

MATH 270B Notes

Ryan Anderson

2024-01-18

Lecture 1/17: QR Decomposition

Want to solve basic matrix equations $Ax = b$. When A is square and invertible, this is trivial - $x = A^{-1}b$. Otherwise, need to solve $x : \min_x \|Ax - b\|_2$.

By way of motivation note that triangular matrices and orthogonal matrices are non square but still very easy to solve. For orthogonal matrices, the inverse is the adjoint, so $Qx = b \rightarrow x = Q^*b$.

The idea then is to write $A = QR$ where Q is orthogonal and R is upper triangular. Then $Ax = b \rightarrow QRx = b \rightarrow Rx = Q^*b$.

To actually obtain this factorization we began as follows: we seek orthonormal vectors q_1, q_2, \dots, q_n such that $\text{span}(q_1, q_2, \dots, q_k) = \text{span}(a_1, a_2, \dots, a_k)$. This is done with the classical Gram-Schmidt algorithm:

$$\begin{aligned} v_j &= a_j \\ \text{for } i &= 1, 2, \dots, j-1 \\ r_{ij} &= q_i^* a_j \\ v_j &= v_j - r_{ij} q_i \\ r_{jj} &= \|v_j\|_2 \\ q_j &= v_j / r_{jj} \end{aligned}$$

Problem is that we can't use classical Gram-Schmidt due to numerical instability. This means that floating point precision problems could ruin the QR factorization. The major change in modified Gram-Schmidt is that we normalize then project away, rather than projecting away and then normalizing.

We want to show existence and uniqueness of the QR decomposition for general matrices.

Theorem: Every $A \in \mathbb{C}^{m \times n}$ has a full (also reduced) QR factorization.

Proof: By the Gram-Schmidt process, we can generate a set of orthogonal vectors q_1, q_2, \dots, q_k with $\text{span}(q_1, q_2, \dots, q_k) = \text{span}(a_1, a_2, \dots, a_k)$. We can also calculate the corresponding entries of R via the G-S process.

Only problem occurs if $r_{jj} = 0$ for some j . In this case, we can just drop the j th column of Q and R and continue the process.

If $A \in \mathbb{C}^{m \times n}$ is full-rank, then we can get a unique reduced QR decomposition.

Example:

```
A <- matrix(c(1,1,0,1,0,1,0,1,1),nrow=3,byrow=T)
A
```

```
##      [,1] [,2] [,3]
## [1,]    1    1    0
## [2,]    1    0    1
```

```
## [3,]    0    1    1
u_1 <- A[,1]
e_1 <- u_1/norm(u_1,type="2")

u_2 <- A[,2] - (A[,2] %*% e_1)*e_1
e_2 <- u_2/norm(u_2,type="2")

u_3 <- A[,3] - (A[,3] %*% e_1)*e_1 - (A[,3] %*% e_2)*e_2
e_3 <- u_3/norm(u_3,type="2")

Q <- matrix(cbind(e_1,e_2,e_3),nrow=3,byrow=F)
R <- matrix(c(A[,1] %*% e_1,A[,2] %*% e_1, A[,3] %*% e_1,0,A[,2] %*% e_2, A[,3] %*% e_2, 0, 0, A[,3] %*% e_3),nrow=3,byrow=F)
A - Q %*% R
```

```
##           [,1]           [,2]           [,3]
## [1,] 2.220446e-16  1.110223e-16 -2.966377e-16
## [2,] 2.220446e-16 -6.816043e-17  1.110223e-16
## [3,] 0.000000e+00  2.220446e-16  1.110223e-16
```

Other uses of the QR inverse include calculating the pseudo-inverse. Recall the Moore-Penrose pseudo-inverse is given as $A^\dagger = (A^T A)^{-1} A^T$.

With the QR decomposition we can write

$$\begin{aligned} A^\dagger &= (QR)^\dagger = ((QR)^T (QR))^{-1} (QR)^T \\ &= (R^T Q^T Q R)^{-1} R^T Q^T \\ &= (R^T R)^{-1} R^T Q^T \\ &= R^{-1} R^{-T} R^T Q^T \\ &= R^{-1} Q^T \end{aligned}$$

If A is square and nonsingular then $A^{-1} = R^{-1} Q^T$.

We can get the absolute value of the determinant of a square matrix out via QR factorization. $\det(A) = \det(Q) * \det(R)$ but Q is unitary so $|\det(Q)| = 1$. Thus $|\det(A)| = |\det(R)| = |\prod_i r_{ii}|$.

We can do experiments to see the use of QR. The pseudoinverse implementation is pretty quick.

```
A <- matrix(rnorm(3000000),nrow=3000,ncol=1000)
x <- rnorm(1000)
b <- A %*% x

start <- Sys.time()
A_pseudo_inv <- solve(t(A) %*% A) %*% t(A)
x_pseudo <- A_pseudo_inv %*% b
end <- Sys.time()
end-start
```

```
## Time difference of 2.498555 secs
```

```
norm(x - x_pseudo,type="2")
```

```
## [1] 1.591807e-13
```

The QR implementation is less quick.

```

start <- Sys.time()
A_qr_factor <- qr(A)
A_Q <- qr.Q(A_qr_factor)
A_R <- qr.R(A_qr_factor)
y <- t(A_Q) %*% b
x_QR <- solve(A_R) %*% y
end <- Sys.time()
end-start

```

```
## Time difference of 5.730671 secs
```

```
norm(x - x_QR,type="2")
```

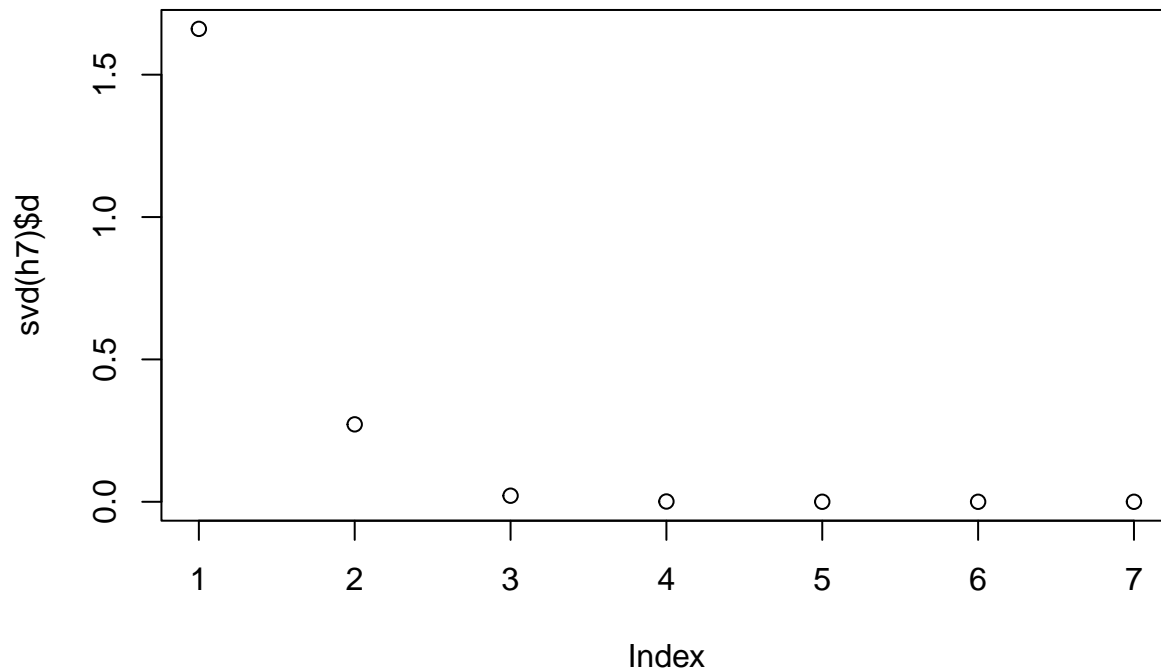
```
## [1] 1.181657e-13
```

If we try to generate a Hilbert matrix we run into problems with precision. Hilbert matrices are cool because they are full-rank but have extremely large condition numbers.

```

hilbert <- function(n) { i <- 1:n; 1 / outer(i - 1, i, "+") }
h7 <- hilbert(7)
plot(svd(h7)$d)

```



Lecture 1/22: Givens Rotations

Givens Rotations

Consider a matrix $G(i, j, \theta)$ given by

$$G(i, j, \theta) = \begin{bmatrix} 1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & 1 \end{bmatrix}$$

As an example consider the expression

$$\begin{bmatrix} \cos(\theta) & \sin(\theta) \\ -\sin(\theta) & \cos(\theta) \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} \sqrt{x^2 + y^2} \\ 0 \end{bmatrix}$$

This indicates that we can perform Givens rotations to zero out the lower triangular portion of a matrix. This is useful for QR factorization.

In particular, we want to perform sequential Givens rotations G_1, G_2, \dots of our initial matrix to yield an upper triangular R and then take the product of their adjoints, which will yield our Q matrix.

