Abstract

3RD CARLO ALBERTO STOCHASTICS WORKSHOP

Jan 10-11, 2014



Via Real Collegio, 30 10024 Moncalieri (TO), Italy

Optimization via information geometry

Giovanni Pistone

de Castro Statistics Initiative, Collegio Carlo Alberto, Moncalieri, Italy

Jan 11, 2014

Plan

- 1. Stochastic Relaxation of an Optimization Problem.
- 2. Natural Gradient
- 3. Second Order Geometry
- Luigi Malagò, Matteo Matteucci, and Giovann Pistone. Towards the geometry of estimation of distribution algorithms based on the exponential family.

In Proceedings of the 11th workshop on Foundations of genetic algorithms, FOGA '11, pages 230–242, New York, NY, USA, 2011. ACM

- Luigi Malagò, Matteo Matteucci, and Giovanni Pistone. Natural gradient, fitness modelling and model selection: A unifying perspective.
 In IEEE Congress on Evolutionary Computation, pages 486–493. IEEE, 2013
- L. Arnold, A. Auger, N. Hansen, and Y. Ollivier. Information-Geometric Optimization Algorithms: A Unifying Picture via Invariance Principles. arXiv:1106.3708, 2011v1; 2013v2

Combinatorial optimization is the maximization of a real function defined on a finite space $f : \Omega \to \mathbb{R}$. This problem is reduced to a continuous optimization problem by considering the relaxed function $F(p) = \mathbb{E}_p[f]$, where p is a positive density in a statistical model \mathcal{M} on Ω . I will present results from a joint work in progress with Luigi Malagò (Università Statale di Milano). This work expands previous works in Optimization by L. Malagò et al., where the geometry of exponential families is show to provide a suitable setting for model based methods in Combinatorial Optimization under a Black Box assumption on the function f. Some basic tools of Algebraic Design of Experiments are used.

Stochastic relaxation

- $(\Omega, \mathcal{F}, \mu)$ (metric) measure space, $\mathcal{P}_{>}$ (strictly) positive probability densities.
- An open statistical model (M, θ, B) is a parameterized subset of P_>, that is M ⊂ P_> and s: M → B, where s is a one-to-one mapping onto an open subset of a Banach space B.
- If f: Ω → ℝ is a bounded continuous function, the mapping
 M ∋ p ↦ E_p[f] is a Stochastic Relaxation SR of f.
- E_p[f] < sup_{ω∈Ω} f(ω) for all p ∈ M if f is not constant, but sup_{p∈M} E_p[f] = sup_{ω∈Ω} f(ω) if there exist a probability measure ν in the weak closure of M · μ whose support is contained in the set of maximizing points of f, {ω ∈ Ω: f(ω) = sup_{ω∈Ω} f(ω)}.
- A SR optimization method is an algorithm to produce a sequence $p_n \in \mathcal{M}, n \in \mathbb{N}$, such that $\lim_{n\to\infty} E_{p_n}[f] = \sup_{\omega\in\Omega} f(\omega)$.
- Such algorithms are best studied in the framework of Information Geometry IG, that is the differential geometry of statistical models.

SR on an exponential family

- The exponential family $q_{\theta} = \exp\left(\sum_{j=1}^{d} \theta_j T_j \psi(\theta)\right) \cdot p$ is a statistical model $\mathcal{M} = \{q_{\theta}\}$ and parameterization $\sigma \colon q_{\theta} \mapsto \theta \in \mathbb{R}^d$.
- 1. $\psi(\theta) = \log (\mathsf{E}_{\rho} [e^{\theta \cdot T}])$ is convex, lower semi-continuous;
 - 2. ψ is analytic on the interior \mathcal{U} the proper domain;
 - 3. $\nabla \psi(\theta) = \mathsf{E}_{\theta}[T]$, Hess $\psi(\theta) = \mathsf{Var}_{\theta}(T)$.
 - U ∋ θ → ∇ψ(θ) = η ∈ N is one-to-one, analytic, monotone;
 N is the interior of the marginal polytope, i.e. the convex set generated by {T(ω): ω ∈ Ω};
- The SR θ → E_θ [f] is well posed iff the border set ∂M contains at least one point of argmax f. A sufficient condition is δ_{T(ω)} ∈ ∂M for all ω ∈ Ω
- The gradient of the SR is

