

## Differentiation

### Continuity

at a point iff  $\lim_{h \rightarrow 0} f(x+h) = f(x)$

Given  $f, g$ , cont. functions, the following functions are cont

|                              |   |
|------------------------------|---|
| $f+g$                        | $ f $   |
| $cf$ for $c \in \mathbb{R}$  | $\exp(f)$   |
| $f^n$ for $n \in \mathbb{N}$ | $f^\alpha$ for $\alpha \in \mathbb{N}, f$ strictly positive |
| $fg$                         | $f/g, g \neq 0$   |
| $f \circ g$                  | $\log f, f$ strictly positive                               |
| $\max(f, g)$                 |   |

### Differentiation in 1D

at a point  $x \iff \lim_{h \rightarrow 0} \frac{f(x+h)-f(x)}{h}$  .uses the same rules as in table above, but not  $\max$  or  $|f|$

Rules for derivative combinations are kinda obv. Extension of quotient rule:

$$\frac{d}{dx} \left( \frac{f}{g^k} \right) = \frac{g \frac{df}{dx} - k f \frac{dg}{dx}}{g^{k+1}}$$

### Taylor's Theorem

For  $k \geq 0$ , and  $f : D \rightarrow \mathbb{R}$  is 'smooth'

$$\begin{aligned} f(x) &= f(x_0) + (x - x_0) \frac{df}{dx}(x_0) + \cdots \\ &+ \frac{(x - x_0)^k}{k!} \frac{d^k f}{dx^k}(x_0) \\ &+ \frac{(x - x_0)^{k+1}}{(k+1)!} \frac{d^{k+1} f}{dx^{k+1}}(\xi) \Leftarrow \text{Lagrange remainder} \\ &\text{for some } \xi \in (x_0, x) \end{aligned}$$

### Partial differentiation

Consider  $f : D \rightarrow \mathbb{R}$ , for  $D \subseteq \mathbb{R}^n$  a partial derivative is, e.g.  $\frac{\partial f}{\partial x}$  or  $\frac{\partial^2 f}{\partial x \partial y}$  (Clairaut's theorem gives us that  $\frac{\partial^2 f}{\partial x \partial y} = \frac{\partial^2 f}{\partial y \partial x}$  if both are cont). Define the Hessian of  $f$  as:

$$\mathbf{H} = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}$$

### n-D differentiation

For a function  $f : D \rightarrow \mathbb{R}$  for  $D \subseteq \mathbb{R}^n$ :

$$\frac{df}{d\mathbf{x}} = \left( \frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right)^T \leftarrow \text{a column vector}$$

Linearity, the product rule and the quotient rule all work as expected from the 1D case.

The chain rule does as well, at least for  $g : \mathbb{R} \rightarrow \mathbb{R}$  with  $f$  as above. But if  $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$  and  $g : \mathbb{R}^n \rightarrow \mathbb{R}$ , then

$$\frac{d}{dx}(g \circ f) = J(f)^T \left( \frac{dg}{dx} \circ f \right)$$

Useful derivatives include:

| Function   | Derivative  |
|--|---|
| $\mathbf{a}^T \mathbf{x}, \mathbf{x}^T \mathbf{a}$ | $\rightarrow \mathbf{a}$  |
| $\mathbf{x}^T \mathbf{x}$                          | $\rightarrow 2\mathbf{x}$   |
| $\mathbf{x}^T \mathbf{A} \mathbf{x}$               | $\rightarrow (\mathbf{A} + \mathbf{A}^T)\mathbf{x} = 2\mathbf{A}\mathbf{x}$ (if $\mathbf{A}$ symmetric) |

### Jacobian

Given  $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$ , we can consider  $f$  as a **column** vector of scalar-valued functions. Then the Jacobian is its 'derivative'

$$J(f) = \begin{bmatrix} \frac{\partial f}{\partial x_1} & \cdots & \frac{\partial f}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

