

# Iterative Proportional Scaling Revisited: A Modern Optimization Perspective

Yiyuan She, Shao Tang

Department of Statistics, Florida State University

# What is IPS?

- ▶ We would like to adjust a given table  $\mathbf{q}$  (say, a table of ones) get a new table  $\boldsymbol{\mu}$  satisfying a set of given margins
- ▶ Consider a **three-way** example. Let  $\{m_{ij+}, m_{+jk}, m_{i+k}\}$  be the pre-specified margins, possibly calculated from  $[m_{ijk}]$
- ▶ Starting with  $\boldsymbol{\mu}^{(0)} = \mathbf{q}$  (often  $\mathbf{1}$ ), IPS performs the following **scaling** operations iteratively for any  $t \geq 0$ :

$$\mu_{ijk}^{(3t+1)} = \mu_{ijk}^{(3t)} \frac{m_{ij+}}{\mu_{ij+}^{(3t)}}, \quad \forall (i, j, k)$$

$$\mu_{ijk}^{(3t+2)} = \mu_{ijk}^{(3t+1)} \frac{m_{i+k}}{\mu_{i+k}^{(3t+1)}}, \quad \forall (i, j, k)$$

$$\mu_{ijk}^{(3t+3)} = \mu_{ijk}^{(3t+2)} \frac{m_{+jk}}{\mu_{+jk}^{(3t+2)}}, \quad \forall (i, j, k)$$

# Who Cares?

- ▶ IPS has widespread applications in **computer science** (Rote & Zachariasen 07), **econometrics** (Lahr & Mesnarad 04), and **mathematics** (Sinkhorn & Knopp 67, Pretzel 80)
- ▶ Matrix raking
- ▶ **Contingency table analysis**: test of association
- ▶ **Log-linear model fitting**: estimation and goodness-of-fit
  - $X^2, G^2$ : expected frequencies. Log-linear:  $\mathbf{X}^T \boldsymbol{\mu} = \mathbf{X}^T \mathbf{y}$

# History

- ▶ Deming & Stephan (40): table adjustment (aiming at  $X^2$ )
- ▶ Ireland & Kullback (68): **KL** entropy subject to constraints
- ▶ Bishop et al (75): log-linear **likelihood** for tables
- ▶ Csiszar (75): **duality** holds though Lagrangian cannot apply
- ▶ **Convergence** studies undergo a long history
  - Brown (59), Birch (63), Fienberg (70), Haberman (74), Bishop et al (75), Csiszar (75), Ruschendorf (95), Pukelsheim & Simeone (09), Kurras (15), just to name a few
- ▶ **Extensions:**
  - Generalized Iterative Scaling (**GIS**) (Darroch & Ratcliff 72). Improved Iterative Scaling (**IIS**) (Pietra et al 97)

# Big Data Era

- ▶ Nowadays, IPS is not as popular as **Newton**-type algorithms, and has become an outdated method

*“primarily of historical interest”* (Agresti 12)

- ▶ The algorithm is pretty motivating in the **big data** era
- ▶ For a 5-way table with 10 levels each, a model including up to three-term interactions has 8,146 parameters
- ▶ Hessian: more than  $10^7$  entries..
- ▶ **First-order** methods? No universal step size.

# Pros and Cons of (Plain) IPS

## Pros

- ▶ No step-size parameter
- ▶ Memory saving
- ▶ Good numerical stability

## Cons

- ▶ No coefficient estimate (or standard errors) but the mean
- ▶ Narrow scope: only applies to tables (MLE, no shrinkage)
- ▶ Slow convergence, though cost-effective per iteration

We would like to study IPS from a modern **optimization** perspective to **improve** and **generalize** this elegant algorithm

# Setting

- ▶ Given a multi-way table, we can introduce dummy variables to form a **binary** design matrix  $\mathbf{X} \in \mathbb{R}^{N \times p}$ .
- ▶ Accordingly, we vectorize the table of counts to  $\mathbf{n} \in \mathbb{R}^N$ .
- ▶ IPS (a general version) can be described by
  1.  $\log \boldsymbol{\mu}^{(0)} \in \log \mathbf{q} + \mathcal{R}(\mathbf{X})$  and  $t \leftarrow 1$ .
  2. With  $\boldsymbol{\mu}^{(t,0)} \leftarrow \boldsymbol{\mu}^{(t-1)}$ ,  $\boldsymbol{\mu}^{(t,j)} = [\mu_i^{(t,j)}]$  ( $1 \leq j \leq p$ ) are obtained sequentially below:

$$\mu_i^{(t,j)} = \begin{cases} \mu_i^{(t,j-1)} \frac{\langle \mathbf{x}_j, \mathbf{n} \rangle}{\langle \mathbf{x}_j, \boldsymbol{\mu}^{(t,j-1)} \rangle} & x_{ij} = 1 \\ \mu_i^{(t,j-1)} & x_{ij} = 0, \end{cases} \quad 1 \leq i \leq N,$$

