

Modeling and Algorithms for Prof Shi's Project

Hongda Li, Yining Zhou, Xiaoping Shi

January 22, 2024

Abstract

We propose some better algorithm for a problem in detecting structure of probability transition matrices from data.

1 Introduction

We describe an optimization problem introduced by Prof Shi and his student Yining. To start we define the following quantities for the optimization problem.

1. $n \in \mathbb{N}$. It denotes the numer of states for a Markov Chain.
2. $p \in \mathbb{R}^{n \times n}$ denotes the probability transition matrix. It's in small case because it's also the variable for the optimization problem. It supports two types of indexing, $p_{i,j}$ for $i, j \in \{1, \dots, n\}$, or p_j with $j \in \{1, \dots, n^2\}$. More on this later.
3. $\eta_{ij} \geq 0$ for $i, j \in \{1, \dots, n\}$ is a parameter of the problem.
4. \hat{p} is the empirically measured probability transition matrix. They are the maximal likelihood estimators for the transition probability in the transition probability matrix.
5. λ is the regularization parameter.
6. \mathbf{C}_n^m is the combinatoric term that counts all possible subset of size $n, n < m$ from a super set of size m .

When p is referred to as a vector we may say $p \in \mathbb{R}^{n^2}$, if it's referred to as the matrix, we will use $p \in \mathbb{R}^{n \times n}$. When indexing p using a tuple, or a single number, it's possible to translate between the two type of indexing scheme using the following bijective map:

$$\begin{aligned}(i, j) &\mapsto k := i \times n + j \\ k &\mapsto (i, j) := (\lfloor k/n \rfloor, \text{mod}(k, n) + 1).\end{aligned}$$

We emphasize, in different programming languages and development environments, the convention of indexing a muti-array using different kind of tuples can be very different. For now we use the above indexing, which is a row major index convention (Like Python).

1.1 Some Mathematical Entities

We introduce some mathematical entities for more context.

1.2 The Optimization Problem

The optimization problem is posed as

$$\operatorname{argmin}_{p \in \mathbb{R}^{n^2}} \left\{ \sum_{i=1}^n \sum_{j=1}^n -\eta_{ij} \log(p_{i,j}) + \lambda \sum_{i=1}^{n^2} \sum_{\substack{j=1 \\ j \neq i}}^{n^2} \frac{1}{2} \frac{|p_i - p_j|}{|\hat{p}_i - \hat{p}_j|} \right\} \quad (1.2.1)$$

$$\text{s.t.: } \left\{ \begin{array}{ll} \sum_{i=1}^n p_{i,j} = 1 & \forall i = 1, \dots, n \\ p_{i,j} \geq 0 & \forall i = 1, \dots, n, j = 1, \dots, n \end{array} \right\} \quad (1.2.2)$$

There are 3 parts to the optimization problem posed above. It has a smooth differentiable function, the sum of $\eta_{ij} \log(p_{i,j})$ but its gradient is not Lipschitz. The has a non-smooth part with $p_i - p_j$ for all indices $1, \dots, n^2$ and $i \neq j$. If $\hat{p}_i - \hat{p}_j = 0$, then we would ignore the term. Finally, it has a linear constraints on all n^2 variables. $p \in \mathbb{R}^{n^2}$ is a vector with the structure $\Delta_n \oplus \Delta_n \oplus \dots \oplus \Delta_n$. Each Δ_n is a probability simplex. It's defined as $\Delta_n = \{x \in \mathbb{R}_+^n : \sum_{i=1}^n x_i = 1\}$. For simplicity we just denote using notation $\Delta_{n \times n} = \Delta_n \oplus \Delta_n \oplus \dots \oplus \Delta_n$.

2 Modeling

The non-smooth part with the absolute value requires some amount of creativity if we were to use common optimization algorithms.

