# VE472 Lecture 3

Jing Liu

UM-SJTU Joint Institute

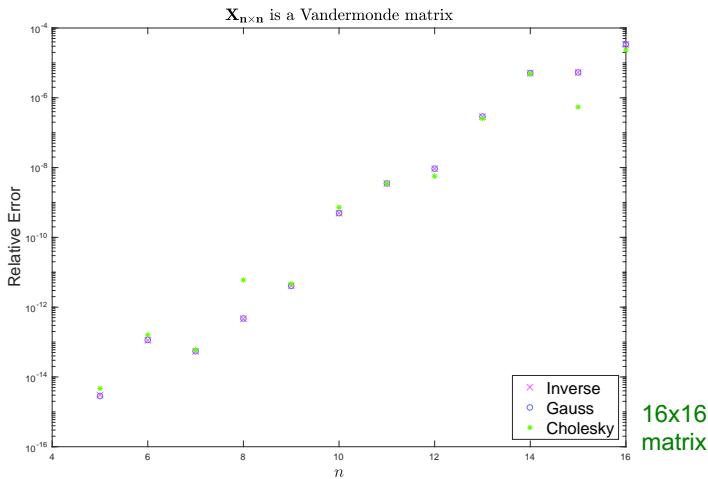Summer

# Random full rank matrix



$\mathbf{X_{n \times 5}}$ is a random matrix of uniform(0,1)

# Vandermonde matrix



$\mathbf{X_{n \times n}}$ is a Vandermonde matrix

16x16 matrix

r o u g h l y   s h a r e   t h e   s a m e   a c c u r a c y.   B u t   t h e y   l o o s e   pr

**Theorem 0.1**

**QR decomposition**

*Let $\mathbf{A}$ be a matrix of $m \times n$ with $m \geq n$. Suppose $\mathbf{A}$ is full rank. Then there is a matrix $\mathbf{Q}$ of $m \times n$ with orthonormal columns, $\mathbf{Q}^T\mathbf{Q} = \mathbf{I}$, and an upper triangular matrix $\mathbf{R}$ of $n \times n$ with positive diagonals $r_{ii} > 0$ such that $\mathbf{A} = \mathbf{QR}$.*
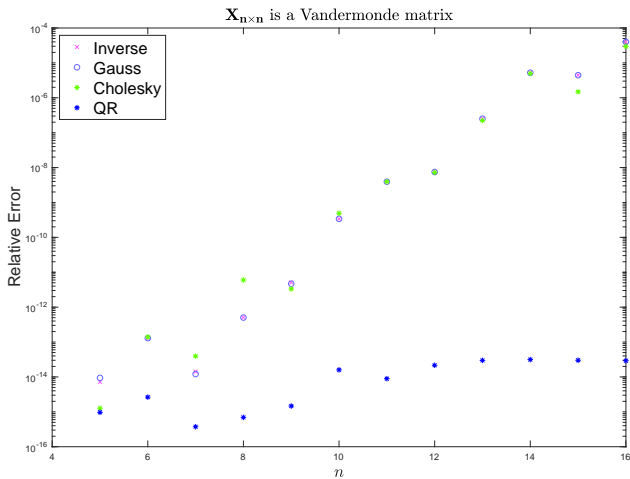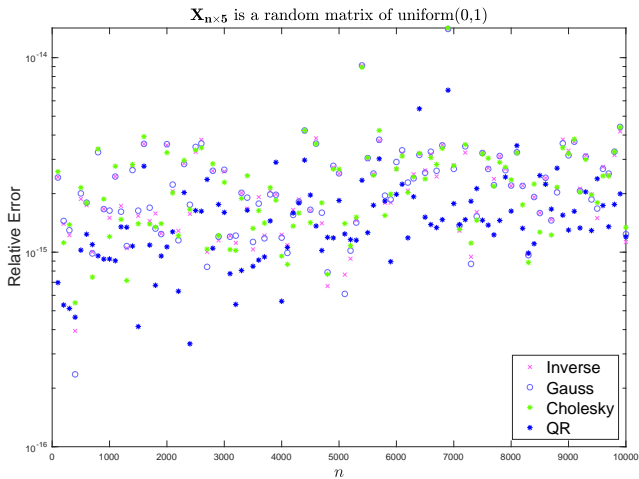


In general, the theorem replies on the idea of projection.

