

Vector spaces, matrices and norms

Vectors and vector spaces

- A vector \underline{x} is an element of a vector space

* In this module, we will consider complex-valued vectors $\underline{x} \in \mathbb{C}^n$, which have the form

$$\underline{x} = [x_1 \ x_2 \ \dots \ x_n]^T, \text{ where } x_i \in \mathbb{C}$$

- A vector space V is a set of vectors that is closed under addition and scalar multiplication

i.e. if $\underline{v}_1, \underline{v}_2 \in V$, then $a\underline{v}_1 + b\underline{v}_2 \in V$ for any scalars a, b

(The scalars are often real/complex, but could also be chosen over some other field)

- A subspace of a vector space is a subset that obeys the rules of a vector space

- The inner product b/w two equal length vectors yields a complex scalar, defined as

$$\underline{x}^H \underline{y} = \sum_{i=1}^n \bar{x}_i y_i$$

* $\underline{x}^H \underline{y} = \overline{\underline{y}^H \underline{x}}$ and $\underline{x}^H \underline{x}$ is real and non-negative ($\underline{x}^H \underline{x} = 0$ iff $\underline{x} = \underline{0}$)

Matrices

- A matrix \underline{A} is a linear operator that performs a linear transformation from one vector to another.

For a matrix $\underline{A}: \mathbb{C}^n \rightarrow \mathbb{C}^m$ where $\underline{A}\underline{x} = \underline{b}$, the vector $\underline{x} \in \mathbb{C}^n$ is mapped to the vector $\underline{b} \in \mathbb{C}^m$

- An $m \times n$ matrix comes from the space \mathbb{C}^{mn} , so we often denote $\underline{A} \in \mathbb{C}^{mn}$

- The product of two matrices $\underline{C} = \underline{A}\underline{B}$ can be expressed using index notation

$$c_{ij} = \sum_k a_{ik} b_{kj} = a_{ik} b_{kj}$$

Treating a column vector or a $n \times 1$ matrix, we can compute the matrix-vector product $\underline{b} = \underline{A}\underline{x}$.

$$b_i = \sum_j A_{ij} x_j = x_1 \begin{bmatrix} a_{11} \\ \vdots \\ a_{1j} \end{bmatrix} + x_2 \begin{bmatrix} a_{21} \\ \vdots \\ a_{2j} \end{bmatrix} + \dots \leftarrow \text{as are columns of } \underline{A}$$

- The conjugate transpose / Hermitian transpose of a matrix is defined as

$$(\underline{A}^H)_{ij} = \bar{a}_{ji}$$

Alternatively, using the transpose, $\underline{A}^H = \overline{\underline{A}^T} = \underline{\bar{A}}^T$ (conjugate transpose = transpose for real-valued matrices)

- A matrix $\underline{M} \in \mathbb{C}^{n \times n}$ is Hermitian if

$$\underline{M}^H = \underline{M}$$

- A Hermitian matrix $\underline{M} \in \mathbb{C}^{n \times n}$ has the following properties:

(i) All eigenvalues are real, i.e. $\lambda \in \mathbb{R}$

If (λ, \underline{x}) is an eigenpair of \underline{M} , then $\underline{M}\underline{x} = \lambda\underline{x}$.

Premultiplying by \underline{x}^H ,

$$\underline{x}^H \underline{M} \underline{x} = \lambda \underline{x}^H \underline{x}$$

Taking the complex conjugate

$$(\underline{x}^H \underline{M} \underline{x})^H = \underline{x}^H \underline{M}^H \underline{x} = \bar{\lambda} \underline{x}^H \underline{x} \quad (\text{use } \underline{M}^H = \underline{M})$$

Comparing the above eqns, they can only hold when $\bar{\lambda} = \lambda$,

i.e. all the eigenvalues are real.

For Personal Use Only -bkwk2

(i) The eigenvectors are orthogonal, i.e. $\underline{x}_i^H \underline{x}_j = 0 \text{ if } i \neq j$.

Consider two eigenpairs of \underline{M} , $(\lambda_1, \underline{x}_1), (\lambda_2, \underline{x}_2)$, i.e. $\underline{M}\underline{x}_1 = \lambda_1 \underline{x}_1$, $\underline{M}\underline{x}_2 = \lambda_2 \underline{x}_2$, where $\lambda_1 \neq \lambda_2$

$$\text{premultipling by } \underline{x}_2^H \underline{M}, \quad \underline{x}_2^H \underline{M} \underline{x}_1 = \lambda_1 \underline{x}_2^H \underline{x}_1 \quad \underline{x}_2^H \underline{M} \underline{x}_2 = \lambda_2 \underline{x}_2^H \underline{x}_2$$

$$\text{Taking the complex conjugate for left egn, } (\underline{x}_2^H \underline{M} \underline{x}_1)^* = \underline{x}_1^H \underline{M} \underline{x}_2 = \lambda_2 \underline{x}_1^H \underline{x}_2 \quad (\text{use } \underline{x}_2^H = \underline{M})$$

Comparing the above eqns, since $\lambda_1 \neq \lambda_2$, they can only hold if $\underline{x}_1^H \underline{x}_2 = 0$

i.e. the eigenvectors are orthogonal.

