- ▶ Numbers on the real line $\mathbb{R}$ cannot all be represented exactly in a computer.
- ▶ Usually, we use *floating point storage*, which uses a fixed finite amount of bit $(0/1)$ storage to represent a subset of the real numbers.

## Double precision floating point numbers (Standard: IEEE 754)

*Double precision* means that we have 64 bits $(0/1)$ at our disposal, split into

| Bits | Interpretation |
|------|----------------|
| 1 | $\text{sign}(x) \in \{-1, +1\}$ |
| 11 | $\text{exponent}(x) \in [-1022, +1023]$ |
| 52 | $\text{fraction}(x) \in [0, 1)$ |

These components combine to define the *computed* version of a number $x$ as

$$\text{comp}(x) = \text{sign} \cdot (1 + \text{fraction}) \cdot 2^{\text{exponent}},$$

and $|x| < 2^{-1022}$ is treated specially.

## Relative storage error

The closest floating point number $\mathrm{comp}(x)$ to a real value $x \not\approx 0$ has bounded *relative error*:

$$\frac{|\mathrm{comp}(x) - x|}{|x|} \lesssim \epsilon_0,$$

where $\epsilon_0$ is called the *machine precision*.

### Machine precision

The machine precision is the smallest number such that $\mathrm{comp}(1 + \epsilon_0) > 1$.
In R the value can be accessed as `.Machine$double.eps`, and is $\epsilon_0 \approx 2.220446 \cdot 10^{-16}$.

This allows a simple model for handling numerical approximation errors:

$$\mathrm{comp}(x) = (1 + \epsilon)x,$$

for some $\epsilon$ such that $|\epsilon| \lesssim \epsilon_0$.

For simple storage of numbers, the error is quite small, but we need to beware of amplifying the error when we perform computations!

# Approximation errors for numerical derivatives

## Finite differences

The derivative of a function $f(\theta)$ can be approximated by the asymmetric finite difference

$$\text{comp}\left\{ \frac{\text{comp}[f(\text{comp}[\theta + h])] - \text{comp}[f(\theta)]}{h} \right\} \approx \frac{f(\theta + h) - f(\theta)}{h} \approx f'(\theta).$$

How small should we choose $h$ in order to minimise the approximation error?

Local polynomial approximations will help in analysing the errors:

## Taylor's theorem

For a function $f : \mathbb{R} \mapsto \mathbb{R}$ with two continuous derivatives near a $\theta \in \mathbb{R}$,

$$f(\theta + h) = f(\theta) + hf'(\theta + t_1 h), \text{ and}$$
$$f(\theta + h) = f(\theta) + hf'(\theta) + \frac{h^2}{2}f''(\theta + t_2 h),$$

for some $t_1 \in [0, 1]$ and $t_2 \in [0, 1]$.

# First order difference approximation error

Assume that $\theta$ is stored exactly.

First, consider the error in computing $f(\theta + h)$:

$$
\begin{aligned}
\text{comp}\{f(\text{comp}[\theta + h])\} &= (1 + \epsilon_f)f([1 + \epsilon_\theta][\theta + h]) \\
&\approx (1 + \epsilon_f)\left[f(\theta + h) + \epsilon_\theta(\theta + h)f'(\theta + t_1 h)\right] \\
&\approx f(\theta + h) + \epsilon_f f(\theta + h) + \epsilon_\theta \theta f'(\theta + t_1 h) + \text{h.o.t.},
\end{aligned}
$$

for some $|\epsilon_f| \lesssim \epsilon_0$ and $|\epsilon_\theta| \lesssim \epsilon_0$.