3.  $\boldsymbol{\mu}^{(t+1)} \leftarrow \boldsymbol{\mu}^{(t,p)}$ . If not converged,  $t \leftarrow t + 1$ , go to Step 2.

# Log-Affine Models

- ▶ Assume  $n_i$  are independent and follow Poisson distributions

$$n_i \sim Poi(\mu_i), \text{ with } \log \boldsymbol{\mu} - \log \mathbf{q} = \mathbf{X}\boldsymbol{\beta} \text{ for some } \boldsymbol{\beta} \in \mathbb{R}^p$$

- ▶ When  $\mathbf{q} = \mathbf{1}$ , it degenerates to a log-linear model
- ▶ The maximum likelihood problem for  $\boldsymbol{\mu}$  is then

$$\begin{aligned} \min_{\boldsymbol{\mu}} l(\boldsymbol{\mu}) &\triangleq -\langle \mathbf{n}, \log \boldsymbol{\mu} \rangle + \langle \mathbf{1}, \boldsymbol{\mu} \rangle \\ \text{s.t. } \boldsymbol{\mu} \in \mathcal{M} &\triangleq \{ \boldsymbol{\mu} \mid \log \boldsymbol{\mu} \in \log \mathbf{q} + \mathcal{R}_{\mathbf{X}}, \boldsymbol{\mu} \succeq \mathbf{0} \}. \end{aligned}$$

- ▶ Table-derived designs are always binary. In general, we could have [binary](#)/[nonnegative](#)/[general](#) designs



## Formulation through **Coefficients**

- ▶ Seen from the initialization requirement and the model description, it is very natural to define the problem via  $\beta$

$$\min_{\beta \in \mathbb{R}^p} l(\beta) = -\langle \mathbf{n}, \mathbf{X}\beta \rangle + \langle \mathbf{q}, \exp(\mathbf{X}\beta) \rangle$$

- ▶ A simple **cyclic coordinate descent** algorithm

$$\beta_j^{(t+1)} \in \arg \min_{\beta_j} l(\beta_1^{(t+1)}, \dots, \beta_{j-1}^{(t+1)}, \beta_j, \beta_{j+1}^{(t)}, \dots, \beta_p^{(t)})$$

- ▶ Initialize through  $\beta$ ; then set  $\mu^{(0)} \leftarrow \mathbf{q} \circ \exp(\mathbf{X}\beta^{(0)})$
- ▶ Mean-update:  $\mu^{(t,j)} \leftarrow \mu^{(t,j-1)} \circ \exp[\mathbf{x}_j(\beta_j^{(t+1)} - \beta_j^{(t)})]$   
(scaling)

- ▶  $\partial l / \partial \beta_j = 0$  gives

$$(\mathbf{x}_j \circ \boldsymbol{\mu}^{(t,j-1)})^T \exp\{\mathbf{x}_j(\beta_j - \beta_j^{(t)})\} - \mathbf{x}_j^T \mathbf{n} = 0.$$

- ▶ With  $x_{ij} = 0$  or  $1$ , we have  $x_{ij} \exp(x_{ij}\beta_j) = x_{ij} \exp(\beta_j)$  or  $\beta_j^{(t+1)} = \beta_j^{(t)} + \log(\frac{\mathbf{x}_j^T \mathbf{n}}{\mathbf{x}_j^T \boldsymbol{\mu}^{(t,j-1)}})$ , and so

$$\boldsymbol{\mu}^{(t,j)} = \boldsymbol{\mu}^{(t,j-1)} \circ \exp\{\mathbf{x}_j \log[(\mathbf{x}_j^T \mathbf{n}) / (\mathbf{x}_j^T \boldsymbol{\mu}^{(t,j-1)})]\}$$

which is exactly IPS!