So  $J_{i,j} = \frac{\partial f_i}{\partial x_j}$

Note, obviously, that  $H(f) = J(\frac{df}{dx})^T$ . Combinations of Jacobian matrices:

$$\begin{aligned} J(f+g) &= J(f) + J(g) \\ J(cf) &= cJ(f) && \text{(for const } c) \\ J(Af) &= AJ(f) && \text{(for const matrix } A) \\ J(f^T g) &= g^T J(f) + f^T J(g) && \text{(dot product rule)} \\ J(fg) &= g \frac{df}{dx}^T + f J(g) && \text{(prod. of scalar, vector)} \\ J(g \circ f) &= (J(g) \circ f) J(f) && \text{(partial chain rule)} \end{aligned}$$

$$\begin{aligned} J(\mathbf{x}) &= \mathbf{I} \\ J(A\mathbf{x}) &= A \end{aligned}$$

### Taylor's Theorem in nD

Given  $f : D \rightarrow \mathbb{R}$  for  $D \subseteq \mathbb{R}^n$ :

$$f(\mathbf{x}) = f(\mathbf{x}_0) + (\mathbf{x} - \mathbf{x}_0)^T \frac{df}{dx}(\mathbf{x}^*) \quad \text{(0th order)}$$

$$\begin{aligned} f(\mathbf{x}) &= f(\mathbf{x}_0) + (\mathbf{x} - \mathbf{x}_0)^T \frac{df}{dx}(\mathbf{x}_0) \quad \text{(1st order)} \\ &\quad + \frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^T H(f)(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}_0) \end{aligned}$$

$$\begin{aligned} f(\mathbf{x}) &= f(\mathbf{x}_0) + (\mathbf{x} - \mathbf{x}_0)^T \frac{df}{dx}(\mathbf{x}_0) \quad \text{(2nd order)} \\ &\quad + \frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^T H(f)(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}_0) + e_3 \end{aligned}$$

## Optimisation

To find:  $\arg \min_{x \in F} f(\mathbf{x})$  for an **objective function**  $f$ , on a **feasible set**  $F$ .

If  $F = \mathbb{R}^n$  then optimisation is **unconstrained**. Otherwise, it is **constrained**, and  $F$  is normally defined as

$$F = \{\mathbf{x} \in \mathbb{R}^n | g_i(\mathbf{x}) = 0, \dots, h_i(\mathbf{x}) \geq 0, \dots\}$$

If  $F = \emptyset$  then the problem is **infeasible/inconsistent**, and if there is no minimum then it is **unbounded**.

### 1D

i.e. for  $f : \mathbb{R} \rightarrow \mathbb{R}$ . First we need  $\frac{df}{dx} = 0$ . We then need to find global minima:

$$\text{global minima} \subseteq \text{local minima} \subseteq \text{stationary points}$$

For constrained optimisation you also need to **check endpoint values**.

If  $k$  is the index of the first non-zero derivative, then

$$x_0 \text{ is local } \begin{cases} \begin{cases} \text{minimum} & \text{if } \frac{d^k f}{dx^k} > 0 \\ \text{maximum} & \text{if } \frac{d^k f}{dx^k} < 0 \end{cases} & \text{if } k \text{ is even} \\ \text{stationary point} & \text{if } k \text{ is odd} \end{cases}$$

Of course, if  $\frac{d^2 f}{dx^2}(x_0) > 0$  then  $x_0$  is immediately a local min.

### nD, unconstrained

i.e. for  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ . First we need  $\frac{df}{d\mathbf{x}} = 0$ .

$$\mathbf{H}(f)(x_0) \text{ is positive definite} \Rightarrow \mathbf{x}_0 \text{ is a minimum}$$

For a symmetric matrix  $\mathbf{A}$

$$\mathbf{A} \text{ is } \begin{cases} \begin{matrix} \text{positive definite} \\ \text{positive semi-definite} \\ \text{negative definite} \\ \text{negative semi-definite} \end{matrix} \end{cases} \text{ if } x^T \mathbf{A} x \begin{cases} > \\ \geq \\ < \\ \leq \end{cases} 0$$

If  $\mathbf{A}$  is 2-by-2 then you can replace  $x^T \mathbf{A} x$  by **eigenvalues** or **pivots** above.