- If the eigenvalues of the matrix $\underline{A} \in \mathbb{C}^{n \times n}$ are λ_i , then

(i) The eigenvalues of the inverse \underline{A}^{-1} are $\frac{1}{\lambda_i}$

$$\underline{A} \underline{x}_i = \lambda_i \underline{x}_i \quad \rightarrow \quad \underline{A}^{-1} \underline{x}_i = \frac{1}{\lambda_i} \underline{x}_i$$

(ii) The eigenvalues of the transpose \underline{A}^T are λ_i . Note that $\det(\underline{A}^T) = \det(\underline{A})$,

$$\det(\underline{A} - \lambda \underline{I}) = 0 \quad ; \quad 0 = \det(\underline{A}^T - \lambda \underline{I}) = \det((\underline{A} - \lambda \underline{I})^T) = \det(\underline{A} - \lambda \underline{I}) \rightarrow \lambda = \lambda_i$$

(iii) The eigenvalues of the conjugate transpose \underline{A}^H are $\bar{\lambda}_i$. Note that $\det(\underline{A}^H) = \det(\underline{A}^T) = \overline{\det(\underline{A})}$

$$\det(\underline{A} - \lambda \underline{I}) = 0 \quad ; \quad 0 = \det(\underline{A}^H - \lambda \underline{I}) = \det((\underline{A} - \lambda \underline{I})^H) = \overline{\det(\underline{A} - \lambda \underline{I})} = \det(\underline{A} - \bar{\lambda} \underline{I}) \rightarrow \lambda = \bar{\lambda}_i$$

- For a Hermitian matrix $\underline{M} \in \mathbb{C}^{n \times n}$, it is positive definite / semi-positive definite if

$$\underline{x}^H \underline{M} \underline{x} > 0 \quad \forall \underline{x} \in \mathbb{C}^n \setminus \{0\}$$

$$\underline{x}^H \underline{M} \underline{x} \geq 0 \quad \forall \underline{x} \in \mathbb{C}^n \setminus \{0\}$$

The eigenvalues of a Hermitian positive definite / semi-positive definite matrix are strictly **real** / **non-negative**.

$$\lambda_i > 0$$

$$\lambda_i \geq 0$$

(say $(\lambda_i, \underline{x}_i)$ is an eigenpair of \underline{M} . $\underline{x}_i^H \underline{M} \underline{x}_i = \underline{x}_i^H (\lambda_i \underline{x}_i) = \lambda_i \underline{x}_i^H \underline{x}_i$, and $\underline{x}_i^H \underline{x}_i$ is real and positive)

→ Note that any Hermitian positive definite matrix is invertible, since $\lambda_i > 0 \rightarrow \det(\underline{M}) = \prod_i \lambda_i \neq 0$

→ For a matrix $\underline{A} \in \mathbb{C}^{n \times n}$, the matrix $\underline{A}^H \underline{A}$ is semi-positive definite, i.e. $\underline{x}^H \underline{A}^H \underline{A} \underline{x} \geq 0$.

$$0 \leq \|\underline{A} \underline{x}\|_2^2 = (\underline{A} \underline{x})^H (\underline{A} \underline{x}) = \underline{x}^H \underline{A}^H \underline{A} \underline{x}$$

so the eigenvalues of $\underline{A}^H \underline{A}$ are non-negative. \leftarrow if the eigenvalues of \underline{A} are λ_i , then the eigenvalues of $\underline{A}^H \underline{A}$ are $\lambda_i \lambda_i^*$

- A matrix $\underline{Q} \in \mathbb{C}^{n \times n}$ is a unitary matrix if $\underline{Q}^H = \underline{Q}^{-1}$ i.e. $\underline{Q}^H \underline{Q} = \underline{Q} \underline{Q}^H = \underline{I}$.

$$\hookrightarrow |\det(\underline{Q})| = 1 \quad (\text{since } \det(\underline{Q}^H \underline{Q}) = \det(\underline{Q}^H) \det(\underline{Q}) = |\det(\underline{Q})|^2 = \det(\underline{I}) = 1)$$

\hookrightarrow Transformation preserves the Euclidean norm of a vector $\|\underline{Q} \underline{x}\|_2 = \|\underline{x}\|_2$ and matrix $\|\underline{Q} \underline{A}\|_2 = \|\underline{A}\|_2$

$$(\|\underline{Q} \underline{x}\|_2^2 = (\underline{Q} \underline{x})^H (\underline{Q} \underline{x}) = \underline{x}^H \underline{Q}^H \underline{Q} \underline{x} = \underline{x}^H \underline{I} \underline{x} = \underline{x}^H \underline{x} = \|\underline{x}\|_2^2)$$

- The rank of a matrix $\underline{A} \in \mathbb{C}^{m \times n}$ is the no. of linearly independent rows or columns, and satisfies

$$\text{rank}(\underline{A}) \leq \min(m, n)$$

It is full rank if $\text{rank}(\underline{A}) = \min(m, n)$ and rank deficient if $\text{rank}(\underline{A}) < \min(m, n)$

- A matrix in which most entries are zero is a sparse matrix — the no. of non-zeroes on each

row of a sparse matrix will be roughly the same and independent of the matrix size.

\hookrightarrow The no. of entries in a dense matrix is n^2

\hookrightarrow The no. of entries in a sparse matrix is $c n$

For Personal Use Only -bkwk2

NORMS

- Denoting a norm of a vector \underline{x} by $\|\underline{x}\|$, a norm is a scalar that satisfies

$$\hookrightarrow \|\underline{x}\| > 0 \text{ when } \underline{x} \neq 0$$

$$\hookrightarrow \|\underline{x}\| = 0 \text{ when } \underline{x} = 0$$

$$\hookrightarrow \|k\underline{x}\| = |k| \|\underline{x}\|$$

$$\hookrightarrow \|\underline{x} + \underline{y}\| \leq \|\underline{x}\| + \|\underline{y}\|$$

* For orthogonal vectors $\underline{x}, \underline{y}$, i.e. $\underline{x}^\top \underline{y} = 0$,

$$\|\underline{x} + \underline{y}\|^2 = \|\underline{x}\|^2 + \|\underline{y}\|^2$$

Vector norms.