- ▶ We get a novel CD characterization of IPS

# Benefits

- ▶ Flexibility in **initialization**: an arbitrary  $\beta^{(0)} \in \mathbb{R}^p$  suffices
- ▶ IPS can produce **parameter** estimates (as well as an estimate of the asymptotic covariance)!
- ▶ No need to scan the data repeatedly to fulfill all margins  
Multiple components of  $\mu$  can be **simultaneously** updated
- ▶ The CD perspective also makes it easy to generalize and improve IPS as well as facilitating theoretical investigations

# IPS for Minimizing $X^2$

- ▶ Deming & Stephan (40) conjectured that IPS minimizes the  $X^2$  statistic, later known to be incorrect
- ▶ Applying CD to  $\min_{\beta \in \mathbb{R}^p} X^2(\beta) \triangleq \sum_i \frac{[n_i - \mu_i(\beta)]^2}{\mu_i(\beta)}$  gives

$$\boldsymbol{\mu}^{(t,j)} = \boldsymbol{\mu}^{(t,j-1)} \circ \exp\left\{\frac{\boldsymbol{x}^j}{2} \log\left[\boldsymbol{x}_j^T (\boldsymbol{n} \circ \boldsymbol{n} \oslash \boldsymbol{\mu}^{(t,j-1)}) / (\boldsymbol{x}_j^T \boldsymbol{\mu}^{(t,j-1)})\right]\right\}$$

- ▶ This is the desired  $X^2$ -version of IPS

# Theoretical Properties

- ▶ From our derivation, it is easy to show for  $\forall t \geq 0, j \in [p]$

$$l(\boldsymbol{\beta}^{(t,j)}) \geq l(\boldsymbol{\beta}^{(t+1,j)})$$

- ▶ With the intercept present, we can further show that the same conclusion holds for  $G^2 \triangleq 2 \sum_i n_i \log(n_i/\mu_i)$ 
  - Many researchers observed that  $G^2$  decreases during IPS
  - This is just a natural outcome of its CD nature
- ▶ Furthermore, the sequence of **iterates** converges to the unique optimal solution, at least **linearly**

# Accelerations

- ▶ Speeding the convergence of IPS is crucial on big tables
- ▶ **Cyclic** update has slow (worst-case) convergence rate
- ▶ One idea: pick the coordinate with the “largest” derivative

$$\text{Gauss-Southwell rule: } j = \arg \max_j |\nabla_j l(\boldsymbol{\beta}^{(t)})|$$

- ▶ But this is not worthy since when the full gradient is available, one can directly update the whole  $\boldsymbol{\beta}$ -vector
- ▶ A powerful idea is to apply **randomization!**
- ▶ In theory, randomized coordinate descent can avoid the worst-case bounds (Nesterov 12). We call it **A-IPS**

## An MM Viewpoint

- ▶ Majorization-Minimization algorithms are popular alternatives to EM algorithms
- ▶ Rather than directly minimizing  $l(\beta)$ , we construct a surrogate function  $g(\beta | \beta^-)$  satisfying

$$g(\beta | \beta^-) \geq l(\beta) , \quad g(\beta | \beta) = l(\beta)$$

- ▶ Let  $\beta^{(t+1)} \in \arg \min_{\beta} g(\beta | \beta^{(t)})$ . Then the function value is guaranteed to be **decreasing**
- ▶ Once getting  $\beta^{(t+1)}$ , update  $\mu$  via **proportional scaling**

$$\mu^{(t+1)} = \mu^{(t)} \circ \exp\left[\sum_j x_j (\beta_j^{(t+1)} - \beta_j^{(t)})\right]$$

- ▶ The problem boils down to the design of proper surrogates

- ▶ Recall the objective function:  $l(\boldsymbol{\beta}) = \sum_i l_i$  with  $l_i = q_i \exp(\tilde{\mathbf{x}}_i^T \boldsymbol{\beta}) - n_i \tilde{\mathbf{x}}_i^T \boldsymbol{\beta}$ .
- ▶ Applying Jensen's inequality gives

$$\begin{aligned}
 l_i &\leq -n_i \tilde{\mathbf{x}}_i^T \boldsymbol{\beta} + \sum_j \alpha_{ij} q_i \exp\left[\frac{x_{ij}}{\alpha_{ij}} (\beta_j - \beta_j^-) + \tilde{\mathbf{x}}_i^T \boldsymbol{\beta}^-\right] \\
 &= -n_i \tilde{\mathbf{x}}_i^T \boldsymbol{\beta} + \sum_j \alpha_{ij} \mu_i^- \exp\left[\frac{x_{ij}}{\alpha_{ij}} (\beta_j - \beta_j^-)\right],
 \end{aligned}$$

where  $\alpha_{ij}$  satisfy  $\alpha_{ij} \geq 0$ ,  $\sum_{j=1}^p \alpha_{ij} = 1$  for all  $i$ .