In the below example, we want to zero out the 2, 1 entry of the matrix, the 5.

$$\begin{bmatrix} 6 & 5 & 0 \\ 5 & 1 & 4 \\ 0 & 4 & 3 \end{bmatrix}$$

We start by calculating G_1 , with $r = \sqrt{6^2 + 5^2} = 7.81, c = 6/7.81 = 0.768, s = -5/7.81 = -0.64$.

$$G_1 = \begin{bmatrix} c & -s & 0 \\ s & c & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

This simplifies to

```
orig_mat <- matrix(c(6,5,0,5,1,4,0,4,3),nrow=3,byrow=T)
G_1 <- matrix(c(0.768,0.64,0,-0.64,0.768,0,0,0,1),nrow=3,byrow=T)
G_1
```

```
##      [,1] [,2] [,3]
## [1,] 0.768 0.640  0
## [2,] -0.640 0.768  0
## [3,] 0.000 0.000  1
```

```
G_1 %%% orig_mat
```

```
##      [,1] [,2] [,3]
## [1,] 7.808000e+00 4.480 2.560
## [2,] 2.220446e-16 -2.432 3.072
## [3,] 0.000000e+00 4.000 3.000
```

We then calculate G_2 .

$$G_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & c & -s \\ 0 & s & c \end{bmatrix}, r = \sqrt{-2.43^2 + 4^2} = 4.68, c = -2.43/r = -0.52, s = -4/r = -0.85$$

Then we get the final matrix.

```
G_2 <- matrix(c(1,0,0,0,-0.52,0.85,0,-0.85,-0.52),nrow=3,byrow=T)
G_2
```

```
##      [,1] [,2] [,3]
## [1,]  1  0.00  0.00
## [2,]  0 -0.52  0.85
## [3,]  0 -0.85 -0.52
```

```
G_2 %%% (G_1 %%% orig_mat)
```

```
##      [,1] [,2] [,3]
## [1,] 7.808000e+00 4.48000 2.56000
## [2,] -1.154632e-16 4.66464 0.95256
## [3,] -1.887379e-16 -0.01280 -4.17120
```

Then we calculate the QR decomposition as

```
Q <- t(G_1) %*% t(G_2)
Q

##      [,1]      [,2]      [,3]
## [1,] 0.768  0.33280  0.5440
## [2,] 0.640 -0.39936 -0.6528
## [3,] 0.000  0.85000 -0.5200

R <- G_2 %*% (G_1 %*% orig_mat)
R

##      [,1]      [,2]      [,3]
## [1,] 7.808000e+00  4.48000  2.56000
## [2,] -1.154632e-16  4.66464  0.95256
## [3,] -1.887379e-16 -0.01280 -4.17120

round(Q %*% R,1)

##      [,1] [,2] [,3]
## [1,]    6    5    0
## [2,]    5    1    4
## [3,]    0    4    3
```

And we can see that $QR = A$ as desired.

Householder Reflections

Householder reflections are another way to perform QR decomposition. A *Householder reflection* is a reflection about a plane that contains the origin.

We can project a vector x onto a hyperplane H by calculating $P_x = (I - \frac{vv^*}{\|v\|^2})x$ where v is a unit vector normal to the hyperplane H .

Since that gets you onto the hyperplane, reflecting across it is given by simply doubling the distance: $F_x = (I - 2\frac{vv^*}{\|v\|^2})x$.

It turns out that reflecting over a particular hyperplane is equivalent to obtaining a vector in the direction of the unit basis vector with magnitude equal to the norm of the vector you want to reflect over. That is, $Fx = \|x\|e_1$. This lets us calculate v as $v = \|x\|e_1 - x$.

Note there are two possible such projections, one that results in a vector $Fx = \|x\|e_1$, and one that gives $F'x = -\|x\|e_1$. For stability reasons we want to choose the reflection that moves x the furthest.

This leads to our algorithm: we sequentially perform Householder reflections to zero out the lower triangular portion of our matrix, moving column-by-column. We can then take the product of the adjoints of the Householder reflections to get our Q matrix. Let $A \in R^{m \times n}$. Then we have

$$\begin{aligned} \text{for } k=1, \dots, n \\ x &= A_{k:m,k} \\ v_k &= \text{sign}(x)\|x\|e_1 + x \\ v_k &= \frac{v_k}{\|v_k\|} \\ A_{k:m,k:n} &= A_{k:m,k:n} - 2v_k(v_k^*A_{k:m,k:n}) \end{aligned}$$

Our R then becomes the final A matrix after the algorithm completes. We don't explicitly compute Q while performing this, but can obtain it. Let Q_1, Q_2, \dots be the reflection operators we generate at each iteration. Then $R = Q_n \dots Q_1 A$, which means $Q_n \dots Q_1 = Q^*$.

Lecture 1/24: Least Squares

Consider the least squares problem $Ax = b$, $A \in \mathbb{C}^{m \times n}$, $m \geq n$ where A is a tall matrix. Clearly there is only an exact solution if $b \in \text{Rng}(A)$, which is uncommon in practice. Instead we solve by taking

$$x_{LS} = \operatorname{argmin}_{x \in \mathbb{C}^n} \|Ax - b\|_2$$

We define the *residual* $r = Ax_{LS} - b$ to be the error using the least squares solution. The geometry of the situation tells us that Ax_{LS} is given by the orthogonal projection of b into $\text{Rng}(A)$.

Theorem: Least Squares Solution Let $A \in \mathbb{C}^{m \times n}$, $m \geq n$, $b \in \mathbb{C}^m$. A vector $x \in \mathbb{C}^n$ minimizes $\|Ax - b\|_2^2 = \|r\|_2^2$ iff $r \perp \text{Rng}(A)$, that is $A^*r = 0$ or $A^*Ax = A^*b$ or $\operatorname{proj}_{\text{Rng}(A)}(b) = Ax$.

The system defined by $A^*Ax = A^*b$ is called the *normal equations*, which are nonsingular iff A is full rank.

Proof For the equivalence of $A^*r = 0$ and $A^*Ax = A^*b$, start by noting $r = b - Ax$. That means $A^*r = 0$ iff $A^*b = A^*Ax$.

For the equivalence of $A^*Ax = A^*b$ and $\operatorname{proj}_{\text{Rng}(A)}(b) = Ax$, we start by noting that $\operatorname{proj}_{\text{Rng}(A)}(b) - Ax = 0$ implies $AA^\dagger b - Ax = 0$. Multiplying through by A^* we get

$$\begin{aligned} A^*AA^\dagger b - A^*Ax &= 0 \\ \Rightarrow A^*b - A^*Ax &= 0 \\ \Rightarrow A^*r &= 0 \end{aligned}$$

We can show that all solutions to the least squares problem are of the form $A^\dagger b + \ker(A)$.

Implementing Least Squares via QR Factorization

For implementation's sake, let A be full-rank. Consider the normal equations $A^*Ax = A^*b$.

The Gram matrix A^*A is square, nonsingular, and positive definite. That means we can take $(A^*A)^{-1}$, which means the solution to the normal equations can be obtained simply by taking $x = (A^*A)^{-1}A^*b = A^\dagger b$. That is, the solution to the normal equations is the pseudoinverse of A when A has full column rank.

To get the pseudoinverse we will proceed by QR factorization.

$$\begin{aligned} A = QR &\Rightarrow A^\dagger = [(QR)^*(QR)]^{-1}(QR)^* \\ &= [R^*Q^*QR]^{-1}R^*Q^* \\ &= [R^*R]^{-1}R^*Q^* \\ &= R^{-1}(R^*)^{-1}R^*Q^* = R^{-1}Q^*. \end{aligned}$$

Thus we can restate the solution to the normal equations as $x = A^\dagger b \Rightarrow x = R^{-1}Q^*b \Rightarrow Rx = Q^*b$, which is now an upper triangular system.

This gives the least squares algorithm via QR factorization:

- 1) Compute $A = QR$.
- 2) Compute Q^*b .
- 3) Solve the upper triangular system $Rx = Q^*b$.

Note this computation is dominated by the first step, taking the factorization. With Householder reflections this takes about mn steps.

Implementing Least Squares via SVD

We can also obtain an algorithm via SVD.

$$\begin{aligned} A = U\Sigma V^* \Rightarrow A^\dagger &= [(U\Sigma V^*)^* U\Sigma V^*]^{-1} (U\Sigma V^*)^* \\ &= [V\Sigma^* U^* U\Sigma V^*]^{-1} V\Sigma^* U^* \\ &= [V\Sigma^2 V^*]^{-1} V\Sigma^* U^* \\ &= V\Sigma^{-1} U^* \end{aligned}$$

This yields $x = A^\dagger b \Rightarrow x = V\Sigma^{-1}U^*b \Rightarrow \Sigma V^*x = U^*b$.

Then our algorithm for least squares via SVD is:

- 1) Compute (reduced) SVD
- 2) Compute U^*b
- 3) Solve $\Sigma w = U^*b$ for w
- 4) Set $x = Vw$

Step 1) dominates the computational runtime here, giving overall $2mn^2 + n^3$. Note for $m \gg n$, SVD and QR are similarly time expensive. When $m \simeq n$, SVD is way more expensive. However, SVD does give a diagonal system, which is trivial, compared to the upper triangular system from QR.

Cholesky factorization is better than both!

Experimenting with Factorizations

We will now experiment with the different factorizations.

```
M = 10000
N = 1000
A = matrix(rnorm(M*N),nrow=M,ncol=N)
x = rnorm(N)
b = A %*% x

start <- Sys.time()
A_qr_factor <- qr(A)
A_Q <- qr.Q(A_qr_factor)
A_R <- qr.R(A_qr_factor)
y <- t(A_Q) %*% b
x_QR <- solve(A_R) %*% y
end <- Sys.time()
end-start
norm(x - x_QR,type="2")

start <- Sys.time()
A_SVD <- svd(A)
u_b <- t(A_SVD$u) %*% b
w <- solve(diag(A_SVD$d)) %*% u_b
x_SVD <- A_SVD$v %*% w
end <- Sys.time()
end-start
norm(x - x_SVD,type="2")
```

Lecture 1/29: Conditioning

Consider a problem $f : X \rightarrow Y$ where we have X data and Y solution. We want to describe *well-conditioned problems*, where small changes in the dataset lead to small changes in the solution.

Let δX denote a small change in the dataset, and let $\delta f = f(X + \delta X) - f(X)$. Then the *absolute condition number* of f at X is given by

$$\hat{\kappa} = \lim_{\delta \rightarrow 0} \sup_{\|\delta X\| \leq \delta} \frac{\|\delta f\|}{\|\delta X\|} = \sup_{\|\delta X\| \leq \delta} \frac{\|\delta f\|}{\|\delta X\|}$$

We also define the *relative condition number*:

$$\kappa = \lim_{\delta \rightarrow 0} \sup_{\|\delta X\| \leq \delta} \frac{\|\delta f\|/\|f\|}{\|\delta X\|/\|X\|} = \sup_{\|\delta X\| \leq \delta} \frac{\|\delta f\|/\|f\|}{\|\delta X\|/\|X\|}$$

Note the similarities to the definition of the derivative. Indeed we can formally describe the relationship between the absolute condition number and the Jacobian of f at X as

$$\delta f \simeq J(X)\delta X \Rightarrow \hat{\kappa} = \|J(X)\|$$

Ex. 1: Computing Eigenvalues Consider two matrices as below

$$\begin{bmatrix} 1 & 1000 \\ 0 & 1 \end{bmatrix}, \begin{bmatrix} 1 & 1000 \\ .001 & 1 \end{bmatrix}$$

The first matrix has eigenvalues 1 and 1, while the second has eigenvalues 0 and 2. This points to the fact that in general, computing eigenvalues is a very poorly conditioned problem.

