













Examples of Model-based Algorithms

Evolutionary computation

EDAs (Larrañaga and Lozano, 2002), DEUM framework (Shakya et al., 2005)

8

Gradient descent

 SGD (Robbins and Monro, 1951), CMA-ES (Hansen and Ostermeier, 2001), NES (Wierstra et al., 2008), SNGD (M. et al., FOGA 2011), IGO (Ollivier et al., 2011),

Boltzmann distribution and Gibbs sampler (Geman and Geman, 1984)

Simulated Annealing and Boltzmann Machines (Aarts and Korst, 1989) The Cross-Entropy method (Rubinstein, 1997)

LP relaxation in pseudo-Boolean optimization (Boros and Hammer, 2001) Methods of Moments (Meziat et al., 2001)



















	18
Stochastic Relaxation of f	
Consider the following optimization problem	
$(P) \qquad \min_{x \in \Omega} \ f(x)$	
We define stochastic relaxation (SR) of f the function	n
$F: p \mapsto \mathbb{E}_p[f]$	
Given a statistical model $\mathcal{M} = \{p(x)\}$, we look for the by generating minimizing sequences $\{p_t\}$ in \mathcal{M} for F	e solution of (P) $T(p)$
Let ξ be a parameterization for \mathcal{M} , i.e., $\mathcal{M} = \{p(x;\xi) $ can be expressed as	$: \xi \in \Xi$ }, the SR
(SR) $\min_{\xi \in \Xi} F(\xi)$	
We move the search to the space of probability distr	ibution
The parameters $\xi \in \Xi$ become the variables of the re	laxed problem
L. Malagò, T. Glasmachers, GECCO, July 13, 2014	

Equivalence of (P) and (SR)	19		
Let us introduce some notation • $x^* \in \Omega^* = \arg \min_{x \in \Omega} f(x)$ the global optima of f • $p_* \in \mathcal{M}^* = \arg \min_{p \in \overline{\mathcal{M}}} F(\xi)$ the global optima of • $\overline{\mathcal{M}}$ the topological closure of \mathcal{M} , i.e., \mathcal{M} together distributions of sequences $\{p_t\} \in \mathcal{M}$	<i>F</i> r all lin	nit	
Candidate solutions for (P) can be sampled by solution	ons of	the (S	R)
Distributions in \mathcal{M}^* have reduced support and for discorresponds to faces of Δ	crete §	2	
(P) and (SR) and equivalent if and only if from a solu can sample points in Ω^* with $\mathbb{P}(X = x^*) = 1$	tion of	(SR)	we
A sufficient condition is the inclusion of the Dirac dist $\overline{\mathcal{M}}$, i.e., there exists a sequence $\{p_t\} \in \mathcal{M}$ such that	ributio	ns δ_{x^*}	in
$\lim_{t \to \infty} F(p_t) = \min_{x \in \Omega} f(x)$			
L. Malagò, T. Glasmachers, GECCO, July 13, 2014			

20 21 The Exponential Family Information Geometry In the following, we consider models in the Exponential Family \mathcal{E} The geometry of statistical models is not Euclidean, and we need tools from differential geometry to properly describe notions such as $p(x,\theta) = \exp\left(\sum_{i=1}^{m} \theta_i T_i(x) - \psi(\theta)\right)$ tangent vectors, shortest paths and distances between distributions Information Geometry (IG), first propose by Amari, consists of the • sufficient statistics $T = (T_1(x), \dots, T_m(x))$ study of statistical models as manifolds of distributions endowed with • natural parameters $\theta = (\theta_1, \dots, \theta_m) \in \Theta$ the Fisher information matrix • log-partition function $\psi(\theta)$ Several statistical models belong to the exponential family, both in the continuous and discrete case L. Malagò, T. Glasmachers, GECCO, July 13, 2014

