

## 9 Lecture 9, Feb 2

### Announcements

- HW1 graded. Feedback:
  - Sketch of solution: <http://hua-zhou.github.io/teaching/biostatm280-2016winter/hw01sol.html>. Please compare to your code carefully and understand why.
  - Q2: *Rounding to even rule* in R.
  - Q5: be familiar with R functions: `crossprod`, `tcrossprod`, `outer`, `row`, `col`, `rowSums`, `colSums`, ...
  - Q6: vectorize code (no looping necessary).
  - Code style. Be professional!
  - Frequent commits in git. I want to see how you developed the solutions.
- HW2 due today @ 11:59PM.
- HW3 posted. Due Feb 11 @ 11:59PM. [http://hua-zhou.github.io/teaching/biostatm280-2016winter/biostat\\_m280\\_2016\\_hw3.pdf](http://hua-zhou.github.io/teaching/biostatm280-2016winter/biostat_m280_2016_hw3.pdf)
- Quiz 2 this Thu Feb 4. In class, closed book.

### Last time

- Linear regression by Cholesky.  $np^2 + p^3/3$  flops for  $\hat{\beta}$  or more for s.e.
- Linear regression by QR  $\mathbf{X} = \mathbf{QR}$ .
- QR by MGS.  $2np^2$  flops to obtain  $\mathbf{Q}_1$  and  $\mathbf{R}$ .
- QR by Householder.  $2np^2 - 2p^3/3$  flops to overwrite  $\mathbf{X}$  by  $\mathbf{R}$  and Householder vectors. It is the algorithm R uses for linear regression.
- QR by Givens.  $3np^2 - p^3$  to overwrite  $\mathbf{X}$  by  $\mathbf{R}$  and Givens rotations.

## Today

- Sweep operator.
- Summary of numerical methods for linear regression.
- Summary of solving linear equations.
- Iterative methods for solving linear equations.

## Sweep operator

Assume  $\mathbf{A} \succeq \mathbf{0}_{n \times n}$ .

- The popular statistical software SAS uses sweep operator for linear regression and matrix inversion.
- Read KL 7.4-7.6 for introduction.

Also see “*A tutorial on the SWEEP operator*” by James H. Goodnight. <http://www.jstor.org/stable/2683825>

- *Sweep* on the  $k$ -th diagonal entry  $a_{kk} \neq 0$  yields  $\hat{\mathbf{A}}$  with entries

$$\begin{aligned}\hat{a}_{kk} &= -\frac{1}{a_{kk}} \\ \hat{a}_{ik} &= \frac{a_{ik}}{a_{kk}} \\ \hat{a}_{kj} &= \frac{a_{kj}}{a_{kk}} \\ \hat{a}_{ij} &= a_{ij} - \frac{a_{ik}a_{kj}}{a_{kk}}, \quad i \neq k, j \neq k.\end{aligned}$$

$n^2$  flops (taking into account of symmetry).

- *Inverse sweep* sends  $\mathbf{A}$  to  $\check{\mathbf{A}}$  with entries

$$\begin{aligned}\check{a}_{kk} &= -\frac{1}{a_{kk}} \\ \check{a}_{ik} &= -\frac{a_{ik}}{a_{kk}} \\ \check{a}_{kj} &= -\frac{a_{kj}}{a_{kk}} \\ \check{a}_{ij} &= a_{ij} - \frac{a_{ik}a_{kj}}{a_{kk}}, \quad i \neq k, j \neq k.\end{aligned}$$

$n^2$  flops (taking into account of symmetry).

