

Exploiting ultra-sparsity in the revised simplex method

Julian Hall

School of Mathematics
University of Edinburgh

`jajhall@ed.ac.uk`

Sparse Days, CERFACS, Toulouse

28 September 2018



THE UNIVERSITY
of EDINBURGH

LP and high performance simplex solvers

$$\text{minimize } f = \mathbf{c}^T \mathbf{x} \quad \text{subject to } \mathbf{Ax} = \mathbf{b} \quad \mathbf{x} \geq \mathbf{0}$$

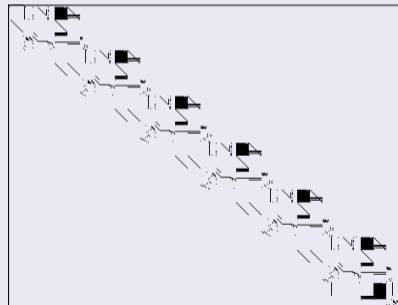
LP and high performance simplex solvers

$$\text{minimize } f = \mathbf{c}^T \mathbf{x} \quad \text{subject to } \mathbf{Ax} = \mathbf{b} \quad \mathbf{x} \geq \mathbf{0}$$

Background

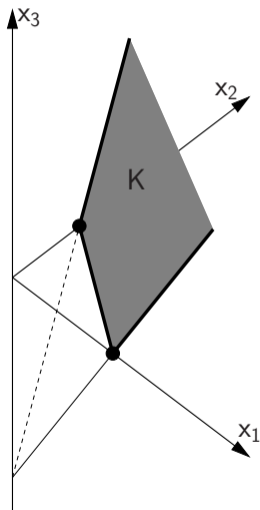
- Fundamental model in optimal decision-making
- Simplex method preferred when solving related problems
- High performance requires
 - Algorithmic tricks
 - Computational tricks

Example



STAIR: 356 rows, 467 columns and 3856 nonzeros

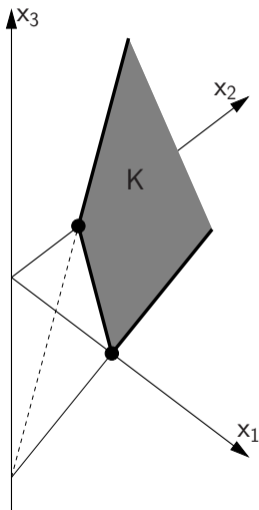
Solving LP problems: Characterizing the feasible region



$$\text{minimize } f = \mathbf{c}^T \mathbf{x} \quad \text{subject to } \mathbf{A}\mathbf{x} = \mathbf{b} \quad \mathbf{x} \geq \mathbf{0}$$

- Optimal solution at a vertex of feasible region K

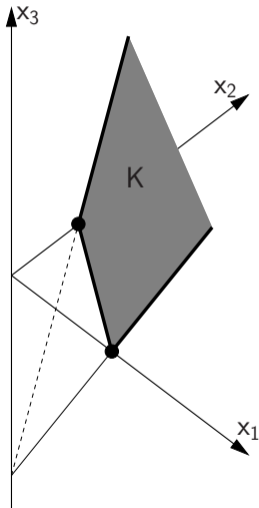
Solving LP problems: Characterizing the feasible region



$$\text{minimize } f = \mathbf{c}^T \mathbf{x} \quad \text{subject to } \mathbf{A}\mathbf{x} = \mathbf{b} \quad \mathbf{x} \geq \mathbf{0}$$

- Optimal solution at a vertex of feasible region K
- Vertex characterised by partition $\mathcal{B} \cup \mathcal{N}$ of index set

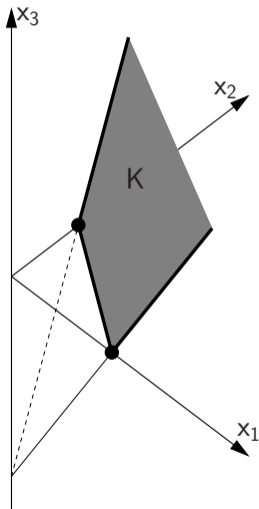
Solving LP problems: Characterizing the feasible region



$$\text{minimize } f = \mathbf{c}^T \mathbf{x} \quad \text{subject to } \mathbf{A}\mathbf{x} = \mathbf{b} \quad \mathbf{x} \geq \mathbf{0}$$