Assume that $f(\cdot)$ is bounded near $\theta$, with $|f(\cdot)| \leq L_0$, $|f'(\cdot)| \leq L_1$, and $|f''(\cdot)| \leq L_2$. Then

$$
\left| \frac{\text{comp}\{f(\text{comp}[\theta + h])\} - \text{comp}\{f(\theta)\}}{h} - \frac{f(\theta + h) - f(\theta)}{h} \right| \lesssim \frac{\epsilon_0(2L_0 + |\theta|L_1)}{h}.
$$

# First order difference approximation error (cont.)

Applying Taylor's theorem to the exact finite difference leads to

$$\left| \frac{f(\theta + h) - f(\theta)}{h} - f'(\theta) \right| = \frac{h}{2} \left| f''(\theta + t_2 h) \right| \leq \frac{h L_2}{2}.$$

Using the triangle inequality ($|E_1 + E_2| \leq |E_1| + |E_2|$), we get the total error bound,

$$\left| \frac{\text{comp}\{f(\text{comp}[\theta + h])\} - \text{comp}\{f(\theta)\}}{h} - f'(\theta) \right| \lesssim \frac{\epsilon_0 (2L_0 + |\theta| L_1)}{h} + \frac{h L_2}{2},$$

which is minimised for $h = \sqrt{\epsilon_0 \frac{4L_0 + 2|\theta| L_1}{L_2}} \propto \epsilon_0^{1/2}$.

Note:
If $\theta$ wasn't stored exactly, the additional error term would be

$$|f'(\text{comp}[\theta]) - f'(\theta)| \approx |f'(\theta) + \epsilon_3 \theta f''(\theta + t_3 h) - f'(\theta)| \lesssim \epsilon_0 |\theta| L_2,$$

which doesn't depend on $h$.

# Optimal stepsize for finite differences

Let $L_k$ be bounds for the $k$:th derivatives around $\theta$. The errors from *floating point cancellation* and *Taylor series truncation* can be bounded and minimised by choosing the *step size* $h$:

▶ Asymmetric first order differences for $f'(\theta)$, using $f(\theta)$ and $f(\theta + h)$, gives the bound

$$\lesssim \frac{\epsilon_0(2L_0 + |\theta|L_1)}{h} + \frac{hL_2}{2}, \quad \text{which is minimised for } h = \sqrt{2\epsilon_0 \frac{2L_0 + |\theta|L_1}{L_2}} \sim \epsilon_0^{1/2}.$$

▶ Symmetric first order differences for $f'(\theta)$, using $f(\theta - h)$ and $f(\theta + h)$, gives the bound
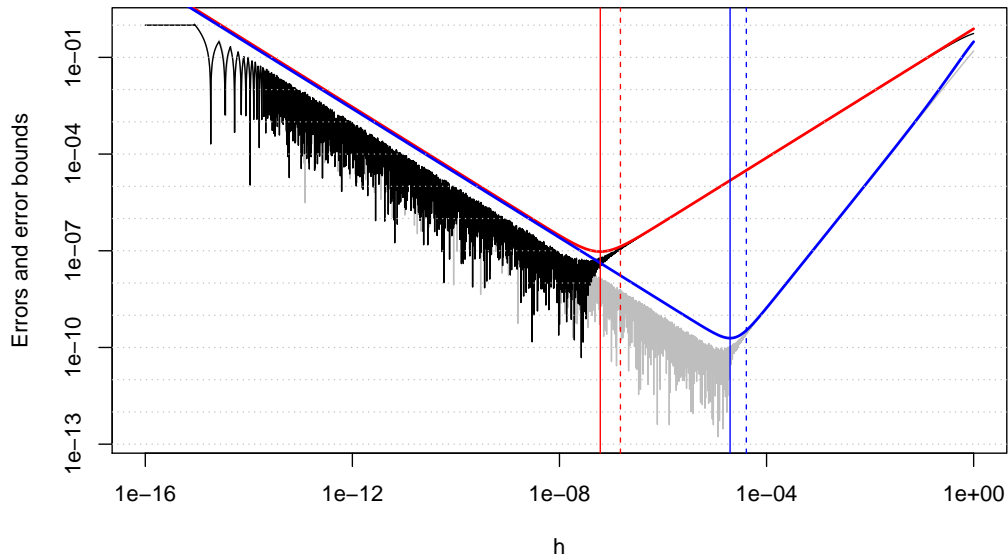
$$\lesssim \frac{\epsilon_0(L_0 + |\theta|L_1)}{h} + \frac{h^2 L_3}{6}, \quad \text{which is minimised for } h = \left(3\epsilon_0 \frac{L_0 + |\theta|L_1}{L_3}\right)^{1/3} \sim \epsilon_0^{1/3}.$$