Also:  $\mathbf{A}$  is positive (semi-)definite  $\Leftrightarrow -\mathbf{A}$  is negative (semi-)definite.

For a positive (semi-)definite matrix  $\mathbf{A}$ , so is :

$$\begin{array}{ll} \mathbf{A} + \mathbf{B} & \text{for } \mathbf{B} \text{ positive semi-definite} \\ c\mathbf{A} & \text{for } c \in \mathbb{R} \\ \mathbf{A}^{-1} & \text{which exists if } \mathbf{A} \text{ is pos def} \\ \mathbf{C}^T \mathbf{A} \mathbf{C} & \text{for full rank } \mathbf{C} \text{ for pos def, otherwise any } \mathbf{C} \end{array}$$

### Convexity

A **set**  $D \subseteq \mathbb{R}^n$  is convex  $\Leftrightarrow \forall \mathbf{x}_1, \mathbf{x}_2 \in D, \forall \alpha \in (0, 1) \quad (1 - \alpha)\mathbf{x}_1 + \alpha\mathbf{x}_2 \in D$ .

A **function**  $f : D \rightarrow \mathbb{R}^n$  on convex  $D \subseteq \mathbb{R}^n$  is convex  $\Leftrightarrow \forall \mathbf{x}_1, \mathbf{x}_2 \in D, \alpha \in (0, 1)$

$$f((1 - \alpha)\mathbf{x}_1 + \alpha\mathbf{x}_2) \leq (1 - \alpha)f(\mathbf{x}_1) + \alpha f(\mathbf{x}_2)$$

(strict convexity defs have strict inequalities)

If we have a convex  $f$ , then

$$\text{global minima} = \text{local minima} = \text{stationary points}$$

But of course in constrained optimisation we still need to **check the boundary**:  
for a **function**  $f : D \rightarrow \mathbb{R}^n$  on convex  $D \subseteq \mathbb{R}^n$  with cont. second-order derivatives

$$f \text{ (-/strictly) convex} \Leftrightarrow \mathbf{H}(f) \text{ positive (semi/-)definite}$$

| Convex       | Concave (strict is *)                           |
|--------------|---|
| $y = ax + b$ | $y = ax + b$                                    |
| * $\exp(x)$  | * $\log x$                                      |
|              | $ x ^p$ for * $p > 1$ (if $p = 1$ only concave) |

For  $f, g$  (strictly) convex functions on convex  $D \subseteq \mathbb{R}^n$ , so are

$$\begin{aligned} & f + g \\ & cf \quad \forall c \in \mathbb{R} \\ & x \mapsto f(\mathbf{A}x + b) \quad (\text{for strict convexity } \mathbf{A} \text{ full rank}) \\ & \max(f, g) \\ & \exp f \end{aligned}$$

### Tricks

If  $g$  is a strictly increasing function,

$$\arg \min_{x \in F} f(x) = \arg \min_{x \in F} g(f(x))$$

e.g. if  $f$  is a product (and a positive function) then  $\log f$  is a sum, and sometimes easier to diff.  
If  $g$  is a 1-1 function,

$$\arg \min_{x \in F} f(x) = g \left( \arg \min_{x \in F} f(g(x)) \right)$$

e.g. to minimise  $\cos(\exp(x))$  we need points where  $\exp(x) = n\pi$ .

### Equality constraints

Use the **Lagrangian**, as it has the same stationary points on the constraint as  $f$ :

$$\Lambda(\mathbf{x}, \lambda_1, \lambda_2 \dots) = f(\mathbf{x}) - \lambda_1 g_1(\mathbf{x}) - \lambda_2 g_2(\mathbf{x}) \dots$$

Then differentiate and eliminate  $\lambda$  to get solutions, then classify them as follows.