- A particular family of norms are known as $l\text{-p}$ norms, defined as

$$\|\underline{x}\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{1/p} \quad \text{where } |\underline{x}| = \sqrt{\operatorname{Re}[\underline{x}]^2 + \operatorname{Im}[\underline{x}]^2}$$

$$\hookrightarrow l\text{-1 norm: } \|\underline{x}\|_1 = |x_1| + \dots + |x_n| = \sum_{i=1}^n |x_i| \quad [\text{taxicab norm}]$$

$$\hookrightarrow l\text{-2 norm: } \|\underline{x}\|_2 = \left(|x_1|^2 + \dots + |x_n|^2 \right)^{1/2} = \left(\sum_{i=1}^n |x_i|^2 \right)^{1/2} = (\underline{x}^\top \underline{x})^{1/2} \quad [\text{Euclidean norm}]$$

$$\hookrightarrow l\text{-}\infty \text{ norm: } \|\underline{x}\|_\infty = \max_i |x_i| \quad [\text{max norm}]$$

- We can define norms that involve a matrix A if A is positive definite.

$$\|\underline{x}\|_A = \sqrt{\underline{x}^\top A \underline{x}} \quad [\text{energy norm}]$$

TF $A \geq I$, we have
+c l-2 norm

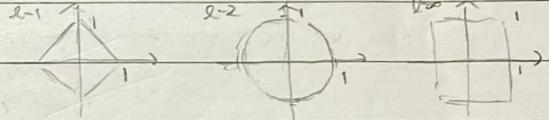
(A being positive definite implies that the energy $\underline{x}^\top A \underline{x}$ is non-negative and only zero if $\underline{x} = 0$)

- It is useful to refer to the unit ball/sphere for a given norm.

$$\{\underline{x} \in V : \|\underline{x}\| \leq 1\}$$

$$\{\underline{x} \in V : \|\underline{x}\| = 1\}$$

The "shape" depends on the chosen norm



Matrix norms induced by vector l-p norms

- The norm of a matrix A is defined as

$$\|A\| = \max_{\underline{x} \in \mathbb{C}^{n \times 1}} \frac{\|A\underline{x}\|}{\|\underline{x}\|}$$

$$\text{or} \quad \|A\| = \max_{\|\underline{x}\|=1} \|\underline{x}^\top A\|$$

i.e. the norm measures the "max amt" by which the matrix A can rescale a vector \underline{x}

- We can express the defn of the operator norm as an inequality

$$\|A\underline{x}\| \leq \|A\| \|\underline{x}\| \quad \forall \underline{x}$$

It follows from this that (i) $\|\underline{x}^\top A\| \leq \|A\| \|\underline{x}\|$ and (ii) $\|\underline{x}^\top (A+B)\| \leq \|A\| \|\underline{x}\| + \|B\| \|\underline{x}\|$

$$(i) \quad \|\underline{x}^\top (A\underline{y})\| \leq \|A\| (\|\underline{x}^\top \underline{y}\|) \leq \|A\| (\|\underline{y}\| \|\underline{x}\|) \rightarrow \frac{\|\underline{x}^\top A\|}{\|\underline{x}\|} \leq \|A\| \|\underline{y}\| \quad \forall \underline{y}$$

since $\|\underline{x}^\top A\| = \max_{\underline{y} \in \mathbb{C}^{n \times 1}} \frac{\|\underline{x}^\top A\underline{y}\|}{\|\underline{y}\|}$, $\|\underline{x}^\top A\| = \frac{\|\underline{x}^\top A\|}{\|\underline{x}\|}$ for a particular \underline{x} $\rightarrow \|\underline{x}^\top A\| \leq \|A\| \|\underline{x}\|$

$$(ii) \quad \|\underline{x}^\top (A+B)\| = \|A\underline{x} + B\underline{x}\| \leq \|A\underline{x}\| + \|B\underline{x}\| \leq \|A\| \|\underline{x}\| + \|B\| \|\underline{x}\| \rightarrow \frac{\|\underline{x}^\top (A+B)\|}{\|\underline{x}\|} \leq \|A\| + \|B\| \|\underline{x}\|$$

since $\|\underline{x}^\top (A+B)\| = \max_{\underline{y} \in \mathbb{C}^{n \times 1}} \frac{\|\underline{x}^\top (A+B)\|}{\|\underline{y}\|}$, $\|\underline{x}^\top (A+B)\| = \frac{\|\underline{x}^\top (A+B)\|}{\|\underline{x}\|}$ for a particular \underline{x} $\rightarrow \|\underline{x}^\top (A+B)\| \leq \|A\| + \|B\| \|\underline{x}\|$

