# How to Smooth Data in Two Way Tables with the Order Restriction

SEIYOUNG CHUNG

ABSTRACT. To smooth a given data in two-way tables with the order restriction, we propose the dual problem and construct an algorithm utilizing the network flows which ends up with the minimum $L1$-norm after a finite number of iterations.

## 1. Introduction

Suppose that an $m \times n$ matrix of data $A = (a_{ij})$ is given and we wish to find $r_1 \leq r_2 \leq \cdots \leq r_m$ and $s_1 \leq s_2 \leq \cdots \leq s_n$ so that

$$(P): \sum_{i=1}^{m} \sum_{j=1}^{n} |a_{ij} - r_i - s_j|$$

is minimized.

This problem rises from smoothing of data in two way tables with the order restriction. This without the order restriction has been discussed in [1], [2], [3] and [4]. Tukey [1] introduces the method of median polish and it is known to converge for physical data, [2]. It decreases the $L1$ norm of the matrix but does not necessarily converges to the minimum $L1$ norm. Kemperman [3] has basically proposed an algorithm suggesting that a flow model be used and Fink [4] has developed an algorithm utilizing a flow, which produces the minimum $L1$ norm. But there is no algorithm for the problem with the order restriction known to us. Taking advantage of the linearity, the dual $(D)$ of the primal $(P)$ is proposed. An algorithm which utilizes network flows is constructed that solves the primal and the dual simultaneously. All the material related to the network and the flow, which is necessary for developing the algorithm, can be found in Rockafellar [5, Chapter 1 and 2].

---

## 2. Dual Problem

In the present section, the dual problem for the minimum problem $(P)$ is proposed and it is given a necessary and sufficient condition for the problems $(P)$ and $(D)$ to exhibit the optimal solutions. The dual is :

$$(D) : \max \sum_{i=1}^{m} \sum_{j=1}^{n} a_{ij} w_{ij}$$

subject to

$(D-1) : w_{ij} \in \{0, \pm 1\}$ for all $(i,j)$

$(D-2) : b_i \geq 0, \quad i = 1, 2, \ldots, m-1$

$\qquad c_j \geq 0, \quad j = 1, 2, \ldots, n-1$

$$(D-3) : \sum_{j=1}^{n} w_{ij} = -b_{i-1} + b_i, \quad i = 1, 2, \ldots, m$$

$$\sum_{i=1}^{m} w_{ij} = -c_{j-1} + c_j, \quad i = 1, 2, \ldots, n$$

$$b_0 = b_m = c_0 = c_n = 0$$

Any sets $\{r_i\}$ and $\{s_j\}$ are called feasible for the primal $(P)$ if they satisfy the order restriction and any sets $\{w_{ij}\}$, $\{b_i\}$ and $\{c_j\}$ are called feasible for the dual $(D)$ if they satisfy the constraints $(D-1)$, $(D-2)$ and $(D-3)$. Define $\text{sgn}(x) = 1$ if $x > 0$ ; $-1$ if $x < 0$ ; $0$ if $x = 0$. Let $W_i = \sum_{j=1}^{n} w_{ij}$ for each $i$ and $W^j = \sum_{i=1}^{m} w_{ij}$ for each $j$. Two conditions, which turn out to be optimal criteria, are defined as:

CONDITION A: For each $(i,j)$, $w_{ij} = \text{sgn}(a_{ij} - r_i - s_j)$ whenever

$$a_{ij} - r_i - s_j \text{ is nonzero.}$$

CONDITION B: $\sum_{i=1}^{m} r_i W_i = 0 = \sum_{j=1}^{n} s_j W^j$.

LEMMA 2.1. *Assume that $d_i \geq 0$ for $i = 1, 2, \ldots, k-1$ and $d_0 = 0 = d_k$. If $x_i = d_i - d_{i-1}$, $i = 1, 2, \ldots, k$ and $y_1 \leq y_2 \leq \cdots \leq y_k$, then $\sum_{i=1}^{k} y_i x_i \leq 0$.*

The above lemma is an easy consequence of mathematical induction.

