732A90 Computational Statistics

Krzysztof Bartoszek (krzysztof.bartoszek@liu.se)

24 I 2019 (P42) Department of Computer and Information Science Linköping University

- Introduction
- Mathematical definition of problem
- 1D optimization
- \bullet kD optimization
- R code examples

Nearly everything is optimization!

- Chemistry
- Physics
 Economics, Industry
- Engineering

BUT EVEN

- Your mobile price planCourse scheduling
- Your lunch choice
- STATISTICS
- Fit parameters to data
 Propose optimal decision

ANY BIOLOGICAL **ORGANISM**

YOU

Industry

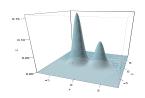
How to produce a cylindrical (WHY?) 0.5L beer can so it requires minimum material?

Given a certain product minimize e.g. material usage, production effort while still meeting consumer requirements.

Economics/Logistics

- Travelling Salesman Problem
- Windmills
- Flight schedule (especially "cheap" airlines)

Statistics Maximize likelihood, model fitting



Maximal likelihood

An i.i.d. sample (X_1,\ldots,X_n) is drawn from a probability distribution $P(X\ominus)$, where \ominus is an unknown parameter so

The joint probability of all the observations is

$$P(X_1, ..., X_n | \Theta) = \prod_{i=1}^n P(X_i | \Theta).$$

Find Θ that maximizes $P(X_1, \dots, X_n | \Theta)$.

Mathematical formulation

The goal is to minimize (maximize)

Objective function: $f(\theta)$ (reproduction, chances of survival, quality of life, cost, profit, likelihood, fit to data)

depending on

Parameters or Unknowns θ (reproduction strategy, resource utilization, consumer choices, height & diameter, production, raw material choice, service times, route, tlight routes/times, parameters)

Mathematical form

$$\min_{\theta \in \Theta} f(\theta) \text{ subject to } \frac{c_i(\theta) = 0, \quad i \in E}{c_i(\theta) \geq 0, \quad i \in I}$$

QUESTION: What should we do if we are interested in

QUESTION: What should we do if the constraints are

- Available environment
- Production: Factories (F₁, F₂), retail outlets (R₁, R₂, R₃), Production: Pactories $\{r_1, r_2\}$, retail outlets $\{n_1, n_2, n_3\}$, cost of shipping $i \to j$: c_{ij} , production a_i per week, requirement b_j per week **to optimize**: x_{ij} amount shipped $i \to j$ per week

$$\begin{aligned} \min_{x_i \in \mathcal{X}_{ij}} \sum_{t_{ij}} c_{tj} x_{ij} & & \text{minimize shipping costs} \\ \sum_{j=1}^{\infty} x_{ij} \leq a_i, j = 1, 2 & & \text{production capacity} \\ \sum_{i=1}^{\infty} x_{ij} \geq b_j, j = 1, 2, 3 & & \text{demand} \end{aligned}$$

 ${\bf Question:}$ What would happen if we drop demand constraint?

 \bullet ML: often no constraints

- Split into pairs/triplets/quadruples
 Think of some human anatomy part/organ:
 What is its function?
 What is its function?
 What could it have been optimized for over the course of
 - Is it still under selection?
 What constraints was and is it under?
- What constraints was and is it under?
 Think of a situation where optimization is needed in your own student/professional/personal/financial situation.
 State the problem in terms of
 Objective function
 Parameters
 Constraints
 Does it have a trivial solution?

- 10 minutes

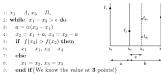
- Constrained optimization
 - Lagrange multipliers, linear programming
 E.g. LASSO
 Not this lecture!
- \bullet Unconstrained optimization

 - Steepest descent
 Newton method
 Quasi-Newton-Methods
 Conjugate gradients
- Why are there different methods?

- Function of a single parameter, find minimum
- What algorithm would you suggest?
- $\begin{tabular}{ll} \bullet & \textbf{Golden-section search} \\ & \textbf{local minimum on } [A,B] & \textbf{interval (constraint)} \\ \end{tabular}$
- Works by narrowing down the search interval with a constant reduction factor
 - $1 \alpha = \frac{\sqrt{5} 1}{2} \approx 0.62$

Ouestion: Does α remind you of something?

