

	Stochastic and combinatorial optimization	Stochastic and combinatorial optimization	Simulated annealing																
<div>EM Algorithm, Stochastic Optimization</div> <div>732A90 Computational Statistics Krzysztof Bartoszek (krzysztof.bartoszek@liu.se)</div> <div>28 II 2018 (A37) Department of Computer and Information Science Linköping University</div>	<ul style="list-style-type: none"><li>So far: Unconstrained optimization<ul style="list-style-type: none"><li>Predictor variables are continuous</li><li>Response function is differentiable</li></ul></li><li>We discussed Steepest descent, Newton, BFGS, CG</li><li>But: predictors can be discrete (scheduling problems, travelling salesman)</li><li>But: outcome can be discrete, noisy or multi-modal</li></ul>	<p>Given a (large) set of states <math>S</math>, find</p> $\min_{s \in S} f(s)$ <ul style="list-style-type: none"><li>Exhaustive search (shortest path algorithm)</li><li>Often exhaustive search is NP-hard (TSP)</li><li>Alternative: stochastic methods random search</li></ul>	<p>Motivation from physics: cooling of metal</p> <ul style="list-style-type: none"><li>Parameters: Energy of metal (decreasing, but not strictly monotonic) Temperature (decreasing)</li><li>Aim: find global minimum energy</li></ul>																
Simulated annealing	Simulated annealing	Simulated annealing: TSP example	Genetic algorithm																
<pre>0. Set <math>k = 1</math> and initialize state <math>s</math>. 1. Compute the temperature <math>T(k)</math>. 2. Set <math>i = 0</math> and <math>j = 0</math>. 3. Generate a new state <math>r</math> and compute <math>\delta f = f(r) - f(s)</math>. 4. Based on <math>\delta f</math>, decide whether to move from state <math>s</math> to state <math>r</math>.    If <math>\delta f \leq 0</math>,      accept state <math>r</math>;    otherwise,      accept state <math>r</math> with a probability <math>P(\delta f, T(k))</math>.    If state <math>r</math> is accepted, set <math>s = r</math> and <math>i = i + 1</math>. 5. If <math>i</math> is equal to the limit for the number of successes at a given temperature,    go to step 1. 6. Set <math>j = j + 1</math>. If <math>j</math> is less than the limit for the number of iterations at    given temperature, go to step 3. 7. If <math>i = 0</math>,    deliver <math>s</math> as the optimum; otherwise,    if <math>k &lt; k_{max}</math>,      set <math>k = k + 1</math> and go to step 1;    otherwise,      issue message that        'algorithm did not converge in <math>k_{max}</math> iterations'.</pre>	<ul style="list-style-type: none"><li><a href="https://www.youtube.com/watch?v=1aq_Fpr4KZc">https://www.youtube.com/watch?v=1aq_Fpr4KZc</a></li><li>Generating new state:<ul style="list-style-type: none"><li>Continuous: choose a new point <math>a</math> (random) distance from the current one</li><li>Discrete: similar or some rearrangement</li></ul></li><li>Selection probability: e.g. <math>\exp(-\delta f(x)/T)</math>: decreasing with <math>f(x)</math>, increasing with <math>T</math></li><li>Temperature function: constant, proportional to <math>k</math>, or<math display="block">T(k+1) = b(k)T(k), \quad b(k) = (\log(k))^{-1}</math></li></ul> <p><b>Remember:</b> A smaller value is better than one on the path to the global minimum! Always keep track of smallest found.</p>	<p>Assume constant temperature</p> <pre>1: Choose initial configuration <math>(Town_1, \dots, Town_n)</math> 2: <math>k = 1</math> 3: <b>while</b> <math>k &lt; k_{max} + 1</math> <b>do</b> 4:   Generate new configuration by rearrangement.       <math>(1, 2, 3, 4, 5, 6, 7, 8, 9) \rightarrow (1, 6, 5, 4, 3, 2, 7, 8, 9)</math>       <math>(1, 2, 3, 4, 5, 6, 7, 8, 9) \rightarrow (1, 7, 8, 2, 3, 4, 5, 6, 9)</math> 5:   Measure difference in path length (<math>\delta f</math>) between old and new configuration 6:   <b>if</b> shorter path found <b>then</b> 7:     accept it 8:   <b>else</b> 9:     accept it with probability <math>P(\delta f)</math> 10:  <b>end if</b> 11:  <math>k++</math> 12: <b>end while</b></pre>	<ul style="list-style-type: none"><li>Inspiration from evolutionary theory: survival of the fittest</li><li>Variables=genotypes</li><li>Observation =organism, characterized by genetic code</li><li>State space=population of organisms</li><li>Objective function=fitness of organism</li></ul> <p>New points are obtained from old points by crossover and mutation, the population only retains the fittest organisms (with better objective function).</p> <p><a href="https://en.wikipedia.org/wiki/List_of_genetic_algorithm_applications">https://en.wikipedia.org/wiki/List_of_genetic_algorithm_applications</a></p>																
Genetic algorithm	Genetic algorithm	Genetic algorithm: TSP example	Genetic algorithm: Mutations																