LEMMA 2.2. *Let* $\{w_{ij}\}$ *be feasible for the dual (D), and* $\{r_i\}$ *and* $\{s_j\}$ *for the primal (P). The following inequality then holds:*

$$\sum_{i=1}^{m}\sum_{j=1}^{n}a_{ij}w_{ij} \leq \sum_{i=1}^{m}\sum_{j=1}^{n}|a_{ij}-r_i-s_j|,$$

*where the equality holds if and only if both Conditions A and B are satisfied.*

PROOF: It follows directly from the constraint $(D-1)$ that the right hand side is greater than or equal to $\sum_{i=1}^{n}\sum_{j=1}^{n}(a_{ij}-r_i-s_i)w_{ij}$ and that they are equal to each other if Condition A is satisfied. Notice that $\sum_{i=1}^{m}\sum_{j=1}^{n}(r_i+s_j)w_{ij} = \sum_{i=1}^{m}r_iW_i + \sum_{j=1}^{n}s_jW^j$. It is the immediate consequence of Lemma 2.1 that $\sum_{i=1}^{m}r_iW_i \leq 0$ and $\sum_{j=1}^{n}s_jW^j \leq 0$, which completes the proof.

In Lemma 2.2, it has been shown that the maximum of the dual is always less than or equal to the minimum of the primal and hence that the feasible solutions are optimal if they are equal. Therefore the sufficiency of Theorem 2.3 stated below is proved. The necessity is proved after constructing an algorithm in section 3.

THEOREM 2.3 (Duality Theorem). *Under the same assumption in Lemma 2.2,* $\{r_i\}$, $\{s_j\}$ *and* $\{w_{ij}\}$ *are optimal if and only if they satisfy both Conditions A and B.*

Assume the necessity is proved. By Lemma 2.2 again, the optimal values then are the same. They hence are dual to each other.

## 3. Algorithm

Our aim is to solve both the primal and the dual simultaneously using a network. Our basic strategy is: start with an obvious feasible solutions, seek improved ones satisfying Condition B and stop when Condition A is also satisfied. To construct an network for this purpose,

let's define the node sets and the arc sets:

$$R = \{R_1, R_2, \ldots, R_m\},$$
$$C = \{C_1, C_2, \ldots, C_n\},$$
$$N = R \cup C,$$
$$Lrc = R \times C,$$
$$Lr = \{(R_{i+1}, R_i) \mid i = 1, 2, \ldots, m-1\},$$
$$Lc = \{(C_j, C_{j+1}) \mid j = 1, 2, \ldots, n-1\},$$
$$L = Lrc \cup Lr \cup Lc.$$

The network being considered here is a network with the node set $N$ and with the arc set $L$. Noticing that $\{w_{ij}\}, \{b_i\}, \{c_j\}, \{r_i\}$ and $\{s_j\}$ may be regarded as the functions on $Lrc$, $Lr$, $Lc$, $R$ and $C$ respectively, they will be denoted by the functions $w$, $b$, $c$, $r$ and $s$ respectively. From now on, $R_i$ and $C_j$ will be denoted by $i$ and $j$ respectively. The following abbreviations also will be used:

$$b(i) \text{ for } b(R_{i+1}, R_i) \text{ and hence for } b_i,$$
$$c(j) \text{ for } c(C_j, C_{j+1}) \text{ and hence for } c_j.$$

The incedence functions for a node set $S$ and a circuit $P$ are defined by:

$$e_S(x) = \begin{cases} 1 & \text{if } x \in S \\ 0 & \text{otherwise,} \end{cases}$$

$$e_P(J) = \begin{cases} 1 & \text{if } J \in P^+ \\ -1 & \text{if } J \in P^- \\ 0 & \text{otherwise.} \end{cases}$$

Without loss of generality, we may assume that $a_{ij} > 0$. For both $\sum_{i=1}^m \sum_{j=1}^n |a_{ij} - r_i - s_j|$ and $\sum_{i=1}^m \sum_{j=1}^n |a_{ij} + x - r_i - s_j|$ assume their minimum at the same $\{r_i\}$ and $\{s_j\}$ for any real $x$.