- Optimal solution at a vertex of feasible region K
- Vertex characterised by partition $\mathcal{B} \cup \mathcal{N}$ of index set with
 - Nonsingular **basis matrix** B

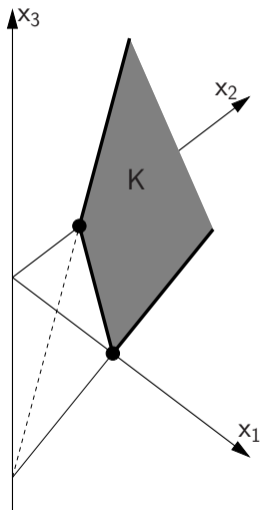
Solving LP problems: Characterizing the feasible region



$$\text{minimize } f = \mathbf{c}^T \mathbf{x} \quad \text{subject to } \mathbf{A}\mathbf{x} = \mathbf{b} \quad \mathbf{x} \geq \mathbf{0}$$

- Optimal solution at a vertex of feasible region K
- Vertex characterised by partition $\mathcal{B} \cup \mathcal{N}$ of index set with
 - Nonsingular **basis matrix** B
 - Equations partitioned as $B\mathbf{x}_B + N\mathbf{x}_N = \mathbf{b}$

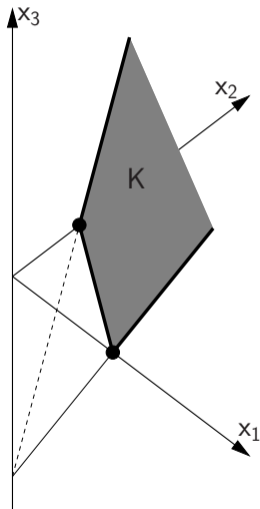
Solving LP problems: Characterizing the feasible region



$$\text{minimize } f = \mathbf{c}^T \mathbf{x} \quad \text{subject to } \mathbf{A}\mathbf{x} = \mathbf{b} \quad \mathbf{x} \geq \mathbf{0}$$

- Optimal solution at a vertex of feasible region K
- Vertex characterised by partition $\mathcal{B} \cup \mathcal{N}$ of index set with
 - Nonsingular **basis matrix** B
 - Equations partitioned as $B\mathbf{x}_B + N\mathbf{x}_N = \mathbf{b}$
 - Vertex given by $\mathbf{x}_N = \mathbf{0}$ and $\mathbf{x}_B = B^{-1}\mathbf{b} = \hat{\mathbf{b}}$

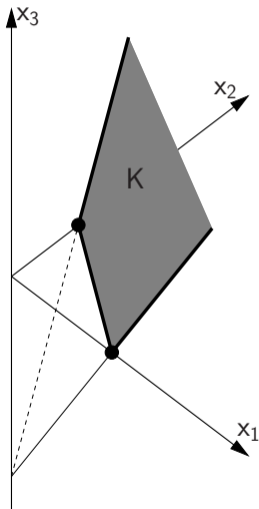
Solving LP problems: Characterizing the feasible region



minimize $f = \mathbf{c}^T \mathbf{x}$ subject to $\mathbf{Ax} = \mathbf{b}$ $\mathbf{x} \geq \mathbf{0}$

- Optimal solution at a vertex of feasible region K
- Vertex characterised by partition $\mathcal{B} \cup \mathcal{N}$ of index set with
 - Nonsingular **basis matrix** B
 - Equations partitioned as $B\mathbf{x}_B + N\mathbf{x}_N = \mathbf{b}$
 - Vertex given by $\mathbf{x}_N = \mathbf{0}$ and $\mathbf{x}_B = B^{-1}\mathbf{b} = \hat{\mathbf{b}}$
 - K given by $\mathbf{x}_B = \hat{\mathbf{b}} - B^{-1}N\mathbf{x}_N$ for some $\mathbf{x}_N \geq \mathbf{0}$

