Statistical Computing (36-350) Lecture 18: Optimization II: Stochastic and Constrained Optimization

Cosma Shalizi

31 October 2012



- Stochastic optimization methods
- Constraints and penalties

OPTIONAL READING (big picture): Francis Spufford, *Red Plenty*; Herbert Simon, *The Sciences of the Artificial* OPTIONAL READING (close up): Bottou and Bosquet, "The Tradeoffs of Large Scale Learning"

Problems with Big Data

Typical statistical objective function, mean-squared error:

$$f(\theta) = \frac{1}{n} \sum_{i=1}^{n} (y_i - m(x_i, \theta))^2$$

Getting a value of f is O(n), ∇f is O(np), **H** is $O(np^2)$ worse still if m slows down with nNot bad when n = 100 or even $n = 10^4$, but if $n = 10^9$ or $n = 10^{12}$ we don't even know which way to move Stochastic Optimization
Constraints and PenaltiesStochastic Gradient Descent
Simulated Annealing



イロト イタト イヨト イヨト

3

Sampling, the Alternative to Sarcastic Gradient Descent

Pick *one* data point *I* at random (uniform on 1:n) Loss there, $(y_I - m(x_I, \theta))^2$, is random, but

$$\mathbb{E}\left[\left(y_{I}-m(x_{I},\theta)\right)^{2}\right]=f(\theta)$$

Generally, if $f(\theta) = n^{-1} \sum_{i=1}^{n} f_i(\theta)$ and f_i are well-behaved,

$$\mathbb{E}[f_{I}(\theta)] = f(\theta)$$

$$\mathbb{E}[\nabla f_{I}(\theta)] = \nabla f(\theta)$$

$$\mathbb{E}\left[\nabla^{2} f_{I}(\theta)\right] = \mathbf{H}(\theta)$$

:. Don't optimize with all the data, optimize with random samples

Stochastic Gradient Descent

Draw lots of one-point samples, let their noise cancel out:

- Start with initial guess θ , learning rate η
- While ((not too tired) and (making adequate progress))
 - At t^{th} iteration, pick random I uniformly
 - $e Set \theta \leftarrow \theta t^{-1} \eta \nabla f_I(\theta)$
- Seturn final θ

Shrinking step-size by 1/t ensures noise in each gradient dies down (Variants: put points in some random order, only check progress after going over each point once, adjust 1/t rate, average a couple of random data points, etc.)

The sample code from the midterm works, though it could be made more efficient

Stochastic Gradient Descent Simulated Annealing

Stochastic Newton's Method

a.k.a. 2nd order stochastic gradient descent

- **0** Start with initial guess θ
- While ((not too tired) and (making adequate progress))
 - At t^{th} iteration, pick uniformly-random I
- Seturn final θ

+ all the Newton-ish tricks to avoid having to recompute the Hessian

Stochastic Gradient Descent Simulated Annealing

Stochastic Gradient Methods

Pros:

- Each iteration is fast (and constant in *n*)
- Never need to hold all data in memory
- Does converge eventually

Cons:

• Noise *does* reduce precision — more iterations to get within ϵ of optimum than non-stochastic GD or Newton

Often low computational cost to get within *statistical* error of the optimum

Simulated Annealing

Use Metropolis to sample from a density $\propto e^{-f(\theta)/T}$ Samples will tend to be near small values of fKeep lowering T as we go along ("cooling", "annealing")

- Set initial θ , T > 0
- While ((not too tired) and (making adequate progress))
 - Proposal: $Z \leftarrow r(\cdot|\theta)$ (e.g., Gaussian noise)
 - Oraw $U \sim \text{Unif}(0, 1)$
 - Acceptance: If $U < e^{-\frac{f(Z)-f(\theta)}{T}}$ then $\theta \leftarrow Z$
 - Reduce T a little
- Seturn final θ

Always moves to lower values of f, sometimes moves to higher No derivatives, works for discrete problems, few guarantees