## First Try

- ▶ Let  $\alpha_{ij} = 1/p$ . Then for all binary designs, we get

$$\boldsymbol{\beta}^{(t+1)} = \boldsymbol{\beta}^{(t)} + \frac{1}{p} \log[(\mathbf{X}^T \mathbf{n}) \oslash (\mathbf{X}^T \boldsymbol{\mu}^{(t)})]$$

$$\boldsymbol{\mu}^{(t+1)} = \boldsymbol{\mu}^{(t)} \circ (\exp\{\frac{1}{p} \mathbf{X} \log[(\mathbf{X}^T \mathbf{n}) \oslash (\mathbf{X}^T \boldsymbol{\mu}^{(t)})]\})$$

- ▶ This MM algorithm resembles IPS, but updates all components synchronously with a smaller stepsize ( $1/p$ )

# Non-negative Designs

- ▶ A perhaps better choice is to set  $\alpha_{ij} = x_{ij}/x_{i+}$ , leading to

$$l_i \leq - \sum_j n_i x_{ij} \beta_j + \sum_j \frac{x_{ij}}{x_{i+}} \mu_i^- \exp[x_{i+}(\beta_j - \beta_j^-)]$$

- ▶ However, the  $g$ -optimization has no a closed-form solution
- ▶ We can majorize  $\exp[x_{i+}(\beta_j - \beta_j^-)]/x_{i+}$  further. For any  $a$ ,  $b$  with  $0 < a \leq b$ ,  $\exp(at)/a - 1/a \leq \exp(bt)/b - 1/b$
- ▶ With  $a = x_{i+}$  and  $b = \max_i x_{i,+} = \|\mathbf{X}\|_\infty \triangleq R$ , we get

$$g_2(\boldsymbol{\beta} \mid \boldsymbol{\beta}^-) = -\mathbf{n}^T \mathbf{X} \boldsymbol{\beta} + \sum_{i,j} x_{ij} \mu_i^- \left\{ \frac{\exp[R(\beta_j - \beta_j^-)]}{R} - \frac{1}{R} + \frac{1}{x_{i+}} \right\}$$

which has an explicit solution for any **nonnegative**  $\mathbf{X}$ !

## Extending the **Generalized Iterative Scaling**

- ▶ The iterates are then given by

$$\boldsymbol{\beta}^{(t+1)} = \boldsymbol{\beta}^{(t)} + \frac{1}{R} \log[(\mathbf{X}^T \mathbf{n}) \oslash (\mathbf{X}^T \boldsymbol{\mu}^{(t)})]$$

$$\boldsymbol{\mu}^{(t+1)} = \boldsymbol{\mu}^{(t)} \circ (\exp\{\frac{1}{R} \mathbf{X} \log[(\mathbf{X}^T \mathbf{n}) \oslash (\mathbf{X}^T \boldsymbol{\mu}^{(t)})]\})$$

- ▶ Larger stepsize than the previous one
- ▶ This takes the celebrated **GIS** (Darroch & Ratcliff 72) as a particular instance
- ▶ In fact, using the MM trick, we can even extend IPS to fit any log-affine models with **arbitrary** designs

# A Reparametrization Trick

- ▶ Assume the intercept is present:  $\mathbf{X} = [\mathbf{1} \ \dot{\mathbf{X}}]$ ,  $\boldsymbol{\beta} = [\beta_0 \ \dot{\boldsymbol{\beta}}^T]^T$
- ▶ Introducing  $\alpha \triangleq \beta_0 + \log \langle \mathbf{q}, \exp(\dot{\mathbf{X}} \dot{\boldsymbol{\beta}}) \rangle$ , or  $\exp(\alpha) = \langle \mathbf{q}, \exp(\beta_0 \mathbf{1} + \dot{\mathbf{X}} \dot{\boldsymbol{\beta}}) \rangle = \langle \mathbf{1}, \boldsymbol{\mu} \rangle$ ,  $l(\boldsymbol{\beta})$  becomes

$$-\langle \mathbf{n}, \dot{\mathbf{X}} \dot{\boldsymbol{\beta}} \rangle + \langle \mathbf{1}, \mathbf{n} \rangle \log \langle \mathbf{q}, \exp(\dot{\mathbf{X}} \dot{\boldsymbol{\beta}}) \rangle + \exp(\alpha) - \langle \mathbf{1}, \mathbf{n} \rangle \alpha$$