Solving LP problems: Optimality conditions



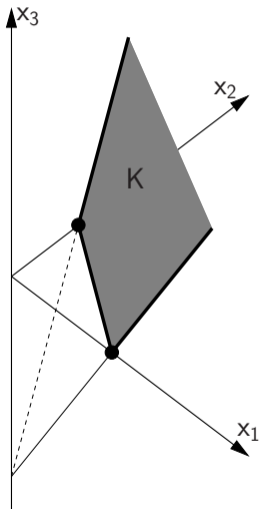
minimize $f = \mathbf{c}^T \mathbf{x}$ subject to $A\mathbf{x} = \mathbf{b}$ $\mathbf{x} \geq \mathbf{0}$

- Objective partitioned according to $\mathcal{B} \cup \mathcal{N}$ as

$$f = \mathbf{c}_B^T \mathbf{x}_B + \mathbf{c}_N^T \mathbf{x}_N = \hat{f} + \hat{\mathbf{c}}_N^T \mathbf{x}_N$$

where $\hat{f} = \mathbf{c}_B^T \hat{\mathbf{b}}$ and $\hat{\mathbf{c}}_N^T = \mathbf{c}_N^T - \mathbf{c}_B^T B^{-1} N$

Solving LP problems: Optimality conditions



minimize $f = \mathbf{c}^T \mathbf{x}$ subject to $\mathbf{Ax} = \mathbf{b}$ $\mathbf{x} \geq \mathbf{0}$

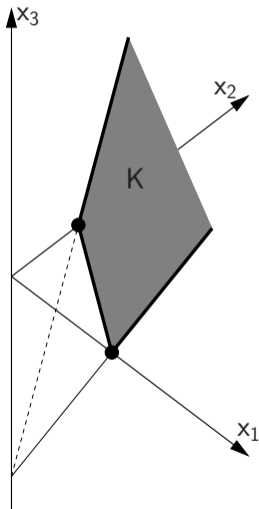
- Objective partitioned according to $\mathcal{B} \cup \mathcal{N}$ as

$$f = \mathbf{c}_B^T \mathbf{x}_B + \mathbf{c}_N^T \mathbf{x}_N = \hat{f} + \hat{\mathbf{c}}_N^T \mathbf{x}_N$$

where $\hat{f} = \mathbf{c}_B^T \hat{\mathbf{b}}$ and $\hat{\mathbf{c}}_N^T = \mathbf{c}_N^T - \mathbf{c}_B^T \mathbf{B}^{-1} \mathbf{N}$

- Partition yields an optimal solution if there is
 - Primal feasibility $\hat{\mathbf{b}} \geq \mathbf{0}$
 - Dual feasibility $\hat{\mathbf{c}}_N \geq \mathbf{0}$

Solving LP problems: Optimality conditions



minimize $f = \mathbf{c}^T \mathbf{x}$ subject to $\mathbf{Ax} = \mathbf{b}$ $\mathbf{x} \geq \mathbf{0}$

- Objective partitioned according to $\mathcal{B} \cup \mathcal{N}$ as

$$f = \mathbf{c}_B^T \mathbf{x}_B + \mathbf{c}_N^T \mathbf{x}_N = \hat{f} + \hat{\mathbf{c}}_N^T \mathbf{x}_N$$

where $\hat{f} = \mathbf{c}_B^T \hat{\mathbf{b}}$ and $\hat{\mathbf{c}}_N^T = \mathbf{c}_N^T - \mathbf{c}_B^T \mathbf{B}^{-1} \mathbf{N}$

- Partition yields an optimal solution if there is
 - Primal feasibility $\hat{\mathbf{b}} \geq \mathbf{0}$
 - Dual feasibility $\hat{\mathbf{c}}_N \geq \mathbf{0}$
- Simplex algorithm adjusts $\mathcal{B} \cup \mathcal{N}$ until optimal