- Check if  $f$  is unbounded (i.e. if there is even a global min)
- evaluate the stationary points to find the min.

### Inequality constraints

The question is if each constraint  $h_i(\mathbf{x})$  is **slack** ( $> 0$ ) or **tight** ( $= 0$ ) Guess (with reasoning) which are so, and solve all options for remaining tight/slack constraints, and check that the ones assumed to be slack are actually slack.

### Numerical integration

#### in 1D

Note  $m = \frac{a+b}{2}$ ,

$$\begin{aligned} \underline{D}_k &= \min_{x \in (a,b)} \frac{d^k f}{dx^k} \\ \overline{D}_K &= \max_{x \in (a,b)} \frac{d^k f}{dx^k} \end{aligned}$$

**Midpoint Rule:**

$$\begin{aligned}
M_1[f, a, b] &= (b - a)f(m) \\
err(M_1)[f, a, b] &\leq -\frac{1}{24}(\mathbf{b} - \mathbf{a})^3 \underline{D_2} \\
M_n[f, a, b] &= \sum_n^{i=1} M_1[f, x_{i-1}, x_i] \\
&= \frac{b-a}{n} \sum_n^{i=1} f\left(\frac{x_{i-1} + x_i}{2}\right) \\
err(M_n)[f, a, b] &= \sum_n^{i=1} err(M_1)[f, x_{i-1}, x_i] \leq -\frac{(b-a)^3}{24n^2} \underline{D_2}
\end{aligned}$$

**Trapezium Rule:**

$$\begin{aligned}
T_1[f, a, b] &= \frac{b-a}{2}(f(a) + f(b)) \\
err(T_1)[f, a, b] &\leq \frac{(\mathbf{b} - \mathbf{a})^3}{12} \overline{D_2} \\
T_n[f, a, b] &= \sum_n^{i=1} T_1[f, x_{i-1}, x_i] \\
&= \frac{b-a}{2n} (f(x_0) + 2f(x_1) + \dots + 2f(x_{n-1}) + f(x_n)) \\
err(T_n)[f, a, b] &\leq -\frac{(b-a)^3}{12n^2} \overline{D_2}
\end{aligned}$$

**Simpson's Rule:**

Note  $n$  must be even.

$$\begin{aligned}
S_2[f, a, b] &= \frac{b-a}{6}(f(a) + 4f(m) + f(b)) \\
err(S_2)[f, a, b] &\leq \frac{(\mathbf{b} - \mathbf{a})^5}{2880} \overline{D_4} \\
S_n[f, a, b] &= \sum_{n/2}^{i=1} S_2[f, x_{2i-2}, x_{2i}] \\
&= \frac{b-a}{3n} (f(x_0) + 4f(x_1) + 2f(x_2) + 4f(x_3) + \dots + 4f(x_{n-1}) + f(x_n)) \\
err(S_n)[f, a, b] &\leq -\frac{(b-a)^5}{180n^4} \overline{D_4}
\end{aligned}$$

Note all these composite methods do about the same amount of work (all eval.  $f$  roughly  $n$  times)

**in nD**

The 1D methods above can be extended (i.e. evaluating  $f$  on grids in  $n$  dimensions) but that means their error converges only at  $O(n^{-2/d})$  and  $O(n^{-4/d})$ , which is very slow. So we use

**Monte Carlo Integration:**

Given  $N$  random vectors  $\mathbf{X}_1, \dots, \mathbf{X}_N$

$$MC_N[f, R] = \text{Area}(R) \frac{1}{N} \sum_i f(\mathbf{X}_i)$$

Useful properties include:

- $E[MC_N[f, R]] = \int_R f(\mathbf{x}) d\mathbf{x}$  (i.e. unbiased)
- $\text{Var}[MC_N[f, R]] = V/N$  for a constant
- $P[|err(MC_N)[f, R]| \geq \epsilon] \rightarrow 0$  as  $N \rightarrow \infty$ , for any  $\epsilon > 0$  (i.e. consistent)
- $err(MC_N)[f, R]$  approaches the Normal distribution  $N(0, \frac{V}{N})$  for large  $N$ .