Ex. 2: Square Root Problems Consider $f : X \rightarrow X$ where $f(x) = \sqrt{X}$. We compute the condition numbers as below by actually taking the Jacobian:

$$\kappa = \frac{\|J(X)\|}{\|f(X)\|/\|X\|} = \frac{\frac{1}{2\sqrt{X}}}{\frac{\sqrt{X}}{|X|}} = \frac{1}{2}$$

Ex. 3: Solving for Polynomial Roots Consider the problem of solving for roots of general polynomials, $p(x) = \sum_i^n c_i x^i$. We have the following theorem:

Let r be a root of $p(x)$, and let $p'(x) = \sum_i^{n-1} (i+1)c_i x^i$. Then

Proof: let $r(c_j) = r, r(c_j + \delta c_j) = \hat{r}$. Then

$$\kappa = \limsup_{\delta \rightarrow 0, \delta c_j \leq \delta} \frac{|\delta c_j|}{|c_j|}$$

A good example of extremely poorly conditioned polynomials for root solving is *Wilkinson's polynomial*, $w(x) = (x-1)(x-2)\dots(x-20)$. It turns out that if the coefficient on the x^{19} term is perturbed by $\simeq 2^{-23}$, the root will change by $\simeq 1$.

Ex. 4: Matrix-vector Multiplication We calculate κ for matrix-vector multiplication as

$$\kappa = \sup_{\delta x} \frac{\|A(x + \delta x) - Ax\|/\|Ax\|}{\|\delta x\|/\|x\|} = \sup_{\delta x} \frac{\|A\delta x\|/\|Ax\|}{\|\delta x\|/\|x\|} = \frac{\|A\|}{\|Ax\|/\|x\|} \leq \|A\|\|A^\dagger\|$$

The final quantity $\|A\|\|A^\dagger\|$ is actually the definition of the condition number of A .

Accuracy and Stability

Consider a problem of computing $f(X)$ from X , with a result in floating point $\tilde{f}(X)$. An algorithm is *accurate* if

$$\frac{\|\tilde{f}(X) - f(X)\|}{\|f(X)\|} = O(\epsilon_{machine})$$

An algorithm is *stable* if $\forall X$ we have

$$\frac{\|\tilde{f}(X) - f(\tilde{X})\|}{\|f(\tilde{X})\|} = O(\epsilon_{machine})$$

for some \tilde{X} such that $\frac{\|\tilde{X} - X\|}{\|X\|} = O(\epsilon_{machine})$.

An algorithm is *backwards stable* if $\forall X$ there exists \tilde{X} with $\tilde{f}(X) = f(\tilde{X})$.

Reading Notes on Stability

Stable algorithms give nearly the right answer to nearly the right question.

Backward stable algorithms give exactly the right answer to nearly the right question.

Stability of Floating Point Arithmetic

We can find the stability of subtraction as follows. We have data $(x_1, x_2) \in \mathbb{C}^2$ and algorithm $\tilde{f}(x_1, x_2) = fl(x_1) -_{FL} fl(x_2)$.

By the axiom of floating point arithmetic, $fl(x_1) = x_1(1 + \epsilon_1)$, $fl(x_2) = x_2(1 + \epsilon_2)$, and $fl(x_1) -_{FL} fl(x_2) = (fl(x_1) - fl(x_2))(1 + \epsilon_3)$ for some $\epsilon_1, \epsilon_2, \epsilon_3 \leq \epsilon_{machine}$.

Combining, we get

$$\begin{aligned} fl(x_1) -_{FL} fl(x_2) &= (fl(x_1) - fl(x_2))(1 + \epsilon_3) \\ &= (x_1(1 + \epsilon_1) - x_2(1 + \epsilon_2))(1 + \epsilon_3) \\ &= x_1(1 + \epsilon_1)(1 + \epsilon_3) - x_2(1 + \epsilon_2)(1 + \epsilon_3) \\ &= x_1(1 + \epsilon_4) - x_2(1 + \epsilon_5) \end{aligned}$$

where $\epsilon_4, \epsilon_5 \leq 2\epsilon_{machine} + O(\epsilon_{machine}^2)$.

This is enough to imply subtraction is backward stable.

Accuracy and Stability

Theorem - let a backward stable algorithm be applied to a problem $f : X \rightarrow Y$ with condition number κ . Then the relative error satisfies

$$\frac{\|\tilde{f}(X) - f(X)\|}{\|f(X)\|} \leq O(\kappa(x)\epsilon_{machine})$$

Stability of QR Factorization QR factorization via Householder reflections is backward stable.

Consider QR factorization via reflections to solve a matrix equation $Ax = b$. We compute $A = QR, y = Q^*b, x = R^{-1}y$.

The first step is backward stable, as we noted above.

The second step is also backward stable, and so we get $(\tilde{Q} + \delta Q)\tilde{y} = b, \|\delta Q\| = \epsilon_{machine}$.

The third step is also backward stable, as we have $(\tilde{R} + \delta R)\tilde{x} = \tilde{y}, \|\delta R\|/\|\tilde{R}\| = O(\epsilon_{machine})$.

Then we can show that the algorithm for solving the matrix equation is backward stable, with $(A + \Delta A)\tilde{x} = b$, $\|\Delta A\|/\|A\| = \epsilon_{machine}$.

Start by concatenating to get

$$\begin{aligned} b &= (\tilde{Q} + \delta Q)(\tilde{R} + \delta R)\tilde{x} \\ &= (\tilde{Q}\tilde{R} + \tilde{Q}\delta R + \delta Q\tilde{R} + \delta Q\delta R)\tilde{x} \\ &= (A + \delta A + \tilde{Q}\delta R + \delta Q\tilde{R} + \delta Q\delta R)\tilde{x} \end{aligned}$$

where the last equality follows from backward stability of QR factorization (hence, $\tilde{Q}\tilde{R} = A + \delta A$). Now we just need to show that $\delta A + \tilde{Q}\delta R + \delta Q\tilde{R} + \delta Q\delta R = \Delta A$ is small relative to A .

We do one trick - using $\tilde{Q}\tilde{R} = A + \delta A$ and unitarity of \tilde{Q} , we get

$$\frac{\|\tilde{R}\|}{\|A\|} \leq \|\tilde{Q}^*\| \frac{\|A + \delta A\|}{\|A\|} = O(1 + \epsilon_{machine})$$

Now we can examine each term:

$$\begin{aligned} \frac{\|\tilde{Q}\delta R\|}{\|A\|} &\leq \|\tilde{Q}\| \frac{\|\delta R\|}{\|\tilde{R}\|} \frac{\|\tilde{R}\|}{\|A\|} = O(\epsilon_{machine})O(1 + \epsilon_{machine}) = O(\epsilon_{machine}) \\ \frac{\|\delta Q\tilde{R}\|}{\|A\|} &\leq \|\delta Q\| \frac{\|\tilde{R}\|}{\|A\|} = \epsilon_{machine}O(1 + \epsilon_{machine}) = O(\epsilon_{machine}) \\ \frac{\|\delta Q\delta R\|}{\|A\|} &\leq \|\delta Q\| \frac{\|\delta R\|}{\|A\|} = O(\epsilon_{machine})O(\epsilon_{machine}) = O(\epsilon_{machine}^2) \end{aligned}$$

Putting this all together, we get $\|\Delta A\|/\|A\| = O(\epsilon_{machine})$.

Lecture 2/5: Stability and Conditioning of Least Squares

Given the least squares problem $Ax = b$, we want to ask the question of how x, y change as A, b change?

Recall the definition of the condition number of a problem $f : X \rightarrow Y$:

$$\kappa(x) = \sup_{\delta x} \frac{\|\delta f(x)\|}{\|f(x)\|} \frac{\|x\|}{\|\delta x\|}$$

Theorem: Let $A \in \mathbb{C}^{m \times n}$ be full rank and consider the solution x to the least squares problem $Ax = b$ and $y = Ax$. Then perturbations in A yield condition on x, y as follows:

$$\begin{bmatrix} b \text{ vs } A & y & x \\ b & \frac{1}{\cos(\theta)} & \frac{\kappa(A)}{\gamma \cos(\theta)} \\ A & \frac{\kappa(A)}{\cos(\theta)} & \kappa(A) + \frac{\kappa(A)^2 \tan(\theta)}{\gamma} \end{bmatrix}$$

Proof: the first few sensitivities are easy to pull out.

$$\begin{aligned} \kappa_{b \rightarrow y} &= \frac{\|J(b)\|}{\|y\|/\|b\|} = \frac{\|P\|}{\|y\|/\|b\|} = \frac{1}{\cos(\theta)} \\ \kappa_{b \rightarrow x} &= \frac{\|J(b)\|}{\|x\|/\|b\|} = \|A^\dagger\| \frac{\|b\|}{\|y\|} \frac{\|y\|}{\|x\|} = \|A^\dagger\| \|A\| \frac{1}{\cos(\theta)} = \frac{\kappa(A)}{\gamma \cos(\theta)} \end{aligned}$$

To calculate the sensitivities with respect to changes in A we have to understand how changes in the colspace of A affect our x, y .

Lecture 2/7: LU Decompositions

Triangular systems are easy to deal with - you can invert via back substitution.

Def: A matrix is *unit lower (upper) triangular* if it is lower (upper) triangular with 1's on the diagonal.

We are concerned with matrices that admit an *LDU decomposition* where we write $A = LDU$ where L is unit lower triangular, D is diagonal, and U is upper triangular.

An *LU decomposition* is a special case of an LDU decomposition where $D = I$.

