

EM Algorithm, Stochastic Optimization

732A90

Computational Statistics

Krzysztof Bartoszek
(krzysztof.bartoszek@liu.se)

7 XII 2020 (Zoom)

Department of Computer and Information Science
Linköping University

Stochastic and combinatorial optimization

- So far: Unconstrained optimization
 - Predictor variables are continuous
 - Response function is differentiable
- We discussed Steepest descent, Newton, BFGS, CG
- But: predictors can be discrete
(scheduling problems, travelling salesman)
- But: outcome can be discrete, noisy or multi-modal

Given a (large) set of states S , find

$$\min_{s \in S} f(s)$$

- Exhaustive search (shortest path algorithm)
- Often exhaustive search is NP-hard (TSP)
- Alternative: stochastic methods
random search

Motivation from physics: cooling of metal

- Parameters:
Energy of metal
(decreasing, but not strictly monotonic)
Temperature (decreasing)
- Aim: find global minimum energy

Simulated annealing

0. Set $k = 1$ and initialize state s .
 1. Compute the temperature $T(k)$.
 2. Set $i = 0$ and $j = 0$.
 3. Generate a new state r and compute $\delta f = f(r) - f(s)$.
 4. Based on δf , decide whether to move from state s to state r .
If $\delta f \leq 0$,
 accept state r ;
otherwise,
 accept state r with a probability $P(\delta f, T(k))$.
If state r is accepted, set $s = r$ and $i = i + 1$.
 5. If i is equal to the limit for the number of successes at a given temperature, go to step 1.
 6. Set $j = j + 1$. If j is less than the limit for the number of iterations at given temperature, go to step 3.
 7. If $i = 0$,
 deliver s as the optimum; otherwise,
 if $k < k_{\max}$,
 set $k = k + 1$ and go to step 1;
 otherwise,
 issue message that
 ‘algorithm did not converge in k_{\max} iterations’.
-

Simulated annealing

- https://www.youtube.com/watch?v=iaq_Fpr4KZc
- Generating new state:
 - Continuous: choose a new point a (random) distance from the current one
 - Discrete: similar or some rearrangement
- Selection probability: e.g $\exp(-\delta f(x)/T)$: decreasing with $f(x)$, increasing with T
- Temperature function: constant, proportional to k , or

$$T(k+1) = b(k)T(k), \quad b(k) = (\log(k))^{-1}$$

Remember: A smaller value is better than one on the path to the global minimum! Always keep track of smallest found.

Simulated annealing: TSP example

Assume constant temperature

- 1: Choose initial configuration $(Town_1, \dots, Town_n)$
- 2: $k = 1$
- 3: **while** $k < k_{max} + 1$ **do**
- 4: Generate new configuration by rearrangement,
 $(1, 2, 3, 4, 5, 6, 7, 8, 9) \rightarrow (1, 6, 5, 4, 3, 2, 7, 8, 9)$
 $(1, 2, 3, 4, 5, 6, 7, 8, 9) \rightarrow (1, 7, 8, 2, 3, 4, 5, 6, 9)$
- 5: Measure difference in path length (δf) between old and new configuration
- 6: **if** shorter path found **then**
- 7: accept it
- 8: **else**
- 9: accept it with probability $P(\delta f)$
- 10: **end if**
- 11: $k++$
- 12: **end while**

Genetic algorithm

- Inspiration from evolutionary theory: survival of the fittest
- Variables=genotypes
- Observation=organism, characterized by genetic code
- State space=population of organisms
- Objective function=fitness of organism

New points are obtained from old points by crossover and mutation, the population only retains the fittest organisms (with better objective function).

https://en.wikipedia.org/wiki/List_of_genetic_algorithm_applications

Genetic algorithm

Encoding points

- 1 Enumerate each element of the state space, S
- 2 Code for observation i is binary representation of i (or something else)

Mutation and recombination rules

Generation k	Generation $k + 1$
----------------	--------------------

Crossover

$$\begin{array}{l} x_i^{(k)} \text{ 11001001} \\ x_j^{(k)} \text{ 00111010} \end{array} \rightarrow x_i^{(k+1)} \text{ 11011010}$$

Inversion

$$x_i^{(k)} \text{ 11101011} \rightarrow x_i^{(k+1)} \text{ 11010111}$$

Mutation

$$x_i^{(k)} \text{ 11101011} \rightarrow x_i^{(k+1)} \text{ 10111011}$$

Clone

$$x_i^{(k)} \text{ 11101011} \rightarrow x_i^{(k+1)} \text{ 11101011}$$

Genetic algorithm

0. Determine a representation of the problem, and define an initial population, $x_1^{(0)}, x_2^{(0)}, \dots, x_n^{(0)}$. Set $k = 0$.
1. Compute the objective function (the “fitness”) for each member of the population, $f(x_i^{(k)})$ and assign probabilities p_i to each item in the population, perhaps proportional to its fitness.
2. Choose (with replacement) a probability sample of size $m \leq n$. This is the reproducing population.
3. Randomly form a new population $x_1^{(k+1)}, x_2^{(k+1)}, \dots, x_n^{(k+1)}$ 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 or pair of individuals.
4. If convergence criteria are met, stop, and deliver $\arg \min_{x_i^{(k+1)}} f(x_i^{(k+1)})$ as the optimum; otherwise, set $k = k + 1$ and go to step 1.