- ▶ It suffices to study

$$\min_{\dot{\boldsymbol{\beta}}} L(\dot{\boldsymbol{\beta}}) \triangleq -\langle \mathbf{n}, \dot{\mathbf{X}} \dot{\boldsymbol{\beta}} \rangle + \langle \mathbf{1}, \mathbf{n} \rangle \log \langle \mathbf{q}, \exp(\dot{\mathbf{X}} \dot{\boldsymbol{\beta}}) \rangle$$

- ▶ This gives a really *nice* form in optimization!

- ▶ Due to the **concavity** of the log function, for any  $\zeta > 0$ , we have  $\log(\zeta x) \leq \zeta x - 1$  (with equality at  $\zeta = 1/x$ )
- ▶  $L(\mathring{\beta}) \leq -\langle \mathbf{1}, \mathbf{n} \rangle \log \zeta - \langle \mathbf{n}, \mathring{X} \mathring{\beta} \rangle + \langle \mathbf{1}, \mathbf{n} \rangle [\zeta \langle \mathbf{q}, \exp(\mathring{X} \mathring{\beta}) \rangle - 1]$
- ▶ Assuming the non-negative setting and applying Jensen's inequality with  $\alpha_{ij} = \mathring{x}_{ij} / \mathring{x}_{i+}$ , we get

$$g_3(\mathring{\beta} \mid \mathring{\beta}^-) = -\langle \mathbf{n}, \mathring{X} \mathring{\beta} \rangle - \langle \mathbf{1}, \mathbf{n} \rangle (\log \zeta + 1) \\ + \zeta \langle \mathbf{1}, \mathbf{n} \rangle \sum_{i \in [N], j \in [p-1]} \frac{\mathring{x}_{ij}}{\mathring{x}_{i+}} \hat{\mu}_i^- \exp[\mathring{x}_{i+} (\mathring{\beta}_j - \mathring{\beta}_j^-)]$$

- ▶ The only choice of  $\zeta$  to guarantee that  $g_3$  is a surrogate is  $\zeta = 1 / \langle \mathbf{q}, \exp(\mathring{X} \mathring{\beta}^-) \rangle$

# Recovering the **Improved Iterative Scaling**

- ▶  $\hat{\beta}_j^{(t+1)}$  is determined by

$$\frac{\exp(\beta_0^{(t)}) \langle \mathbf{1}, \mathbf{n} \rangle}{\langle \mathbf{1}, \boldsymbol{\mu}^{(t)} \rangle} \sum_i \hat{x}_{ij} \hat{\mu}_i^{(t)} \exp[\hat{x}_{i+} (\hat{\beta}_j^{(t+1)} - \hat{\beta}_j^{(t)})] = \sum_i n_i \hat{x}_{ij}$$

- ▶ Interestingly, this MM algorithm fully restores the **IIS** (Pietra et al 97), if we add a mean-normalization and an intercept update
- ▶ These two steps are however unnecessary in our algorithm

# Quadratic Surrogates

- ▶ The re-parameterized form allows for linearization

$$Q(\dot{\boldsymbol{\beta}} \mid \dot{\boldsymbol{\beta}}^-) = L(\dot{\boldsymbol{\beta}}^-) + \langle \dot{\boldsymbol{\beta}} - \dot{\boldsymbol{\beta}}^-, \nabla_{\dot{\boldsymbol{\beta}}} L(\dot{\boldsymbol{\beta}}^-) \rangle + \frac{1}{2} (\dot{\boldsymbol{\beta}} - \dot{\boldsymbol{\beta}}^-)^T \mathbf{W} (\dot{\boldsymbol{\beta}} - \dot{\boldsymbol{\beta}}^-)$$

- ▶  $Q$  is a valid surrogate function provided that

$$\langle \mathbf{1}, \mathbf{n} \rangle \dot{\mathbf{X}}^T \left[ \frac{\text{diag}(\dot{\boldsymbol{\mu}})}{\langle \mathbf{1}, \dot{\boldsymbol{\mu}} \rangle} - \frac{\dot{\boldsymbol{\mu}} \dot{\boldsymbol{\mu}}^T}{\langle \mathbf{1}, \dot{\boldsymbol{\mu}} \rangle^2} \right] \dot{\mathbf{X}} \preceq \mathbf{W}, \quad \forall \dot{\boldsymbol{\mu}} \succeq \mathbf{0}$$