- To quantify the "size" of the original vector \underline{x} and the transformed vector $A\underline{x}$, we need to choose a norm

$$l\text{-1 norm: } \|A\|_1 = \max_{\underline{x} \in \mathbb{C}^{n \times 1}} \frac{\|A\underline{x}\|_1}{\|\underline{x}\|_1} = \max_{\underline{x}} \sum_{j=1}^n |a_{j1}| \quad \begin{matrix} \text{l-1 norm of the column of } A \\ \text{in the mat. } l\text{-1 norm} \end{matrix}$$

$$l\text{-2 norm: } \|A\|_2 = \max_{\underline{x} \in \mathbb{C}^{n \times 1}} \frac{\|A\underline{x}\|_2}{\|\underline{x}\|_2} = \sqrt{\lambda_{\max}(A^\top A)} \quad \begin{matrix} \text{some of the largest eigenvalue of } A^\top A \\ \text{or largest singular value of } A \end{matrix}$$

$$l\text{-}\infty \text{ norm: } \|A\|_\infty = \max_{\underline{x} \in \mathbb{C}^{n \times 1}} \frac{\|A\underline{x}\|_\infty}{\|\underline{x}\|_\infty} = \max_{\underline{x}} \sum_{j=1}^n |a_{j1}| \quad \begin{matrix} \text{l-1 norm of the row of } A \\ \text{in the mat. } l\text{-2 norm} \end{matrix}$$

* If A is Hermitian, the eigenvalues of $A^\top A = A A^\top$ are the square of eigenvalues of A , so $\|A\|_2 = \lambda_{\max}(A)$

For Personal Use Only -bkwk2

- For the ℓ_2 norm of a matrix A , $\|A\|_2$,

$$\|A\|_2^2 = \max_{\|\mathbf{x}\|=1} \frac{\|\mathbf{Ax}\|^2}{\|\mathbf{x}\|^2} = \max_{\|\mathbf{x}\|=1} \frac{(\mathbf{Ax})^H(\mathbf{Ax})}{\mathbf{x}^H \mathbf{x}} = \max_{\|\mathbf{x}\|=1} \frac{\sum_{i=1}^m \sum_{j=1}^n a_{ij}^2 x_j^2}{\mathbf{x}^H \mathbf{x}} = \frac{\sum_{i=1}^m \lambda_{\max}(A^H A)_i}{\|\mathbf{x}\|^2} = \lambda_{\max}(A^H A)$$

$$\therefore \|A\|_2 = \sqrt{\lambda_{\max}(A^H A)}$$

• Note $\|A^H A\|_2 = \frac{1}{\lambda_{\min}(A^H A)}$ since $\lambda(A^H A) = \frac{1}{\lambda(A^H A)}$, $\rightarrow \lambda_{\max}(A^H A) = \frac{1}{\lambda_{\min}(A^H A)}$

- For the ℓ_1 and ℓ_∞ norm of a matrix, $\|A\|_1$, $\|A\|_\infty$,

recall that the operator norm $\|A\|_1$ is the min. value of C that satisfies $\|\mathbf{Ax}\|_1 \leq C\|\mathbf{x}\|_1$, and see $\|\mathbf{x}\|_1 = \sum_{j=1}^n |x_j|$

(i) ℓ_1 norm: $\|A\|_1 = \left\| \sum_{j=1}^n x_j a_{ij} \right\|_1 \leq \sum_{j=1}^n \|x_j a_{ij}\|_1 = \sum_{j=1}^n |x_j| \|a_{ij}\|_1 \leq \max_{j \in \{1, n\}} \|a_{ij}\|_1$

choose a vector \mathbf{x} s.t. $x_i = \begin{cases} 1 & \text{if } i=j \\ 0 & \text{otherwise} \end{cases}$, then $\|A\|_1 = \max_{j \in \{1, n\}} \|a_{ij}\|_1 \rightarrow \|A\|_1 = \max_{j \in \{1, n\}} \|a_{ij}\|_1$

(ii) ℓ_∞ norm: $\|A\|_\infty = \max_{i \in \{1, m\}} \left\| \sum_{j=1}^n x_j a_{ij} \right\|_\infty \leq \max_{i \in \{1, m\}} \sum_{j=1}^n |x_j| |a_{ij}| \leq \max_{i \in \{1, m\}} \sum_{j=1}^n |a_{ij}|$

choose a vector \mathbf{x} s.t. $x_i = \text{sign}(a_{mi})$, then $\|A\|_\infty = \max_{i \in \{1, m\}} \sum_{j=1}^n |a_{ij}| \rightarrow \|A\|_\infty = \max_{i \in \{1, m\}} \sum_{j=1}^n |a_{ij}|$

• Note that the vector and matrix ℓ_2 norms are invariant under rotation, i.e. for a unitary matrix Q ,

$$\|Q\mathbf{x}\|_2 = \|\mathbf{x}\|_2$$

$$\|Q\mathbf{A}\|_2 = \|\mathbf{A}\|_2$$

Frobenius norm

- The Frobenius norm of a matrix A is defined as

$$\|A\|_F = \sqrt{\sum_{i,j} (a_{ij})^2}$$

It can be shown that the Frobenius norm is also given by the following.

$$\|A\|_F = \text{trace}(A^H A) = \sqrt{\sum_{i=1}^{m \times n} \sigma_i^2(A)}$$

• Note that the Frobenius norm of a matrix is invariant under rotation, i.e. for a unitary matrix Q ,

$$\|Q\mathbf{A}\|_F = \|A\|_F$$

Equivalence of norms

- choosing a norm depends on what we want to measure

↳ e.g.: ℓ_1 norm for probabilities, ℓ_2 norm for distances

- On all the vector spaces that we consider, all norms are equivalent, i.e. \exists const. $c_1 > 0, c_2 > 0$ s.t.

$$c_1 \|\mathbf{x}\|_X \leq \|\mathbf{x}\|_Y \leq c_2 \|\mathbf{x}\|_X \quad \forall \mathbf{x} \in V$$

For Personal Use Only -bkwk2

Stability and condition number

condition number.