Q: Why is this theorem useful in terms of dealing with a big data matrix $\mathbf{X}$?
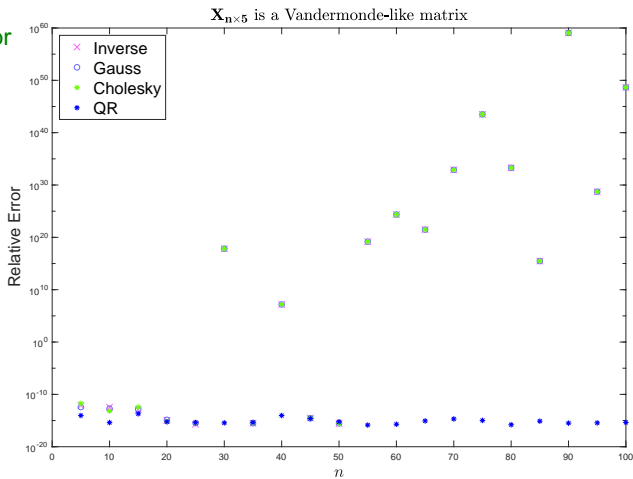
see notes

# Vandermonde matrix



$\mathbf{X_{n \times n}}$ is a Vandermonde matrix

# Random full rank matrix



$\mathbf{X_{n \times 5}}$ is a random matrix of uniform(0,1)

# Vandermonde-like matrix

- QR is no question the slowest!

```
>> clear all
>> n = 1000; X = rand(n, 100); y = randn(n, 1);
tic; for i = 1:10000
    XX = transpose(X)*X; yy = transpose(X)*y;
    [L, U] = lu(XX); bhat = U\(L\yy);
end; toc;
tic; for i = 1:10000
    XX = transpose(X)*X; yy = transpose(X)*y;
    C = chol(XX, 'lower');
    bhat =  transpose(C)\(C\yy);
end; toc;
tic ;for i = 1:10000
    [Q, R] = qr(X); bhat =  R\(transpose(Q)*y);
end; toc;
```

```
Elapsed time is 3.266948 seconds.
Elapsed time is 2.457167 seconds.
Elapsed time is 44.321296 seconds.      stable but slow
```

Theorem 0.2 (Singular Value Decomposition)

*Let $\mathbf{A}$ be a rank $k$ matrix of $m \times n$ with $m \geq n$, then we have $\sigma_1 \geq \cdots \geq \sigma_k > 0$*

k linearly independent columns

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^{\mathrm{T}}$$

eigenvectors of X^TX

$$= \left[\begin{array}{ccc|ccc} \mathbf{u}_1 & \cdots & \mathbf{u}_k & \mathbf{u}_{k+1} & \cdots & \mathbf{u}_m \end{array}\right] \left[\begin{array}{ccc|c} \sigma_1 & \cdots & 0 & \\ \vdots & \ddots & \vdots & \mathbf{0}_{k \times (n-k)} \\ 0 & \cdots & \sigma_k & \\ \hline & \mathbf{0}_{(m-k) \times k} & & \mathbf{0}_{(m-k) \times (n-k)} \end{array}\right] \left[\begin{array}{c} \mathbf{v}_1^{\mathrm{T}} \\ \vdots \\ \mathbf{v}_k^{\mathrm{T}} \\ \hline \mathbf{v}_{k+1}^{\mathrm{T}} \\ \vdots \\ \mathbf{v}_n^{\mathrm{T}} \end{array}\right]$$

*where $\mathbf{U}$ and $\mathbf{V}$ are orthogonal matrices of size $m \times m$ and $n \times n$, respectively.*

Q: Why is this theorem useful in terms of dealing with a big data matrix $\mathbf{X}$?