- $\check{\mathbf{A}} = \mathbf{A}$ .
- Successively sweeping all diagonal entries of  $\mathbf{A}$  yields  $-\mathbf{A}^{-1}$ .
- Exercise: invert a  $2 \times 2$  matrix, say  $\mathbf{A} = \begin{pmatrix} 4 & 3 \\ 3 & 2 \end{pmatrix}$  on paper using sweep operator.
- Block form of sweep: Let the symmetric matrix  $\mathbf{A}$  be partitioned as  $\mathbf{A} = \begin{pmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{pmatrix}$ . If possible, sweep on the diagonal entries of  $\mathbf{A}_{11}$  yields

$$\hat{\mathbf{A}} = \begin{pmatrix} -\mathbf{A}_{11}^{-1} & \mathbf{A}_{11}^{-1}\mathbf{A}_{12} \\ \mathbf{A}_{21}\mathbf{A}_{11}^{-1} & \mathbf{A}_{22} - \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{A}_{12} \end{pmatrix}.$$

Order does *not* matter. The block  $\mathbf{A}_{22} - \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{A}_{12}$  is recognized as the *Schur complement* of  $\mathbf{A}_{11}$ .

- Pd and determinant:  $\mathbf{A}$  is pd if and only if each diagonal entry can be swept in succession and is positive until it is swept. When a diagonal entry of a pd matrix  $\mathbf{A}$  is swept, it becomes negative and remains negative thereafter. Taking the product of diagonal entries just before each is swept yields the determinant of  $\mathbf{A}$ .

- Linear regression by sweep. Sweep on  $\begin{pmatrix} \mathbf{X}^\top \mathbf{X} & \mathbf{X}^\top \mathbf{y} \\ \mathbf{y}^\top \mathbf{X} & \mathbf{y}^\top \mathbf{y} \end{pmatrix}$  yields

$$\begin{pmatrix} -(\mathbf{X}^\top \mathbf{X})^{-1} & (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \\ \mathbf{y}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} & \mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \end{pmatrix} = \begin{pmatrix} -\frac{1}{\sigma^2} \text{Var}(\hat{\boldsymbol{\beta}}) & \hat{\boldsymbol{\beta}} \\ \hat{\boldsymbol{\beta}}^\top & \|\mathbf{y} - \hat{\mathbf{y}}\|_2^2 \end{pmatrix}.$$

In total  $np^2 + p^3$  flops.

- Sweep is useful for stepwise regression, (conditional) multivariate normal density calculation, MANOVA, ...
- LAPACK does *not* implement sweep operator. It is statisticians' secret tool.
- Warning: the `sweep()` function in R has nothing to do with the sweep operator here.
- Demo code: <http://hua-zhou.github.io/teaching/biostatm280-2016winter/sweep.html>

## Summary of linear regression: Table on KL p105

Method	Flops	Remarks	Software	Stability
Sweep	$np^2 + p^3$	$(\mathbf{X}^\top \mathbf{X})^{-1}$ available	SAS	less stable
Cholesky	$np^2 + p^3/3$			less stable
QR by Householder	$2np^2 - 2p^3/3$		R	stable
QR by MGS	$2np^2$	$\mathbf{Q}_1$ available		most stable

Table 1: Numerical methods for linear regression. In order of stability.

Remarks:

- When  $n \gg p$ , sweep and Cholesky are twice faster than QR and need less space.
- Sweep and Cholesky is based on the Gram matrix  $\mathbf{X}^T \mathbf{X}$ , which can be dynamically updated with incoming data. They can easily handle huge  $n$ , moderate  $p$  data sets that cannot fit into memory.
- QR methods are more stable and produce numerically more accurate solution.
- Although sweep is slower than Cholesky, it yields standard errors and so on.
- Sweep is useful for stepwise regression, multivariate normal calculation, and numerous other statistical applications.
- MGS appears slower than Householder, but it yields  $\mathbf{Q}_1$ .

There is simply no such thing as a universal ‘gold standard’ when it comes to algorithms.

anonymous reviewer

## Summary of solving linear equations

Consider linear system  $\mathbf{A}\mathbf{x} = \mathbf{b}$ .

- We now know some good numerical methods for the least squares problem, which is essentially “solving” an overdetermined system (a tall  $\mathbf{A}$ ).

- Table 2 compares the flops of some methods (in order of stability) for solving a square (unstructured)  $\mathbf{A} \in \mathbb{R}^{n \times n}$ .

Method	Flops	Stability
Gaussian elimination	$2n^3/3$	less stable
QR by Householder	$4n^3/3$	
QR by MGS	$n^3$	
SVD	$6n^3$	most stable

Table 2: Flops of different numerical methods for  $n \times n$  square linear systems, assuming availability of the right hand side at time of decomposition.