22

Characterization of the Tangent Space of
$$\mathcal{E}$$

The one dimensional model
 $p(\theta) = \exp\{\theta T - \psi(\theta)\}$
is a curve in the manifold, with tangent vector
 $T - \frac{d}{d\theta}\psi(\theta)$
On the other side, given a vector field, at each p we have a vector
 $U(p)$ tangent to some curve, we obtain a differential equation
 $\frac{d}{d\theta}\log p(\theta) = U(p),$

whose solution is a one dimensional model in $\ensuremath{\mathcal{E}}$

L. Malagò, T. Glasmachers, GECCO, July 13, 2014

A statistical model can be modeled as a manifold of distributions by introducing an affine chart in *p* such that any density *q* is locally represented w.r.t. to the reference measure *p*, i.e., $\frac{q}{p} - 1$ The tangent space at each point *p* is defined by $T_p = \{v : \mathbb{E}_p[v] = 0\}$ Consider a curve $p(\theta)$ such that p(0) = p, then $\frac{\dot{p}(\theta)}{p} \in T_p$ In a moving coordinate system, tangent (velocity) vectors in $T_{p(\theta)}$ to the curve are given by logarithmic derivatives $\frac{\dot{p}(\theta)}{p(\theta)} = \frac{d}{d\theta} \log p(\theta)$

Characterization of the Tangent Space of \mathcal{E}

Geometry of the Exponential Family In case of a finite sample space \mathcal{X} , we have $p(x;\theta) = \exp\left(\sum_{i=1}^{k} \theta_i T_i(x) - \psi(\theta)\right) \quad \theta \in \mathbb{R}^k$ and $\mathsf{T}_{\theta} = \left\{ v : v = \sum_{i=1}^{k} a_i (T_i(x) - \mathbb{E}_{\theta}[T_i]), a_i \in \mathbb{R} \right\}$

Since $\nabla \mathbb{E}_{\theta}[f] = \operatorname{Cov}_{\theta}(f, T)$, if $f \in \mathsf{T}_p$ the steepest direction is given by $f - \mathbb{E}_{\theta}[f]$ otherwise we take the projection of f onto T_p

24

$$\hat{f} = \sum_{i=1}^{k} \hat{a}_i (T_i(x) - \mathbb{E}_{\theta}[T_i]),$$

and obtain \hat{f} by solving a system of linear equations



















The Log-Likelihood Trick Assume the form $W_f(\xi) = \mathbb{E}_{x \sim P_{\xi}} \Big[w(f(x)) \Big].$
$\nabla_{\xi} W_{f}(\xi) = \nabla_{\xi} \mathbb{E}_{x \sim P_{\xi}} [w(f(x))]$ $= \nabla_{\xi} \int_{\Omega} w(f(x)) p(x \xi) dx$ $= \int_{\Omega} w(f(x)) \cdot \nabla_{\xi} p(x \xi) dx$ $= \int_{\Omega} w(f(x)) \cdot \nabla_{\xi} p(x \xi) \cdot \frac{p(x \xi)}{p(x \xi)} dx$ $= \int_{\Omega} w(f(x)) \cdot \frac{\nabla_{\xi} p(x \xi)}{p(x \xi)} \cdot p(x \xi) dx$ $= \mathbb{E}_{x \sim P_{\xi}} [w(f(x)) \cdot \nabla_{\xi} \log (p(x \xi))]$ The gradient of the expectation can be written as the expectation of a weighted gradient of the log likelihood

The Log-Likel	ihood Tri	ck			35		
$\nabla_{\xi} W$	$V_f(\xi) = \mathbb{E}_{\mathfrak{g}}$	$x \sim P_{\xi} \left[w \right]$	$(f(x))\cdot abla_{\xi}$	$\log(p(x \xi))$;))]		
The expect	ted value	can be	e estimate	d efficientl	у.		
Its Monte 0	Carlo esti	mate re	eads:				
$ abla_{\xi} W_f(\xi)$	$(\xi) \approx \frac{1}{N}_{x_1},$	$\sum_{\dots,x_N \sim P_i}$	$w(f(x_i))$	$\cdot \nabla_{\xi} \log(p)$	$(x_i \xi))$		
 Note: neith gradient of 	her the grant f .	adient	of W_f nor	its approx	imation r	equire	e the
L. Malaç	jò, T. Glasmache	rs, GECCO,	July 13, 2014				