Given $A = LDU$, we can solve $Ax = b$ by solving $LDUx = b$. We can solve $Ly = b$ by forward substitution, $Dz = y$ by scaling, and $Ux = z$ by back substitution.

We can compute the LU decomposition by Gaussian elimination.

Ex. Consider the matrix

$$\begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 4 \\ 3 & 9 & 27 \end{bmatrix}$$

Proceed one step at a time. First let's try to zero out the (3,1) entry, so we want to take $R_3 - 3R_1 \rightarrow R_3$. We can write that in matrix form as

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ -3 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 4 \\ 3 & 9 & 27 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 4 \\ 0 & 6 & 24 \end{bmatrix}$$

Note that our 1st elimination matrix was unit lower triangular. In fact its inverse will be as well, as the inverse is given by just flipping the sign on the (3,1) entry. Hence we also have

$$\begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 4 \\ 3 & 9 & 27 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 3 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 4 \\ 0 & 6 & 24 \end{bmatrix}$$

Now say we want to zero out the (2,1) entry. We can do this by taking $R_2 - R_1 \rightarrow R_2$. We can write this as

$$\begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 4 \\ 0 & 6 & 24 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & 3 \\ 0 & 6 & 24 \end{bmatrix}$$

Now we can see where this is going. The product of all these individual elimination matrices will give us our L matrix, and the product of all the modified matrices will give us our U matrix. That is, we now have

$$\begin{aligned} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 4 \\ 3 & 9 & 27 \end{bmatrix} &= \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 3 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & 3 \\ 0 & 6 & 24 \end{bmatrix} \\ &= \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 3 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & 3 \\ 0 & 6 & 24 \end{bmatrix} \end{aligned}$$

Lastly, we just need to eliminate the (3,2) entry. We can do this by taking $R_3 - 6R_2 \rightarrow R_3$. We can write this as

$$\begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & 3 \\ 0 & 6 & 24 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 6 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & 3 \\ 0 & 0 & 6 \end{bmatrix}$$

Combining all the above, we get

$$\begin{aligned} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 4 \\ 3 & 9 & 27 \end{bmatrix} &= \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 3 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 6 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & 3 \\ 0 & 0 & 6 \end{bmatrix} \\ &= \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 3 & 6 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & 3 \\ 0 & 0 & 6 \end{bmatrix} \\ A &= LU \end{aligned}$$

Algorithmically, we can write this as follows. In Gaussian elimination for LU , let A be $m \times m$. Then we have

$$\begin{aligned} U &= A, L = I \\ \text{for } k &= 1 \text{ to } m - 1 \\ \text{for } j &= k + 1 \text{ to } m \\ L_{j,k} &= U_{j,k}/U_{k,k} \\ U_{j,k:m} &= U_{j,k:m} - L_{j,k}U_{k,k:m} \end{aligned}$$

Stability of LU Decomposition

Consider the matrix

$$A = \begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix} \simeq \begin{bmatrix} 10^{-20} & 1 \\ 1 & 1 \end{bmatrix}$$

The LU decomposition of A is

$$L = \begin{bmatrix} 1 & 0 \\ 10^{20} & 1 \end{bmatrix}, U = \begin{bmatrix} 10^{-20} & 1 \\ 0 & 1 - 10^{20} \end{bmatrix}$$

Consider perturbations of the decomposition matrices so we get

$$\tilde{L} = \begin{bmatrix} 1 & 0 \\ 10^{20} & 1 \end{bmatrix}, \tilde{U} = \begin{bmatrix} 10^{-20} & 1 \\ 0 & -10^{20} \end{bmatrix} \Rightarrow \tilde{L}\tilde{U} = \begin{bmatrix} 10^{-20} & 1 \\ 1 & 0 \end{bmatrix}$$

Hence even if we can produce L, U stably, may not mean $Ax = b$ is solved stably. This points to the fact that the LU factorization for solving $Ax = b$ is not backward stable.

We can get around this by modifying the algorithm to better choose the way we eliminated entries. The choice of such eliminations is called choosing a *pivot*.

Lecture 2/14: Cholesky Factorization

Recall that a real symmetric matrix M is *positive definite* if for all non-zero vectors x , we have $x^T M x > 0$. A real symmetric matrix is *positive semi-definite* if for all non-zero vectors x , we have $x^T M x \geq 0$.

Note that we can write $x^T A x$ as

$$x^T A x = \sum_{i=1}^n \sum_{j=1}^n A_{ij} x_i x_j = \sum_{i=1}^n A_{ii} x_i^2 + 2 \sum_{i>j} A_{ij} x_i x_j$$

Equivalently, a real symmetric matrix M is positive definite if all its eigenvalues are positive, and positive semi-definite if all its eigenvalues are non-negative.

Recall the *Gram matrix* of a matrix B is $B^T B$. Every Gram matrix is PSD, and a Gram matrix $A = B^T B$ is positive definite if $x^T A x = \|Bx\|^2 > 0$ which is true when B is non-singular.

Cholesky Factorization

Every positive definite matrix A has a *Cholesky factorization* of the form $A = R^T R$ where R is upper triangular with positive diagonal elements. R is the *Cholesky factor* of A .

The complexity of computing R is $O(n^3/3)$. R is sort of like the “square root” of A .

We can write the Cholesky factorization as follows. Let A be a positive definite matrix. Then we have

$$\begin{aligned} A &= R^T R \\ \begin{bmatrix} A_{11} & A_{1,2:n} \\ A_{2:n,1} & A_{2:n,2:n} \end{bmatrix} &= \begin{bmatrix} R_{11} & 0 \\ R_{1,2:n}^T & R_{2:n,2:n}^T \end{bmatrix} \begin{bmatrix} R_{11} & R_{1,2:n} \\ 0 & R_{2:n,2:n} \end{bmatrix} \\ &= \begin{bmatrix} R_{11}^2 & R_{11} R_{1,2:n} \\ R_{1,2:n}^T R_{11} & R_{1,2:n}^T R_{1,2:n} + R_{2:n,2:n}^T R_{2:n,2:n} \end{bmatrix} \end{aligned}$$

Hence, we get $R_{11} = \sqrt{A_{11}}$ and $R_{1,2:n} = \frac{A_{1,2:n}}{R_{11}}$ and $A_{2:n,2:n} - R_{1,2:n}^T R_{1,2:n} = R_{2:n,2:n}^T R_{2:n,2:n}$, and we can recursively apply this to the bottom right submatrix.

Ex: Let $A = \begin{bmatrix} 25 & 15 & -5 \\ 15 & 18 & 0 \\ -5 & 0 & 11 \end{bmatrix}$. Then we have

$$\begin{aligned} \begin{bmatrix} 25 & 15 & -5 \\ 15 & 18 & 0 \\ -5 & 0 & 11 \end{bmatrix} &= \begin{bmatrix} R_{11} & 0 & 0 \\ R_{21} & R_{22} & 0 \\ R_{31} & R_{32} & R_{33} \end{bmatrix} \begin{bmatrix} R_{11} & R_{21} & R_{31} \\ 0 & R_{22} & R_{32} \\ 0 & 0 & R_{33} \end{bmatrix} \\ &= \begin{bmatrix} 5 & 0 & 0 \\ 3 & R_{22} & 0 \\ -1 & R_{23} & R_{33} \end{bmatrix} \begin{bmatrix} 5 & 3 & -1 \\ 0 & R_{22} & R_{23} \\ 0 & 0 & R_{33} \end{bmatrix} \end{aligned}$$

After completing the first row/column, we start with the bottom right submatrix.

$$\begin{aligned} A_{2:n,2:n} - R_{1,2:n}^T R_{1,2:n} &= R_{2:n,2:n}^T R_{2:n,2:n} \\ \begin{bmatrix} 18 & 0 \\ 0 & 11 \end{bmatrix} - \begin{bmatrix} 3 & -1 \end{bmatrix} \begin{bmatrix} 3 \\ -1 \end{bmatrix} &= \begin{bmatrix} R_{22} & 0 \\ R_{23} & R_{33} \end{bmatrix} \begin{bmatrix} R_{22} & R_{23} \\ 0 & R_{33} \end{bmatrix} \\ \begin{bmatrix} 9 & 3 \\ 3 & 10 \end{bmatrix} &= \begin{bmatrix} R_{22}^2 & R_{22} R_{23} \\ R_{22} R_{23} & R_{23}^2 + R_{33}^2 \end{bmatrix} \\ &\Rightarrow R_{22} = 3, R_{23} = 1, R_{33} = \sqrt{10 - 1} = 3 \end{aligned}$$

Thus we have $R = \begin{bmatrix} 5 & 3 & -1 \\ 0 & 3 & 1 \\ 0 & 0 & 3 \end{bmatrix}$.

Lecture 2/21: Eigenvalue Problems

To start note that the matrix $A = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$ is not diagonalizable. That means, there does not exist B, D such that $A = BDB^{-1}$.

A series of definitions is useful. For a matrix A , a pair (x, λ) is an *eigenpair* of A if $Ax = \lambda x$ for some scalar λ . The scalar λ is called the *eigenvalue* of A corresponding to the eigenvector x . The *eigenspace* is the space on which A acts like scalar multiplication: the set of eigenvectors corresponding to an eigenvalue λ is denoted E_λ . Note $AE_\lambda \subset E_\lambda$. The *spectrum* of A is the set of all eigenvalues of A .

Thm: The determinant of a matrix A is the product of its eigenvalues, and the trace of A is the sum of its eigenvalues.

A factorization of a square matrix A into a product $A = X\Lambda X^{-1}$ where Λ is diagonal and X is invertible is called an *eigenvalue factorization* of A . The columns of X are the eigenvectors of A , and the diagonal elements of Λ are the eigenvalues of A .

Every matrix has an SVD, but not every matrix has an eigenvalue factorization. For example, the matrix $A = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$ does not have an eigenvalue factorization. However, it does have an SVD.

Characteristic Polynomials and Eigenvalues

The characteristic polynomial of a matrix $A \in \mathbb{C}^{m \times m}$ is defined as $p_A(z) = \det(zI - A)$. The characteristic polynomial is a monic polynomial of degree n in λ .