- Solve an underdetermined system (a flat  $\mathbf{A} \in \mathbb{R}^{m \times n}$  of full row rank) by QR – version 1. First compute QR on  $\mathbf{A}^T$

$$\mathbf{A}^T = \mathbf{Q}\mathbf{R} = \mathbf{Q} \begin{pmatrix} \mathbf{R}_1 \\ \mathbf{0}_{(n-m) \times m} \end{pmatrix}.$$

Then  $\mathbf{A}\mathbf{x} = \mathbf{b}$  becomes

$$(\mathbf{Q}\mathbf{R})^T \mathbf{x} = \begin{pmatrix} \mathbf{R}_1^T & \mathbf{0}_{m \times (n-m)} \end{pmatrix} \begin{pmatrix} \mathbf{z}_1 \\ \mathbf{z}_2 \end{pmatrix} = \mathbf{b},$$

where

$$\mathbf{Q}^T \mathbf{x} = \begin{pmatrix} \mathbf{z}_1 \\ \mathbf{z}_2 \end{pmatrix}.$$

$\mathbf{z}_1$  is determined from  $\mathbf{R}_1^T \mathbf{z}_1 = \mathbf{b}$ . If we take  $\mathbf{z}_2 = \mathbf{0}_{n-m}$ , then we obtain the minimum norm solution (why?).

- Solve an underdetermined system (a flat  $\mathbf{A} \in \mathbb{R}^{m \times n}$  of full row rank) by QR – version 2. First compute QR with column pivoting on  $\mathbf{A}$

$$\mathbf{A}\mathbf{\Pi} = \mathbf{Q} \begin{pmatrix} \mathbf{R}_1 & \mathbf{R}_2 \end{pmatrix},$$

where  $\mathbf{R}_1 \in \mathbb{R}^{m \times m}$  is upper triangular and  $\mathbf{R}_2 \in \mathbb{R}^{m \times (n-m)}$ . Thus  $\mathbf{A}\mathbf{x} = \mathbf{b}$  transforms to

$$\begin{pmatrix} \mathbf{R}_1 & \mathbf{R}_2 \end{pmatrix} \begin{pmatrix} \mathbf{z}_1 \\ \mathbf{z}_2 \end{pmatrix} = \mathbf{Q}^T \mathbf{b},$$

where

$$\mathbf{\Pi}^T \mathbf{x} = \begin{pmatrix} z_1 \\ z_2 \end{pmatrix}.$$

One solution is obtained by  $\mathbf{z}_1 = \mathbf{R}_1^{-1} \mathbf{Q}^T \mathbf{b}$  and  $\mathbf{z}_2 = \mathbf{0}_{n-m}$ . It is not guaranteed to be of minimum norm.

## Condition number for linear equations (matrix inversion)

- Assume  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is nonsingular and consider the system of linear equation  $\mathbf{Ax} = \mathbf{b}$ . The solution is  $\mathbf{x} = \mathbf{A}^{-1} \mathbf{b}$ . We want to know how the solution changes with a small perturbation of the input  $\mathbf{b}$  (or  $\mathbf{A}$ ).

- Let  $\tilde{\mathbf{b}} = \mathbf{b} + \Delta \mathbf{b}$ . Then  $\tilde{\mathbf{x}} = \mathbf{A}^{-1}(\mathbf{b} + \Delta \mathbf{b}) = \mathbf{x} + \Delta \mathbf{x}$ . Thus

$$\|\Delta \mathbf{x}\| = \|\mathbf{A}^{-1} \Delta \mathbf{b}\| \leq \|\mathbf{A}^{-1}\| \|\Delta \mathbf{b}\|.$$

Because  $\mathbf{b} = \mathbf{Ax}$ ,  $\frac{1}{\|\mathbf{x}\|} \leq \|\mathbf{A}\| \frac{1}{\|\mathbf{b}\|}$ . This results

$$\frac{\|\Delta \mathbf{x}\|}{\|\mathbf{x}\|} \leq \|\mathbf{A}\| \|\mathbf{A}^{-1}\| \frac{\|\Delta \mathbf{b}\|}{\|\mathbf{b}\|}.$$