So the error (the standard error  $\sqrt{\frac{V}{N}}$ ) decreases as  $O(N^{-1/2})$  regardless of dimension, which is slower than, e.g. Simpson's but much better if  $d$  is large.

To use this, we need to generate random vectors in  $R$ , so the simplest way to do that is choose an  $S$  such that  $R \subseteq S$  and random vectors in  $S$  can be easily generated, and then discard vectors in  $S \setminus R$ . If necessary, we can calculate  $\text{Area}(R)$  as  $\text{Area}(S) \times$  the proportion of accepted samples from  $S$ .

Techniques to improve this: **importance sampling**, **control variates** ( $\int_R f(\mathbf{x}) d\mathbf{x} = \int_R (f(\mathbf{x}) - \hat{f}(\mathbf{x})) d\mathbf{x} + \int_R \hat{f}(\mathbf{x}) d\mathbf{x}$ ), **stratified sampling**

## Accuracy

For  $u \in \mathbb{R}$ :  
 $\tilde{u} - u$  is the **error**  
 $|\tilde{u} - u|$  is the **absolute error**  
 $|\tilde{u} - u|/|u|$  is the **relative error**

If  $u \in \mathbb{R}^n$  replace  $|x|$  with the Euclidean norm  $\|x\|$ .

## Error

**Rounding error** only roughly 15 sig digits can be stored

**Machine epsilon**,  $\epsilon$  the smallest quantity st. every  $x \in \mathbb{R}$  can be stored as  $\tilde{x}$ , a floating point with relative error  $\leq \epsilon$

For double precision  $\epsilon \approx 2^{-53} \approx 1.1 \times 10^{-16}$

Most methods run with accuracy to within a few multiples of  $\epsilon$ . Others will run to  $\sqrt{\epsilon}$  which is bigger than  $\epsilon$ .

**Catastrophic cancellation** occurs if  $a \approx b$ , then small relative errors in  $a$  or  $b$  can cause large relative errors in  $a - b$ . Avoid this by rewriting formulae to avoid subtraction of this kind

**Truncation error** error when approximating an infinite procedure with a finite one (e.g. a Taylor poly vs the orig. function)

**Roundoff error** is introduced by approximating a real number with a stored (floating point) number

**Propagation** is when a small error is magnified by the following procedures (often called **instability**), e.g. solving  $A\mathbf{x} = \mathbf{b}$

## Rates of convergence

Given an algorithm that seeks  $\mathbf{x}^*$  with iterative guesses  $\mathbf{x}_1, \mathbf{x}_2, \dots$ , the absolute error at the  $n$ -th step is  $\epsilon_n = x_n - x^*$ .

|  |                     |     |   |   |
|--|---------------------|-----|---|---|
| If $ \epsilon_n  \rightarrow 0$ , convergence is | linear              | (is | $\frac{ \epsilon_{n+1} }{ \epsilon_n } \rightarrow \alpha \in (0, 1)$                   | $c^n$ for $c < 1$<br>(throughout)<br>$1/n$<br>$c^{n^2}$<br>$c^{2^n}$ is order-2 |
|  | sublinear           |     | $\frac{ \epsilon_{n+1} }{ \epsilon_n } \rightarrow 1$                                   |   |
|  | logarithmic         |     | $\frac{ \epsilon_{n+2} - \epsilon_{n+1} }{ \epsilon_{n+1} - \epsilon_n } \rightarrow 1$ |   |
|  | sublinear)          |     |   |   |
|  | superlinear         |     | $\frac{ \epsilon_{n+1} }{ \epsilon_n } \rightarrow 0$                                   |   |
|  | order-q ( $q > 1$ ) |     | $\frac{ \epsilon_{n+1} }{ \epsilon_n ^q} \rightarrow \alpha > 0$                        |   |

Order-2 convergence is often called quadratic. It is very useful when comparing algorithms to know that:

If  $A_n$  converges with order  $q$ , then  $A_{2n}$  converges with order  $q^2$ .