Thm: λ is an eigenvalue of A iff it is a root of the characteristic polynomial. If $A \in \mathbb{C}^{m \times m}$ and you count with algebraic multiplicity, then there are m eigenvalues not necessarily distinct.

We distinguish two notions of multiplicity. The *algebraic multiplicity* of an eigenvalue λ is the number of times λ appears as a root of the characteristic polynomial. The *geometric multiplicity* of an eigenvalue λ is the dimension of the eigenspace E_λ .

Ex: let $A = \begin{bmatrix} 2 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 2 \end{bmatrix}, B = \begin{bmatrix} 2 & 1 & 0 \\ 0 & 2 & 1 \\ 0 & 0 & 2 \end{bmatrix}$. $p_A(z) = (z - 2)^3, p_B(z) = (z - 2)^3$. The algebraic multiplicity of 2 is 3 for both A and B .

However, note that the eigenvectors of A are just the coordinate vectors e_1, e_2, e_3 . This means that E_λ has dimension 3, and hence λ has geometric multiplicity 3.

For B the only eigenvector is e_1 , so the geometric multiplicity of λ is 1.

Thm: The geometric multiplicity of an eigenvalue is less than or equal to its algebraic multiplicity. Moreover, if X is nonsingular then A and XAX^{-1} have the same eigenvalues and the same multiplicities.

Diagonalization and Unitary Diagonalization

An eigenvalue whose algebraic and geometric multiplicities are equal is called *simple*. An eigenvalue whose algebraic multiplicity is greater than its geometric multiplicity is called *defective*. A matrix with any defective eigenvalues is called *defective* (otherwise *non-defective*).

Theorem: A matrix $A \in \mathbb{C}^{m \times m}$ is non-defective iff it has an eigendecomposition $A = X\Lambda X^{-1}$ where X is nonsingular and Λ is diagonal.

Furthermore, we might get lucky and have X be unitary, in which case the eigendecomposition is given as $A = X\Lambda X^*$. Then we say that A is *unitarily diagonalizable*.

A *Schur factorization* of A is given by $A = QTQ^*$ where Q is unitary and T is upper triangular. Since A, T are similar, and since upper triangular matrices contain their eigenvalues on the diagonal, the diagonal elements of T are the eigenvalues of A .

Theorem: Every square matrix has a Schur factorization.

Now note that if $A = QTQ^*$ is normal, i.e. $AA^* = A^*A$, then T is also normal. Since T is normal and upper triangular, it is diagonal. Thus, if A is normal, then the Schur form gives us unitary diagonalization!

Theorem: A matrix A is normal iff it is unitarily diagonalizable.

To summarize:

1. Every matrix A has an SVD $A = U\Sigma V^*$.
2. Every square matrix A has a unitary triangularization (Schur factorization) $A = QTQ^*$.
3. Every non-defective matrix A has a diagonalization (resp. eigendecomposition) $A = X\Lambda X^{-1}$.
4. Every normal matrix A has a unitary diagonalization $A = Q\Lambda Q^*$.

These diagonalizations are useful for finding eigenvalues. Recall that root-finding is ill-conditioned, so we want to avoid calculating eigenvalues by solving the characteristic polynomial.

Iterative Methods for Eigenvalue Computation

When computing Schur form iteratively we end up with intermediate matrices whose product is our final matrix. This resembles compiling Householder reflections into a final QR factorization.

Because we need to be careful about pivoting, we want to have our intermediate steps end in *Hessenberg form*, with zeroes below the first subdiagonal.

Consider $A = \begin{bmatrix} a_{11} & \dots & a_{1m} \\ \vdots & \ddots & \vdots \\ a_{m1} & \dots & a_{mm} \end{bmatrix}$. Then apply our Householder reflections to get

$$\begin{bmatrix} a_{11} & \dots & a_{1m} \\ \vdots & \ddots & \vdots \\ a_{m1} & \dots & a_{mm} \end{bmatrix} \xrightarrow{Q_1^*} \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1m} \\ a'_{21} & a'_{22} & \dots & a'_{2m} \\ 0 & a_{32} & \dots & a_{3m} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & a_{mm} \end{bmatrix} \xrightarrow{Q_1} \begin{bmatrix} a'_{11} & a'_{12} & \dots & a'_{1m} \\ a'_{21} & a'_{22} & \dots & a'_{2m} \\ 0 & a_{32} & \dots & a_{3m} \\ 0 & 0 & a_{43} & \dots \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & a_{mm} \end{bmatrix}$$

Then we iterate on the sub-blocks. That is, we first compute the Householder \hat{P}_1 which sends $(a_{2:m,1})$ to e_1 and let $P_1 = \begin{bmatrix} 1 & 0 \\ 0 & \hat{P}_1 \end{bmatrix}$. Then we compute $P_1 A P_1^*$ and get

$$\begin{bmatrix} a'_{11} & a'_{12} & \dots & a'_{1m} \\ a'_{21} & a'_{22} & \dots & a'_{2m} \\ 0 & a_{32} & \dots & a_{3m} \\ 0 & a_{42} & a_{43} & \dots \\ \vdots & \vdots & \ddots & \vdots \\ 0 & a_{m3} & \dots & a_{mm} \end{bmatrix}$$

Then compute \hat{P}_2 to send $(a_{3:m,2})$ to e_2 and let $P_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & \hat{P}_2 \end{bmatrix}$.

Once complete, we get $P_n \dots P_1 A P_1^{-1} \dots P_n^{-1} = H$, where H is Hessenberg. Then let $Q = P_n \dots P_1$ and $A = Q H Q^*$.

Lecture 2/26: Hessenberg Reductions

The algorithm for Hessenberg reduction via Householder reflections is as follows:

```
A <- matrix(rnorm(n*n), n,n)
H <- A
V <- vector(mode = "list", length = n-2)

# Householder transformation
for (k in 1:(n-2)) {
  v <- H[(k+1):n, k]
  sgn <- sign(v[1])
  if (sgn == 0) sgn <- 1
  v[1] <- v[1] + sgn * norm(v,type="2")
}
```

```

    if (norm(v,type="2") != 0) v <- v / norm(v,type="2")

    H[(k+1):n,k:n] <- H[(k+1):n, k:n] - 2 * v %*% (t(v) %*% H[(k+1):n,k:n])
    H[ , (k+1):n] <- H[ , (k+1):n] - (2 * (H[ , (k+1):n] %*% v)) %*% t(v)
    V[[k]] <- v
  }
  Q <- diag(nrow=n)
  for (j in (n-2):1) {
    Q[(j+1):n, ] <- Q[(j+1):n, ] - (2 * V[[j]]) %*% (t(V[[j]]) %*% Q[(j+1):n, ])
  }

  list(H = H, P = Q)

```

This Householder method reduces any matrix to Hessenberg form which is nearly triangular in $\frac{10}{3}n^3$ flops. If A is Hermitian, then the Hessenberg form is tridiagonal, and this allows us to do computations in only $\frac{4}{3}n^3$ flops.

What's critical here is that the Hessenberg form of a matrix has the same eigenvalues as the original matrix.

```

library(pracma)
A_diag <- diag(1:10)
A_hess_1 <- hessenberg(A_diag)
list(round(A_hess_1$H), round(A_hess_1$P))

## [[1]]
##      [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9] [,10]
## [1,]    1    0    0    0    0    0    0    0    0    0
## [2,]    0    2    0    0    0    0    0    0    0    0
## [3,]    0    0    3    0    0    0    0    0    0    0
## [4,]    0    0    0    4    0    0    0    0    0    0
## [5,]    0    0    0    0    5    0    0    0    0    0
## [6,]    0    0    0    0    0    6    0    0    0    0
## [7,]    0    0    0    0    0    0    7    0    0    0
## [8,]    0    0    0    0    0    0    0    8    0    0
## [9,]    0    0    0    0    0    0    0    0    9    0
## [10,]   0    0    0    0    0    0    0    0    0   10
##
## [[2]]
##      [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9] [,10]
## [1,]    1    0    0    0    0    0    0    0    0    0
## [2,]    0    1    0    0    0    0    0    0    0    0
## [3,]    0    0    1    0    0    0    0    0    0    0
## [4,]    0    0    0    1    0    0    0    0    0    0
## [5,]    0    0    0    0    1    0    0    0    0    0
## [6,]    0    0    0    0    0    1    0    0    0    0
## [7,]    0    0    0    0    0    0    1    0    0    0
## [8,]    0    0    0    0    0    0    0    1    0    0
## [9,]    0    0    0    0    0    0    0    0    1    0
## [10,]   0    0    0    0    0    0    0    0    0    1

```

```

Q_unit <- randortho(10,type="orthonormal")
A_trans <- Q_unit %*% A_diag %*% t(Q_unit)
A_hess_2 <- hessenberg(A_trans)
list(round(A_hess_2$H), round(A_hess_2$P))

```

```
## [[1]]
```



```

##      [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9] [,10]
## [1,]    4   -2    0    0    0    0    0    0    0    0
## [2,]   -2    6   -2    0    0    0    0    0    0    0
## [3,]    0   -2    7    2    0    0    0    0    0    0
## [4,]    0    0    2    6    2    0    0    0    0    0
## [5,]    0    0    0    2    5   -2    0    0    0    0
## [6,]    0    0    0    0   -2    4    1    0    0    0
## [7,]    0    0    0    0    0    1    6    2    0    0
## [8,]    0    0    0    0    0    0    2    6   -2    0
## [9,]    0    0    0    0    0    0    0   -2    7    1
## [10,]   0    0    0    0    0    0    0    0    1    5
##
## [[2]]
##      [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9] [,10]
## [1,]    1    0    0    0    0    0    0    0    0    0
## [2,]    0    0    1    0    0    0    1    0    0    1
## [3,]    0   -1    0    1    0    0    0    0    0    0
## [4,]    0    0    0    0    0    0    0    0    0    0
## [5,]    0    0    1    0    0    0    0    0   -1    0
## [6,]    0    0    0    0    0    0   -1    0    0    0
## [7,]    0    0    0    0    0    0    0   -1    0    0
## [8,]    0    0    0   -1    0    1    0    0    0    0
## [9,]    0    0    0    0    1    0    0    0    0   -1
## [10,]   0    1    0    0    0    1    0    0    0    0

```

Other Eigenvalue methods

Other method for eigenvalue computation include the Jacobi algorithm, the divide-and-conquer algorithm, and the bisection algorithm.