- ▶ An obvious choice is  $\mathbf{W} = (\langle \mathbf{1}, \mathbf{n} \rangle \|\dot{\mathbf{X}}\|_2^2 / 2) \mathbf{I}$
- ▶ A finer one:

$$\mathbf{W} = \dot{\mathbf{X}}^T (\langle \mathbf{1}, \mathbf{n} \rangle \mathbf{I} - \mathbf{1} \mathbf{1}^T) \dot{\mathbf{X}} / 2$$

- ▶ Fancier techniques like Nesterov's **second acceleration** can be applied to speed the convergence

- ▶ Initialization:  $t \leftarrow 0$ ,  $\dot{\beta} \leftarrow \dot{\beta}^{(0)}$ ,  $\dot{\eta}^{(0)} \leftarrow \dot{\beta}^{(0)}$ ,  $\theta_0 = 1$

- ▶ The momentum trick to guarantee faster convergence:

1.  $\dot{\alpha}^{(t)} \leftarrow (1 - \theta_t)\dot{\beta}^{(t)} + \theta_t\dot{\eta}^{(t)}$
2.  $\dot{\eta}^{(t+1)} \leftarrow \dot{\eta}^{(t)} - \theta_t^{-1}\mathbf{W}^{-1}\{-\dot{\mathbf{X}}^T \mathbf{n} + \frac{\langle \mathbf{1}, \mathbf{n} \rangle}{\langle \mathbf{q}, \exp(\dot{\mathbf{X}}\dot{\alpha}^{(t)}) \rangle} \dot{\mathbf{X}}^T [\mathbf{q} \circ \exp(\dot{\mathbf{X}}\dot{\alpha}^{(t)})]\}$
3.  $\dot{\beta}^{(t+1)} \leftarrow (1 - \theta_t)\dot{\beta}^{(t)} + \theta_t\dot{\eta}^{(t+1)}$
4.  $\dot{\mu}^{(t+1)} \leftarrow \dot{\mu}^{(t)} \circ \exp[\dot{\mathbf{X}}(\dot{\beta}^{(t+1)} - \dot{\beta}^{(t)})]$
5.  $\theta_{t+1} \leftarrow (\sqrt{\theta_t^4 + 4\theta_t^2} - \theta_t^2)/2$
6.  $t \leftarrow t + 1$



# Divide-and-Conquer for Big Data

- ▶ The CD idea can be extended further to BCD
- ▶ We recommend **randomized** BCD. But the key is not to pick a block at random, but to perform **random blocking**
- ▶ This **B-IPS** algorithm can combine the virtues of IPS and Newton-like algorithms and is extremely scalable
  - Typically we adopt block sizes  $\approx 200$

# Shrinkage Estimation

- ▶ We use the  $\ell_1$ -penalized contingency **table** analysis as an example (the  $\ell_2$ -penalty is much easier to handle)
- ▶ To identify significant association terms, we can solve

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^p} l_1(\boldsymbol{\beta}) \triangleq -\langle \mathbf{n}, \mathbf{X}\boldsymbol{\beta} \rangle + \langle \mathbf{q}, \exp(\mathbf{X}\boldsymbol{\beta}) \rangle + \sum_{j=1}^p \lambda_j |\beta_j|,$$

where the penalty is not imposed on the intercept ( $\lambda_1 = 0$ )

- ▶ IPS/CD:  $\min_{\beta_j \in \mathbb{R}} l_1(\beta_0^-, \dots, \beta_{j-1}^-, \beta_j, \beta_{j+1}^-, \dots, \beta_{p-1}^-)$

# Poisson ‘Soft-Thresholding’ Operator

- ▶ Define

$$\Theta(\mathbf{x}, \mathbf{n}, \mathbf{q}; \lambda) = \begin{cases} \log \frac{\langle \mathbf{x}, \mathbf{n} \rangle - \lambda \operatorname{sgn} \langle \mathbf{x}, \mathbf{n} - \mathbf{q} \rangle}{\langle \mathbf{x}, \mathbf{q} \rangle}, & |\langle \mathbf{x}, \mathbf{n} - \mathbf{q} \rangle| \geq \lambda \\ \mathbf{0}, & |\langle \mathbf{x}, \mathbf{n} - \mathbf{q} \rangle| < \lambda \end{cases}$$