Primal simplex algorithm

Assume $\hat{\mathbf{b}} \geq \mathbf{0}$ Seek $\hat{\mathbf{c}}_N \geq \mathbf{0}$

Scan $\hat{c}_j < 0$ for q to leave \mathcal{N}

Scan $\hat{b}_i / \hat{a}_{iq} > 0$ for p to leave \mathcal{B}

	\mathcal{N}		RHS
\mathcal{B}	$\hat{\mathbf{a}}_q$		$\hat{\mathbf{b}}$
	\hat{a}_{pq}	$\hat{\mathbf{a}}_p^T$	\hat{b}_p
	\hat{c}_q	$\hat{\mathbf{c}}_N^T$	

Primal simplex algorithm

Assume $\hat{\mathbf{b}} \geq \mathbf{0}$ Seek $\hat{\mathbf{c}}_N \geq \mathbf{0}$

Scan $\hat{c}_j < 0$ for q to leave \mathcal{N}

Scan $\hat{b}_i / \hat{a}_{iq} > 0$ for p to leave \mathcal{B}

Update: Exchange p and q between \mathcal{B} and \mathcal{N}

Update $\hat{\mathbf{b}} := \hat{\mathbf{b}} - \alpha_P \hat{\mathbf{a}}_q$ $\alpha_P = \hat{b}_p / \hat{a}_{pq}$

Update $\hat{\mathbf{c}}_N^T := \hat{\mathbf{c}}_N^T + \alpha_D \hat{\mathbf{a}}_p^T$ $\alpha_D = -\hat{c}_q / \hat{a}_{pq}$

	\mathcal{N}		RHS
\mathcal{B}	$\hat{\mathbf{a}}_q$		$\hat{\mathbf{b}}$
	\hat{a}_{pq}	$\hat{\mathbf{a}}_p^T$	\hat{b}_p
	\hat{c}_q	$\hat{\mathbf{c}}_N^T$	

Primal simplex algorithm

Assume $\hat{\mathbf{b}} \geq \mathbf{0}$ Seek $\hat{\mathbf{c}}_N \geq \mathbf{0}$

Scan $\hat{c}_j < 0$ for q to leave \mathcal{N}

Scan $\hat{b}_i / \hat{a}_{iq} > 0$ for p to leave \mathcal{B}

Update: Exchange p and q between \mathcal{B} and \mathcal{N}

Update $\hat{\mathbf{b}} := \hat{\mathbf{b}} - \alpha_P \hat{\mathbf{a}}_q$ $\alpha_P = \hat{b}_p / \hat{a}_{pq}$

Update $\hat{\mathbf{c}}_N^T := \hat{\mathbf{c}}_N^T + \alpha_D \hat{\mathbf{a}}_p^T$ $\alpha_D = -\hat{c}_q / \hat{a}_{pq}$

Data required

- Pivotal row $\hat{\mathbf{a}}_p^T = \mathbf{e}_p^T B^{-1} N$
- Pivotal column $\hat{\mathbf{a}}_q = B^{-1} \mathbf{a}_q$

	\mathcal{N}		RHS
\mathcal{B}	$\hat{\mathbf{a}}_q$		$\hat{\mathbf{b}}$
	\hat{a}_{pq}	$\hat{\mathbf{a}}_p^T$	\hat{b}_p
	\hat{c}_q	$\hat{\mathbf{c}}_N^T$	

Primal simplex algorithm

Assume $\hat{\mathbf{b}} \geq \mathbf{0}$ Seek $\hat{\mathbf{c}}_N \geq \mathbf{0}$

Scan $\hat{c}_j < 0$ for q to leave \mathcal{N}

Scan $\hat{b}_i / \hat{a}_{iq} > 0$ for p to leave \mathcal{B}

Update: Exchange p and q between \mathcal{B} and \mathcal{N}

Update $\hat{\mathbf{b}} := \hat{\mathbf{b}} - \alpha_P \hat{\mathbf{a}}_q$ $\alpha_P = \hat{b}_p / \hat{a}_{pq}$

Update $\hat{\mathbf{c}}_N^T := \hat{\mathbf{c}}_N^T + \alpha_D \hat{\mathbf{a}}_p^T$ $\alpha_D = -\hat{c}_q / \hat{a}_{pq}$