▶ 2nd order differences for $f''(\theta)$, using $f(\theta - h)$, $f(\theta + h)$, and $f(\theta)$, gives the bound
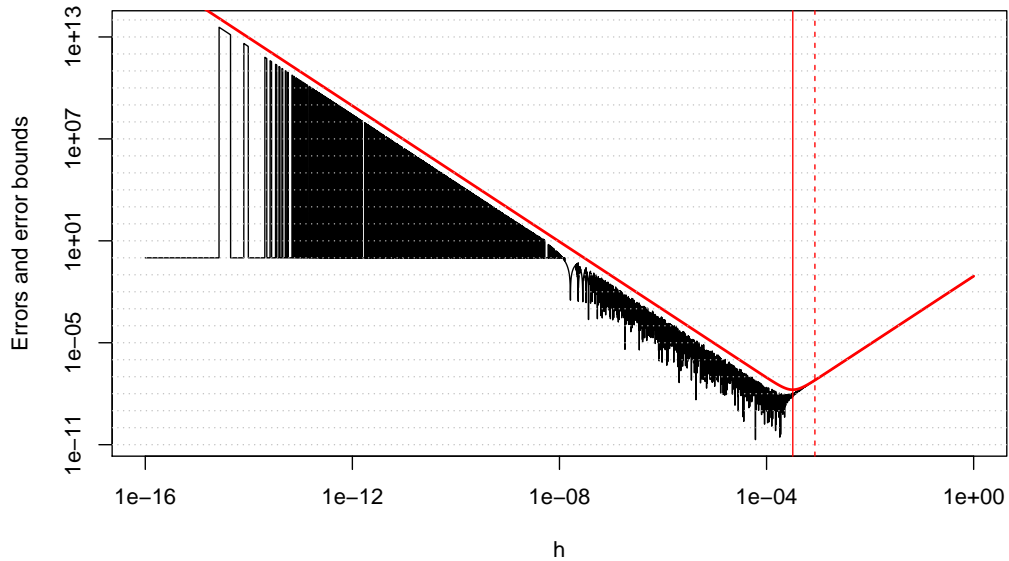
$$\lesssim \frac{\epsilon_0(4L_0 + 2|\theta|L_1)}{h^2} + \frac{h^2 L_4}{12}, \quad \text{which is minimised for } h = \left(24\epsilon_0 \frac{2L_0 + |\theta|L_1}{L_4}\right)^{1/4} \sim \epsilon_0^{1/4}.$$

▶ Approximate rule of thumb: Plugin $L_k \equiv 1$ and a representative $|\theta|$ value.

**Asymmetric (black/red), and symmetric (grey/blue) 1st order difference errors**

**2nd order difference errors**

Errors and error bounds

h

# Numerics for Least Squares estimation of linear models

- In matrix/vector form, we write linear models for
  - observations $\boldsymbol{y} = \left[y_1, \ldots, y_n\right]^\top : n \times 1$,
  - explanatory/predictor covariates $\boldsymbol{X} = \left[\boldsymbol{x}_1, \ldots, \boldsymbol{x}_p\right] : n \times p,\ n > p$,
  - parameters $\boldsymbol{\beta} = \left[\beta_1, \ldots, \beta_p\right]^\top : p \times 1$, and
  - observation noise $\boldsymbol{e} = \left[e_1, \ldots, e_n\right]^\top : n \times 1$,

  as

  $$\boldsymbol{y} = \sum_{k=1}^{p} \boldsymbol{x}_k \beta_k + \boldsymbol{e} = \boldsymbol{X}\boldsymbol{\beta} + \boldsymbol{e}.$$