Jacobi Algorithm For symmetric matrices, we want to repeatedly diagonalize small submatrices of A until A is diagonal, which obviously gives you the eigenvalues! The trick here is to use Givens rotations to diagonalize 2×2 submatrices. That is, for a submatrix $\begin{bmatrix} a & b \\ c & d \end{bmatrix}$, we want to find a Givens rotation G such that $G^T \begin{bmatrix} a & b \\ c & d \end{bmatrix} G = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix}$.

We start by trying $G = \begin{bmatrix} c & s \\ -s & c \end{bmatrix}$, where $c = \cos \theta$ and $s = \sin \theta$, $\tan \theta = \frac{2d}{b-a}$. We repeatedly apply this to the 2×2 submatrices, usually beginning with the submatrix with largest off-diagonal entries.

Lecture 2/28: QR Algorithm

Let A be real and symmetric. Hence it has real eigenvalues and orthonormal eigenvectors, and is diagonalizable by its Schur form.

The (reduced) QR factorization of A is $A = QR$, where Q is orthogonal and R is upper triangular.

Power Iteration as a Motivation

One jumping off point is the power iteration method for calculating eigenvalues. We start with a random unit vector $v^{(0)}$ and repeatedly apply A and renormalize. In the limit, $v^{(k)} \rightarrow v_1$, the eigenvector corresponding to the largest eigenvalue.

Consider an extension of power iteration where we begin with a matrix $V^{(0)} = (v_1^{(0)}, \dots, v_n^{(0)})$ and repeatedly apply A and renormalize. In the limit, $V^{(k)} \rightarrow V$, where the columns of V are the eigenvectors of A . We

assume that there is a gap between the first and second eigenvalues, and the rate of convergence is determined by the max of the eigengaps.

QR Algorithm

The basic idea going forward is that we want to use the similarity transform enabled by the QR factorization. We start with $A = QR$, then triangularize via $R = Q^*A$, then right multiply by Q again to get $RQ = Q^*AQ$. We then iterate on this process.

In particular, we get the QR algorithm via the following steps:

1. Start with $A_0 = A$.
2. Compute the QR factorization $A_k = Q_k R_k$.
3. Set $A_{k+1} = R_k Q_k$.
4. Repeat until A_k is upper triangular.

Since A_k is always similar to A , it has the same eigenvalues.

```
A <- matrix(c(1,2,3,4),nrow=2,byrow=T)
for(i in 1:100) {
  Q <- qr.Q(qr(A))
  R <- qr.R(qr(A))
  A <- R %*% Q
}
```

```
A
##           [,1]      [,2]
## [1,]  5.372281e+00 -1.0000000
## [2,]  1.071841e-115 -0.3722813
```

```
Q
```

```
##           [,1]      [,2]
## [1,] -1.000000e+00 -2.879117e-115
## [2,] -2.879117e-115  1.000000e+00
```

Practical Modifications to the QR Algorithm

The QR algorithm is not practical for large matrices, but there are some modifications that make it more practical.

We first reduce A to tridiagonal form via Householder reflections. Then we apply the QR algorithm to the tridiagonal matrix. We also pick a shift μ_k to accelerate convergence.

With these modifications in mind, the QR algorithm becomes:

1. Start with the tridiagonalization of A , $Q_0 A Q_0^T = T_0$.
2. Pick a shift μ_k (one easy choice is the final diagonal element of our current iterate $\mu_k = A_{mm}^{(k-1)}$).
3. Compute the QR factorization of $A^{(k-1)} - \mu_k I$: $A^{(k)} - \mu_k I = Q_k R_k$.
4. Recombine the factors to get $A^{(k)} = R_k Q_k + \mu_k I$.
5. If any off-diagonal elements are small, we can set them to zero and continue.
6. Reapply the QR algorithm to the resulting submatrices.

Analysis of QR Algorithm

QR works because it is secretly just a simultaneous iteration of the power iteration method.

Thm: Let the pure QR algorithm be applied to a real symmetric matrix A with distinct eigenvalues. Then the iterates $A^{(k)}$ converge to a diagonal matrix with the eigenvalues of A on the diagonal at a rate given by $\max(\lambda_{k+1}/\lambda_k)$, and the iterates Q_k converge to the matrix of eigenvectors of A .

Lecture 3/4: Jacobi Algorithm

Consider the off-diagonal sum of a matrix A , $\text{off}(A) = \sum_{i \neq j} |A_{ij}|$.

Theorem: Let A be a real symmetric matrix. Let $A^{(d)}$ be the d -th iterate of the Jacobi algorithm. Then

$$\text{off}(A^{(d)}) \leq \left(1 - \frac{2}{N^2 - N}\right)^d \text{off}(A)$$

Bisection Algorithm

Let $A^{(1)}, \dots, A^{(m)}$ be the principal upper-left submatrices of a real symmetric matrix A . Each submatrix $A^{(k)}$ has eigenvalues $\lambda_1^{(k)} \geq \dots \geq \lambda_n^{(k)}$. One surprising result is that the eigenvalues of the submatrices *interlace*:

$$\lambda_j^{(k+1)} < \lambda_j^{(k)} < \lambda_{j+1}^{(k+1)}$$

Ex. Let A be a tridiagonal matrix such as

$$A = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 2 & 1 \\ 0 & 0 & 1 & -1 \end{bmatrix}$$

The determinants of the submatrices are $\det(A^{(1)}) = 1, \det(A^{(2)}) = -1, \det(A^{(3)}) = -3, \det(A^{(4)}) = 4$.

This sequence of determinants of the principal submatrices is called the *Sturm sequence* of a matrix. The number of sign changes in the Sturm sequence gives the number of negative eigenvalues of the matrix.

The interesting extension here is to consider shifts of A by multiples of the identity. We can use the Sturm sequence to find the number of eigenvalues not only if they're negative, but as it happens in any interval!

Lecture 3/6: Arnoldi Iteration

The goal of *Arnoldi iteration* is to obtain a Hessenberg form of a not-necessarily-symmetric matrix $A \in \mathbb{C}^{m \times m}$. That is, we want $A = QHQ^*$, where H is upper Hessenberg and Q is unitary.

Start by forming the first $n < m$ columns of Q and the $(n+1) \times n$ upper-left submatrix of H , which should also be Hessenberg:

$$Q_n = [q_1 \quad \dots \quad q_n], \tilde{H}_n = \begin{bmatrix} h_{11} & h_{12} & \dots & h_{1n} \\ h_{21} & h_{22} & \dots & h_{2n} \\ 0 & h_{32} & \dots & h_{3n} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & h_{n+1,n} \end{bmatrix}$$

Note that this gives $AQ_n = Q_{n+1}\tilde{H}_n$. More explicitly, we have

$$AQ_n = h_{1n}q_1 + h_{2n}q_2 + \dots + h_{n+1,n}q_{n+1} \Rightarrow \langle AQ_n, q_i \rangle = h_{in}$$

Algorithm for Arnoldi Iteration

The above suggests the following algorithm for calculating Q, H and $A = QHQ^*$:

1. Choose b arbitrarily and set $q_1 = b / \|b\|_2$.
2. For $k = 1, \dots, n$:
 - a. $v = Aq_k$
 - b. For $j = 1, \dots, k$:
 - i. $h_{jk} = q_j^* v$

```

ii. $v = v - h_{\{jk\}}q_j$
c. $h_{\{k+1,k\}} = \|v\|_2$
d. If $h_{\{k+1,k\}} = 0$, stop.
e. $q_{\{k+1\}} = v/h_{\{k+1,k\}}$

```

Krylov Subspaces

Given A, b , we define the n -th Krylov subspace as $\mathbb{K}_n(A, b) = \text{span}\{b, Ab, A^2b, \dots, A^{n-1}b\}$. The Arnoldi iteration is a way to find an orthonormal basis for the Krylov subspace.

Note that the above equation we used to express Aq_n in terms of q_i leads to an observation about Krylov subspaces: in particular, q_i generated by the Arnoldi iteration are orthonormal (generated similarly to Gram-Schmidt process) and form a basis for $\mathbb{K}_n(A, b)$.

We define the *Krylov matrix* as

$$K_n = [b \quad Ab \quad \dots \quad A^{n-1}b]$$

The QR factorization of K_n is $K_n = Q_n R_n$ where Q_n is as above, the first n columns of the Arnoldi iteration.

Note that as n increases, the Krylov subspaces become better approximations to the eigenspaces of A .

Theorem: The matrices Q_n generated by Arnoldi iteration are (reduced) QR factors of $K_n = Q_n R_n$. The Hessenberg matrices H_n are the projections $H_n = Q_n^* A Q_n$, that is, the orthogonal projection of A onto \mathbb{K}_n . The iterates of Arnoldi satisfy the relation $AQ_n = Q_{n+1}H_n$.

Further Observations on Krylov Subspaces

Let $x \in \mathbb{K}_n$, so $x = c_0b + c_1Ab + \dots + c_{n-1}A^{n-1}b$. We would prefer to work with polynomials here, so let $x = q(A)b$, $q(z) = c_0 + c_1z + \dots + c_{n-1}z^{n-1}$.

Consider \mathbb{P}^n , the set of all monic polynomials of degree n . The *Arnoldi/Lanczos approximation problem* is to find $p^n \in \mathbb{P}_n$ such that $\|p^n(A)b\|_2$ is minimized. Note that $p^n(A)b = A^n b + q(A)b = A^n b + x$, $x \in \mathbb{K}_n$.

Theorem: If K_n is rank n , then the Arnoldi approximation problem has a unique solution p^n , which is the characteristic polynomial of H_n .

The sketch of the proof is as follows. We seek y such that $\|A^n b - Q_n y\|$ is minimal. This is just a least squares problem where we find the point in \mathbb{K}_n closest to $A^n b$.