Continuous Optimization with NES

37

- In the context of evolution strategies such a scheme was first proposed by Wierstra et al. in 2008.
- Original Natural Evolution Strategies (NES) approach: apply SNGD to optimize expected fitness $W_f(\xi) = \mathbb{E}_{x \sim P_{\xi}}[f(x)]$ with multi-variate Gaussian search distributions $\mathcal{N}(m, C)$.
- SNGD with Gaussians is rather easy since the Fisher matrix is a closed form term:

$$\mathcal{I}_{i,j} = \frac{\partial m^T}{\partial \xi_i} C^{-1} \frac{\partial m}{\partial \xi_j} + \frac{1}{2} \operatorname{tr} \left(C^{-1} \frac{\partial C}{\partial \xi_i} C^{-1} \frac{\partial C}{\partial \xi_j} \right)$$

 Practical versions of NES apply many performance enhancing techniques like rank-based utilities and non-uniform learning rates that complement the SNGD approach.

L. Malagò, T. Glasmachers, GECCO, July 13, 2014

Natural Gradients for Gaussian Distributions

• In \mathbb{R}^d we start with the Gaussian density

L. Malagò, T. Glasmachers, GECCO, July 13, 2014

$$p(x|\xi) = \frac{1}{\sqrt{(2\pi)^d \det(C)}} \exp\left(-\frac{1}{2}(x-m)^T C^{-1}(x-m)\right)$$

Its natural logarithm is

$$\log(p(x|\xi)) = -\frac{d}{2}\log(2\pi) - \frac{1}{2}\operatorname{tr}(\log(C)) - \frac{1}{2}(x-m)^{T}C^{-1}(x-m) .$$

38

- It is often beneficial to represent the covariance matrix with a factor: $C = AA^T$. For $x \sim \mathcal{N}(m, C)$ we introduce the transformed sample $z = A^{-1}(x m) \sim \mathcal{N}(0, I)$.
- In tailored coordinates

$$(m', A') = \left(m + A\delta, A\left(\mathbf{I} + \frac{1}{2}M\right)\right)$$

centered to the current distribution (m, A) the Fisher matrix w.r.t. the local parameters $\xi = (\delta, M)$ becomes the identity.

_

Natural Gradients for Gaussian Distributions • The (natural) gradient of the log density at $(\delta, M) = 0$ is $\tilde{\nabla}_{\delta} \log(p(x|\xi)) = A^{-1}(x-m) = z$ $\tilde{\nabla}_{M} \log(p(x|\xi)) = \frac{1}{2} (A(x-m)(x-m)^{T}A - I) = \frac{1}{2} (zz^{T} - I)$ • The stochastic (natural) gradient of W_{f} becomes $G_{\delta}(\xi) = \frac{1}{N} \sum_{i=1}^{N} f(x_{i}) \cdot z_{i}$ $G_{M}(\xi) = \frac{1}{2N} \sum_{i=1}^{N} f(x_{i}) \cdot (z_{i}z_{i}^{T} - I)$ • Tricks of the trade: replace "raw fitness" with "rank-based utility weights" $f(x_{1}) \leq \cdots \leq f(x_{N}) \rightarrow u_{1} \geq \cdots \geq u_{N}$ to achieve better invariance and faster convergence.









SNGD with Exponential Family

- Now assume bitstrings $\Omega = \{0, 1\}^n$.
- Then the probability simplex Δ and hence the statistical manifold Θ is finite dimensional.

44

46

- The sufficient statistics $T_i(x)$ are (square free) monomials.
- Each monomial characterizes a subset of bits the dependencies of which can be modeled.
- If the chosen model contains all interactions of variables in *f* then there is only one (global) optimum of W_f. The natural gradient will lead us there.

Dynamic Stochastic Relaxation

- Expected fitness W_f(ξ) = E_{x∼Pξ}[f(x)] is only one possible stochastic relaxation of f.
- We have already seen the generalization $W_f(\xi) = \mathbb{E}_{x \sim P_{\xi}}[w(f(x))]$ with a transformation w.