2.1 Representation of the Non-smooth Part

Claim 2.1.1. The nonsmooth objective can be model as

$$\lambda \sum_{i=1}^{n^2} \sum_{j>i}^{n^2} \frac{|p_i - p_j|}{|\hat{p}_i - \hat{p}_j|} = \|Cp\|_1 \quad \text{where } C \text{ is } \mathbf{C}_2^{n^2} \text{ by } n^2.$$

The one norm represents the summation of absolute values. The transformation Cp is linear transformation and it's a vector of length $\mathbf{C}_2^{n^2}$. The vector is long and it has a dimension of $(1/2)(1 + n^2)n^2$. Each term inside of the summation is a row of the matrix C . Each row of matrix C has exactly 2 non-zero elements in it. Suppose that $i \in \{1, \dots, \mathbf{C}_2^{n^2}\}$ denoting the index for a specific row of matrix C denotes the index of a specific row, and $j \in \{1, \dots, n^2\}$ denotes a specific column of matrix C . Mathematically describing the matrix is difficult, but it can be algorithmically defined. A sparse matrix format can be described by as a mapping from (i, j) , the set of all indices to the element in the vector. The following [algorithm 1](#) construct such a mapping.

Algorithm 1 Matrix Make Algorithm

```
1: Let  $C$  be a  $\mathbf{C}_2^{n^2}$  by  $n^2$  zero matrix.
2: for  $i = 1, \dots, n^2$  do
3:   for  $j = 1, \dots, n^2$  do
4:     if  $|\hat{p}_i - \hat{p}_j| == 0$  then
5:       continue
6:     end if
7:      $C[i \times n^2 + j, i] := \lambda/|\hat{p}_i - \hat{p}_j|$ 
8:      $C[i \times n^2 + j, j] := -\lambda/|\hat{p}_i - \hat{p}_j|$ 
9:   end for
10: end for
```

Remark 2.1.1. In practice, we should use sparse matrix data format such as SCR, SCC in programming languages.

2.2 Modeling it for Sequential Quadratic Programming

To start, refer to [Sequential Qaudratic Programming Wikipedia](#) for brief overview about what sequential quadratic programming problem is. To use sequential programming, we model the non-smooth $\|Cp\|_1$ parts of the objective function as a linear constraints. Define $u \in \mathbb{R}_+^{C_2^{n^2}}$ then the below problem is equivalent to the original formulation:

$$\underset{p, u}{\operatorname{argmin}} \left\{ \sum_{i=1}^n \sum_{j=1}^n -\eta_{ij} \log(p_{i,j}) + \sum_{i=1}^{C_2^{n^2}} u_k \right\} \quad (2.2.1)$$

$$\text{s.t: } \begin{cases} -u \leq Cp \leq u \\ p \in \Delta_{n \times n} \\ u \in \mathbb{R}_+^{C_2^{n^2}} \end{cases} \quad (2.2.2)$$

This is a Non-linear programming problem and it has a convex objective. Common NLP packages in programming languages can solve this efficiently. However it's potentially possible that these solvers are not adapated to huge sparse matrix C that has special structure to it.

Remark 2.2.1. To formulate into a linear programming with some relaxations, consider that the non-linear objective $-\eta_{ij} \log(p_{i,j})$ can be discretized.

2.3 Modeling it for Operator Splitting

Operator splitting method aims for objective function of the type $f + g$ where f, g are both convex function and they are proximal friendly. And $\operatorname{ri}.\operatorname{dom}(f) \cap \operatorname{ri}.\operatorname{dom}(g) \neq \emptyset$. Recall that proximal operator of function f is the resolvent operator on the subgradient $(I + \partial f)^{-1}$. Subgradient is just a generalize type of gradient that can handle continuous function that is not necessarily differentiable. To model the original formulation in such a form, we introduce these quantities:

$$1. f_1(x_1) : \mathbb{R}^{n^2} \mapsto \mathbb{R} := \sum_{i=1}^{n^2} -\eta_{i,j}(\log(p_{i,j}))$$

$$2. f_2(x_2) : \mathbb{R}^{C_2^{n^2}} \mapsto \mathbb{R} := \|x_2\|_1.$$

$$3. f_3(x_3) : \mathbb{R}^{n^2} \mapsto \bar{\mathbb{R}} := \delta_{\Delta_{n \times n}}(x_3)$$

The above three functions represent the three parts of the summed objective.