## Dual argorithm 3.1.

Initially, $w \equiv 0$, $b \equiv 0$, $c \equiv 0$, $r \equiv 0$ and $s \equiv 0$.

Step 1: Given any $w$, $r$ and $s$, set

$$UL = \{(i,j) \in Lrc \mid r(i) + s(j) < a_{ij}, \quad w(i,j) = 0\}.$$

If $UL = \emptyset$, then stop. The current $w$, $b$, $c$, $r$ and $s$ are optimal.

Step 2: Given $w$, $b$, $c$, $r$ and $s$, paint the network:

1) Any arc $(i,j)$ in $Lrc$ is painted;
red if $w(i,j) = \text{sgn}[a_{ij} - r(i) - s(j)]$ whenever $a_{ij} \neq r(i) - s(j)$,
black if $r(i) + s(j) = a_{ij}$, $w(i,j) = 1$,
white if $[r(i) + s(j) < a_{ij}, w(i,j) = 0]$ or
$\qquad [r(i) + s(j) = a_{ij}, w(i,j) = -1]$,
green if $r(i) + s(j) = a_{ij}$, $w(i,j) = 0$.

2) Any arc $(i+1, i)$ in $Lr$ is painted;
red if $r(i+1) > r(i)$, $b(i) = 0$,
white if $r(i+1) = r(i)$, $b(i) = 0$,
green if $r(i+1) = r(i)$, $b(i) > 0$.

3) Any arc $(j, j+1)$ in $Lc$ is painted;
red if $s(j+1) > s(j)$, $c(j) = 0$,
white if $s(j+1) = s(j)$, $c(j) = 0$,
green if $s(j+1) = s(j)$, $c(j) = 0$.

Step 3: Select $(i*, j*)$ in $UL$ and apply PNA(Painted Network Algorithm) with $N^+ = \{j*\}$ and with $N^- = \{i*\}$. The same arc $(i*, j*)$ should be selected as long as it is still in $UL$.

1) If PNA ends up with a path $P$, let

$$w' = w + e_P,$$
$$b' = b + e_P,$$
$$c' = c + e_P,$$

and go to Step 1.

2) If PNA ends up with a cut $Q = [S, N - S]$, define

$$\delta = \min \begin{cases} r(i) + s(j) - a_{ij} & \text{for } (i,j) \in Q^+ \text{ with } w(i,j) = -1 \\ a_{ij} - r(i) - s(j) & \text{for } (i,j) \in Q^- \text{ with } w(i,j) = 0 \text{ or } 1 \\ r(i+1) - r(i) & \text{for } (i+1, i) \in Q^+ \\ s(j+1) - s(j) & \text{for } (j, j+1) \in Q^+, \end{cases}$$

and set

$$r' = r + \delta e_{N-S},$$
$$s' = s - \delta e_{N-S}.$$

Go to Step 1.

In Propositions 3.2 through 3.5, we will show that the functions $w$, $b$, $c$, $r$ and $s$ remain feasible and Condition B is satisfied after each iteration. Notice that $\delta$ in Step 3 is positive, which follows from the painting conditions.

PROPOSITION 3.2. *New $w'$, $b'$ and $c'$ remain feasible after each iteration.*

PROOF: Since there occurs a change in them only after a circuit $P$, it suffices to prove in such a case. Let $(i,j) \in Lrc$ be in the circuit $P$. If $(i,j)$ is green, then $w(i,j) = 0$ and hence $w'(i,j) = 1$ if $(i,j) \in P^+$ ; $-1$ if $(i,j) \in P^-$. If it is black, then $w(i,j) = 1$ and $w'(i,j) = 0$. If it is white, then $w(i,j) = 0$ or $-1$ and $w'(i,j) = 1$ or $0$. Since the color of any arc in a circuit is green, black or white, the constraint $(D-1)$ holds for $w'$. By the same manner we can show that $b' \geq 0$ and $c' \geq 0$.