The solution is $Q_n^* p^n(A)b = 0$. Since $A = QHQ^*$, we have $Q_n^* p^n(QHQ^*)b = Q_n^* Q p^n(H) Q^* b = 0$. By the Cayley-Hamilton theorem, we are done.

Other notes:

1. \mathbb{P}^n is invariant to scalar translation, so the Arnoldi method is also translation invariant.
2. The Arnoldi method is also invariant to scaling, so if we transform $A \rightarrow \sigma A$, then the eigenvalues of H_n become multiplied by σ (these are the *Ritz values*).
3. Ritz values are invariant under unitary similarity transforms.

HW 3/11

Question Eigenvalue Problems 1

To calculate the eigenvalues of a matrix $A = D + ww^T$, where D is a diagonal matrix and w is a column vector, we can use the following approach: Let's denote the diagonal elements of D as d_1, d_2, \dots, d_n , and the elements of the column vector w as w_1, w_2, \dots, w_n . Step 1: Find the characteristic polynomial of A . The characteristic polynomial of A is given by $\det(A - \lambda I)$, where I is the identity matrix and λ represents the eigenvalues. We can expand the determinant as follows:

$$\det(A - \lambda I) = \det(D + ww^T - \lambda I) = \det(D - \lambda I + ww^T) = \det(D - \lambda I) \det(I + (D - \lambda I)^{-1} ww^T)$$

Since D is a diagonal matrix, $(D - \lambda I)^{-1}$ is also a diagonal matrix with entries $(d_i - \lambda)^{-1}$. Step 2: Simplify the second determinant factor. Let's denote the second determinant factor as $f(\lambda)$:

$$f(\lambda) = \det(I + (D - \lambda I)^{-1} w w^T)$$

Using the Matrix Determinant Lemma, we can simplify $f(\lambda)$ as:

$$f(\lambda) = 1 + w^T (D - \lambda I)^{-1} w$$

Substituting the diagonal elements, we get:

$$f(\lambda) = 1 + \sum_{i=1}^n \frac{w_i^2}{d_i - \lambda}$$

Step 3: Find the characteristic polynomial by combining the two determinant factors. The characteristic polynomial of A is given by:

$$\det(A - \lambda I) = \prod_{i=1}^n (d_i - \lambda) \left(1 + \sum_{j=1}^n \frac{w_j^2}{d_j - \lambda} \right)$$

Step 4: Find the eigenvalues by solving the characteristic equation. The eigenvalues of A are the values of λ that satisfy the characteristic equation:

$$\prod_{i=1}^n (d_i - \lambda) \left(1 + \sum_{j=1}^n \frac{w_j^2}{d_j - \lambda} \right) = 0$$

This equation can be solved numerically or analytically, depending on the specific values of d_i and w_i . In summary, the eigenvalues of the matrix $A = D + w w^T$, where D is a diagonal matrix and w is a column vector, can be obtained by solving the characteristic equation involving the diagonal elements of D and the elements of w .

Question Eigenvalue Problems 2: Problem 30.5

Calculating e-values via Jacobi iteration.

```
rows <- 4
A_mat <- matrix(rnorm(rows^2), nrow=rows)
A <- A_mat + t(A_mat)
n <- nrow(A)
U <- diag(n)
max_iter <- 20000
tol <- 1e-8
off_diag_entries <- rep(0, max_iter)
for (iter in 1:max_iter) {
  A_off_diag <- A - diag(diag(A))
  off_diag_entries[iter] <- sum(A_off_diag)
  max_off_diag <- which(abs(A_off_diag) == max(abs(A_off_diag)), arr.ind = TRUE)

  if (A_off_diag[max_off_diag] < tol) {
    break
  }
}

i <- max_off_diag[1]
j <- max_off_diag[2]
```

```

phi <- atan(2 * A[i, j] / (A[j, j] - A[i, i])) / 2
c <- cos(phi)
s <- sin(phi)

P <- diag(n)
P[i, i] <- c
P[j, j] <- c
P[i, j] <- -s
P[j, i] <- s

A <- P %*% A %*% t(P)
U <- U %*% P
}

```

```

## Warning in if (A_off_diag[max_off_diag] < tol) {: the condition has length > 1
## and only the first element will be used

```

```

## Warning in if (A_off_diag[max_off_diag] < tol) {: the condition has length > 1
## and only the first element will be used

```

```

## Warning in if (A_off_diag[max_off_diag] < tol) {: the condition has length > 1
## and only the first element will be used

```

```

## Warning in if (A_off_diag[max_off_diag] < tol) {: the condition has length > 1
## and only the first element will be used

```

```

eigenvalues <- diag(A)
eigenvectors <- U

```

```

list(eigenvalues = eigenvalues, eigenvectors = eigenvectors)

```

```

## $eigenvalues
## [1] -0.04450106  2.19214568 -4.02569123  3.26747197
##
## $eigenvectors
##           [,1]      [,2]      [,3]      [,4]
## [1,]  0.8042380  0.0000000  0.2732838 -0.5277473
## [2,] -0.2037353  0.7739592  0.5995658  0.0000000
## [3,] -0.2490113 -0.6332355  0.7328071  0.0000000
## [4,]  0.4996865  0.0000000  0.1697958  0.8494014

```

```

det(A) - prod(eigenvalues)

```

```

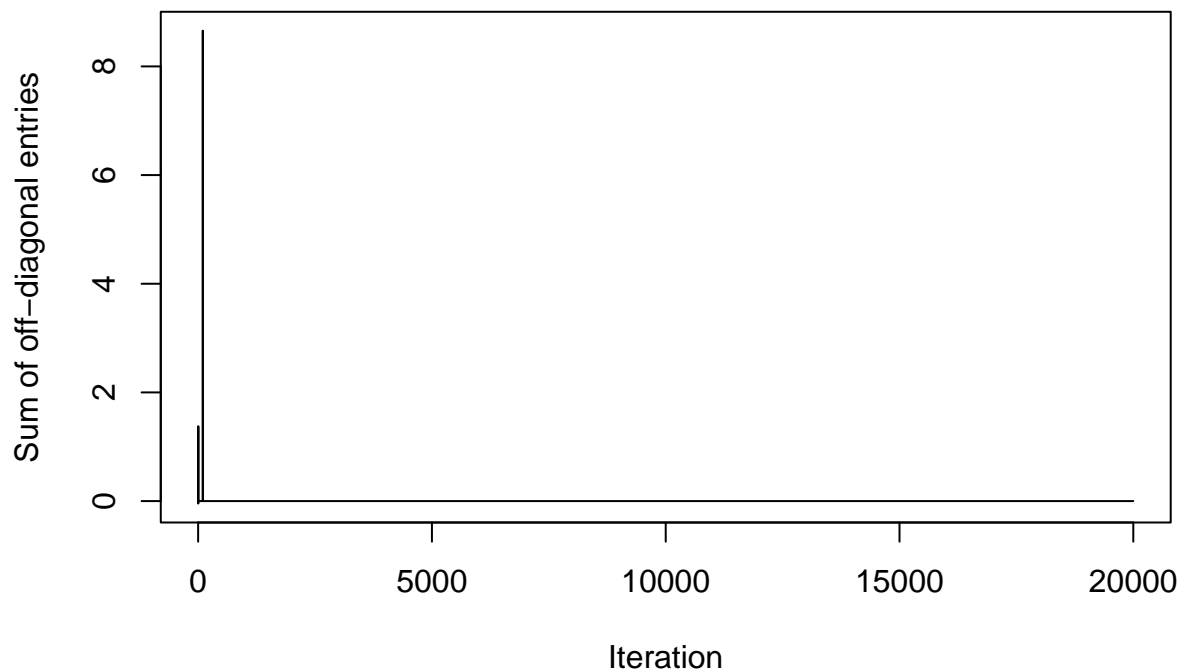
## [1] 12.9001

```

```

plot(1:max_iter, off_diag_entries, type = "l", xlab = "Iteration", ylab = "Sum of off-diagonal entries").