프 에 에 프 어

agrange Multipliers nequality Constraints and Penalties

Maximizing a multinomial likelihood

I roll dice *n* times; $n_1, \dots n_6$ count the outcomes Likelihood and log-likelihood:

$$L(\theta_{1}, \theta_{2}, \theta_{3}, \theta_{4}, \theta_{5}, \theta_{6}) = \frac{n!}{n_{1}!n_{2}!n_{3}!n_{4}!n_{5}!n_{6}!} \prod_{i=1}^{6} \theta_{i}^{n_{i}}$$
$$\ell(\theta_{1}, \theta_{2}, \theta_{3}, \theta_{4}, \theta_{5}, \theta_{6}) = \log \frac{n!}{n_{1}!n_{2}!n_{3}!n_{4}!n_{5}!n_{6}!} + \sum_{i=1}^{6} n_{i}\log \theta_{i}$$

Optimize by taking the derivative and setting to zero:

$$\frac{\partial \ell}{\partial \theta_1} = \frac{n_1}{\theta_1} = 0$$

$$\therefore \theta_1 = \infty$$

or $n_1 = 0$

Stochastic Optimization Constraints and Penalties Inequality Constraints and Penalties

We forgot that $\sum_{i=1}^{6} \theta_i = 1$ We could use the constraint to eliminate one of the variables

$$\theta_6 = 1 - \sum_{i=1}^5 \theta_i$$

Then solve the equations

$$\frac{\partial \ell}{\partial \theta_i} = \frac{n_1}{\theta_i} - \frac{n_6}{1 - \sum_{j=1}^5 \theta_j} = 0$$

BUT eliminating a variable with the constraint is usually messy

Lagrange Multipliers Inequality Constraints and Penalties

Lagrange Multipliers

$$g(\theta) = c \iff g(\theta) - c = 0$$

Lagrangian:

$$\mathcal{L}(\theta, \lambda) = f(\theta) - \lambda(g(\theta) - c)$$

= f when the constraint is satisfied Now do *unconstrained* minimization over θ and λ :

$$\begin{split} \nabla_{\theta} \mathcal{L} |_{\theta^{*}, \lambda^{*}} &= \nabla f(\theta^{*}) - \lambda^{*} \nabla g(\theta^{*}) = \mathbf{0} \\ \left. \frac{\partial \mathcal{L}}{\partial \lambda} \right|_{\theta^{*}, \lambda^{*}} &= g(\theta^{*}) - c = \mathbf{0} \end{split}$$

optimizing Lagrange multiplier λ enforces constraint More constraints, more multipliers Try the dice again:

$$\begin{aligned} \mathscr{L} &= \log \frac{n!}{\prod_{i} n_{i}!} + \sum_{i=1}^{6} n_{i} \log(\theta_{i}) - \lambda \left(\sum_{i=1}^{6} \theta_{i} - 1 \right) \\ \frac{\partial \mathscr{L}}{\partial \theta_{i}} \bigg|_{\theta_{i} = \theta_{i}^{*}} &= \frac{n_{i}}{\theta_{i}^{*}} - \lambda^{*} = 0 \\ \frac{n_{i}}{\lambda^{*}} &= \theta_{i}^{*} \\ \sum_{i=1}^{6} \frac{n_{i}}{\lambda^{*}} &= \sum_{i=1}^{6} \theta_{i}^{*} = 1 \\ \lambda^{*} &= \sum_{i=1}^{6} n_{i} \Rightarrow \theta_{i}^{*} = \frac{n_{i}}{\sum_{i=1}^{6} n_{i}} \end{aligned}$$

・ロト ・四ト ・ヨト ・ヨト

2

Thinking About the Lagrange Multipliers

Constrained minimized is generally higher than the unconstrained Changing the constraint level *c* changes θ^* , $f(\theta^*)$

$$\begin{aligned} \frac{\partial f(\theta^*)}{\partial c} &= \frac{\partial \mathscr{L}(\theta^*, \lambda^*)}{\partial c} \\ &= \left[\nabla f(\theta^*) - \lambda^* \nabla g(\theta^*) \right] \frac{\partial \theta^*}{\partial c} - \left[g(\theta^*) - c \right] \frac{\partial \lambda^*}{\partial c} + \lambda^* = \lambda^* \end{aligned}$$

 λ^* = Rate of change in optimal value as the constraint is relaxed λ^* = "Shadow price": How much would you pay for minute change in the level of the constraint

Inequality Constraints

What about an *inequality* constraint?

$$b(\theta) \le d \iff b(\theta) - d \le \mathbf{0}$$

The region where the constraint is satisfied is the **feasible set** *Roughly* two cases:

- Unconstrained optimum is inside the feasible set ⇒ constraint is inactive
- Optimum is outside feasible set; constraint is active, binds or bites; *constrained* optimum is usually on the boundary

Add a Lagrange multiplier; $\lambda \neq 0 \Leftrightarrow$ constraint binds

Mathematical Programming

Older than computer programming...