- The least squares estimate of $\boldsymbol{\beta}$ is *in theory* given by $\widehat{\boldsymbol{\beta}} = (\boldsymbol{X}^\top \boldsymbol{X})^{-1} \boldsymbol{X}^\top \boldsymbol{y}$.
- We will first analyse the floating point errors involved in $\boldsymbol{X}^\top \operatorname{comp}(\boldsymbol{y})$ and $(\boldsymbol{X}^\top \boldsymbol{X})^{-1} \operatorname{comp}(\boldsymbol{X}^\top \boldsymbol{y})$.
- Then we will develop a method for computing $\widehat{\boldsymbol{\beta}}$ that does *not* involve inverting any matrices, and avoids unneccesarily amplifying the floating point errors.
- We need some analytical tools!

## Singular Value Decomposition (SVD, `svd(A)` in R)

For any rectangular matrix $\boldsymbol{A} \neq \boldsymbol{0}$, $n \times m$, there exist matrices $\boldsymbol{U} : n \times p$, $\boldsymbol{V} : m \times p$, and $\boldsymbol{D} = \mathrm{diag}\{d_1, \ldots, d_p\}$ such that

$\boldsymbol{U}^\top \boldsymbol{U} = \boldsymbol{I}_p$, i.e. the column vectors in $\boldsymbol{U} = \begin{bmatrix} \boldsymbol{u}_1, \ldots, \boldsymbol{u}_p \end{bmatrix}$ are orthonormal,

$\boldsymbol{V}^\top \boldsymbol{V} = \boldsymbol{I}_p$, i.e. the column vectors in $\boldsymbol{V} = \begin{bmatrix} \boldsymbol{v}_1, \ldots, \boldsymbol{v}_p \end{bmatrix}$ are orthonormal, and

$$d_1 \geq d_2 \geq \cdots \geq d_r > 0 = d_{r+1} = \cdots = d_p,$$

for some $r \leq p = \min(n, m)$, and

$$\boldsymbol{A} = \boldsymbol{U}\boldsymbol{D}\boldsymbol{V}^\top = \sum_{k=1}^{r} d_k \boldsymbol{u}_k \boldsymbol{v}_k^\top.$$

The value $r$ is the *rank* of $A$. A matrix with $\mathrm{rank}(A) = p$ is said to have *full rank*.

The column vectors of $\boldsymbol{U}$ and $\boldsymbol{V}$ are called the left and right *singular vectors* of $\boldsymbol{A}$, and the corresponding $d_1, \ldots, d_p$ are the *singular values* of $\boldsymbol{A}$.

The SVD algorithm is also used for *dimension reduction* in *Principal Component Analysis*.

The *complete* SVD expands $U$ and $V$ to square and completely orthogonal matrices $\begin{bmatrix} U & U_0 \end{bmatrix}$ and $\begin{bmatrix} V & V_0 \end{bmatrix}$ as needed, and expands $D$ with zeros.

For a *tall* matrix with full rank, $n > m = p = r$. The complete decomposition is

$$A = \begin{bmatrix} U & U_0 \end{bmatrix} \begin{bmatrix} D \\ 0 \end{bmatrix} V^\top$$

and

$$U^\top U = I_p, \qquad U_0^\top U_0 = I_{n-p}, \qquad U^\top U_0 = 0,$$
$$UU^\top + U_0 U_0^\top = I_n,$$
$$V^\top V = VV^\top = I_m$$

These identities allow us to split a vector $w \in \mathbb{R}^n$ into a component $w_A$ that lies within the vector space spanned by the columns of $A$, and a component $w_0$ orthogonal to $A$:

$$w = UU^\top w + U_0 U_0^\top w = w_A + w_0, \qquad A^\top w_A = A^\top w, \qquad A^\top w_0 = 0.$$

## Vector length amplification

Multiplying with $A^\top$ leads to different vector length amplification depending on in which direction $w$ is pointing:

$$\|A^\top w\| = \sqrt{w^\top A A^\top w} = \sqrt{w^\top U D V^\top V D U^\top w} = \sqrt{w^\top U D^2 U^\top w}$$
$$= \|D U^\top w\| = \|D U^\top w_A\|$$