Golden section (minimization)



10: end while



 $732 \texttt{A90_ComputationalStatisticsVT2019_Lecture02codeSlide16.R}$

f has to be **UNIMODAL**

Find

$$\min_{\vec{x} \in \mathbb{R}^n} f(\vec{x})$$

Using (known, or numerically evaluated)

Gradient
$$\nabla f(\vec{x}) = \left(\frac{\partial f(\vec{x})}{\partial x_1}, \dots, \frac{\partial f(\vec{x})}{\partial x_n}\right)^T$$

Hessian
$$\nabla^2 f(\vec{x}) = \left[\frac{\partial^2 f(\vec{x})}{\partial x_i \partial x_j} \right]_{i,i=1}^n$$

- Provide a (good) starting point \vec{x}_0 , $\vec{x} = \vec{x}_0$
- Choose a direction \vec{p} (||p|| 1) and step size a
- $\bullet \ \text{Move to} \ \vec{x} := \vec{x} + \alpha \vec{p}$
- ${\color{red} \bullet}$ Repeat step 2 until convergence

Taylor's theorem

$$f(\vec{x} + a\vec{p}) = f(\vec{x}) + \alpha \vec{p} \cdot \nabla f(\vec{x}) + o(\alpha^2)$$

$$\vec{p}$$
 s.t. $\vec{p}^T \cdot \nabla f(\vec{x}) < 0$ is a $descent$ direction.

 ${\bf Steepest} \,\, {\rm descent} \,\, {\rm is} \,\,$

$$\vec{p} = -(\bigtriangledown f(\vec{x})) / \|\bigtriangledown f(\vec{x})\|$$

- \bullet Expensive way: find the global minimum in direction \vec{p}
- Trade-off way: find a decrease which is sufficient

BACKTRACKING

- 1: Choose (large) $\alpha_0 > 0, \, \rho \in (0,1), \, c \in (0,1),$
- 2: $\alpha = \alpha_0$ 3: repeat
- 4: $\alpha = \rho \alpha$ 5: $\operatorname{until} f(\vec{x} + \alpha \vec{p}) \leq f(\vec{x}) + c\alpha \vec{p}^T \nabla f(\vec{x})$

- Newton-Raphson method
- Hessian ignored in steepest descent
- If f is quadratic

$$f(\vec{p}) = \frac{1}{2}\vec{p}^T \mathbf{A} \vec{p}' + \vec{b}^T \vec{p}' + c,$$
n

then minimum

$$p^* = \mathbf{A}^{-1} \vec{b}$$
.

 \bullet Taylor expansion of f

$$f(\vec{x} + a\vec{p}) = f(\vec{x}) - \alpha \vec{p}^T \cdot \bigtriangledown f(\vec{x}) + \frac{\alpha^2}{2} \vec{p}^T \bigtriangledown^2 f(\vec{x}) \vec{p} + o(\alpha^3)$$

• $x := x + \alpha \vec{p}$ where

$$\vec{p} = -\left(\bigtriangledown^2 f(\vec{x})\right)^{-1} \bigtriangledown f(\vec{x})$$

- $(\nabla^2 f(\vec{x}))^{-1}$ is expensive to compute, there are quicker approaches, e.g. Cholesky decomposition
- \bullet Hessian should be positive definite for \vec{p} to be a descent direction (if not see book)
- Memory expensive need to store $O(n^2)$ elements

BUT

• Method converges quickly esp. near optimum

- \bullet Compute an approximation to the Hessian, B, that will allow for efficient choice of $\vec{p}.$
- SECANT CONDITION: (quasi-Newton condition)

$$\mathbf{B}_{k+1} (\vec{x}_{k-1} - \vec{x}_k) = \nabla f(\vec{x}_{k+1}) - \nabla f(\vec{x}_k)$$

- BFGS Algorithm
- $\begin{aligned} & \text{BFGS Algorithm} \\ & \text{1: Choose } \mathbf{B}_0 > 0, \vec{x}_0, k = 0 \\ & \text{2: repeat} \\ & \text{3: } \vec{p}_k \text{ is solution of } \mathbf{B}_k \vec{p}_k = \nabla f(\vec{x}_k) \\ & \text{5: } \vec{x}_{k+1} + \vec{x}_k + \alpha_k \vec{p}_k \\ & \text{6: calculate } \mathbf{B}_{k-1} \text{ {\{nex slide\}}} \\ & \text{7: } k = k+1 \end{aligned}$

How to compute B_{k+1} ?