 $\nabla(\boldsymbol{\theta} \mapsto \mathsf{E}_{\boldsymbol{\theta}}[f]) = (\mathsf{Cov}_{\theta}(f, T_1), \dots, \mathsf{Cov}_{\theta}(f, T_d))$

which suggests to take the θ -MS approximation of f on Span (T_1, \ldots, T_d) as direction of steepest ascent.

• This ideas prompt for a systematic development of the geometric picture of statistical models.

EDA

• Estimation of Distribution is a model-based optimization algorithm.

Input: N, M	▷ population size, selected population size
Input: $\mathcal{M} = \{p(x; \xi)\}$	> parametric model
1: $t \leftarrow 0$	
2: $\mathcal{P}^t = \text{INITRANDOM}($)
3: repeat	
4: $\mathcal{P}_s^t = \text{Selection}$	$(\mathcal{P}^t, M) \qquad \qquad \triangleright \text{ select } M \text{ samples}$
5: $\xi^{t+1} = \text{ESTIMAT}$	$ON(\mathcal{P}_s^t, \mathcal{M})$ \triangleright opt. model selection
6: $\mathcal{P}^{t+1} = \text{SAMPLEI}$	$\Re(\xi^{t+1}, N)$ $\triangleright N$ samples
7: $t \leftarrow t+1$	
8: until StoppingCri	ΓERIA()

Information Geometry

- Shun-ichi Amari and Hiroshi Nagaoka. Methods of information geometry. American Mathematical Society, Providence, RI, 2000. Translated from the 1993 Japanese original by Daishi Harada
- Giovanni Pistone and Carlo Sempi. An infinite-dimensional geometric structure on the space of all the probability measures equivalent to a given one. Ann. Statist., 23(5):1543–1561, October 1995
- Paolo Gibilisco and Giovanni Pistone. Connections on non-parametric statistical manifolds by Orlicz space geometry. IDAOP. 1(2):325–347, 1998
- Giovanni Pistone and Maria Piera Rogantin. The exponential statistical manifold: mean parameters, orthogonality and space transformations. Bernoulli, 5(4):721–760, August 1999
- Alberto Cena. Geometric structures on the non-parametric statistical manifold. PhD thesis, Dottorato in Matematica, Università di Milano, 2002
- Alberto Cena and Giovanni Pistone. Exponential statistical manifold. Ann. Inst. Statist. Math., 59(1):27–56, 2007
- Luigi Malagò. On the geometry of optimization based on the exponential family relaxation. PhD thesis, Politecnico di Milano, 2012
- Giovanni Pistone. Nonparametric information geometry. In Frank Nielsen and Freedeèric Barbaresco, editors, *Geometric Science of Information*, number 8085 in LNCS, pages 5–36, Berlin Heidelberg, 2013. Springer-Verlag. First International Conference, GSI 2013 Paris, France, August 28–30, 2013 Proceedings

SNGD

• Stochastic Natural Gradient Descent is a ST algorithm that requires the estimation of a gradient.

• Input: N, λ	population size, learning rate
Optional: M	\triangleright selected population size (default $M = N$)
1: $t \leftarrow 0$	
2: $\theta^t \leftarrow (0, \dots, 0)$	uniform distribution
3: $\mathcal{P}^t \leftarrow \text{InitRando}$	M() ▷ random initial population
4: repeat	
5: $\mathcal{P}_s^t = \text{Selection}$	$DN(\mathcal{P}^t, M)$ \triangleright opt. select M samples
6: $\widehat{\nabla}\mathbb{E}[f] \leftarrow \widehat{\mathrm{Cov}}(f)$	$(T, T_i)_{i=1}^m$ \triangleright empirical covariances
7: $\widehat{I} \leftarrow [\widehat{Cov}(T_i, T_i)]$	$[i,j]_{i,j=1}^m agestarrow \{T_i(x)\}$ may be learned
8: $\theta^{t+1} \leftarrow \theta^t - \lambda \hat{I}$	$\hat{\nabla} \mathbb{E}[f]$
9: $\mathcal{P}^{t+1} \leftarrow \text{Gibbs}$	SAMPLER (θ^{t+1}, N) $\triangleright N$ samples
10: $t \leftarrow t+1$	
11: until StoppingCi	RITERIA()