- $\kappa(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^{-1}\|$  is called the *condition number for inversion*. It depends on the matrix norm being used.  $\kappa_p$  means condition number defined from matrix- $p$  norm.
- Large condition number means “bad”.
- Some useful facts

$$\begin{aligned} \kappa(\mathbf{A}) &= \kappa(\mathbf{A}^{-1}) \\ \kappa(c\mathbf{A}) &= \kappa(\mathbf{A}) \\ \kappa(\mathbf{A}) &\geq 1 \\ \kappa_1(\mathbf{A}) &= \kappa_\infty(\mathbf{A}^\top) \\ \kappa_2(\mathbf{A}) &= \kappa_2(\mathbf{A}^\top) = \frac{\sigma_1(\mathbf{A})}{\sigma_n(\mathbf{A})} \\ \kappa_2(\mathbf{A}^\top \mathbf{A}) &= \frac{\lambda_1(\mathbf{A}^\top \mathbf{A})}{\lambda_n(\mathbf{A}^\top \mathbf{A})} = \kappa_2^2(\mathbf{A}) \geq \kappa_2(\mathbf{A}). \end{aligned}$$

The last fact says that the condition number of  $\mathbf{A}^\top \mathbf{A}$  can be much larger than that of  $\mathbf{A}$ .

- The smallest singular value  $\sigma_n$  indicates the “distance to the trouble”.
- Condition number for the least squares problem is more complicated. Roughly speaking, the method based on normal equation (Cholesky, sweep) has a condition depending on  $\kappa_2(\mathbf{X})^2$ . QR for a “small residuals” least squares problem has a condition depending on  $\kappa_2(\mathbf{X})$ .
- Numerically, consider the simple case

$$\mathbf{X} = \begin{pmatrix} 1 & 1 \\ 10^{-3} & 0 \\ 0 & 10^{-3} \end{pmatrix}.$$

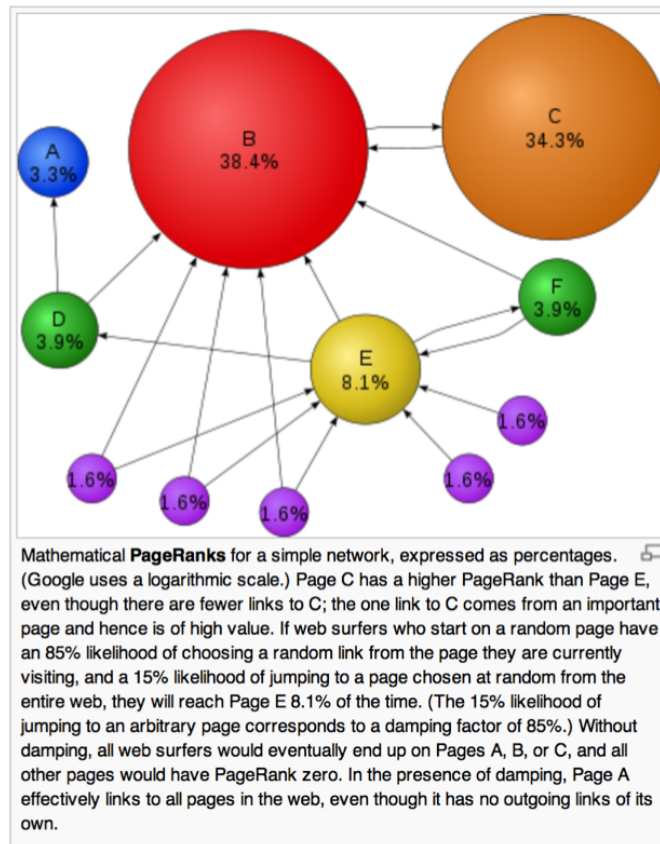
Forming normal equation yields a singular Gramian matrix

$$\mathbf{X}^\top \mathbf{X} = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}$$

if executed with a precision of 6 digits.