	\mathcal{N}		RHS
\mathcal{B}	$\hat{\mathbf{a}}_q$		$\hat{\mathbf{b}}$
	\hat{a}_{pq}	$\hat{\mathbf{a}}_p^T$	\hat{b}_p
	\hat{c}_q	$\hat{\mathbf{c}}_N^T$	

Computation

Pivotal row via $B^T \pi_p = \mathbf{e}_p$ **BTRAN** and $\hat{\mathbf{a}}_p^T = \pi_p^T N$ **PRICE**

Pivotal column via $B \hat{\mathbf{a}}_q = \mathbf{a}_q$ **FTRAN** Represent B^{-1} **INVERT**

Update B^{-1} exploiting $\bar{B} = B + (\mathbf{a}_q - B\mathbf{e}_p)\mathbf{e}_p^T$ **UPDATE-BASIS**

Recall: major computational components

- **BTRAN**: Form $\boldsymbol{\pi}_p = B^{-T} \mathbf{e}_p$
- **PRICE**: Form $\hat{\mathbf{a}}_p^T = \boldsymbol{\pi}_p^T N$
- **FTRAN**: Form $\hat{\mathbf{a}}_q = B^{-1} \mathbf{a}_q$

Simplex method: Hyper-sparsity

Recall: major computational components

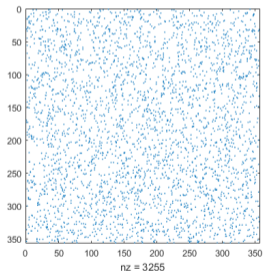
- **BTRAN**: Form $\pi_p = B^{-T} \mathbf{e}_p$
- **PRICE**: Form $\hat{\mathbf{a}}_p^T = \pi_p^T N$
- **FTRAN**: Form $\hat{\mathbf{a}}_q = B^{-1} \mathbf{a}_q$

Phenomenon of hyper-sparsity

- Vectors \mathbf{e}_p and \mathbf{a}_q are sparse
 - In **BTRAN**, π_p is a row of B^{-1}
 - In **FTRAN**, $\hat{\mathbf{a}}_q$ is a linear combination of a few columns of B^{-1}
- Results π_p and $\hat{\mathbf{a}}_q$ may be sparse—because B^{-1} is sparse
 - In **PRICE**, $\hat{\mathbf{a}}_p^T$ is a linear combination of a few rows of N

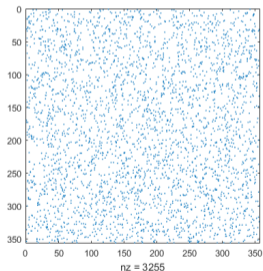
Inverse of a sparse matrix and solution of $Bx = r$

Random sparse matrix B
 $m = 356$ and density 2.5%

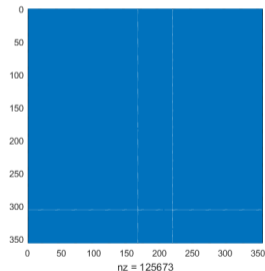


Inverse of a sparse matrix and solution of $Bx = r$

Random sparse matrix B
 $m = 356$ and density 2.5%

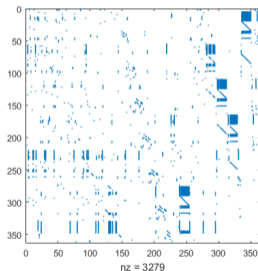


B^{-1} has density of 99%



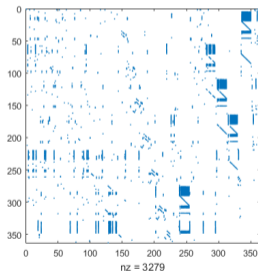
Inverse of a sparse matrix and solution of $Bx = r$