Optimize $f(\theta)$ subject to $g(\theta) = c$ and $h(\theta) \le d$

"Give us the best deal on f, keeping in mind that we've only got d to spend, and the books have to balance" Linear programming (Kantorovich, 1938)

- f, b both linear in θ
- θ^* always at a corner of the feasible set

Lagrange Multipliers Inequality Constraints and Penalties

Back to the Factory

Constraints:

$$40(cars) + 60(trucks) \leq 1600$$

 $1(cars) + 3(trucks) \leq 70$

Revenue: \$13k/car, \$27k/truck The feasible region:

4



Lagrange Multipliers Inequality Constraints and Penalties

Back to the Factory

Constraints:

$$\begin{array}{rcl} 40(cars) + 60(trucks) &\leq & 1600\\ 1(cars) + 3(trucks) &\leq & 70 \end{array}$$

Revenue: \$13k/car, \$27k/truck The feasible region:

4



Barrier Methods

(a.k.a. "interior point", "central path", etc.) Having constraints switch on and off abruptly is annoying especially with gradient methods Fix $\mu > 0$ and try minimizing

 $f(\theta) - \mu \log(d - h(\theta))$

"pushes away" from the barrier — more and more weakly as $\mu \rightarrow 0$

- Initial θ in feasible set, initial μ
- While ((not too tired) and (making adequate progress))
 - Minimize $f(\theta) \mu \log(d h(\theta))$
 - **2** Reduce μ
- Seturn final θ

イロト イタト イヨト イヨト

Constraints vs. Penalties

$$\underset{\theta: b(\theta) \le d}{\operatorname{argmin}} f(\theta) \iff \underset{\theta, \lambda}{\operatorname{argmin}} f(\theta) - \lambda(b(\theta) - d)$$

 $\begin{array}{ll} d \text{ doesn't matter for doing the second minimization over } \theta \\ \text{Constrained optimization} & \Leftrightarrow & \text{Penalized optimization} \\ \text{Constraint level } d & \Leftrightarrow & \text{Penalty factor } \lambda \end{array}$

Statistical Applications of Penalization

Minimize MSE of linear function $\beta \cdot x$: ordinary least squares regression

penalty on length of coefficient vector, $\sum \beta_j^2$: ridge regression stabilizes estimate when data are noisy, p > n, collinear penalty on sum of coefficients, $\sum |\beta_j|$: lasso stability + drive small coefficients to 0 ("sparsity") Minimize MSE of function + penalty on curvature : spline fit smooth regressions w/o assuming specific form

Smoothing over time, space, other relations

e.g., social or genetic ties

Usually decide on penalty factor/constraint level by trying to predict out of sample

R implementation

In penalty form, just chose λ and modify your objective function constrOptim implements the barrier method Try this:

```
factory <- matrix(c(40,1,60,3),nrow=2,
    dimnames=list(c("labor","steel"),c("car","truck")))
available <- c(1600,70); names(available) <- rownames(factory)
prices <- c(car=13,truck=27)
revenue <- function(output) { return(-output %*% prices) }
plan <- constrOptim(theta=c(5,5),f=revenue,grad=NULL,
    ui=-factory,ci=-available,method="Nelder-Mead")
plan$par
```

only works with constraints like $\mathbf{u}\theta \ge c$, so minus signs

Summary

- Stochastic optimization methods use probability in the search
 - Stochastic gradient descent samples the data; gives up precision for speed
 - Simulated annealing randomly moves against the objective function; escapes local minima
- Constraints are usually part of optimization
 - Constrained optimum generally at the boundary of feasible set
 - Lagrange multipliers turn constrained problems into unconstrained ones
 - Multipliers are prices: trade-off between tightening constraint and worsening optimal value