- SVD is faster than QR, but a magnitude slower than LU or Cholesky .

```
>> n = 1000; X = rand(n, 100); y = randn(n, 1);
tic;for i = 1:10000
    XX = transpose(X)*X; yy = transpose(X)*y;
    C = chol(XX, 'lower');
    bhat = transpose(C)\(C\yy);
end; toc;
tic;for i = 1:10000
    [Q, R] = qr(X); bhat = R\(transpose(Q)*y);
end; toc;
tic;for i = 1:10000
    [U, S, V] = svd(X, 'econ');    compact svd
    s = diag(S); s = 1./s;  inverse of S
    bhat = V*diag(s)*transpose(U)*y;
end; toc;
```

```
Elapsed time is 2.521359 seconds.
Elapsed time is 45.891408 seconds.
Elapsed time is 19.551785 seconds.    faster than QR, slower
                                       than LU
```
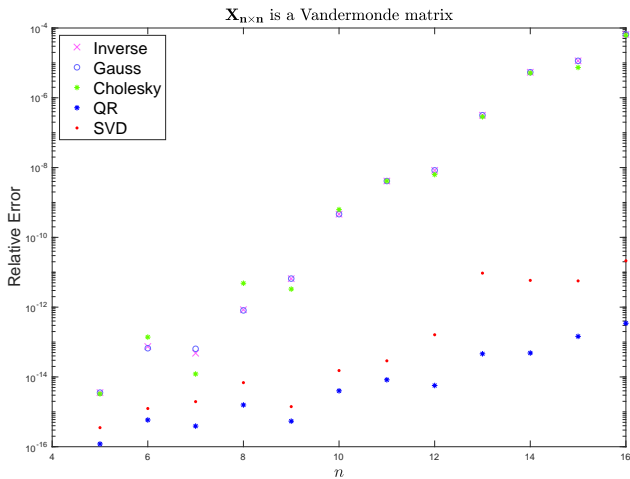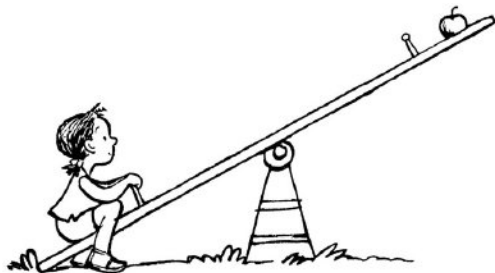
# Vandermonde matrix



$\mathbf{X_{n \times n}}$ is a Vandermonde matrix

- Recall the stability of linear model can be studied using the eigenvalues of

$$\mathbf{X}^{\mathrm{T}}\mathbf{X}$$

  having small eigenvalues indicate columns are nearly linearly dependent.



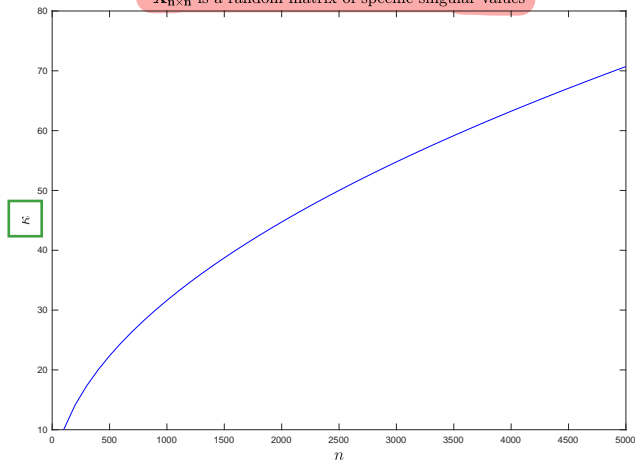- Using SVD gives us the stability of our model given the data as a byproduct.

$$\kappa = \frac{\max \sigma_i}{\min \sigma_i}$$ samllest non zero sv

# Increasing $\kappa$ as $n$ increases by construction



true random

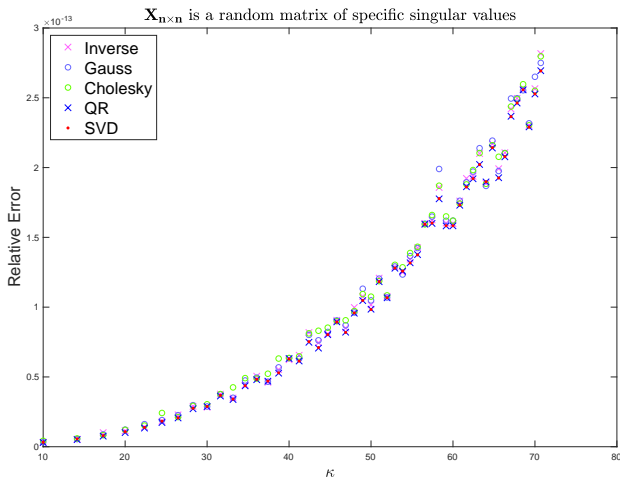$X_{n \times n}$ is a random matrix of specific singular values