L. Malagò, T. Glasmachers, GECCO, July 13, 2014

- In IGO the weight function depends on the *f*-quantile under the current distribution P_{ξ_0} : $w(f(x)) = \tilde{w}(q_{\xi_0}^{-1}(f(x)))$, where $q_{\xi_0} : [0,1] \to \mathbb{R}$ encodes the quantiles of the distribution of f(x), $x \sim P_{\xi_0}$.
- E.g., $q_{\xi_0}(1/2)$ is the median of *f*-values, and $\tilde{w}(q) = 1$ for q < 1/2and $\tilde{w}(q) = 0$ for $q \ge 1/2$ encodes truncation (selection): only the better half of the distribution enters the update equation.
- The dynamic choice of $w = w(\xi_0)$ rescales f to a locally relevant range. It emphasizes local improvements relative to the current f distribution.

```
L. Malagò, T. Glasmachers, GECCO, July 13, 2014
```



L. Malagò, T. Glasmachers, GECCO, July 13, 2014

Dynamic Stochastic Relaxation

- Benefit 1: W_f becomes invariant under rank-preserving transformations of fitness values.
- Benefit 2: the rank-based utility weights of NES are obtained automatically in a principled manner.
- Drawback: the objective function $W_{f}^{\xi_{0}}(\xi)$ becomes dependent on the current distribution $\xi = \xi_{0}$

47

• This means that the following situation may exist in principle:

 $W_{f}^{\xi_{1}}(\xi_{2}) > W_{f}^{\xi_{1}}(\xi_{1})$ $W_{f}^{\xi_{2}}(\xi_{3}) > W_{f}^{\xi_{2}}(\xi_{2})$ $W_{f}^{\xi_{3}}(\xi_{1}) > W_{f}^{\xi_{3}}(\xi_{3})$

and the "optimization" turns around in circles...

Provably, in important special cases this does not happen.

Vector Field, ODE, and Flow

- Algorithmic and hence discrete time optimization is only one possibility.
- The natural gradient $\tilde{\nabla}W_f(\xi)$ defines the vector field $V:\Xi \to T\Xi$ via $V(\xi) = \tilde{\nabla}W_f(\xi)$.
- This vector field is naturally associated with the differential equation $\dot{\gamma}(t) = V(\gamma(t))$ with solution curves $\gamma : \mathbb{R} \to \Xi$.
- The solution curves are collected in the flow $\varphi(\xi, t)$.

L. Malagò, T. Glasmachers, GECCO, July 13, 2014

- Note 1: just like the natural gradient itself this flow is deterministic. This is achieved in the limit of infinite samples in the MC approximation, corresponding to infinite population size.
- Note 2: in each point the flow moves tangential to the vector field. This corresponds to re-evaluating the gradient after an infinitesimal step, or to an infinitesimal leaning rate in the gradient descent procedure.

SGD Algorithms

- An SGD algorithm is a two-fold approximation of the flow:
 - it discretizes time and performs Euler steps,
 - it relies on a stochastic gradient based on sampling.
- If the IGO objective is in use then such algorithms have been termed IGO algorithms.

49

• NES is a rather pure example of an IGO algorithm.

L. Malagò, T. Glasmachers, GECCO, July 13, 2014

- Surprisingly many established algorithm fall into this category or have a close connection to SGD (or IGO) algorithms
- This derivation of randomized direct search algorithms opens a new perspective: the EA is an approximation of a much more simple and beautiful mathematical object, namely the deterministic, time continuous gradient flow, thought of as an idealized EA.

50

48

Connection to CMA-ES

- Modern CMA-ES is a versatile algorithm that comes in many variants, e.g., with mirrored sampling, a restart strategy, a noise handling mechanism, an elitist variant, constraint handling, ...
- Even the basic reference algorithm is rather complex.
- It updates the distribution mean with global weighted recombination:

L. Malagò, T. Glasmachers, GECCO, July 13, 2014

$$m \leftarrow \sum_{i=1}^{\mu} w_i \cdot x_i$$
.