## Rootfinding

### Termination

Forward error bounds:

|                                    |                                  |
|------------------------------------|----------------------------------|
| $ x_n - x_{n-1}  < tol$            | a small change in abs. error     |
| $ x_n - x_{n-1}  < tol x_n $       | a small change in relative error |
| $ x_n - x_{n-1}  < tol(1 +  x_n )$ | measures both abs. and rel       |

OR Backward error bound:  $|f(x_n)| < tol$  function close to zero (backward error bound)

Practical bounds: (must include)  $n = N$  iteration ran out of time  
a step was ill-defined so an error must be raised

### Interval bisection

Given  $m_n = \frac{a_n + b_n}{2}$ ,

If  $f(a_n)f(m_n) < 0$  set  $(a_{n+1}, b_{n+1})$  to be  $(a_n, m_n)$ , otherwise  $(m_n, b_n)$ .

It is slow, but guarantees convergence given a bracket (use either maths or test increasing brackets of the form  $(-2^k, 2^k)$ ), but this method does not generalise to higher dimensions.

Convergence is linear, as  $\epsilon_{n+1} = \frac{1}{2}\epsilon_n$ , and we need  $n = \lceil \log \frac{b_0 - a_0}{2tol} \rceil$

### Newton's method in 1D

Uses 1st-order Taylor poly as an approximation, and finds the root of that.

$$x_{n+1} = x_n - \frac{f(x_n)}{\frac{df}{dx}(x_n)}$$

### Problems

- $\frac{df}{dx}(x_n) = 0$  means the iteration is undefined
- it can diverge (if  $\frac{df}{dx}$  is close/eq. to 0)
- it can converge to a different root
- it can also get stuck in a loop

### Error

By the same Taylor expansion,

$$x^* - x_n \approx \frac{f(x_n)}{\frac{df}{dx}(x_n)} = x_{n+1} - x_n$$

So we can make one more iteration to get an **a posteriori error estimate**.

Given

$$A_I(c) = \frac{\max_{\beta \in I} \left| \frac{d^2 f}{dx^2}(\beta) \right|}{\min_{\alpha \in I} \left| \frac{d^2 f}{dx^2}(\alpha) \right|}$$

- For an interval  $I = (x^* - c, x^*)$ , if  $x_0 \in I$  and  $\frac{cA_I(c)}{2} < 1$ , Newton's method converges at least quadratically to  $x^*$ .
- Given a bracket for  $x^*$  of the form  $(m - c, m + c)$ , Newton's method converges at least quadratically starting at  $m$ , if  $\frac{cA_I(c)}{2} < 1$  for the interval  $J = (m - 2c, m + 2c)$ .
- Such an interval  $J$  exists if  $\frac{df}{dx}(x^*) \neq 0$ .

If  $f$  has a very small derivative at the root (so  $f$  is **ill-conditioned**) then there is possibility for significant roundoff error.

### Secant method, 1D

Uses **linear interpolation** and solves that approximation.

$$x_{n+1} = x_n - \frac{f(x_n)(x_n - x_{n-1})}{f(x_n) - f(x_{n-1})}$$

Only needs 1 eval. of  $f$  per step, as we can remember  $f(x_{n-1})$ .

Can fail if  $f(x_n) \approx f(x_{n-1})$ , and a useful termination condition is  $|f(x_n) - f(x_{n-1})| < \text{tol}|f(x_n)|$  to avoid this.