Optimal B for LP problem stair
 $m = 356$ and density 2.5%

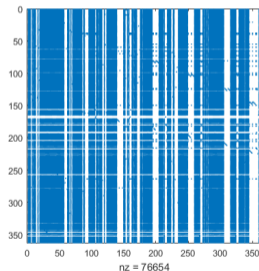


Inverse of a sparse matrix and solution of $Bx = r$

Optimal B for LP problem stair
 $m = 356$ and density 2.5%



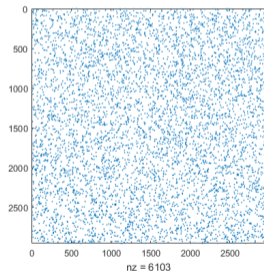
B^{-1} has density of 58%, so $B^{-1}r$ is typically dense



Inverse of a sparse matrix and solution of $Bx = r$

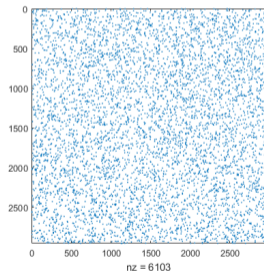
Random matrix

$m = 2953$ and density 0.07%

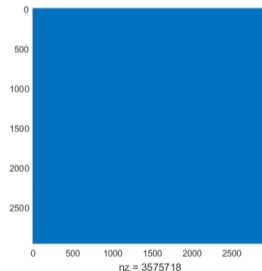


Inverse of a sparse matrix and solution of $Bx = r$

Random matrix
 $m = 2953$ and density 0.07%

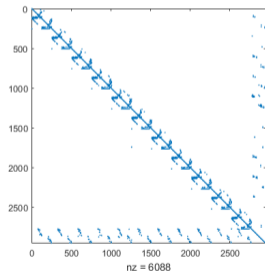


B^{-1} has density of 100%



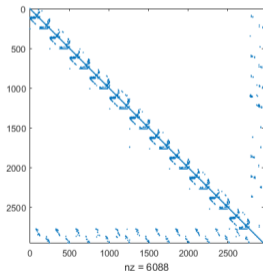
Inverse of a sparse matrix and solution of $Bx = r$

Optimal B for LP problem pds-02
 $m = 2953$ and density 0.07%

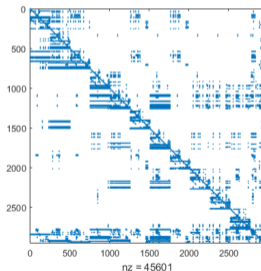


Inverse of a sparse matrix and solution of $Bx = r$

Optimal B for LP problem pds-02
 $m = 2953$ and density 0.07%



B^{-1} has density of 0.52%, so $B^{-1}r$ is typically **sparse**—when r is sparse



FTRAN: Form $\hat{\mathbf{a}}_q = B^{-1}\mathbf{a}_q$

- Triangular solves with $B = LU$
- Few columns of L and U are used to form $\hat{\mathbf{a}}_q$ Gilbert and Peierls (1988)

FTRAN: Form $\hat{\mathbf{a}}_q = B^{-1}\mathbf{a}_q$

- Triangular solves with $B = LU$
- Few columns of L and U are used to form $\hat{\mathbf{a}}_q$ Gilbert and Peierls (1988)

BTRAN: Form $\pi_p = B^{-T}\mathbf{e}_p$

- Transposed triangular solves
- Store L and U row-wise (also)
- Use FTRAN code

FTRAN: Form $\hat{\mathbf{a}}_q = B^{-1}\mathbf{a}_q$

- Triangular solves with $B = LU$
- Few columns of L and U are used to form $\hat{\mathbf{a}}_q$ Gilbert and Peierls (1988)

BTRAN: Form $\pi_p = B^{-T}\mathbf{e}_p$

- Transposed triangular solves
- Store L and U row-wise (also)
- Use FTRAN code

PRICE: Form $\hat{\mathbf{a}}_p^T = \pi_p^T N$

- Hyper-sparsity: π_p^T is sparse
- Store N row-wise
- Form $\hat{\mathbf{a}}_p^T$ as a combination of rows of N for nonzeros in π_p^T

H and McKinnon (1998–2005)
COAP best paper prize (2005)