- In R, the `kappa()` function (wrapper of `DTRCON` in LAPACK and `DTRCO` in LINPACK) computes or approximates (default) the condition number of a matrix.
- In regression problems, standardizing the predictors could improve the condition. See demo on the Longley data <http://hua-zhou.github.io/teaching/biostatm280-2016winter/longleycond.html>.
- In design of experiments (DoE), people favor orthogonal design. Why?

## Iterative method for solving linear equations: introduction



- Direct method (flops fixed *a priori*) vs iterative methods:
  - Direct method (GE/LU, Cholesky, QR, SVD): good for dense, small or moderate sized, unstructured  $\mathbf{A}$
  - Iterative methods (Jacobi, Gauss-Seidel, SOR, conjugate-gradient, GMRES): good for large, sparse, or structured linear system, parallel computing, warm start
- PageRank (HW3):
  - $\mathbf{A} \in \{0, 1\}^{n \times n}$  the connectivity matrix with entries

$$a_{ij} = \begin{cases} 1 & \text{if page } i \text{ links to page } j \\ 0 & \text{otherwise} \end{cases}.$$

$n \approx 10^9$  in Feb 2016.

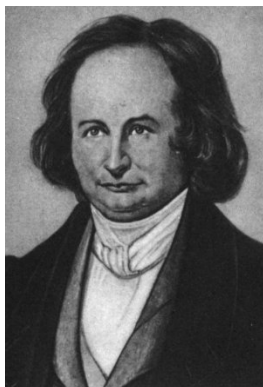


- $r_i = \sum_j a_{ij}$  is the out-degree of page  $i$ .
- Larry Page: Imagine a random surfer wandering on internet according to following rules:
  - \* From a page  $i$  with  $r_i > 0$ 
    - with probability  $p$ , (s)he randomly chooses a link on page  $i$  (uniformly) and follows that link to the next page
    - with probability  $1 - p$ , (s)he randomly chooses one page from the set of all  $n$  pages (uniformly) and proceeds to that page
  - \* From a page  $i$  with  $r_i = 0$  (a dangling page), (s)he randomly chooses one page from the set of all  $n$  pages (uniformly) and proceeds to that page

The process defines a Markov chain on the space of  $n$  pages. Stationary distribution of this Markov chain gives the ranks (probabilities) of each page.

- Stationary distribution is the top left eigenvector of the transition matrix  $\mathbf{P}$  corresponding to eigenvalue 1. Equivalently it can be cast as a linear equation.
- Largest matrix computation in world (?).
- GE/LU will take  $2 \times (10^9)^3 / 3 / 10^{12} \approx 6.66 \times 10^{14}$  seconds  $\approx 2 \times 10^7$  years on a tera-flop supercomputer!
- Iterative methods come to the rescue.

## Jacobi method



Solve linear system  $\mathbf{A}\mathbf{x} = \mathbf{b}$ .

- Jacobi iteration:

$$x_i^{(t+1)} = \frac{b_i - \sum_{j=1}^{i-1} a_{ij}x_j^{(t)} - \sum_{j=i+1}^n a_{ij}x_j^{(t)}}{a_{ii}}.$$

- Requires non-zero diagonal element!
- One round costs  $2n^2$  flops with an unstructured  $\mathbf{A}$ . Gain over GE/LU if converges in  $o(n)$  iterations. Saving is huge for *sparse* or *structured*  $\mathbf{A}$ . By structured, we mean the matrix-vector multiplication  $\mathbf{A}\mathbf{v}$  is fast.
- Splitting:  $\mathbf{A} = \mathbf{L} + \mathbf{D} + \mathbf{U}$ .
- Jacobi:  $\mathbf{D}\mathbf{x}^{(t+1)} = -(\mathbf{L} + \mathbf{U})\mathbf{x}^{(t)} + \mathbf{b}$ , i.e.,

$$\mathbf{x}^{(t+1)} = -\mathbf{D}^{-1}(\mathbf{L} + \mathbf{U})\mathbf{x}^{(t)} + \mathbf{D}^{-1}\mathbf{b}.$$