### Error

Same conditions as Newton's Method, to get order- $\phi$  convergence, for  $\phi$  between 1.5 and 1.7. Remember that two steps of the secant method have order  $\phi^2 > 2$  convergence, which is faster than Newton's method (so if Newton is slower than two secant steps, use this instead).

### Newton's method in dD

We are solving  $\mathbf{f}(\mathbf{x}) = \mathbf{0}$  for  $\mathbf{f} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ . To get this method split  $\mathbf{f}$  into  $d$  functions  $\mathbb{R}^d \rightarrow \mathbb{R}$  and apply the 1st order Taylor poly, to get:

$$\mathbf{x}_{n+1} = \mathbf{x}_n - (\mathbf{J}(\mathbf{f})(\mathbf{x}_n))^{-1} \mathbf{f}(\mathbf{x}_n)$$

or, which is about 3 times **faster**, solve for  $\Delta \mathbf{x}$  in  $(\mathbf{J}(\mathbf{f})(\mathbf{x}_n))\Delta \mathbf{x} = -\mathbf{f}(\mathbf{x}_n)$  and set  $\mathbf{x}_{n+1} = \mathbf{x}_n + \Delta \mathbf{x}$

If  $\mathbf{J}(\mathbf{f})(\mathbf{x}_n)$  is singular, there may not be a solution, so we must stop. The cost per iteration is  $O(d^3)$  to solve for  $\Delta \mathbf{x}$  and  $O(d^2)$  to evaluate the Jacobian.

We can make this less fragile if we insist that  $\|\mathbf{f}(\mathbf{x}_{n+1})\| < \|\mathbf{f}(\mathbf{x}_n)\|$  ( $\dagger$ ), and if it is not true, set  $\mathbf{x}_{n+1} = \mathbf{x}_n + \lambda \Delta \mathbf{x}$  for  $\lambda \in (0, 1)$ . Damping with small  $\lambda$  will always satisfy ( $\dagger$ ), but may only be linear convergence.

### Quasi-Newton methods

Essentially we want to approximate  $\mathbf{J}_n$ , with  $\hat{\mathbf{J}}_n$ . This must satisfy the secant equation:  $(\mathbf{y}_r = \mathbf{f}(\mathbf{x}_r))$

$$\mathbf{y}_n = \mathbf{y}_{n-1} + \hat{\mathbf{J}}_n(\mathbf{x}_n - \mathbf{x}_{n-1})$$

And also must not change for vectors orthogonal to  $\mathbf{x}_n - \mathbf{x}_{n-1}$  (Trid).

A starting approximate Jacobian is needed (either  $-\mathbf{I}$  or  $\mathbf{J}(\mathbf{f})(\mathbf{x}_0)$ ), and we can optimise further by tracking the inverse of the approximation and only updating that.

Termination is similar to Newton's, and it may need more steps, but they will likely be faster.



## Numerical optimisation

Similar termination conditions to rootfinding:

- $|x_n - x_{n-1}| < tol(1 + |x_n|)$  - loc. of approx min didn't change much (might signify step length too short, but good for Newton-like methods)
- $|f(x_n) - f(x_{n-1})| < tol(1 + |f(x_n)|)$  - val. of approx min didn't change much (dangerous if fun. flat far from min)
- $|\frac{df}{dx}(x_n)| < tol$  - when gradient close to zero (only applicable if methods computes gradient)
- $n = N$  or ill-defined steps.

**Warning:** The relative error in **loc. of min** is  $O(\sqrt{\epsilon}) > O(\epsilon)$  (from Taylor's rule if relative error in value is less than  $\epsilon$ ).

### Golden section search, 1D:

A **bracket** is a triple  $(a, b, c)$  with  $a < b < c$ , and  $f(a) > f(b) < f(c)$ , which must contain at least 1 local minimum.

Refinement process:

- choose  $z \in (a, c)$  and find the appropriate triple out of  $a, b, c, z$  that is a bracket.
- Choose  $z$  using the golden ratio: if  $(a, b)$  is the larger of  $(a, b)$  and  $(b, c)$ ,  $b - z = \phi(b - a)$ . Also  $z - a = c - b$ .
- If the previous bracket followed the same ratio, then we have  $\phi$  as the golden ratio.