A toy example I

- $f(x_1, x_2) = a_0 + a_1x_1 + a_2x_2 + a_{12}x_1x_2, x_1, x_2 = \pm 1, a_0, a_1, a_2, a_{12} \in \mathbb{R}.$
- f is a real random variable on the sample space $\Omega = \{+1, -1\}^2$ with the uniform probability λ .
- $X_1, X_2: \Omega \to \pm 1$ generate an orthonormal basis $1, X_1, X_2, X_1X_2$ of $L^2(\Omega, \lambda)$ and f is the general form of a real random variable on such a space.
- $\mathcal{P}_{>}$ is the open simplex of positive densities on (Ω, λ) , \mathcal{E} is a statistical model.

The relaxed mapping $F: \mathcal{E}_{\rightarrow} \mathbb{R}$,

 $F(p) = \mathsf{E}_{p}[f] = a_{0} + a_{1}\mathsf{E}_{p}[X_{1}] + a_{2}\mathsf{E}_{p}[X_{2}] + a_{12}\mathsf{E}_{p}[X_{1}X_{2}]$

is strictly bounded by the maximum of $f, \ \mathsf{E}_{\rho}\left[f\right] < \max_{x \in \Omega} \ \text{if} \ f \ \text{is not}$ constant.

Model, border



A toy example II

- We are looking for a sequence p_n , $n \in \mathbb{N}$, such that $\mathsf{E}_{p_n}[f] \to \max_{x \in \Omega} f(x)$ as $n \to \infty$.
- The existence of such a sequence in a nontrivial condition for the model $\ensuremath{\mathcal{E}}.$

This condition is satisfied by the independence model, when we can write

$$F(\eta^1, \eta^2) = a_0 + a_1 \eta^1 + a_2 \eta^2 + a_{12} \eta^1 \eta^2, \quad \eta^i = \mathsf{E}_p[X_i],$$

- Not all functions on $\mathcal{P}_>$ are extremised on the vertices. For example the entropy, or the polarization measure

$$p\mapsto \sum_{\omega\in\Omega}p(\omega)^2(1-p(\omega))$$

Border: polarization



$$F(\eta_1,\eta_2)=\eta_1+2\eta_2+3\eta_1\eta_2$$



 $abla F(\eta_1, \eta_2) = (1 + 3\eta_2, 2 + 3\eta_1)$



F and ∇F



Integral curves



∇F and integral curves

Issues/Actions



Issues

issue 1 There are critical points in the interior of $[-1,+1]^2$. issue 2 The gradient points in the right direction outside $[-1,+1]^2$.

Actions

action 1 Order according decreasing effects. action 2 Use a modified gradient.

$Y_1 = X_1 X_2, Y_2 = X_2, Y_1 Y_2 = X_1$

			1	X_1	X_2	X_1X_2	Y_1	Y_2	$Y_1 Y_2$
1	1	1	1	1	1	1	1	1	1
2	-1	1	1	-1	1	-1	-1	1	-1
3	1	-1	1	1	-1	-1	-1	-1	1
4	-1	-1	1	-1	-1	1	1	-1	-1

Same function

 $X_1 + 2X_2 + 3X_1X_2 = 3Y_1 + 2Y_2 + Y_1Y_2$

New model

 $Y_1 = X_1 X_2$, $Y_2 = X_2$ independent

- Emanuele Corsano, Davide Cucci, Luigi Malagò, and Matteo Matteucci. Implicit model selection based on variable transformations in estimation of distribution.
 In LION, pages 360–365, 2012
- Davide Cucci, Luigi Malagò, and Matteo Matteucci. Variable transformations in estimation of distribution algorithms.
 In PPSN (1), pages 428–437, 2012