 For the purpose of simplified analysis we drop cumulative step size adaptation and consider only the so-called rank-μ update

$$C \leftarrow (1 - \gamma_C) \cdot C + \gamma_C \cdot \sum_{i=1}^{\mu} w_i \cdot (x_i - m)(x_i - m)^T$$
.

Solution to CMA-ES • Both equations can be written as updates $m \leftarrow m + \gamma_m \cdot \sum_{i=1}^N w_i \cdot (x_i - m)$ $C \leftarrow C + \gamma_C \cdot \sum_{i=1}^N w_i \cdot ((x_i - m)(x_i - m)^T - C)$ with fixed learning rate $\gamma_m = 1$ and $w_i = 0$ for $i > \mu$. • The change of coordinates $C = AA^T$, x = Az + m reveals: $m \leftarrow m + \gamma_m \cdot A \cdot \sum_{i=1}^N w_i \cdot z_i$ $C \leftarrow C + \gamma_C \cdot A \cdot \left(\sum_{i=1}^N w_i \cdot (z_i z_i^T - I)\right) \cdot A^T$ • This is essentially the IGO/NES SGD update (see Akimoto 2010). L Malagh, T. Glasmachers, GECCO, July 13, 2014

Connection to Maximum Likelihood Estimation

The CMA-ES update equations can be written as

$$m \leftarrow (1 - \gamma_m) \cdot m + \gamma_m \cdot \hat{m}_{\mathsf{ML}} ,$$

$$C \leftarrow (1 - \gamma_C) \cdot C + \gamma_C \cdot C_{\mathsf{ML}}$$
.

52

54

- Here $\hat{m}_{ML} = \sum_{i=1}^{N} w_i \cdot x_i$ is the weighted Maximum Likelihood (ML) estimator of m.
- The term

$$\hat{C}_{\mathsf{ML}} = \sum_{i=1}^{\mu} w_i (x_i - m) (x_i - m)^T$$

is the weighted ML estimator of C, provided that m remains fixed.

- For a learning rate of $\gamma_m = \gamma_C = 1$ we obtain $m \leftarrow \hat{m}_{\text{ML}}$ and $C \leftarrow \hat{C}_{\text{ML}}$.
- Due to the large learning rate this works in practice only for a large enough sample size μ.

L. Malagò, T. Glasmachers, GECCO, July 13, 2014

Convergence Results

- Convergence has been analyzed for many optimization algorithms. For most realistic evolutionary algorithm this is a hard task.
- One of the benefits of the gradient flow (as an idealization of actual algorithms) is that it is much easier to analyze.
- Question: Do all trajectories converge to the optimum?
- More formally, let δ_{x^*} denote the Dirac peak over an (isolated) optimum $x^* \in \Omega$. Does it hold $\lim_{t\to\infty} P_{\phi(\xi,t)} = \delta_{x^*}$ for all initial conditions ξ ?
- Other notions of convergence are possible, e.g., convergence of expected fitness.
- Note: convergence of the flow does not directly imply convergence of stochastic approximate algorithms!

L. Malagò, T. Glasmachers, GECCO, July 13, 2014



 This is because the fact that the optimization will stop at all is meaningless in practice; instead we need to know (roughly) how long it takes.

L. Malagò, T. Glasmachers, GECCO, July 13, 2014

Representability: The Marginal Polytope

Convergence of the gradient flow to the optimum has a hidden prerequisite:

55

- The statistical model must be sufficiently rich to focus the probability mass on an optimum, or more exactly, on a subset of the optimal set $\Omega^* \subset \Omega$.
- In general this property is a non-trivial prerequisite for optimization with a stochastic relaxation.
- As such it is not a prerequisite for convergence of the flow to an optimal distribution within (the closure of) the statistical model.
- However, it is a prerequisite for convergence to an optimal distribution, possibly outside the family, and hence to an optimum of the original problem $\min_{x \in \Omega} f(x)$.