Noting that for any node $i \in R$ in the circuit $P$, there exist exactly two arcs in $P$, say $L_1$ and $L_2$, using it as their end node, we can classify all the cases into three different cases according to where the arcs $L_1$ and $L_2$ belong to:

Case 1: Both arcs are in $Lrc$,

Case 2: Both arcs are in $Lr$,

Case 3: One is in $Lr$ and the other in $Lrc$.

In Case 1, $w$ at $L_1$ is increased by one and $w$ at $L_2$ decreased by 1 or vice versa and hence $W_i$ in $(D-3)$ is unchanged. But $b(i)$ and $b(i-1)$ remain unchanged since $(i+1,i)$ and $(i,i-1)$ should not be in $P$. Now the first part in $(D-3)$ holds. In Case 2, $b(i)$ and $b(i-1)$ are increased by 1 if both arcs are in $P^+$ and decreased by 1 if in $P^-$. But $W_i$ is the same as before since no arcs in $Lrc$ connected to the node $i$ is in $P$, and the first part of $(D-3)$ holds. For Case 3, the similar argument can be employed, and the first part in $(D-3)$ holds for any case. We can employ the same argument to show that the second part holds. Since they are feasible at the beginning, mathematical induction completes the proof.

PROPOSITION 3.3. *The order restriction is satisfied by new $r'$ and $s'$.*

PROOF: Initially the restriction is satisfied. Assume this is the case before a certain iteration. Notice that there occurs a change in $r$ and $s$ after a cut $Q = [S, N - S]$ and that $r' = r$ and $s' = s$ except on $N - S$. For any $(i+1, i)$ in $Q^-$, $r'(i+1) = r(i+1) + \delta$ but $r'(i) = r(i)$. For any $(i+1, i)$ in $Q^+$, $r'(i+1) = r(i+1)$ but $r'(i) = r(i) + \delta$. From the definition of $\delta$ in Step 3, it follows that $\delta \leq r(i+1) - r(i)$. In all we conclude that $r'(i+1) \leq r'(i)$ for any $i$. By the same manner it follows that the restriction holds for new $s'$.

PROPOSITION 3.4. *Condition B holds after each iteration.*

PROOF: After a circuit $P$, there occurs a change in $w$ only at the arcs in $P$ and no change in $r$. For any node $i$ in $P$, we have three different cases as in the proof of Proposition 3.2. Let's use the same classification there. For Case 1 and Case 2, we have shown in the proof of Proposition 3.2 that $W$ remains unchanged. In Case 3, there should exist the subpath $\overline{P}$ of $P$ using $i$ as its end node and consisting of only nodes in $R$. For we can go forward or backward from $R$ to $C$ only through the arcs in $Lrc$. Let $i_1$ be the other end node of $\overline{P}$. Then $W_i' = W_i + 1$ and $W_{i_1}' = W_{i_1} - 1$ or vice versa. But $r$ is constant on the set of nodes in $\overline{P}$ since all the arcs in $\overline{P}$ are green or white. In all, we've shown that $\sum_{i=1}^{m} r_i W_i' = 0$ holds after a circuit. Now assume that the outcome is a cut $Q = [S, N - S]$. Then $r' = r$ on $R \cap S$ and $r' = r + \delta$ on $R - S$. Because $w$ experiences no change after a cut, we have:

$$\sum_{i=1}^{m} r_i' W_i' = \sum_{i=1}^{m} r_i W_i + \delta \sum_{i \in R-S} W_i = \delta \sum_{i \in R-S} W_i.$$

The constraint $(D - 3)$ can be rewritten as $\sum_{i=1}^{k} W_i = b_k$ for $k = 1, 2, \ldots, m$. So $\sum_{i \in R-S} W_i = b_m = 0$ if $R \cap S = \emptyset$. We thus may assume that $R \cap S$ is not empty. Let $\{i_0 + 1, i_0 + 2, \ldots, i_0 + k\}$ be in $R \cap S$ for some integer $k \geq 0$. Any arc $(i+1, i)$ in $Lr \cap Q$ is red or white and hence $b(i) = 0$. Therefore $b(i_0) = 0 = b(i_0 + k)$, which implies that $\sum_{i=i_0+1}^{i_0+k} W_i = 0$ and that $\sum_{i \in R \cap S} W_i = 0$. Now $\sum_{i=1}^{m} W_i = 0$ shows that $\sum_{i \in R-S} W_i = 0$. The same argument can be employed to prove the second part in Condition B.