- The condition no. of an invertible matrix A is defined by

$$K(A) = \|A\| \|A^{-1}\|$$

Note that $K(A) \geq 1$ since $K(A) = \|A\| \|A^{-1}\| \geq \|A\| \|A^{-1}\| = \|I\| = 1$

- For the ℓ_2 norm, $\|A\|_2 = \sqrt{\lambda_{\max}(A^T A)}$ and $\|A^{-1}\|_2 = \frac{1}{\sqrt{\lambda_{\min}(A^T A)}}$, so

$$K_2(A) = \|A\|_2 \|A^{-1}\|_2 = \frac{\sqrt{\lambda_{\max}(A^T A)}}{\sqrt{\lambda_{\min}(A^T A)}}$$

If A is Hermitian, then $K_2(A) = \frac{\lambda_{\max}(A)}{\lambda_{\min}(A)}$

- For small K , the shape dilates (stretch equally in all dir.)

For large K , the shape stretches in a particular dir. (in the limit collapses into a line)

Error bounds and stability.

- Consider the problem where we have an error in making the observation b

$$\underline{A} \underline{x} = \underline{b}$$

where \underline{b} is the error in the observation and $\underline{f}x$ is the consequent error in the soln.

$$\underline{A} \underline{x} = \underline{b} : \|A\| \|\underline{x}\| \geq \|\underline{A} \underline{x}\| = \|\underline{b}\| \rightarrow \frac{1}{\|\underline{x}\|} \leq \frac{\|\underline{A}\|}{\|\underline{b}\|} \quad [1]$$

$$\underline{A} \underline{f} \underline{x} = \underline{f} \underline{b} : \underline{f} \underline{x} = \underline{A}^{-1} \underline{f} \underline{b} \rightarrow \|\underline{f} \underline{x}\| = \|\underline{A}^{-1} \underline{f} \underline{b}\| \leq \|\underline{A}^{-1}\| \|\underline{f} \underline{b}\| \quad [2]$$

$$\text{Combining [1], [2], } \frac{\|\underline{f} \underline{x}\|}{\|\underline{x}\|} \leq \|\underline{A}\| \|\underline{A}^{-1}\| \frac{\|\underline{f} \underline{b}\|}{\|\underline{b}\|} = K(A) \frac{\|\underline{f} \underline{b}\|}{\|\underline{b}\|}$$

i.e. for large condition no. $K(A)$, we can expect small errors in \underline{b} to cause large errors in \underline{x} .

- Consider the problem where we have an error in the matrix A (e.g. floating pt. error in LU decomposition)

$$(A + \delta A)(x + \delta x) = b$$

where δA is the error in the matrix and δx is the consequent error in the soln.

$$(A + \delta A)(x + \delta x) = \underline{A} \underline{x} + \delta A \underline{x} + \delta A (\underline{x} + \delta x) = \underline{b} \rightarrow \delta x = \underline{A}^{-1} \delta A (\underline{x} + \delta x)$$

$$\therefore \|\delta x\| = \|\underline{A}^{-1} \delta A (\underline{x} + \delta x)\| \leq \|\underline{A}^{-1}\| \|\delta A\| \|\underline{x} + \delta x\| \rightarrow \frac{\|\delta x\|}{\|\underline{x} + \delta x\|} \leq \|\underline{A}^{-1}\| \|\delta A\| \frac{\|\underline{b}\|}{\|\underline{b}\|} = K(A) \frac{\|\delta A\|}{\|\underline{b}\|}$$

i.e. for large condition no. $K(A)$, we can expect small errors in A to cause large errors in x .

- A matrix w/ a large condition no. $K(A)$ is said to be ill-conditioned (\neq small determinant)

* Conditioning shouldn't be confused w/ the determinant (i.e. large condition no. \neq small determinant).

$$\hookrightarrow A = \begin{bmatrix} 1 & -1 & \cdots & -1 \\ 0 & 1 & \cdots & -1 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{bmatrix} \quad \det(A) = 1, \quad K_2(A) = n^{2^{n-1}}$$

$$\hookrightarrow D = \text{diag}(10^{-1}, 10^{-2}, \dots, 10^{-n}) \quad \det(D) = 10^{-n}, \quad K_2(D) = 1$$

For Personal Use Only -bkwk2

Interpolation and least squares methods

Polynomial interpolation

- If we have n data pts in a 2D space (x, y) , we can usually fit a polynomial w/ n coeffs

e.g. given n data pts $f = [f_0(x_0, y_0) \dots f_n(x_n, y_n)]^T$, we can interpolate the data pts w/ a polynomial of the form $f(x, y) = c_0 + c_1 x + c_2 y + c_3 xy + c_4 x^2 + c_5 y^2$

For each data pt $f(x_i, y_i)$, we rely on eqn, so we can write our system of eqns in matrix form

$$\begin{bmatrix} 1 & x_0 & y_0 & x_0^2 & y_0^2 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_n & y_n & x_n^2 & y_n^2 \end{bmatrix} \begin{bmatrix} c_0 \\ c_1 \\ c_2 \\ c_3 \\ c_4 \\ c_5 \end{bmatrix} = \begin{bmatrix} f_0 \\ \vdots \\ f_n \end{bmatrix}$$

- In general, there are a few issues w/ solving for c , i.e. finding the interpolating polynomial:

↳ We can only solve the matrix eqn when A is full rank, i.e. data pts are not colinear

↳ The matrix A is notoriously ill-conditioned, and the condition no. $K(A)$ grows w/ polynomial degree n .

↳ If we choose equally spaced interpolation pts, we get the Runge effect - large amplitude oscillations near the ends of the domain, which become more severe at high polynomial degrees.