Representability: The Marginal Polytope

Assume an exponential family

$$p(x|\xi) = \exp\left(\sum_{i=1}^{k} \xi_i T_i(x) - \psi(\xi)\right)$$

56

of distributions. Then this property has a geometric interpretation in distribution space.

- Distribution are restricted to the *marginal polytope*, which is the convex hull of the canonical statistics $T_i(x)$.
- For discrete search spaces Ω it is a sub-polytope of the probability simplex.

L. Malagò, T. Glasmachers, GECCO, July 13, 2014

• The optimum of the stochastically relaxed problem describes an optimum of $f: \Omega \to \mathbb{R}$ iff a subset of $S \subset \Omega^*$ corresponds to an *exposed face* W_f of the marginal polytope, i.e., if $S = T^{-1}(W_f) \subset \Omega^*$.

Representability: The Marginal Polytope

- 57
- Consider the family of Gaussian distributions $\mathcal{N}(m, \sigma^2)$ with adaptive mean m and step size σ .
- The Dirac distribution δ_x of a point $x \in \mathbb{R}^d$ is obtained in the limit

$$\delta_x = \lim_{\sigma \to 0} \mathcal{N}(x, \sigma^2) \quad .$$

- Hence the probability mass can be focused to an arbitrarily small neighborhood of an optimal point $x^* \in \Omega^* \subset \mathbb{R}^d$.
- As a consequence convergence is possible also for all supersets of distributions, i.e., multi-variate Gaussians.















Summary

 Information geometry provides update equations for optimization from first principles.

68

70

- This often amounts to stochastic natural gradient descent applied to a stochastically relaxed problem.
- This is an approximation to an optimization flow.

L. Malagò, T. Glasmachers, GECCO, July 13, 2014

- It can help the analysis of existing algorithms like CMA-ES.
- It is a generic design principle for optimization algorithms on any search space and any family of search distributions.
- Dedicated algorithms such as NES were built on this principle.
- Algorithms respecting the information geometry of their search distributions are among the top performers.

Summary

 However, randomized algorithms deviate widely from the idealized, deterministic gradient flow. These stochastic effects are outside the framework, they must be controlled by other means.

69

- Information geometric tools must be augmented with "orthogonal" tools for control of stochastic effects—together they provide a modern perspective on EA research.
- The same problem decomposition is a promising route for theoretical analysis: the gradient flow is becoming a well-investigated object, while more traditional tools (Markov chain analysis, etc.) may be necessary to connect it to real EAs.

L. Malagò, T. Glasmachers, GECCO, July 13, 2014

References 1 Ollivier et al. Information-Geometric Optimization Algorithms: A Unifying Picture via Invariance Principles. arXiv:1106.3708, 2011. 2 Wierstra et al. Natural Evolution Strategies. Congress on Evolutionary Computation (CEC), 2008. 3 Glasmachers et al. Exponential Natural Evolution Strategies. Genetic and Evolutionary Computation Conference (GECCO), 2010, 4 Akimoto et al. Bidirectional relation between CMA evolution strategies and natural evolution strategies. Parallel Problem Solving from Nature (PPSN) XI 2010. 5 Akimoto et al. Convergence of the Continuous Time Trajectory of Isotropic Evolution Strategies on Monotonic C2-composite Functions. Parallel Problem Solving from Nature (PPSN) XII, 2012. 6 Glasmachers. Convergence of the IGO-Flow of Isotropic Gaussian Distributions on Convex Quadratic Problems. Parallel Problem Solving from Nature (PPSN) XII, 2012. 7 Malagò et al. Natural Gradient, Fitness Modelling and Model Selection: A Unifying Perspective. Congress on Evolutionary Computation (CEC), 2013. 8 Malagò et al. Towards the geometry of estimation of distribution algorithms based on the exponential family. FOGA, 2011 9 Beyer. Convergence Analysis of Evolutionary Algorithms which are Based on the Paradigm of Information Geometry. Personal communication L. Malagò, T. Glasmachers, GECCO, July 13, 2014