- ▶ Then, an  $\ell_1$ -IPS algorithm for tables can be developed (similar to the fast CD lasso algorithm)

$$\beta_j = \Theta(\mathbf{x}, \mathbf{n}, \mathbf{q} \circ \exp(\mathbf{X}_{\setminus j} \boldsymbol{\beta}_{\setminus j}^-); \lambda), \quad j = 1, \dots, p$$

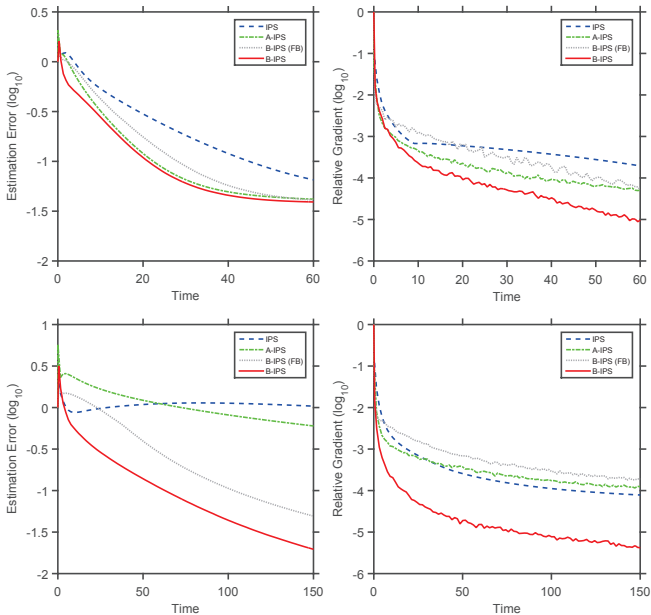


Figure: Performance comparison in two (small) table settings ( $p = 523$ ,  $N = 100K$ )

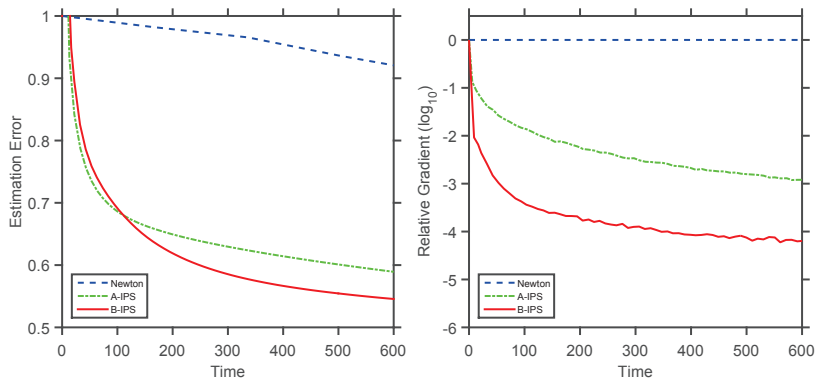


Figure: Large table setting ( $p = 8,146$ ,  $N = 100K$ ).

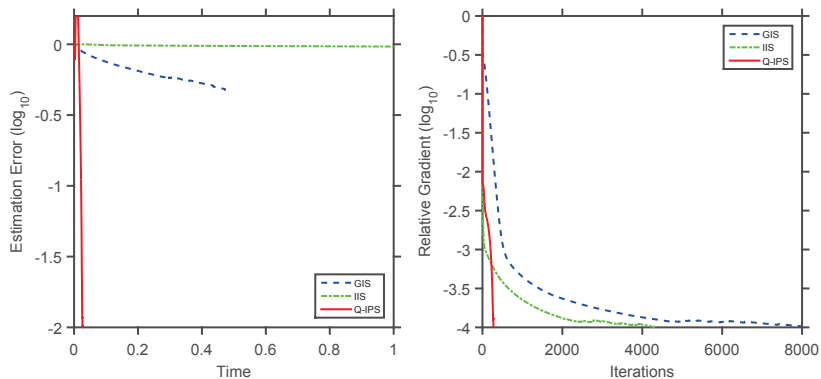


Figure: Nonnegative design (small):  $N = 1000, p = 100$ .