```



Question Arnoldi Iteration 1

Generate a 10x10 full rank matrix and apply Arnoldi to get a Hessenberg form, then do the same for rank-5 matrix.

```
A_mat <- matrix(rnorm(10^2),nrow=10)
A <- A_mat + t(A_mat)
n <- nrow(A)
b <- rnorm(n)
```

```
q <- matrix(0, nrow=n+1, ncol=n)
H <- matrix(0, nrow=n+1, ncol=n)
q[1,] <- b / sqrt(sum(b^2))
for(i in 1:10) {
  v <- A %*% q[i,]
  for(j in 1:i) {
    H[i,j] <- q[i,] %*% v
    v <- v - H[i,j] * q[i,]
  }
  H[i+1,i] <- sqrt(sum(v^2))
  if(H[i+1,i] < 1e-8) {
    break
  }
  q[i+1,] <- v / H[i+1,i]
}
list(q,round(H))
```

```
## [[1]]
##           [,1]      [,2]      [,3]      [,4]      [,5]      [,6]
## [1,] 0.1564715 -0.29762095 -0.46668107 -0.342079840 -0.29162731 0.142388468
## [2,] 0.2858651 0.52046863 -0.08676443 0.130188640 -0.65273125 0.056860844
## [3,] 0.3813461 -0.05802480 -0.51642973 0.174159072 0.07536823 -0.150745232
## [4,] 0.4345946 0.47164651 -0.13378407 0.136311573 -0.41515256 -0.008200014
```

```
## [5,] 0.4762687 -0.05183329 -0.51241772 0.135392115 0.26372442 -0.019014643
## [6,] 0.4943464 0.40437640 -0.10657312 0.072733776 -0.27386669 0.119442503
## [7,] 0.5022691 -0.10410033 -0.46341977 0.081245517 0.38224653 0.099528150
## [8,] 0.4982643 0.35371945 -0.06176727 0.026110271 -0.19292612 0.200464667
## [9,] 0.4968893 -0.13885353 -0.42657879 0.048641673 0.43909467 0.155859128
## [10,] 0.4911798 0.32873718 -0.03605212 0.003658302 -0.15755131 0.235216888
## [11,] 0.4905672 -0.15496747 -0.40860428 0.034009830 0.46249755 0.178693032
##      [,7]      [,8]      [,9]     [,10]
## [1,] 0.13758753 0.623375495 0.1463327 0.133667674
## [2,] -0.10378242 0.006487267 -0.1351589 -0.405640662
## [3,] 0.21931864 0.593431788 0.3499784 0.055008488
## [4,] -0.09351794 -0.121436568 -0.2231949 -0.553598037
## [5,] 0.27057252 0.479917439 0.3407860 0.001031574
## [6,] -0.04987018 -0.241711314 -0.2180654 -0.614589240
## [7,] 0.30289760 0.373362831 0.3574915 -0.023769023
## [8,] -0.03134473 -0.325012056 -0.2083504 -0.628240793
## [9,] 0.31403522 0.315343696 0.3656350 -0.025303938
## [10,] -0.02535023 -0.363699053 -0.2042931 -0.628137276
## [11,] 0.31758420 0.290075528 0.3685381 -0.023574025
##
## [[2]]
##      [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9] [,10]
## [1,] -1    0    0    0    0    0    0    0    0    0
## [2,] 0     0    0    0    0    0    0    0    0    0
## [3,] 0     0    0    0    0    0    0    0    0    0
## [4,] 0     0    0    0    0    0    0    0    0    0
## [5,] 1     0    0    0    0    0    0    0    0    0
## [6,] 1     0    0    0    0    0    0    0    0    0
## [7,] 1     0    0    0    0    0    0    0    0    0
## [8,] 1     0    0    0    0    0    0    0    0    0
## [9,] 1     0    0    0    0    0    0    0    0    0
## [10,] 1    0    0    0    0    0    0    0    0    0
## [11,] 0     0    0    0    0    0    0    0    0    8
```

```
# rank 5 10x10 matrix
A_mat <- matrix(rnorm(10^2),nrow=10)
A <- A_mat + t(A_mat)
A <- qr.Q(qr(A)) %*% diag(c(1,2,3,4,5,0,0,0,0,0)) %*% qr.Q(qr(A))
n <- nrow(A)
b <- rnorm(n)
q <- matrix(0, nrow=n+1, ncol=n)
H <- matrix(0, nrow=n+1, ncol=n)
q[1,] <- b / sqrt(sum(b^2))
for(i in 1:10) {
  v <- A %*% q[i,]
  for(j in 1:i) {
    H[i,j] <- q[i,] %*% v
    v <- v - H[i,j] * q[i,]
  }
  H[i+1,i] <- sqrt(sum(v^2))
  if(H[i+1,i] < 1e-8) {
    break
  }
  q[i+1,] <- v / H[i+1,i]
```



```
}
list(q,round(H))
```

```
## [[1]]
##           [,1]      [,2]      [,3]      [,4]      [,5]      [,6]
## [1,] -0.09188447 -0.2511217 -0.11264649 -0.49873883  0.62084585 -0.05759405
## [2,] -0.41023418 -0.4428584  0.09142340 -0.20504634 -0.26719933 -0.28913119
## [3,]  0.01604561  0.2490338  0.02191136  0.46534806 -0.25786843 -0.33653186
## [4,]  0.01766526  0.0377262  0.20650397 -0.52130233 -0.02292175 -0.07711026
## [5,] -0.26217332 -0.5671032  0.11196616  0.09356634  0.47012998 -0.36575724
## [6,]  0.23995884  0.1820970 -0.15184010  0.33091210  0.43885237  0.07137299
## [7,]  0.16469198  0.2166528 -0.15403788 -0.38342608 -0.12860805  0.44401797
## [8,] -0.29505352 -0.2507315  0.10759943 -0.15738350 -0.52469114 -0.09982841
## [9,]  0.05248942  0.3152221  0.06267945  0.64179051 -0.46780544 -0.24777713
## [10,]  0.20109099  0.3106947  0.18805989 -0.37185664 -0.03829614  0.02478890
## [11,] -0.27816476 -0.5763093  0.16311920 -0.13597279  0.42176994 -0.33088609
##           [,7]      [,8]      [,9]      [,10]
## [1,]  0.24046903 -0.23547154 -0.320178712 -0.2500244
## [2,]  0.25175377  0.39025216 -0.153078600  0.4371270
## [3,]  0.29460957 -0.41084346 -0.521599950  0.1152924
## [4,]  0.76156959 -0.08485718  0.174283837  0.2446006
## [5,] -0.07781600 -0.19855341 -0.137412106  0.4113473
## [6,] -0.36664715 -0.45657982 -0.222333929 -0.4320753
## [7,] -0.15725899  0.41835745  0.226193538 -0.5390603
## [8,]  0.13795005  0.58031214  0.125646231  0.3960859
## [9,]  0.09816882 -0.25873140 -0.272495231  0.2258052
## [10,]  0.73378043 -0.27418425  0.239458745  0.1278868
## [11,]  0.12077728 -0.11039981 -0.009485581  0.4807855
##
## [[2]]
##           [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9] [,10]
## [1,]      0    0    0    0    0    0    0    0    0    0
## [2,]     -1    0    0    0    0    0    0    0    0    0
## [3,]     -1    0    0    0    0    0    0    0    0    0
## [4,]      0    0    0    0    0    0    0    0    0    0
## [5,]     -1    0    0    0    0    0    0    0    0    0
## [6,]     -1    0    0    0    0    0    0    0    0    0
## [7,]     -1    0    0    0    0    0    0    0    0    0
## [8,]     -1    0    0    0    0    0    0    0    0    0
## [9,]     -1    0    0    0    0    0    0    0    0    0
## [10,]    -1    0    0    0    0    0    0    0    0    0
## [11,]      0    0    0    0    0    0    0    0    0    3
```

Lecture 3/11: GMRES

The method of generalized minimal residuals (GMRES) can be used for solving eigenvalue problems.

GMRES

Recall that the n th Krylov subspace $\mathbb{K}_n(A, b) = \text{span}(b, Ab, \dots, A^{n-1}b)$ is the subspace spanned by the first n powers of the matrix A applied to the vector b .

GMRES is an iterative method that finds a vector $x_n \in \mathbb{K}_n$ that minimizes the L2 norm of the residual

$r_n = b - Ax_n$. More generally, consider AK_n ,

$$AK_n = [Ab \quad A^2b \quad \dots \quad A^n b]$$

and look for c s.t. $\|AK_n c - b\|$ is minimized. Then let $x_n = K_n c$.

This is not numerically stable in general, as the columns of AK_n are likely to be correlated due to the fact that we're approaching the eigenvalues of A in later columns! Instead, use a sequence of matrices Q_n whose columns span \mathbb{K}^n . Then we want to find y such that $\|AQ_n y - b\|$ is minimized, and let $x_n = Q_n y$.

Lecture 3/13: GMRES

Let A be a real symmetric positive definite matrix. Our goal is to solve $Ax = b$ for solution x_* .

Consider the quadratic form $f(x) = \frac{1}{2}x^T Ax - b^T x + c$. It will turn out that this quadratic form is minimized at the solution x_* .

Note that thanks to the positive definiteness of A that the image of $f(x)$ has a nice geometry - it is a paraboloid bowl (obviously convex). Moreover, take $f'(x) = Ax - b$, which we will call the *residual*.

Let $e_i = x_i - x_*$ be the error of each iteration and $r_i = b - Ax_i = -f'(x_i)$ be the residual.

Steepest Descent

With the *steepest descent* method, we start at an arbitrary guess x_0 and take a series of steps x_1, x_2, \dots such that each step is in the direction of the negative gradient of $f(x)$ at the current point. We terminate when we're close enough to a solution, or when the difference between successive steps is small enough.

The steepest descent method is given by the iteration

$$x_{i+1} = x_i - \alpha_i f'(x_i)$$

where α_i is the step size.

We choose α_i via *line search* which yields $\alpha_i = \arg \min_{\alpha} f(x_i - \alpha f'(x_i)) = \frac{r_i^T r_i}{r_i^T A r_i}$. This is the optimal step size in the sense that it minimizes $f(x_{i+1})$.

However, steepest descent can be inefficient because it doesn't take into account all the geometry. We can make a transformation of the space to get closer to the minimum in each iteration.

Conjugate Gradient Method

Let two vectors u, v be *conjugate* if $\langle u, Av \rangle = \langle v, Au \rangle = 0$. Since A is symmetric and pos. def., this defines an inner product $\langle u, v \rangle_A = u^T A v$. Thus if $\langle u, v \rangle_A = 0$, then u and v are *A-orthogonal*.

Let $A \in \mathbb{R}^{n \times n}$, $P = \{p_1, \dots, p_n\}$ be a set of *A-orthogonal* vectors. P is a basis for \mathbb{R}^n and moreover we have

$$\begin{aligned} x_* &= \sum_{i=1}^n \alpha_i p_i \\ Ax_* &= \sum_{i=1}^n \alpha_i A p_i \\ p_k^T Ax_* &= \sum_{i=1}^n \alpha_i p_k^T A p_i = \sum_{i=1}^n \alpha_i \langle p_k, p_i \rangle_A \\ p_k^T b &= \alpha_k \langle p_k, p_k \rangle_A \end{aligned}$$

so that $\alpha_k = \frac{p_k^T b}{p_k^T A p_k}$.

This motivates the *conjugate gradient method*: find a sequence of conjugate vectors (via something like Gram-Schmidt), then compute the coefficients α_i and update the solution in n steps.

The conjugate gradient method is given by the iteration

```
x_0 = rand(n)
r_0 = b - Ax_0
p_0 = r_0
for i in 1:n {
    alpha_i = (r_i^T r_i) / (p_i^T A p_i)
    x_{i+1} = x_i + alpha_i p_i
    r_{i+1} = r_i - alpha_i A p_i
    beta_{i+1} = (r_{i+1}^T r_{i+1}) / (r_i^T r_i)
    p_{i+1} = r_{i+1} + beta_{i+1} p_i
}
```

The conjugate gradient method is guaranteed to converge in n steps if A is symmetric and positive definite.