The norm is minimised when $w_A = u_p \|w_A\|$ and maximised when $w_A = u_1 \|w_A\|$, so that

$$d_p \|w_A\| \le \|A^\top w\| \le d_1 \|w_A\|.$$

Each component of $\mathrm{comp}(w)$ has a relative error $\epsilon_i$, $|\epsilon_i| \lesssim \epsilon_0$.

$$\|\mathrm{comp}(w) - w\| = \sqrt{\sum_{i=1}^n (\epsilon_i w_i)^2} \le \max_{i \in \{1, \ldots, n\}} (|\epsilon_i|) \sqrt{\sum_{i=1}^n w_i^2} \lesssim \epsilon_0 \|w\|.$$

For $\|[\mathrm{comp}(w) - w]_A\|$, we'll ignore the additional $\sqrt{n}$ term from $\lesssim \epsilon_0 (\|w_A\| + \sqrt{n})$.

# Numerical matrix multiplication error

## Numerical error propagation: condition numbers

Let $\boldsymbol{X} = \boldsymbol{U}\boldsymbol{D}\boldsymbol{V}^\top$ be the SVD of $\boldsymbol{X}$, and we will assume $\boldsymbol{X}$ has full rank. Given that $\boldsymbol{X}^\top \boldsymbol{y} \neq \boldsymbol{0}$, the relative error of $\boldsymbol{X}^\top \boldsymbol{y}$ is

$$\frac{\|\boldsymbol{X}^\top \operatorname{comp}(\boldsymbol{y}) - \boldsymbol{X}^\top \boldsymbol{y}\|}{\|\boldsymbol{X}^\top \boldsymbol{y}\|} = \frac{\|\boldsymbol{X}^\top [\operatorname{comp}(\boldsymbol{y}) - \boldsymbol{y}]\|}{\|\boldsymbol{X}^\top \boldsymbol{y}\|} \leq \frac{d_1 \|[\operatorname{comp}(\boldsymbol{y}) - \boldsymbol{y}]_U\|}{d_p \|\boldsymbol{y}_U\|} \lesssim \epsilon_0 \kappa(\boldsymbol{X}),$$

where $\kappa(\boldsymbol{X}) = d_1/d_p$ is the *condition number* of $\boldsymbol{X}$.

For the second step of $\widehat{\boldsymbol{\beta}} = (\boldsymbol{X}^\top \boldsymbol{X})^{-1} \operatorname{comp}(\boldsymbol{X}^\top \boldsymbol{y})$,

$$(\boldsymbol{X}^\top \boldsymbol{X})^{-1} = (\boldsymbol{V}\boldsymbol{D}\boldsymbol{U}^\top \boldsymbol{U}\boldsymbol{D}\boldsymbol{V}^\top)^{-1} = (\boldsymbol{V}\boldsymbol{D}^2\boldsymbol{V}^\top)^{-1} = \boldsymbol{V}\boldsymbol{D}^{-2}\boldsymbol{V}^\top$$

This looks like another SVD, except for a reverse order of the singular values, $1/d_p^2, \ldots, 1/d_1^2$. The same bounds as before apply, so that the condition number is $\kappa[(\boldsymbol{X}^\top \boldsymbol{X})^{-1}] = d_1^2/d_p^2 = \kappa(\boldsymbol{X})^2$. The resulting relative computational error is

$$\frac{\|(\boldsymbol{X}^\top \boldsymbol{X})^{-1} \operatorname{comp}(\boldsymbol{X}^\top \boldsymbol{y}) - (\boldsymbol{X}^\top \boldsymbol{X})^{-1} \boldsymbol{X}^\top \boldsymbol{y}\|}{\|(\boldsymbol{X}^\top \boldsymbol{X})^{-1} \boldsymbol{X}^\top \boldsymbol{y}\|} \lesssim \epsilon_0 \kappa(\boldsymbol{X})^2.$$