- We want \mathbf{B}_{k+1} and \mathbf{B}_k to be close to each other
 - $\label{eq:basic_bound} \begin{aligned} & \underset{\mathbf{B}}{\text{min}} & \|\mathbf{B} \mathbf{B}_k\| \\ & s.t. & \mathbf{B} \mathbf{B}^T, \text{ secant condition} \end{aligned}$
- $\bullet \ \vec{y}_k = \bigtriangledown f(\vec{x}_{k+1}) \bigtriangledown f(\vec{x}_k), \ \vec{s}_k = \vec{x}_{k+1} \vec{x}_k$

$$\mathbf{B}_{k+1} = \mathbf{B}_k - \frac{\mathbf{B}_k \vec{y}_k \vec{y}_k^T \mathbf{B}_k}{y_k^T \mathbf{B}_k \dot{y}_k} + \frac{\vec{s}_k \vec{s}_k^T}{\dot{y}_k \vec{s}_k^T}$$

- Closed form Sherman-Morrison formula for \mathbf{B}_{6+1}^{-1}
- We have to store \mathbf{B}_k^{-1}

- BGFS: Broyden-Fletcher-Goldfarb-Shanno
- More iterations than Newton's method (uses approximation)
- Each iteration quicker, no numeric inversion
- Good for large scale problems
- Choice of B₀?

$$f(\vec{x}) = \frac{1}{2}\vec{x}^T \mathbf{A}\vec{x} - \vec{b}^T\vec{x}$$

for A symmetric positive definite.

$$\bigvee f(\vec{x}) = \mathbf{A}\vec{x} - \vec{b} = r(\vec{x})$$

Two vectors \vec{p} and \vec{q} are conjugate with respect to \bf{A} if

$$\vec{p}^T \mathbf{A} \vec{q} = 0$$

IDEA: p' and q' are orthogonal w.r.t. to an inner product associated with A. Use this to find a basis that will allow for easy finding of \vec{x}

Conjugate Gradient method

- $\vec{p}_0 = \vec{r}_0$
- $\blacksquare \vec{p}_{k+1} = -\vec{r}_k + \beta_{k+1} \vec{p}_k$
- Conjugate condition has to be satisfied so

$$\beta_{k+1} = \frac{\vec{r}_k^T \mathbf{A} \vec{p}_{k-1}}{\vec{p}_L^T \mathbf{A} \vec{p}_k}$$

Exercise: check this

 \bullet Convergence in $\dim(\mathbf{A})$ steps (or unless cutoff for \vec{r}_k)

- \bullet If $f(\cdot)$ general, use $\bigtriangledown f(\cdot)$ instead of $r(\cdot)$

- 1: Choose $\vec{x}_0, \vec{p}_0 = -\nabla f(\vec{x}_0), k = 0$ 2: while $\nabla f(x_k) \neq \vec{0}$ do

 3: find suitable α_k {aud now update step}

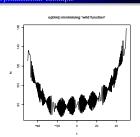
 4: $\vec{x}_1 + \vec{x}_2 + \alpha_2 \vec{p}_1$ {aud now update step}

 5: $\beta_{k+1} = (\vec{y}_1 + \vec{y}_2 + (\vec{x}_{k+1})) / (\nabla^2 f(\vec{x}_k)) \nabla f(\vec{x}_k)$ [Fletcher-Rewes update, other possible)

 6: $\vec{p}_{k+1} = -\nabla f(\vec{x}_{k+1}) + \beta_{k+1} \vec{p}_k$ 7: k = k + 1

- 8: end while

- Local minimum convergence
- But this is true of all methods that cannot "jump out" of descent path
- Faster than steepest descent
- Slower than Newton and Quasi-Newton but significantly less memory



 $732 \texttt{A} 90 _\texttt{ComputationalStatisticsVT2019} _\texttt{Lecture02codeSlide31.R}$

- \bullet Optimization is everywhere
- Numerical methods for finding minimum
- 1D: Golden section (unimodal), optimize()
- kD: choose step size and direction (gradient), optim()