Using a different set of basis functions

- So far, we have considered the monomial basis $\{1, x, \dots, x^{n-1}\}$, so we have polynomials of the form

$$f = c_{n-1} x^{n-1} + c_{n-2} x^{n-2} + \dots + c_1 x + c_0$$

where we would solve for the coefficients c by solving the matrix eqn.

$$\begin{bmatrix} 1 & x_0 & \dots & x_0^{n-2} & x_0^{n-1} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & x_n & \dots & x_n^{n-2} & x_n^{n-1} \end{bmatrix} \begin{bmatrix} c_0 \\ c_1 \\ \vdots \\ c_{n-2} \\ c_{n-1} \end{bmatrix} = \begin{bmatrix} f_0 \\ f_1 \\ \vdots \\ f_{n-2} \\ f_{n-1} \end{bmatrix}$$

- Instead, we could consider the Legendre polynomial basis $\{P_0(x), P_1(x), \dots, P_{n-1}(x)\}$, where

$$(n+1)P_{n+1}(x) = (2n+1)xP_n(x) - nP_{n-1}(x), \quad \text{and} \quad P_0 = 1, P_1 = x$$

Note that the Legendre polynomials form an orthogonal basis on the interval $[-1, 1]$, i.e.

$$\langle P_m, P_n \rangle = \int_{-1}^1 P_m(x) P_n(x) dx = 0 \quad \text{if } m \neq n$$

Therefore, we have polynomials of the form

$$f = \alpha_{n-1} P_{n-1}(x) + \alpha_{n-2} P_{n-2}(x) + \dots + \alpha_1 P_1(x) + \alpha_0 P_0(x)$$

where we would solve for the coefficients α by solving the matrix eqn.

$$\begin{bmatrix} P_0(x_0) & P_1(x_0) & \dots & P_{n-2}(x_0) & P_{n-1}(x_0) \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ P_0(x_{n-1}) & P_1(x_{n-1}) & \dots & P_{n-2}(x_{n-1}) & P_{n-1}(x_{n-1}) \end{bmatrix} \begin{bmatrix} \alpha_0 \\ \alpha_1 \\ \vdots \\ \alpha_{n-2} \\ \alpha_{n-1} \end{bmatrix} = \begin{bmatrix} f_0 \\ f_1 \\ \vdots \\ f_{n-2} \\ f_{n-1} \end{bmatrix}$$

- Using the Legendre polynomials instead of monomials as the basis significantly reduces the condition no. $K(A)$ of the Vandermonde matrix A .

For Personal Use Only -bkwk2

Interpolation w/ non-equispaced points.

- The Lebesgue const. $\Lambda(T)$ bounds the interpolation error. For a function f , where p_n is the best approx of f among all polynomials of degree n and p_0 is the degree n polynomial interpolant of f ,

$$\|f - p_0\| \leq (\Lambda(T) + 1) \|f - p_n\|$$

Note that the Lebesgue const $\Lambda(T)$ depends on the choice of interpolation pt's T .

- For a good polynomial interpolant p_0 , we want the Lebesgue const $\Lambda(T)$ to be as small as possible.
- By using non-equispaced sampling pts - roots of orthogonal polynomials in particular (they tend to cluster close to the ends of the interval), we can mitigate the Runge effect (i.e. small $\Lambda(T)$)
- The Lebesgue const $\Lambda(T)$ grows exponentially in n for equispaced pts but only logarithmically in n for Chebyshev pts (roots of Chebyshev polynomials, defined as $T_n(x) = \cos(n \arccos(x))$, $x \in [-1, 1]$).

Least squares method for data fitting of over-determined systems.

- To find a sol'n to the matrix eqn $\underline{A}\underline{x} = \underline{b}$, the vector \underline{b} must lie in the column space of \underline{A} . If \underline{A} is a mxn matrix w/ $m > n$ (skinny matrix), in general \underline{b} does not lie in the column space of \underline{A} .
- For an approx sol'n $\underline{\hat{x}}$, we can define the residual vector $\underline{\Sigma}$ as

$$\underline{\Sigma} = \underline{A}\underline{\hat{x}} - \underline{b}$$

We want to choose $\underline{\hat{x}}$ s.t. the norm of the residual vector $\|\underline{\Sigma}\|$ is min. for some norm, i.e.

$$\underline{\hat{x}} = \underset{\underline{x} \in \mathbb{C}^n}{\operatorname{arg\,min}} \|\underline{\Sigma}(\underline{x})\| = \underset{\underline{x} \in \mathbb{C}^n}{\operatorname{arg\,min}} \|\underline{A}\underline{x} - \underline{b}\|$$

- If we use the L_2 norm, we seek $\underline{\hat{x}} = \underset{\underline{x} \in \mathbb{C}^n}{\operatorname{arg\,min}} \|\underline{A}\underline{x} - \underline{b}\|_2$. Let $r(\underline{x}) = \|\underline{A}\underline{x} - \underline{b}\|_2^2$, we have

$$r(\underline{x}) = \|\underline{A}\underline{x} - \underline{b}\|_2^2 = \sum_{i=1}^m \left| \sum_{j=1}^n (a_{ij} \underline{x}_j) - b_i \right|^2 \rightarrow \text{least squares method}$$