$F(\zeta_1,\zeta_2)=3\zeta_1+2\zeta_2+\zeta_1\zeta_2$







Natural gradient

 $E_{\theta} [(X_1 + 2X_2 + 3X_1X_2)X_1] = 1 + 2\eta_1\eta_2 + 3\eta_2$ $E_{\theta} [X_1 + 2X_2 + 3 + X_3] E_{\theta} [X_1] = \eta_1^2 + 2\eta_1\eta_2 + 3\eta_1^2\eta_2$ $Cov_{\theta} (X_1 + 2X_2 + 3X_1X_2, X_1) = (1 - \eta_1^2)(1 + 3\eta_2)$

 $E_{\theta} [(X_1 + 2X_2 + 3X_1X_2)X_2] = \eta_1\eta_2 + 2 + 3\eta_1$ $E_{\theta} [X_1 + 2X_2 + 3 + X_3] E_{\theta} [X_2] = \eta_1\eta_2 + 2\eta_2^2 + 3\eta_1\eta_2^2$ $Cov_{\theta} (X_1 + 2X_2 + 3X_1X_2, X_2) = (1 - \eta_2^2)(2 + 3\eta_1)$

$$egin{aligned} \widetilde{
abla} F(\eta_1,\eta_2) &= ig[(1-\eta_1^2)(1+3\eta_2),(1-\eta_2^2)(2+3\eta_1)ig] \ &=
abla F(\eta_1,\eta_2) igg[egin{aligned} 1-\eta_1^2 & 0 \ 0 & 1-\eta_2^2 \end{bmatrix} \end{aligned}$$





$$\widetilde{
abla} F(\eta_1,\eta_2) = ig((1-\eta_1^2)(1+3\eta_2),(1-\eta_2^2)(2+3\eta_1)ig)$$











IGO

- Other relaxation are possible.
- If λ is a positive density, the function

$$\tilde{f}: \Omega \ni x \mapsto \mathsf{P}_{\lambda}[f \le f(x)]$$

has the same maximum as f, and

$$\widehat{F} \colon \mathcal{M} \ni p \mapsto \mathsf{E}_p\left[\widehat{f}\right]$$

is a relaxation of \hat{f} .

•

- \hat{f} is invariant under monotone transformations of the values of f and it is a stochastic ordering measure..
- L. Arnold, A. Auger, N. Hansen, and Y. Ollivier. Information-Geometric Optimization Algorithms: A Unifying Picture via Invariance Principles. arXiv:1106.3708, 2011v1; 2013v2

IGO:
$$F(\eta_1, \eta_2) = \frac{5}{8} + \frac{1}{8}\eta_2 + \frac{1}{4}\eta_1\eta_2$$



Charts II

• The centering is a parallel transport on the tangent spaces of the manifold:

$${}^{e}\mathbb{U}^{p}\colon\mathcal{V}\ni U\mapsto U-\mathsf{E}_{p}\left[U\right]\in{}^{e}\mathbb{U}^{p}\mathcal{V}.$$

• $\theta = I_{\mathcal{B}}^{-1}(p) \mathsf{E}_{p}[U^{e} \mathbb{U}^{p}X]$, where

$$I_{\mathcal{B}}(p) = \left[\mathsf{Cov}_{p}\left(X_{i}, X_{j}\right)\right]_{ij} = \mathsf{E}_{p}\left[XX'\right] - \mathsf{E}_{p}\left[X\right]\mathsf{E}_{p}\left[X'\right] \quad (1)$$

is the Fisher matrix of the basis $\mathcal{B} = \{X_1, \ldots, X_m\}$.