## Gauss-Seidel



- Gauss-Seidel iteration:

$$x_i^{(t+1)} = \frac{b_i - \sum_{j=1}^{i-1} a_{ij}x_j^{(t+1)} - \sum_{j=i+1}^n a_{ij}x_j^{(t)}}{a_{ii}}.$$

- With splitting,  $(\mathbf{D} + \mathbf{L})\mathbf{x}^{(t+1)} = -\mathbf{U}\mathbf{x}^{(t)} + \mathbf{b}$ , i.e.,

$$\mathbf{x}^{(t+1)} = -(\mathbf{D} + \mathbf{L})^{-1}\mathbf{U}\mathbf{x}^{(t)} + (\mathbf{D} + \mathbf{L})^{-1}\mathbf{b}.$$

- GS converges for any  $\mathbf{x}^{(0)}$  for symmetric and pd  $\mathbf{A}$ .
- Convergence rate of Gauss-Seidel is the spectral radius of the  $(\mathbf{D} + \mathbf{L})^{-1}\mathbf{U}$ .
- Comparing Jacobi and GS, Jacobi is particularly attractive for parallel computing.

## Successive over-relaxation (SOR)

- SOR:  $x_i^{(t+1)} = \omega(b_i - \sum_{j=1}^{i-1} a_{ij}x_j^{(t+1)} - \sum_{j=i+1}^n a_{ij}x_j^{(t)})/a_{ii} + (1 - \omega)x_i^{(t)}$ , i.e.,

$$\mathbf{x}^{(t+1)} = (\mathbf{D} + \omega\mathbf{L})^{-1}[(1 - \omega)\mathbf{D} - \omega\mathbf{U}]\mathbf{x}^{(t)} + (\mathbf{D} + \omega\mathbf{L})^{-1}(\mathbf{D} + \mathbf{L})^{-1}\omega\mathbf{b}.$$

- Need to pick  $\omega \in [0, 1]$  beforehand, with the goal of improving convergence rate.

## Conjugate gradient method

Solving  $\mathbf{Ax} = \mathbf{b}$  is equivalent to minimizing the quadratic function  $\frac{1}{2}\mathbf{x}^T\mathbf{Ax} - \mathbf{b}^T\mathbf{x}$ . To do later, when we study optimization. Conjugate gradient and its variants are the top-notch iterative methods for solving huge, structured linear systems.

**Table 1. Kershaw's results for a fusion problem.**

Method	Number of iterations
Gauss Seidel	208,000
Block successive overrelaxation methods	765
Incomplete Cholesky conjugate gradients	25

## A list of “easy” linear systems

Consider  $\mathbf{Ax} = \mathbf{b}$ ,  $\mathbf{A} \in \mathbb{R}^{n \times n}$ . Or, consider matrix inverse (if you want).  $\mathbf{A}$  can be huge. Keep massive data in mind: 1000 Genome Project, Netflix, Google PageRank, finance, spatial statistics, ... We should be alert to many easy linear systems. *Don't waste computing resources by bad choices of algorithms!*

- Diagonal:  $n$  flops.

- Tridiagonal/banded: Band LU, band Cholesky, ... roughly  $O(n)$  flops
- Triangular:  $n^2$  flops
- Block diagonal: Suppose  $n = \sum_i n_i$ .  $(\sum_i n_i)^3$  vs  $\sum_i n_i^3$ .
- Kronecker product:  $(\mathbf{A} \otimes \mathbf{B})^{-1} = \mathbf{A}^{-1} \otimes \mathbf{B}^{-1}$ ,  $(\mathbf{C}^\top \otimes \mathbf{A})\text{vec}\mathbf{B} = \text{vec}(\mathbf{ABC})$  fits iterative method.
- Sparsity: iterative method, or sparse matrix decomposition.  
Remark: Probably the easiest recognizable structure. Familiarize yourself with the sparse matrix computation tools in JULIA, MATLAB, R (**Matrix** package), MKL (sparse BLAS), ... as much as possible.
- Easy plus low rank:  $\mathbf{U} \in \mathbb{R}^{n \times m}$ ,  $\mathbf{V} \in \mathbb{R}^{n \times m}$ ,  $m \ll n$ . Woodbury formula

$$(\mathbf{A} + \mathbf{UV}^\top)^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{U}(\mathbf{I}_m + \mathbf{V}^\top\mathbf{A}^{-1}\mathbf{U})^{-1}\mathbf{V}^\top\mathbf{A}^{-1}.$$

Keep HW2 Q5 in mind.