Extending and minimizing $r(\underline{x})$,

$$r(\underline{x}) = \|\underline{A}\underline{x} - \underline{b}\|_2^2 = (\underline{A}\underline{x} - \underline{b})^H (\underline{A}\underline{x} - \underline{b}) = \underline{x}^H \underline{A}^H \underline{A} \underline{x} - \underline{x}^H \underline{A}^H \underline{b} - \underline{b}^H \underline{A} \underline{x} + \underline{b}^H \underline{b}$$

$$\therefore \frac{d r(\underline{x})}{d \underline{x}} = 2 \underline{x}^H \underline{A}^H \underline{A} - 2 \underline{A}^H \underline{b} = 0 \rightarrow \underline{A}^H \underline{A} \underline{x} = \underline{A}^H \underline{b}$$

Rearranging the least squares problem, we have

$$\underline{\hat{x}} = (\underline{A}^H \underline{A})^{-1} \underline{A}^H \underline{b} = \underline{A}^+ \underline{b}$$

where $\underline{A}^+ = (\underline{A}^H \underline{A})^{-1} \underline{A}^H$ is the pseudoinverse / Penrose-Moore inverse.

- a The inverse $(\underline{A}^H \underline{A})^{-1}$ can be computed when \underline{A} is full rank

- To solve the least squares problem $\underline{A}^H \underline{A} \underline{\hat{x}} = \underline{A}^H \underline{b}$, if \underline{A} is full rank, we can apply LU factorisation or Cholesky factorisation to $\underline{A}^H \underline{A}$, but $\underline{A}^H \underline{A}$ is notoriously ill-conditioned

- We can use more specialised methods like QR factorisation or SVD to mitigate the ill-conditioning of $\underline{A}^H \underline{A}$ (QR may have accuracy/stability problems and SVD is expensive)

- * In practice, do not write your own least-squares solver → use specialised functions/libraries

For Personal Use Only -bkwk2

Iterative methods for linear systems

Direct and iterative methods.

- To solve a system of the form $\underline{A}\underline{x} = \underline{b}$, solns can be categorised as direct or iterative.

① Direct methods.

- A direct method computes a soln to $\underline{A}\underline{x} = \underline{b}$ in a known/predictable no. of operations.
W/o round-off errors, the soln is exact. (e.g. LU factorisation)
- In general, direct methods cannot handle singular eqns.
- For reasonably conditioned systems, direct solvers can be robust.
- Direct solvers do not scale well on parallel computers as substitution steps are serial.

② Iterative methods.

- An iterative method ~~seeks~~ an approx'dn via a series of iterations, and is usually terminated when the error/residual falls below a prescribed threshold.
- Some iterative methods can solve singular eqns.
- Iterations can be terminated early if high accuracy is not reqd.
- Iterative methods use less memory and are suited to large parallel computers.

Power iteration for eigenvector w/ largest eigenvalue

- we can use power iteration to find the eigenvector associated w/ the largest absolute eigenvalue.

- A vector $\underline{x} \in \mathbb{C}^n$ can be expressed in terms of the n eigenvectors of an $n \times n$ matrix.

$$\underline{x} = \sum_{i=1}^n \alpha_i \underline{u}_i \quad \text{where } \alpha_i \in \mathbb{C},$$

If we multiply \underline{x} repeatedly by \underline{A} ,

$$\underline{A}^k \underline{x} = \sum_{i=1}^n \alpha_i \lambda_i^k \underline{u}_i$$

∴ if the largest eigenvalue is distinct, the resulting vector will be aligned w/ the eigenvector of the largest eigenvalue (given \underline{x} has a component in the dir. of the eigenvector)

- similar methods inc. the inverse power method $\underline{A}^T \underline{x} = \frac{1}{\lambda} \underline{x} \rightarrow$ eigenvector w/ smallest eigenvalue

and the shifted inverse power method $(\underline{A} - \underline{\alpha} \underline{I})^{-1} = \frac{1}{\lambda - \alpha} \underline{x} \rightarrow$ eigenvector w/ eigenvalue closest to α .

* These methods can converge slowly

For Personal Use Only -bkwk2

Rayleigh quotient for eigenvalues

- we can use the Rayleigh quotient $R(\underline{A}, \underline{x})$ to find the eigenvalue λ^* correspondingly to the given estimate of an eigenvector \underline{x} .

- Let \underline{x} be an estimate of an eigenvector of the matrix \underline{A} , so $\underline{A}\underline{x} \approx \lambda^*\underline{x}$.

To find an estimate of the corresponding eigenvalue λ^* , we consider the minimisation problem in ℓ_2 norm,

$$\lambda^* = \underset{\lambda \in \mathbb{C}}{\operatorname{arg\,min}} \|\underline{A}\underline{x} - \lambda \underline{x}\|_2$$

This is minimised when λ^* is equal the Rayleigh quotient $R(\underline{A}, \underline{x})$,

$$R(\underline{A}, \underline{x}) = \frac{\underline{x}^\top \underline{A} \underline{x}}{\underline{x}^\top \underline{x}}$$

- To prove this, consider any $\mu \in \mathbb{C}$,

$$\begin{aligned} \|\underline{A}\underline{x} - \mu \underline{x}\|_2^2 &= \underline{x}^\top \underline{A}^\top \underline{A} \underline{x} - \mu \underline{x}^\top \underline{A}^\top \underline{A} \underline{x} - \bar{\mu} \underline{x}^\top \underline{A} \underline{x} + \mu^2 \underline{x}^\top \underline{x} = \|\underline{A} \underline{x}\|_2^2 + \|\underline{x}\|_2^2 (\mu^2 - \mu \bar{\mu} - \bar{\mu} \mu) \\ &= \|\underline{A} \underline{x}\|_2^2 + \|\underline{x}\|_2^2 ((\mu - \bar{\mu})(\bar{\mu} - \mu) - |\mu|^2) \geq \|\underline{A} \underline{x}\|_2^2 - |\mu|^2 \|\underline{x}\|_2^2 = \|\underline{A} \underline{x} - \mu \underline{x}\|_2^2 \end{aligned}$$

- The Rayleigh quotient $R(\underline{A}, \underline{x})$ is often defined for Hermitian matrices only, where it is real-valued.