This is linear convergence, with one eval. of  $f$  per iteration - the bracket shrinks by  $\approx 0.618$  per iter. A starting bracket is required - e.g using  $(0, (\phi + 1)^n, (\phi + 1)^{n+1})$

### Gradient descent in d dimensions

$$\mathbf{x}_{n+1} = \mathbf{x}_n + \alpha_n \mathbf{d}_n$$

Criteria for  $\mathbf{d}_n$  and  $\alpha_n$ : (Note  $\mathbf{g}_n = \frac{df}{dx}(\mathbf{x}_n)$ )

1.  $\mathbf{d}_n$  is a **descent direction**, so  $\mathbf{g}_n^T \mathbf{d}_n < 0$
2.  $\mathbf{d}_n, \mathbf{d}_{n-1}, \dots$  do not veer wildly
3.  $\alpha_n$  seeks to roughly minimise  $f(\mathbf{x}_n + \alpha_n \mathbf{d}_n)$
4. the method does not stop - i.e.  $\sum_1^\infty \alpha_n = \infty$

Options for  $\mathbf{d}_n$ :

- **coordinate gradient descent:** cycle through unit vectors  $\pm e_i$ . Fast, needs no derivatives, but naive
- **gradient descent:**  $\mathbf{d}_n = -\mathbf{g}_n$ . Satisfies conditions (indeed best sol.)
- **conjugate grad descent:** directions that satisfy  $\mathbf{d}_n^T \mathbf{H}(f)(\mathbf{x}_n) \mathbf{d}_m = 0$  for  $m < n$
- in machine learning when obj. fun is an average loss, use only part of the data to find  $\mathbf{d}_n$  (deterministic/random).

Options for  $\alpha_n$ :

- Constant: simple, but often slow progress or not actually descent
- **backtracking**: choose an initial step length  $a'_n$  and multiply by  $0 < \rho < 1$  until it is suitable.
- Good starting step length:  $\alpha'_n = \alpha_{n-1} \frac{\mathbf{g}_{n-1}^T \mathbf{d}_{n-1}}{\mathbf{g}_n^T \mathbf{d}_n}$
- **Armijo rule** for suitable  $\alpha_n$ :  $f(\mathbf{x}_n + \alpha_n \mathbf{d}_n) < f(\mathbf{x}_n) + \sigma \alpha_n \mathbf{g}_n^T \mathbf{d}_n$  for (roughly)  $\sigma \in (10^{-4}, 10^{-1})$ .

Since we have the gradient  $\mathbf{g}_n$  it is sensible to use it in the termination conditions (e.g. as  $\|\mathbf{g}_n\| < \text{tol}(1 + \|\mathbf{g}_0\|)$ )

### Newton methods in d dimensions

Use a 2nd order Taylor expansion and minimise that:

$$\mathbf{x}_{n+1} = \mathbf{x}_n - (\mathbf{H}(f)(\mathbf{x}_n))^{-1} \mathbf{g}_n$$

Much more work, as Hessian calculation is  $O(d^2)$  and  $O(d^3)$  to solve for  $\Delta \mathbf{x}$  for the inverse, which is undefined if  $\mathbf{H}(f)(\mathbf{x}_n)$  is singular.

Also the Newton direction isn't necessarily a descent direction. (it might find a maximum).

If  $f$  has a well-behaved local minimum  $\mathbf{x}^*$   $\mathbf{H}(f)(\mathbf{x}_n)$  should be pos def near  $\mathbf{x}^*$  so Newton will be quadratic conv. (not always true, e.g. if Hess only semi-def).

Can try to fix the Hessian by adding  $\lambda \mathbf{I}$  for some  $\lambda$ .

### Quasi-Newton methods

Approximate the the Hessian somehow.