Hyper-sparsity: PRICE

- PRICE forms $\hat{\mathbf{a}}_p^T = \boldsymbol{\pi}_p^T N$

Hyper-sparsity: PRICE

- PRICE forms $\hat{\mathbf{a}}_p^T = \boldsymbol{\pi}_p^T N$
- Simplifying notation, this is

$$\mathbf{y}^T = \boldsymbol{\pi}^T A$$

Hyper-sparsity: PRICE

- PRICE forms $\hat{\mathbf{a}}_p^T = \boldsymbol{\pi}_p^T N$
- Simplifying notation, this is

$$\mathbf{y}^T = \boldsymbol{\pi}^T A = \sum_{\pi_i \neq 0} \pi_i \times \mathbf{a}_i^T \quad \text{where} \quad A = \begin{bmatrix} \mathbf{a}_1^T \\ \mathbf{a}_2^T \\ \vdots \\ \mathbf{a}_m^T \end{bmatrix}$$

Hyper-sparsity: PRICE

- PRICE forms $\hat{\mathbf{a}}_p^T = \boldsymbol{\pi}_p^T N$

- Simplifying notation, this is

$$\mathbf{y}^T = \boldsymbol{\pi}^T A = \sum_{\pi_i \neq 0} \pi_i \times \mathbf{a}_i^T \quad \text{where} \quad A = \begin{bmatrix} \mathbf{a}_1^T \\ \mathbf{a}_2^T \\ \vdots \\ \mathbf{a}_m^T \end{bmatrix}$$

- In detail

- Start with zero full-length vector \mathbf{y}^T

Hyper-sparsity: PRICE

- PRICE forms $\widehat{\mathbf{a}}_p^T = \boldsymbol{\pi}_p^T N$

- Simplifying notation, this is

$$\mathbf{y}^T = \boldsymbol{\pi}^T A = \sum_{\pi_i \neq 0} \pi_i \times \mathbf{a}_i^T \quad \text{where} \quad A = \begin{bmatrix} \mathbf{a}_1^T \\ \mathbf{a}_2^T \\ \vdots \\ \mathbf{a}_m^T \end{bmatrix}$$

- In detail

- Start with zero full-length vector \mathbf{y}^T
- For each nonzero π_i

$$\mathbf{y}^T := \mathbf{y}^T + \pi_i \mathbf{a}_i^T$$

Hyper-sparsity: PRICE

- PRICE forms $\hat{\mathbf{a}}_p^T = \boldsymbol{\pi}_p^T N$

- Simplifying notation, this is

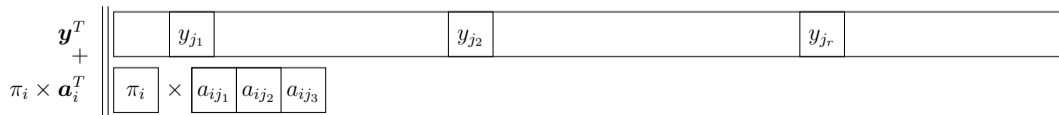
$$\mathbf{y}^T = \boldsymbol{\pi}^T A = \sum_{\pi_i \neq 0} \pi_i \times \mathbf{a}_i^T \quad \text{where} \quad A = \begin{bmatrix} \mathbf{a}_1^T \\ \mathbf{a}_2^T \\ \vdots \\ \mathbf{a}_m^T \end{bmatrix}$$

- In detail

- Start with zero full-length vector \mathbf{y}^T
- For each nonzero π_i

$$\mathbf{y}^T := \mathbf{y}^T + \pi_i \mathbf{a}_i^T$$

- Pass through each entry in (packed) sparse row \mathbf{a}_i^T



PRICE: What more can be done?

How sparse is π ?

- Problem-dependent!
- Extreme case
 - $m = 1,143,250$
 - Average π has 82 nonzeros
 - One nonzero for every 14,000 entries!

PRICE: What more can be done?

How sparse is π ?

- Problem-dependent!
- Extreme case
 - $m = 1,143,250$
 - Average π has 82 nonzeros
 - One nonzero for every 14,000 entries!

