(\psi) = \frac{1}{2^{n/2}\Gamma(n/2)} \psi^{n/2-1} \exp\left(-\frac{\psi}{2}\right) \quad (13)$$

**general case:**  $\mathcal{H} = \mathcal{U}\mathcal{D}\mathcal{U}^{-1}$ ,  $\mathcal{D} = \mathbf{diag}[\Delta_1, \dots, \Delta_n]$

$$F_{\mathcal{H}\chi^2}(\psi) = \int_0^\infty \frac{2}{\pi} \frac{\sin \frac{t\psi}{2}}{t} \cos\left(-t\psi + \frac{1}{2} \sum_{j=1}^n \tan^{-1} 2\Delta_j t\right) \times \prod_{j=1}^n (1 + \Delta_j^2 t^2)^{-\frac{1}{4}} dt, \quad (14)$$

## approximation for the general case

$$F_{\tau\chi^2}(\psi) = \frac{\Upsilon^\eta}{\Gamma(\eta)} \int_0^\psi t^{\eta-1} \exp(-\Upsilon t) dt \quad (15)$$

$$f_{\tau\chi^2}(\psi) = \frac{\Upsilon^\eta}{\Gamma(\eta)} \psi^{\eta-1} \exp(-\Upsilon\psi) \quad (16)$$

$\Upsilon$  and  $\eta$  account for the first two moments of  $\vec{z}^T \mathcal{H} \vec{z}$ :

$$\Upsilon = \frac{1}{2} \frac{\sum_{i=1}^n \Delta_i}{\sum_{i=1}^n \Delta_i^2}, \quad \eta = \frac{1}{2} \frac{(\sum_{i=1}^n \Delta_i)^2}{\sum_{i=1}^n \Delta_i^2} \quad (17)$$

approximation for the general case

$$F_{\tau\chi^2}(\psi) = \frac{\Upsilon^\eta}{\Gamma(\eta)} \int_0^\psi t^{\eta-1} \exp(-\Upsilon t) dt \quad (15)$$

$$f_{\tau\chi^2}(\psi) = \frac{\Upsilon^\eta}{\Gamma(\eta)} \psi^{\eta-1} \exp(-\Upsilon\psi) \quad (16)$$

$\Upsilon$  and  $\eta$  account for the first two moments of  $\vec{z}^T \mathcal{H} \vec{z}$ :

$$\Upsilon = \frac{1}{2} \frac{\sum_{i=1}^n \Delta_i}{\sum_{i=1}^n \Delta_i^2}, \quad \eta = \frac{1}{2} \frac{(\sum_{i=1}^n \Delta_i)^2}{\sum_{i=1}^n \Delta_i^2} \quad (17)$$

Accuracy depends on the eigenvalues'  $\{\Delta_i\}$  *standard deviation*.