Charts I

- On the finite sample space Ω , $\#\Omega = n$, consider a set of random variables $\mathcal{B} = \{X_1, \ldots, X_m\}$ such that $\sum_J \alpha_j X_j$ is constant only if the α_j 's are zero, which implies, in turn, the linear independence of \mathcal{B} .
- We define $\mathcal{V} = \text{Span}(X_1, \dots, X_m)$ and

$$\mathcal{E}_{\mathcal{V}}(p) = \left\{ q \in \mathcal{P}_{>} \colon q \propto \mathrm{e}^{U} p, U \in \mathcal{V}
ight\}.$$

As *E_V*(*q*) = *E_V*(*p*) if, and only if, *q* ∈ *E_V*(*p*) = *E_V*, each choice of a specific reference *p* is the chart centered at *p*

$$\sigma_{p} \colon \exp\left(\sum_{j} \theta^{j \ e} \mathbb{U}^{p} X_{j} - \psi_{p}(\theta)\right) \cdot p \mapsto \theta,$$

where
$${}^{e}\mathbb{U}^{p}$$
 is the centering at p .

Gradients I

Given a function $\phi: \mathcal{E}_{\mathcal{V}} \to \mathbb{R}$ let $\phi_p = \phi \circ e_p$, $e_p = \sigma_p^{-1}$, its representation in the chart centered at p:

$$\begin{array}{c} \mathcal{E}_{\mathcal{V}} \xrightarrow{\phi} \mathbb{R} \\ \stackrel{e_p}{\longrightarrow} \\ \mathbb{R}^m \end{array}$$

The derivative of $\boldsymbol{\theta} \mapsto \phi_p(\boldsymbol{\theta})$ at $\boldsymbol{\theta} = \boldsymbol{0}$ along $\boldsymbol{\alpha} \in \mathbb{R}^m$ is

$$\nabla \phi_{\rho}(\mathbf{0})\boldsymbol{\alpha} = \nabla \phi_{\rho}(\mathbf{0})I_{\mathcal{B}}^{-1}(p)I_{\mathcal{B}}(p)\boldsymbol{\alpha} = \\ \left(I_{\mathcal{B}}^{-1}(p)\nabla \phi_{\rho}(\mathbf{0})'\right)'I_{\mathcal{B}}(p)\boldsymbol{\alpha} = g_{\rho}(I_{\mathcal{B}}^{-1}(p)\nabla \phi_{\rho}(\mathbf{0})',\boldsymbol{\alpha}).$$

The mapping $\widetilde{\nabla}\phi \colon p \mapsto l_{\mathcal{B}}^{-1}(p)(\nabla \phi_p(\mathbf{0}))' \in \mathbb{R}^m$ is Amari's natural gradient.

Gradients II



$\dot{\sigma}_{p} \circ \nabla \phi(p) = I_{\mathcal{B}}^{-1} \nabla \phi_{p}(\mathbf{0}) = \widetilde{\nabla} \phi(p)$

Levi-Civita connection

If $D_Y V$ is the vector field on $\mathcal{E}_{\mathcal{V}}$ whose value at p has coordinates

$$\dot{\sigma}_{p}(D_{Y}V(p)) = dV_{p}(\mathbf{0})\alpha + \frac{1}{2}I_{\mathcal{B}}^{-1}(p)\left(dI_{\mathcal{B},p}(\mathbf{0})\alpha\right)V_{p}(\mathbf{0}), \quad \alpha = \dot{\sigma}_{p}(Y(p)),$$

then

$$D_Y g(V, W) = g(D_Y V, W) + g(V, D_Y W)$$

$$\dot{\sigma}_p (D_W V(p) - D_V W(p)) = \dot{\sigma}[V, W](p),$$

i.e. $D_Y V$ is the metric covariant derivative.

- See VIII, §4 of Serge Lang. Differential and Riemannian manifolds, volume 160 of Graduate Texts in Mathematics.
 Springer-Verlag, New York, third edition, 1995,
- §5.3.2 of P.-A. Absil, R. Mahony, and R. Sepulchre. Optimization algorithms on matrix manifolds. Princeton University Press, Princeton, NJ, 2008.
 With a foreword by Paul Van Dooren.
- L. Malagò, G. Pistone, work in progress on the Newton method.