Genetic algorithm: TSP example

Encoding and crossover

- Encode tours as A_1, \dots, A_n but

Parent 1: FAB|ECGD Parent 2: DEA|CGBF

Child: FAB|CGBF Child: DEA|ECGD

Instead

- 1 Remove FAB from DEACGBF \rightarrow DECG.
Child becomes FABDECG.
- 2 Second child will be by taking prefix from Parent 2:
DEAFBCG

Genetic algorithm: Mutations

- If a population is small and only crossover: the input domain becomes limited and may converge to a local minimum.
- Large initial populations are computationally heavy.
- Mutations allow one to explore more of S : jump out of local minimum.
- In TSP: mutation move a city in the tour to another position.
- Reproduction: Among m tours selected at step 2, two best are selected for reproduction, two worst replaced by children.
- If m is large, some tours might never be parents, global solution may be missed. Random chance of reproduction?
- Mutation probability is usually small (unless you want to jump wildly)

EM algorithm

Fundamental algorithm of computational statistics!

Model depends on the data which are observed (known) **Y** and **latent** (unobserved) data **Z**.

The data's (**both Y's and Z's**) distribution depends on some parameters θ .

AIM: Find MLE of θ .

- All data is known: Apply unconstrained optimization (discussed in Lecture 2)
- Unobserved data
 - **Sometimes** it is possible to look at the marginal distribution of the observed data.
 - Otherwise: **EM algorithm**

EM algorithm

Let

$$Q(\theta, \theta^k) = \int \log p(\mathbf{Y}, \mathbf{z}|\theta) p(\mathbf{z}|\mathbf{Y}, \theta^k) d\mathbf{z} = \mathbb{E} \left[\log \text{lik}(\theta|\mathbf{Y}, \mathbf{Z}) | \theta^k, \mathbf{Y} \right]$$

- 1: $k = 0, \theta^0 = \theta^0$
- 2: **while** Convergence not attained **and** $k < k_{max} + 1$ **do**
- 3: **E-step:** Derive $Q(\theta, \theta^k)$
- 4: **M-step:** $\theta^{k+1} = \operatorname{argmax}_{\theta} Q(\theta, \theta^k)$
- 5: $k++$
- 6: **end while**

Example: Normal data with missing values (but here analytical approach is also possible)

732A90_ComputationalStatisticsHT2020_Lecture06codeSlide15.R

EM algorithm: R

```
> Y<-rnorm(100)
> Y[sample(1:length(Y),20,replace=FALSE)]<-NA
> EM.Norm(Y,0.0001,100)
[1] 1.0000 0.1000 -997.5705
[1] 0.1341894 1.3227095 -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
> mean(Y,na.rm=TRUE)
[1] -0.08226328
> var(Y,na.rm=TRUE)
[1] 1.411775
```

Notice: can be done by studying marginal distribution of observed data.

EM algorithm: Applications

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)

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

Direct MLE leads to numerical problems.

Introduce latent class variables and use EM.

EM algorithm: Gaussian mixtures

1. Initialize the means μ_k , covariances Σ_k and mixing coefficients π_k , and evaluate the initial value of the log likelihood.
2. **E step.** Evaluate the responsibilities using the current parameter values

$$\gamma(z_{nk}) = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \mu_j, \Sigma_j)}. \quad (9.23)$$

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

$$\mu_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n \quad (9.24)$$

$$\Sigma_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \mu_k^{\text{new}}) (\mathbf{x}_n - \mu_k^{\text{new}})^T \quad (9.25)$$

$$\pi_k^{\text{new}} = \frac{N_k}{N} \quad (9.26)$$

where

$$N_k = \sum_{n=1}^N \gamma(z_{nk}). \quad (9.27)$$

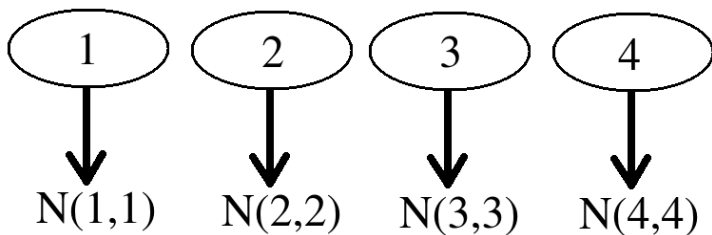
4. Evaluate the log likelihood

$$\ln p(\mathbf{X} | \mu, \Sigma, \pi) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k) \right\} \quad (9.28)$$

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

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

Gaussian mixtures: example



$$P(1) = P(2) = P(3) = P(4) = 0.25$$

- ➊ draw class $Z \in \{1, 2, 3, 4\}$ uniformly
- ➋ draw normal distribution $\mathcal{N}(Z, Z)$ with density $\phi_{Z,Z}(\cdot)$

We can write the mixture density as

$$f(x) = 0.25\phi_{1,1}(x) + 0.25\phi_{2,2}(x) + 0.25\phi_{3,3}(x) + 0.25\phi_{4,4}(x).$$

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)