<p>Encoding points</p> <ul style="list-style-type: none"><li>Enumerate each element of the state space, <math>S</math></li><li>Code for observation <math>i</math> is binary representation of <math>i</math> (or something else)</li></ul> <p>Mutation and recombination rules</p> <table><tr><th>Generation <math>k</math></th><th>Generation <math>k + 1</math></th></tr><tr><td><math>x_1^{(k)}</math> 11001001</td><td>Crossover <math>x_1^{(k+1)}</math> 11011030</td></tr><tr><td><math>x_2^{(k)}</math> 00111010</td><td><math>x_2^{(k+1)}</math> 11011030</td></tr><tr><td><math>x_3^{(k)}</math> 11101011</td><td>Inversion <math>x_3^{(k+1)}</math> 11011031</td></tr><tr><td><math>x_4^{(k)}</math> 11101011</td><td>Mutation <math>x_4^{(k+1)}</math> 10111031</td></tr><tr><td><math>x_5^{(k)}</math> 11101011</td><td>Cycle <math>x_5^{(k+1)}</math> 11101031</td></tr></table>	Generation $k$	Generation $k + 1$	$x_1^{(k)}$ 11001001	Crossover $x_1^{(k+1)}$ 11011030	$x_2^{(k)}$ 00111010	$x_2^{(k+1)}$ 11011030	$x_3^{(k)}$ 11101011	Inversion $x_3^{(k+1)}$ 11011031	$x_4^{(k)}$ 11101011	Mutation $x_4^{(k+1)}$ 10111031	$x_5^{(k)}$ 11101011	Cycle $x_5^{(k+1)}$ 11101031	<pre>0. Determine a representation of the problem, and define an initial population, <math>x_1^{(0)}, x_2^{(0)}, \dots, x_n^{(0)}</math>. Set <math>k = 0</math>. 1. Compute the objective function (the "fitness") for each member of the population, <math>f(x_i^{(k)})</math> and assign probabilities <math>p_i</math> to each item in the population, perhaps proportional to its fitness. 2. Choose (with replacement) a probability sample of size <math>m \leq n</math>. This is the reproducing population. 3. Randomly form a new population <math>x_1^{(k+1)}, x_2^{(k+1)}, \dots, x_n^{(k+1)}</math> from the reproducing population, using various mutation and recombination rules (see Table 6.2). This may be done using random selection of the rule for each individual of pair of individuals. 4. If convergence criteria are met, stop, and deliver <math>\arg\min_{x_i^{(k+1)}} f(x_i^{(k+1)})</math> as the optimum; otherwise, set <math>k = k + 1</math> and go to step 1.</pre>	<p>Encoding and crossover</p> <ul style="list-style-type: none"><li>Encode tours as <math>A_1, \dots, A_n</math> but</li></ul> <table><tr><td>Parent 1: FAB/ECGD</td><td>Parent 2: DEA/CGBF</td></tr><tr><td>Child: FAB/CGBF</td><td>Child: DEA/ECGD</td></tr></table> <p>Instead</p> <ul style="list-style-type: none"><li>Remove FAB from DEACGBF <math>\rightarrow</math> DECG. Child becomes FABDECG.</li><li>Second child will be by taking prefix from Parent 2: DEAFBCG</li></ul>	Parent 1: FAB/ECGD	Parent 2: DEA/CGBF	Child: FAB/CGBF	Child: DEA/ECGD	<ul style="list-style-type: none"><li>If a population is small and only crossover, the input domain becomes limited and may converge to a local minimum.</li><li>Large initial populations are computationally heavy.</li><li>Mutations allow one to explore more of <math>S</math>: jump out of local minimum.</li><li>In TSP: mutation move a city in the tour to another position.</li><li>Reproduction: Among <math>m</math> tours selected at step 2, two best are selected for reproduction, two worst replaced by children.</li><li>If <math>m</math> is large, some tours might never be parents, global solution may be missed. Random chance of reproduction?</li><li>Mutation probability is usually small (unless you want to jump wildly)</li></ul>
Generation $k$	Generation $k + 1$																		
$x_1^{(k)}$ 11001001	Crossover $x_1^{(k+1)}$ 11011030																		
$x_2^{(k)}$ 00111010	$x_2^{(k+1)}$ 11011030																		
$x_3^{(k)}$ 11101011	Inversion $x_3^{(k+1)}$ 11011031																		
$x_4^{(k)}$ 11101011	Mutation $x_4^{(k+1)}$ 10111031																		
$x_5^{(k)}$ 11101011	Cycle $x_5^{(k+1)}$ 11101031																		
Parent 1: FAB/ECGD	Parent 2: DEA/CGBF																		
Child: FAB/CGBF	Child: DEA/ECGD																		
EM algorithm	EM algorithm	EM algorithm: R	EM algorithm: R																
<p><b>Fundamental algorithm</b> of computational statistics!</p> <p>Model depends on the data which are observed (known) <math>\mathbf{Y}</math> and <b>latent</b> (unobserved) data <math>\mathbf{Z}</math>.</p> <p>The data's (both <math>\mathbf{Y}</math>'s and <math>\mathbf{Z}</math>'s) distribution depends on some parameters <math>\theta</math>.</p> <p><b>AIM:</b> Find MLE of <math>\theta</math>.</p> <ul style="list-style-type: none"><li>All data is known: Apply unconstrained optimization (discussed in Lecture 2)</li><li>Unobserved data<ul style="list-style-type: none"><li>Sometimes it is possible to look at the marginal distribution of the observed data.</li><li>Otherwise: <b>EM algorithm</b></li></ul></li></ul>	<p>Let</p> $Q(\theta, \theta^k) = \int \log p(\mathbf{Y}, \mathbf{z}   \theta) p(\mathbf{z}   \mathbf{Y}, \theta^k) d\mathbf{z} = \mathbb{E} \left[ \log \text{lik}(\theta   \mathbf{Y}, \mathbf{Z})   \theta^k, \mathbf{Y} \right]$ <pre>1: <math>k = 0, \theta^0 = \theta^0</math> 2: <b>while</b> Convergence not attained <b>and</b> <math>k &lt; k_{max} + 1</math> <b>do</b> 3:   <b>E-step:</b> Derive <math>Q(\theta, \theta^k)</math> 4:   <b>M-step:</b> <math>\theta^{k+1} = \arg\max_{\theta} Q(\theta, \theta^k)</math> 5:   <math>k++</math> 6: <b>end while</b></pre> <p><b>Example:</b> Normal data with missing values (but here analytical approach is also possible)</p>	732A90_ComputationalStatisticsVT2019Lecture06codeSlide15.R	<pre>&gt; Y&lt;-rnorm(100) &gt; Y[sample(1:length(Y),20,replace=FALSE)]&lt;-NA &gt; EM.Norm(Y,0.0001,100) [1] 1.0000 0.1000 -997.5705 [1] 0.1341894 1.3227895 -128.2789837 [1] -0.03897274 1.38734070 -126.86036252 [1] -0.07360517 1.39307050 -126.80801589 [1] -0.08053165 1.39392861 -126.80593837 [1] -0.08191695 1.39408871 -126.80585537 &gt; mean(Y,na.rm=TRUE) [1] -0.08226328 &gt; var(Y,na.rm=TRUE) [1] 1.411775</pre> <p>Notice: can be done by studying marginal distribution of observed data.</p>																