Cache inefficiency

- Read nonzeros of π - in cache anyway?
- Read row \mathbf{a}_i^T - contiguous
- Read entries in \mathbf{y}^T - in cache?

$$\mathbf{y}^T + \pi_i \times \mathbf{a}_i^T$$

The diagram illustrates the dot product of a sparse vector \mathbf{y}^T and a row \mathbf{a}_i^T . The vector \mathbf{y}^T is represented as a horizontal bar with segments: white, green (labeled y_{j_1}), red, white, green (labeled y_{j_2}), red, white, green (labeled y_{j_r}), red, white. Below it, the row \mathbf{a}_i^T is shown as a box containing π_i multiplied by three boxes containing a_{ij_1} , a_{ij_2} , and a_{ij_3} .

Aim: Avoid full-length double \mathbf{y}^T

- Pack nonzeros in double \mathbf{v}
- Store full-length array of pointers \mathbf{p} into \mathbf{v}

Aim: Avoid full-length double \mathbf{y}^T

- Pack nonzeros in double \mathbf{v}
- Store full-length array of pointers \mathbf{p} into \mathbf{v}
 - Up to 255 pointers using `unsigned char` - 8 times smaller
 - Up to 65535 pointers using `unsigned short` - 4 times smaller
 - Switch to standard hyper-sparse technique if \mathbf{y}^T fills in too much
- Hope that \mathbf{v} fits into cache

Aim: Avoid full-length double \mathbf{y}^T

- Pack nonzeros in double \mathbf{v}
- Store full-length array of pointers \mathbf{p} into \mathbf{v}
 - Up to 255 pointers using `unsigned char` - 8 times smaller
 - Up to 65535 pointers using `unsigned short` - 4 times smaller
 - Switch to standard hyper-sparse technique if \mathbf{y}^T fills in too much
- Hope that \mathbf{v} fits into cache

Does it work?

- Yes! (Poirrier: ISMP 2018)
- No! (Mészáros: FICO)

My experiments showed no performance improvement on amenable problems

Aim: Avoid any full-length \mathbf{y}^T !

- Store nonzeros and indices in C++ `map<int, double>`
- Use C++ utilities to find nonzeros
- Search cost is $\ln |\mathbf{y}|$

Aim: Avoid any full-length \mathbf{y}^T !

- Store nonzeros and indices in C++ `map<int, double>`
- Use C++ utilities to find nonzeros
- Search cost is $\ln |\mathbf{y}|$

Does it work?

My experiments showed no performance improvement on amenable problems

Why doesn't it work

Example problem

- $m = 1,143,250$; $n = 1,367,843$
- Average π has 82 nonzeros
- Average $\mathbf{y}^T = \pi^T A$ has 175 nonzeros

CPU cache structure

- L1d cache: 32K
- L1i cache: 32K
- L2 cache: 256K
- L3 cache: 8192K

Why doesn't it work

Example problem

- $m = 1,143,250$; $n = 1,367,843$
- Average π has 82 nonzeros
- Average $\mathbf{y}^T = \pi^T A$ has 175 nonzeros

CPU cache structure

- L1d cache: 32K
- L1i cache: 32K
- L2 cache: 256K
- L3 cache: 8192K

How might it work?

- Better implementation
- Better understanding of memory access, cache occupancy and CPU activity
- Force \mathbf{y} pointer array to remain in cache
- Ideas?

Ultra-sparsity

- Exists
- Leads to cache inefficiency
- Hard to exploit!

Ultra-sparsity

- Exists
- Leads to cache inefficiency
- Hard to exploit!

Slides:

<http://www.maths.ed.ac.uk/hall/SparseDays18>

Code:

<https://github.com/ERGO-Code/HiGHS>



J. A. J. Hall and K. I. M. McKinnon.

Hyper-sparsity in the revised simplex method and how to exploit it.

Computational Optimization and Applications,
32(3):259–283, December 2005.



Q. Huangfu and J. A. J. Hall.

Parallelizing the dual revised simplex method.

Mathematical Programming Computation, 10(1):119–142,
2018.