Let $f(x_1, x_2, x_3) := f_1(x_1) + f_2(x_2) + f_3(x_3)$. However, they share different variables x_1, x_2, x_3 , from different dimension. We want g to represents the constraints that $x_2 = Cp$, and $x_1 = x_3 = p$. Simplifying: $x_2 = Cx_1, x_1 = x_3$. This is a linear system of the form

$$\begin{bmatrix} C & -I & \mathbf{0} \\ I & \mathbf{0} & -I \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \mathbf{0}.$$

Simply denote the above as $Ax = \mathbf{0}$, then g is just a non-smooth function that happens to be an indicator function of a convex set. The convex set is the set of all solutions to the linear system. Therefore setting $g = \delta_{Ax=\mathbf{0}}(x)$, to be the indicator of the linear constraints, we had a complete equivalent representation of the original form of the problem.

The matrix A is a $(C_2^{n^2} + n^2) \times (3n^2)$ matrix. It's still a sparse matrix. Conviently, if we use \otimes then the above matrix admit kronecker product representation:

2.4 Another Method that is not Cursed by the Dimension

...

2.5 A Quick Tour into Operator Splitting Method

Operator splitting method is relatively new and doesn't have a mature toolings in software. Using the method requires the use of the Proximal Mapping, which can be non-trivial to implement for even simple functions. Read section [1, section 2.7] of Ryu's Textbook for more information about the algorithm and operator splitting in general.

3 Implementation, Software, and Toolings

The software for Sequential Quadratic programming is relatively mature. Whether they support sparse matrix inversion depends on the developers. The software operator splitting method is less common. To use Operator splitting method such as Douglas Rachford (Equivalently the ADMM via a the Dual), Peachman Rachford, it's required to compute the proximal operator of functions. The following theorem is relevant to efficiently compute the proximal mapping for our model.

Theorem 1 (Separable Proximal Theorem). Suppose that function $f(x)$ can be written as the Euclidean cross product $(f_1(x_1), f_2(x_2), \dots, f_m(x_m))$, and we assume that each of x_1, x_2, \dots, x_n are vector in $\mathbb{R}^{k_1}, \mathbb{R}^{k_2}, \dots, \mathbb{R}^{k_m}$ such that $x = \bigoplus_{i=1}^m x_i$, then we can write the proximal gradient operator of f in parallel form which is given as:

$$\text{prox}_{f,\lambda}((x_1, x_2, \dots, x_m)) = (\text{prox}_{f_1,\lambda}(x_1), \text{prox}_{f_2,\lambda}(x_2), \dots, \text{prox}_{f_m,\lambda}(x_m)).$$

In general, if there exists any permutation of the list of x_1, x_2, \dots, x_m , the parallelization property of prox will be come applicable.

Recall objective $f(x_1, x_2, x_3) = f_1(x_1) + f_2(x_2) + f_3(x_3)$ from section 2.1. It has the applicable form and hence we have:

$$\text{prox}_{\lambda f}((x_1, x_2, x_3)) = \text{prox}_{\lambda f_1}(x_1) \oplus \text{prox}_{\lambda f_2}(x_2) \oplus \text{prox}_{\lambda f_3}(x_3).$$

Next, we proceed to write each f_i for $i = 1, 2, 3$ into separable sum as well. f_1, f_2 are trivial to do however $\delta_{\Delta_{n \times n}}$ requires some thought. Consider $x_3 = \bigoplus_{i=1}^n x_{3,i}$, so that the vector is the direct product, then

$$\delta_{\Delta_{n \times n}}(x_3) = \sum_{i=1}^n \delta_{\Delta_n}(x_{3,i}).$$

The indicator function $\delta_{\Delta_{n \times n}}$ acts on groups of n elements in the vector x_3 individually and separately, hence it admits the above separable summation form. We now list the proximal mapping for each of the atomic operations: $\lambda \log(\cdot)$, $\lambda |\cdot|$, δ_{Δ_n} and $\delta_{Ax=0}(\cdot)$. Most of them can be found in the literatures.