# Increasing relative error

虽然在增大但
是k<100依然
可以接受

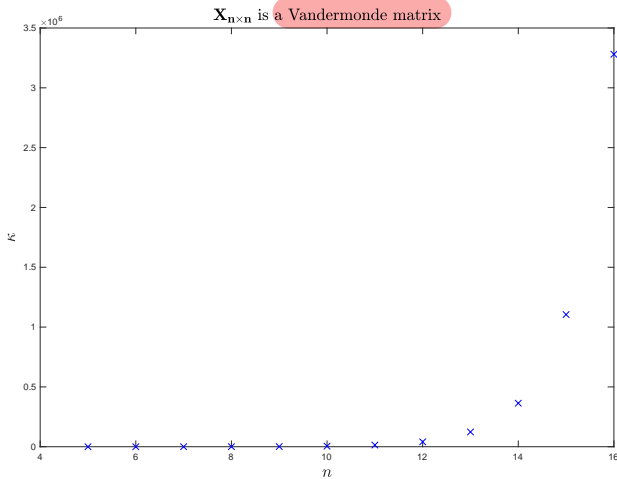

$\mathbf{X_{n \times n}}$ is a random matrix of specific singular values

(Legend: Inverse, Gauss, Cholesky, QR, SVD)

Relative Error vs $\kappa$

# Really big $\kappa$

$\mathbf{X_{n \times n}}$ is a Vandermonde matrix

# Really big $\kappa$ really big relative error



$\mathbf{X_{n \times n}}$ is a Vandermonde matrix

# Decreasing $\kappa$ as $n$ increases for a fixed $k$

fix number of columns, k
decreses with incresing n



$\mathbf{X_{n \times 5}}$ is a random matrix of uniform(0,1)

# Only having really big $\kappa$ is problematic



$\mathbf{X_{n \times 5}}$ is a random matrix of uniform(0,1)

small k

$\mathbf{X}_{5000 \times k}$ is a random matrix of uniform(0,1)

big column number

# A more complex data is more problematic than a large data



$X_{5000 \times k}$ is a random matrix of uniform(0,1)

- So far we have considered the data matrix $\mathbf{X}$ that is full rank when solving

$$\underset{\mathbf{b} \in \mathbb{R}^{k+1}}{\arg\min} \|\mathbf{y} - \mathbf{X}\mathbf{b}\|^2$$

- Q: What happens when $\mathbf{A}$ is rank deficient or "close" to being rank deficient?

- Such problems often arise with big data , e.g. extracting signals from noisy data, digital image restoration as well as big prediction or classification.

### Theorem 0.3

Let $\sigma_{min} > 0$ denote the smallest singular value of $\mathbf{X}$ and $\mathbf{U\Sigma V}^{\mathrm{T}}$ be the singular value decomposition of $\mathbf{X}$. If $\mathbf{b}$ minimises $\|\mathbf{y} - \mathbf{X}\mathbf{b}\|$ , then

can have large norm
change 'y' a bit, 'b' may change a lot $\qquad \|\mathbf{b}\| \geq \dfrac{|\mathbf{u}^{\mathrm{T}}\mathbf{b}|}{\sigma_{min}}$ ← u^Ty

where $\mathbf{u}$ is the last column of $\mathbf{U}$. Furthermore changing $\mathbf{y}$ to $\mathbf{y} + \delta\mathbf{y}$ can induce a change of $\delta\mathbf{b}$ to $\mathbf{b}$, where $\|\delta\mathbf{b}\|$ is as large as $\|\delta\mathbf{y}\| / \sigma_{min}$.

Q: What does the last theorem tell us if $\mathbf{A}$ is nearly rank deficient?

### Theorem 0.4

Let $\mathbf{X}$ be a rank $r$ matrix of $n \times (k+1)$ with $n \geq k+1$. If $r < k+1$, then there is an ~~$n-r$~~ k+1-r dimensional set of vectors $\mathbf{b} \in \mathbb{R}^{k+1}$ that solve the following

$$\underset{\mathbf{b} \in \mathbb{R}^{k+1}}{\arg\min} \|\mathbf{y} - \mathbf{X}\mathbf{b}\|^2 \qquad where \quad \mathbf{y} \in \mathbb{R}^n$$

Furthermore, the singular value decomposition of $\mathbf{X}$ can be written as

$$\mathbf{A} = \begin{bmatrix} \mathbf{U}_1 & \mathbf{U}_2 \end{bmatrix} \begin{bmatrix} \mathbf{\Sigma}_{r \times r} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{V}_1 & \mathbf{V}_2 \end{bmatrix}^{\mathrm{T}}, \quad where \ \mathbf{U}_1 \ and \ \mathbf{V}_1 \ have \ r \ columns,$$

and all the minimisers take the following form

$$\mathbf{b} = \mathbf{V}_1 \mathbf{\Sigma}_{r \times r}^{-1} \mathbf{U}_1^{\mathrm{T}} \mathbf{y} + \mathbf{V}_2 \mathbf{z}, \qquad for \ any \ \mathbf{z} \in \mathbb{R}^r.$$

Q: Which of those minimisers shall we use as the best $\mathbf{b}$?