# Numerical example

```r
n <- 100000
X <- cbind(1, (1:n), (1:n)^2) # 3 columns, (badly!) defining a quadratic regression curve model

## Compute the condition number
s <- svd(X) # Gives a list(u, d, v) for U, D=diag(d), and V
print(condition_number <- max(s$d) / min(s$d))

## [1] 13416843918

## Relative error bound for X'y
.Machine$double.eps * condition_number

## [1] 2.979138e-06

## Relative error bound for (X'X)^(-1)
.Machine$double.eps * condition_number^2

## [1] 39970.63
```

# Numerical solves

Let's investigate a case without measurement noise, so that the solution in a perfect would would be exactly equal to the true parameters:

```
beta_true <- c(1, 1, 1)
y <- X %*% beta_true # '%*%' is the matrix multiplication operator in R
```

A direct solve of the linear system $X^\top X \widehat{\beta} = X^\top y$ (the so called *normal equations*) fails:

```
beta1 <- solve(t(X) %*% X, t(X) %*% y)

## Error in solve.default(t(X) %*% X, t(X) %*% y):  system is computationally singular:
reciprocal condition number = 5.5549e-21
```

Using the computed SVD, $\widehat{\beta} = V D^{-1} U^\top y$:

```
beta2 <- s$v %*% ((t(s$u) %*% y) / s$d)
vec_norm(beta2 - beta_true) / vec_norm(beta_true) # You'll define vec_norm() in Lab 5!

## [1] 6.450106e-05
```

The error is a factor $\sim 20$ larger than our estimated bound
$\epsilon_0 \kappa(X) = \epsilon_0 \kappa(V D^{-1} U^\top) = 2.9791378 \times 10^{-6}$.
Note: This is reasonable. What issues did our analysis ignore?

# QR decomposition

We introduced the SVD mainly to help with our theoretical analysis, but it is expensive to compute.

For least squares problems, the main alternative is the following method:

### QR decomposition (here only for tall matrices)

For any square or tall matrix $A : n \times m$, $n \geq m$, there exist matrices $Q : n \times m$ and $R : m \times m$, such that $Q^\top Q = I_m$ and $R$ is upper triangular, and $A = QR$.

The least squares solution based on $X = QR$ becomes

$$\widehat{\beta} = (R^\top Q^\top QR)^{-1} R^\top Q^\top y = R^{-\top} Q^\top y \qquad \text{(if } X \text{ has full rank, } R \text{ is invertible)}$$

i.e. one matrix multiplication with $\kappa(Q) = 1$ and one triangular linear system solve, $\kappa(R) = \kappa(X)$.

```
beta3 <- qr.solve(X, y)
vec_norm(beta3 - beta_true) / vec_norm(beta_true)

## [1] 3.277137e-05
```

About half the error of the SVD method. What about the speed?

# Computational cost

- ▶ For large models, computational speed and memory usage are vital issues.
- ▶ Choosing the right algorithm can mean the difference of waiting for a few seconds and waiting for several weeks!

```
n <- 100000
m <- 50
X <- matrix(rnorm(n * m), n, m)
rbind(
  system.time({ svd(X) }),
  system.time({ qr(X) })
)

##      user.self sys.self elapsed user.child sys.child
## [1,]     0.853    0.040   0.893          0         0
## [2,]     0.284    0.032   0.316          0         0
```

Both methods take $\propto n^3$ operations to compute for a wide range of least squares problems, but QR is consistently around a factor 3 faster than SVD.
In R, the `lm()` function uses QR decomposition internally.

# Summary

▶ If we're not careful, the finite computer representation error may be amplified by computations.

▶ Theoretically correct formulas are not necessarily appropriate to compute directly as written.

▶ Numerical matrix decomposition methods and method parameters chosen by minimising error bounds can help minimise numerical errors.

▶ Large condition numbers may need manual intervention; can we formulate a linear statistical model in more than one way? (See Lab 5 for examples of simple condition number reduction methods.)