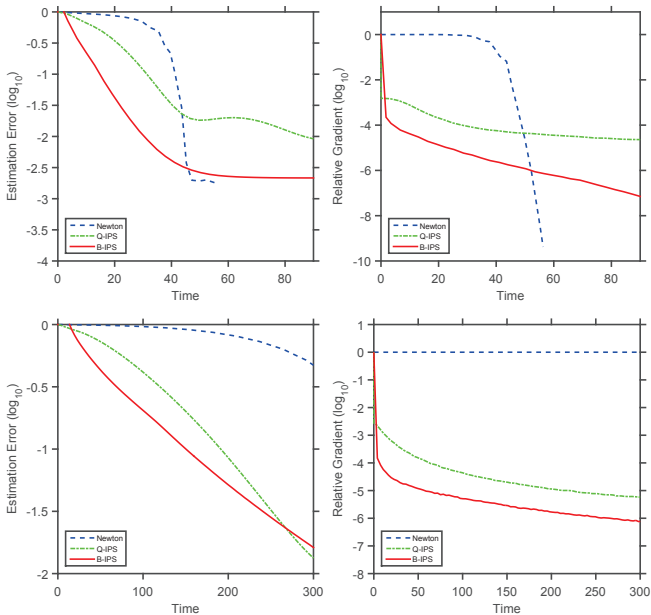


Figure: Large non-negative designs ( $p = 2,000, p = 4,000$ )

► General designs

	$p = 1,000$ tol = 1e-6		$p = 4,000$ tol = 1e-6		$p = 10,000$ $T_{\max} = 600$	
	Time	estErr	Time	estErr	relGrad	estErr
Newton	35.2	0.12	528.5	0.1	—	—
L-BFGS	35.6	2.3	276.5	2.6	7622	$9.3e + 2$
B-IPS	16.8	3.2	116.5	23.2	5.70	$1.0e + 2$



## $\ell_1$ Solution Path on Bank Marketing Data

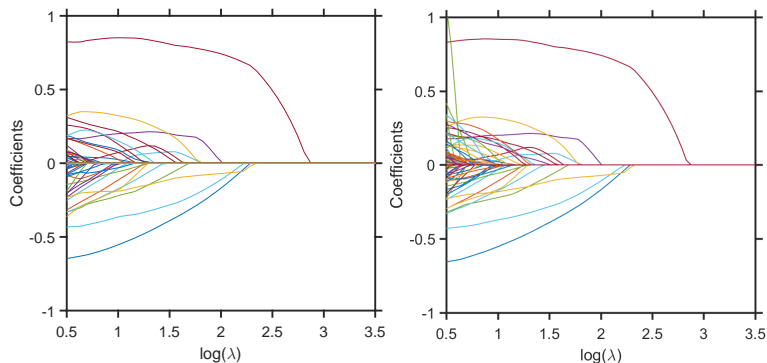


Figure:  $\ell_1$  solution path for two-way ( $p = 739$ ) and three-way ( $p = 5,784$ ) homogeneous association models on the left and right panel, respectively.

## $\ell_1$ Selected Variables

Main	month = mar (+) month = oct (+) education = univ (+) default = unkown (-) loan = yes (-)	month = may (-) month = nov (-) poutcome = success (+) contact = telephone (-)
Two-way	(poutcome = non-exist) * month = mar (+) (poutcome = non-exist) * month = oct (+) (poutcome = non-exist) * education = univ (+) (poutcome = non-exist) * contact = telephone (-) (poutcome = non-exist) * housing = unkown (-)	
Three-way	(job = technician) * (education = professional) * (marital = married) (+) (poutcome = non-exist) * education = univ * (marital = single) (+)	

Table: Selected variables with signs of coefficients in parentheses.

# Bank Marketing Data Analysis

- ▶ `month` plays an important role in the marketing campaign.
- ▶ Contact communication type being cellular and the outcome of previous marketing campaign being successful are favorable factors for successful subscription.
- ▶ All selected two-way interactions involve `poutcome = non-exist`.
- ▶ The three-way terms selected show no collinearity with the other selected variables.

# Summary

- ▶ We revisited the conventional IPS from a modern optimization perspective
- ▶ Based on a coordinate descent characterization, we showed that IPS can deliver estimates of **coefficients**
- ▶ From a majorization-minimization standpoint, IPS can be substantially **generalized** to handle any log-affine models
- ▶ IPS can be accelerated in various ways to provide extremely **scalable** algorithms in big data applications
- ▶ Our techniques apply to **shrinkage** estimation on count data