PROPOSITION 3.5. *If $UL = \emptyset$, then $w$, $r$ and $s$ are optimal.*

Proposition 3.5 is simply a corollary of the sufficiency of Theorem 2.3 and Propositions 3.2, 3.3 and 3.4. Now we will prove the finiteness of the algorithm.

PROPOSITION 3.6. *New $UL$ is a subset of old $UL$. Furthermore new $UL$ is a proper subset of old $UL$ after a circuit.*

PROOF: Assume that PNA ends up with a circuit $P$. Let any arc $(i, j) \in Lrc$ be not in old $UL$. If the arc $(i, j)$ is in $P$, it is green, white or black. For any color it is that $r(i) + s(j) = a_{ij}$. Because no change occurs in $r$ and $s$ after a circuit, the node $(i, j)$ is not in new $UL$. Since $w$ is unchanged on the set of arcs outside $P$, the first part holds after a circuit. The second part follows from the fact that the arc $(i*, j*)$ is in $P^+$ and hence $w(i*, j*)$ becomes 1. For the case of a cut $Q = [S, N - S]$, assume that $(i, j)$ in $Lrc$ is in new $UL$ but not in old $UL$. Since there occurs no change in $w$ after a cut, $w(i, j) = 0$ and $(i, j)$ is green, which in turn shows that both nodes are either in $S$ or in $N - S$. It is easy to show that $r'(i) + s'(j) = a_{ij}$ for any case. This contradicts to the fact that $(i, j)$ is in new $UL$.

PROPOSITION 3.7. *The algorithm is finite.*

PROOF: Because of Propostion 3.6, it suffices to show that whenever PNA ends up with a cut $Q = [S, N - S]$ at a certain iteration, either

(a) New $UL$ is a proper subset of old $UL$

or

(b) The outcome of PNA is a circuit after a finite number of iterations in which there occur only cuts.

If $\delta = a_{ij} - r(i) - s(j)$ for any $(i, j) \in Q^-$ with $w(i, j) = 0$, then $a_{ij} = r'(i) + s'(j)$ and (a) occurs. Assume this is not the case and let $\delta$ attain its minimum at an arc $L_0$. It is easy to show that the arc $L_0$ becomes green or white if it is in $Q^+$ and black if in $Q^-$. The arc $(i*, j*)$ is still in $UL$ and PNA is supposed to be initiated with the same $N^+ = \{j*\}$ and $N^- = \{i*\}$. If we have another $Q' = [S', N - S']$, then $S$ is a proper subset of $S'$. For there occurs no change in $r$ and $s$ on $S$, the colors of the arcs inside $S$ are unchanged after the cut $Q$

and finally one end node of the arc $L_0$ comes in $S'$. But there are only finite number of nodes, so there must be a circuit after a finite number of cuts in a row, which estableishes the proof.

We have shown that the algorithm produces the optimal solutions and the optimal values are the same. Therefore the necessity of theorem 2.3 is proved by lemma 2.2.

## REFERENCES

1. J.W. Tukey, "Exploratory Data Analysis," Addison-Wesley, 1977.
2. A.F. Siegel, *Low Median and Least Absolute Residual Analysis of Two-Way Tables*, J. Amer. Stat. Assoc. (1983), 371–374.
3. J.H.B. Kemperman, *Least Absolute Value and Median Polish in Inequality in Statistics and Probability*, IMS Lecture Note Series 5 (1984), 84–103.
4. A.M. Fink, *How to Polish off Median Polish*, ISU Math. Dept., Ames. Iowa (to appear).
5. R.T. Rockafellar, "Network flows and monotropic optimization," Wiley-Interscience, 1984.

Department of Mathematics
Chungnam National University
Taejon, 302-764, Korea