- Easy plus border: For  $\mathbf{A}$  pd and  $\mathbf{V}$  full row rank,

$$\begin{pmatrix} \mathbf{A} & \mathbf{V}^\top \\ \mathbf{V} & \mathbf{0} \end{pmatrix}^{-1} = \begin{pmatrix} \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{V}^\top(\mathbf{VA}^{-1}\mathbf{V}^\top)^{-1}\mathbf{VA}^{-1} & \mathbf{A}^{-1}\mathbf{V}^\top(\mathbf{VA}^{-1}\mathbf{V}^\top)^{-1} \\ (\mathbf{VA}^{-1}\mathbf{V}^\top)^{-1}\mathbf{VA}^{-1} & -(\mathbf{VA}^{-1}\mathbf{V}^\top)^{-1} \end{pmatrix}.$$

- Orthogonal:  $n^2$  flops *at most*. Permutation matrix, Householder matrix, Jacobi matrix, ... take less.
- Toeplitz systems:

$$\mathbf{T} = \begin{pmatrix} r_0 & r_1 & r_2 & r_3 \\ r_{-1} & r_0 & r_1 & r_2 \\ r_{-2} & r_{-1} & r_0 & r_1 \\ r_{-3} & r_{-2} & r_{-1} & r_0 \end{pmatrix}.$$

$\mathbf{T}\mathbf{x} = \mathbf{b}$ , where  $\mathbf{T}$  is pd and Toeplitz, can be solved in  $O(n^2)$  flops. Durbin algorithm (Yule-Walker equation), Levinson algorithm (general  $\mathbf{b}$ ), Trench algorithm (inverse). These matrices occur in auto-regressive models and econometrics.

- Circulant systems: Toeplitz matrix with wraparound

$$C(\mathbf{z}) = \begin{pmatrix} z_0 & z_4 & z_3 & z_2 & z_1 \\ z_1 & z_0 & z_4 & z_3 & z_2 \\ z_2 & z_1 & z_0 & z_4 & z_3 \\ z_3 & z_2 & z_1 & z_0 & z_4 \\ z_4 & z_3 & z_2 & z_1 & z_0 \end{pmatrix},$$

FFT type algorithms: DCT (discrete cosine transform) and DST (discrete sine transform).

- Vandermonde matrix: such as in interpolation and approximation problems

$$\mathbf{V}(x_0, \dots, x_n) = \begin{pmatrix} 1 & 1 & \cdots & 1 \\ x_0 & x_1 & \cdots & x_n \\ \vdots & \vdots & & \vdots \\ x_0^n & x_1^n & \cdots & x_n^n \end{pmatrix}.$$

$\mathbf{V}\mathbf{x} = \mathbf{b}$  or  $\mathbf{V}^T\mathbf{x} = \mathbf{b}$  can be solved in  $O(n^2)$  flops.

- Cauchy-like matrices:

$$\mathbf{\Omega}\mathbf{A} - \mathbf{A}\mathbf{\Lambda} = \mathbf{R}\mathbf{S}^T,$$

where  $\mathbf{\Omega} = \text{diag}(\omega_1, \dots, \omega_n)$  and  $\mathbf{\Lambda} = (\lambda_1, \dots, \lambda_n)$ .  $O(n)$  flops for LU and QR.

- Structured-rank problems: semiseparable matrices (LU and QR takes  $O(n)$  flops), quasiseparable matrices, ...
- Fast multiple method (FMM) for kernel matrix.
- ...

Other computations such as matrix-vector multiplication with these “easy” matrices are typically fast too.

Bottom line: Don’t blindly use `solve()`.