**Mixture models**  $Z$  is a latent variable,  $P(Z = k) = \pi_k$

- Mixed data comes from different sources (e.g. for regression, classification)
- Clustering
  - 1 Density in each cluster is normally distributed.
  - 2 Cluster label is latent (we do not know what are the chances an observation is from the given cluster)

$$p(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x|\mu_k, \Sigma_k) \quad (\text{informally})$$

Direct MLE leads to numerical problems.  
Introduce latent class variables and use EM.

1. Initialize the means  $\mu_0$ , covariances  $\Sigma_0$  and mixing coefficients  $\pi_0$ , and evaluate the initial value of the log likelihood.
2. **E step.** Evaluate the responsibilities using the current parameter values

$$r_i^{(t+1)} = \frac{\pi_k \mathcal{N}(x_i|\mu_k, \Sigma_k)}{\sum_{k=1}^K \pi_k \mathcal{N}(x_i|\mu_k, \Sigma_k)} \tag{9.23}$$

3. **M step.** Re-estimate the parameters using the current responsibilities

$$\mu_k^{(t+1)} = \frac{1}{N_k} \sum_{i=1}^N r_i^{(t+1)} x_i \tag{9.24}$$

$$\Sigma_k^{(t+1)} = \frac{1}{N_k} \sum_{i=1}^N r_i^{(t+1)} (x_i - \mu_k^{(t+1)})(x_i - \mu_k^{(t+1)})^T \tag{9.25}$$

$$\pi_k^{(t+1)} = \frac{N_k}{N} \tag{9.26}$$

where 
$$N_k = \sum_{i=1}^N r_i^{(t+1)} \tag{9.27}$$

4. Evaluate the log likelihood 
$$\ln p(X|\mu, \Sigma, \pi) = \sum_{i=1}^N \ln \left( \sum_{k=1}^K r_i^{(t+1)} \pi_k \mathcal{N}(x_i|\mu_k, \Sigma_k) \right) \tag{9.28}$$

and check for convergence of either the parameters or the log likelihood. If the convergence criterion is not satisfied return to step 2.

Source: Andrew Ng, Machine Learning

$$E x_{nk} = \gamma(x_{nk})$$

- Random walk over the state space in search of minimum
- 1 Follow decreasing path
  - 2 **BUT** with a certain probability go to higher values, to avoid local minima traps.
  - 3 **Never forget** best found conformation!
  - 4 Simulated annealing, Genetic algorithm, **EM algorithm**, Stochastic gradient descent (see 2016 slides)