$$\begin{aligned} \text{prox}_{-\lambda \log(\cdot)}(x) &= \frac{x + \sqrt{x^2 + 4\lambda}}{2} \\ \text{prox}_{\lambda |\cdot|} &= \text{sign}(x) \max(|x| - \lambda, 0) \\ \text{prox}_{\delta_{\Delta_n}}(x) &= \Pi_{\Delta_n} \\ \text{prox}_{\delta_{\{Ax=0\}}}(x) &= \Pi_{Ax=0}. \end{aligned}$$

3.1 Projection onto Probability Simplex and Sparse Linear System

Projecting onto a probability simplex and a linear system that is sparse requires a bit more attention for software implementation sake. We introduce the following theorem from Beck textbook to assist with the projection onto a probability simplex.

Theorem 2 (Projecting onto Hyperplane Box Intersections). Let $C \subseteq \mathbb{R}^n$ to be non-empty and defined as $C := \text{Box}[l, u] \cap H_{a,b}^-$, where $\text{Box}[l, u] = \{x \in \mathbb{R}^n : l \leq x \leq u\}$. We assume that $l \in [-\infty, \infty)$ and $u \in (-\infty, \infty]$. And $H_{a,b}^- = \{x \in \mathbb{R}^n : \langle a, x \rangle = b\}$, where $a \in \mathbb{R}^n \setminus \{0\}$. Then the projection onto the set is equivalent to the following conditions

$$\begin{aligned} \Pi_C(x) &= \Pi_{\text{Box}[l, u]}(x - \mu a) \\ \text{where } u \in \mathbb{R} \text{ solves: } &\langle a, \Pi_{\text{Box}[l, u]}(x - ua) \rangle = b. \end{aligned}$$

Proof. See [2, theorem 6.27] by Beck for more information. \square

Remark 3.1.1. The equation for u has a continuous mapping $\mathbb{R} \mapsto \mathbb{R}$, it can be solved efficiently using method such as the bisection method.

Theorem 3 (Projecting onto the Probability Simplex). The projection onto the simplex Δ_n can be computed via

$$\begin{aligned}\Pi_{\Delta_n}(x) &= \Pi_{\mathbb{R}_+}(x - \mu a) \\ \text{where } \mu \in \mathbb{R} \text{ solves: } \langle \mathbf{1}, \Pi_{\mathbb{R}_+}(x - \mu \mathbf{1}) \rangle &= 1.\end{aligned}$$

Proof. Use the previous theorem 2, observe that u can be a vector of infinity and hence setting \mathbb{R}_+ to be the box, and the probability simplex constraints $\sum_{i=1}^n x_i = 1$ is the hyperplane $H_{a,b}^-$ with $a = \mathbf{1}, b = 1$. \square

Theorem 4 (Projection onto Affine Spaces). Define the set $C = \{x \in \mathbb{R}^n : Ax = b\}$, assuming that $A \in \mathbb{R}^{m \times n}$ is full rank, then the projection onto the set is computed via

$$\Pi_C(x) = x - (A^T A)^{-1} A(Ax - b).$$

Remark 3.1.2. In the case when A is not full rank, do a reduced QR decomposition of $A = QR$ where R is $k \times n$ with $k < m$. In the case of our application, the matrix A is a full rank matrix. The matrix $(A^T A)^{-1}$ in our case can be very slow to invert since it's $\mathbf{C}_2^{n^2} \times \mathbf{C}_2^{n^2}$. One way to make it easier is to use sparse matrix invert method. In our case the matrix we are inverting is sparse and symmetric, it can be solved efficiently say for example,

3.2 The Software for These Algorithms

option **SLSQP** in “scipy.optimize.minimize” performs sequential quadratic programming. I am not sure whether they would support efficient sparse matrix inverse.

References

- [1] E. K. Ryu and W. Yin, *Large-Scale Convex Optimization: Algorithms & Analyses via Monotone Operators*. Cambridge: Cambridge University Press, 2022. [Online]. Available: <https://large-scale-book.mathopt.com/>
- [2] A. Beck, *First-Order Methods in Optimization* / *SIAM Publications Library*, ser. MOS-SIAM Series in Optimization. SIAM. [Online]. Available: <https://epubs.siam.org/doi/book/10.1137/1.9781611974997>