If also has other properties, inc. second-order accuracy for estimating eigenvalue
eigenvector error $O(\epsilon)$, eigenvalue error $O(\epsilon)$

stationary methods for $\underline{A}\underline{x} = \underline{b}$

- consider decomposing the matrix operator s.t. $\underline{A} = \underline{N} - \underline{P}$. The matrix eqn. $\underline{A}\underline{x} = \underline{b}$ becomes,

$$(\underline{N} - \underline{P})\underline{x} = \underline{b} \quad \rightarrow \quad \underline{N}\underline{x} = \underline{b} + \underline{P}\underline{x}$$

We can compute an approx soln \underline{x}_K by

$$\underline{N}\underline{x}_K = \underline{b} + \underline{P}\underline{x}_K$$

where \underline{x}_K is the most recent known estimate of the soln

- we keep iterating for the new estimate \underline{x}_{K+1} and hopefully this converges to the exact soln.

- The key is to split & s.t. eqns of the form $\underline{N}\underline{x} = \underline{f}$ are easy (cheap) to solve, examples inc.

↳ Richardson : $\underline{N} = \underline{I}$

↳ Jacobi : $\underline{N} = \operatorname{diag}(\underline{A})$.

↳ Gauss-Seidel : $\underline{N} = L(\underline{A})$ (lower triangular part of \underline{A} , inc. the diagonal).

- Defining the error at the K th iterate, as $e_K = \underline{x}_{\text{exact}} - \underline{x}_K$, we have,

$$\underline{N}(\underline{x}_{\text{exact}} - \underline{x}_K) = (\underline{N} - \underline{P})\underline{x}_{\text{exact}} + \underline{P}(\underline{x}_{\text{exact}} - \underline{x}_K) \rightarrow \underline{N}e_K = \underline{P}e_K.$$

$$\therefore e_{K+1} = \underline{N}^{-1} \underline{P} e_K$$

Expressing the initial error vector e_0 as a linear combination of the eigenvectors of $\underline{N}^{-1}\underline{P}$,

$$e_0 = c_1 \underline{u}_1 + \dots + c_n \underline{u}_n$$

The error e_K is therefore

$$e_K = (\underline{N}^{-1}\underline{P})^K e_0 = c_1 \lambda_1^K \underline{u}_1 + \dots + c_n \lambda_n^K \underline{u}_n$$

- For convergence, we req. $\underline{P} \rightarrow 0$ as K increases, which is true if the absolute value of every eigenvalue of $\underline{N}^{-1}\underline{P}$ is less than one — i.e. the spectral radius of $\underline{N}^{-1}\underline{P}$ is less than 1, $|\rho(\underline{N}^{-1}\underline{P})| < 1$.

* The smaller $p < 1$, the faster the convergence.

For Personal Use Only -bkwk2

Conjugate gradient (CG) method

- The CG method is a Krylov subspace method for solving $\Delta x = b$, where Δ is Hermitian positive definite, and is more powerful than the stationary methods outlined above.
- The CG method is often presented as a direct method, but is applied in practice usually as an iterative method.

① CG method as a direct method.

- consider $\Delta x = b$, where Δ is a non-Hermitian positive-definite matrix.

Consider that we have a set P of n non-zero vectors $P = \{p_1, \dots, p_m\}$ that are A -conjugate, i.e.,

$$p_i^H \Delta p_j = 0 \quad \text{if } i \neq j.$$

Note that since Δ is positive definite, $p_i^H \Delta p_i > 0$

- Using P as a basis for the soln, $x = \sum_{i=0}^m \alpha_i p_i$, so we have

$$b = \Delta x = \sum_{i=1}^m \alpha_i \Delta p_i$$

Pre-multiplying by P^H

$$P^H b = P^H \Delta x = \sum_{i=1}^m \alpha_i p_i^H \Delta p_i = \alpha_1 p_1^H \Delta p_1$$

We can therefore solve for all α_i and get $x = \sum_{i=1}^m \alpha_i p_i$, where α_i is given by

$$\alpha_i = \frac{p_i^H b}{p_i^H \Delta p_i}$$

- To generate the A-conjugate set P , we can simply pick n linearly independent vectors and apply the Gram-Schmidt process to build P , but this is not practical for large problems build + store
P too expensive

② CG method as an iterative method.

- The CG algorithm involves computing elements of P using a shift difference eqn. Moreover, we may not even need all of P - we can get an accurate soln w/ a few elements of P .

- The CG algorithm only involves matrix-vector products and dot products, and only req. storing a few work vectors (low space complexity)

- When solving $\Delta x = b$ where Δ is a non-Hermitian matrix:

↳ The error $e_k = x_{\text{exact}} - x_k$ is monotone in the A -norm. ($\|e_k\|_A = \sqrt{x_k^H \Delta x_k}$), i.e. $\|e_{k+1}\|_A \leq \|e_k\|_A$

and converges to 0, i.e. $\|e_k\|_A$ for some $k \leq n$. (in exact arithmetic)

↳ The no. of iterations req. to solve x exactly is equal to the no. of distinct eigenvalues of Δ (so $k \leq n$)

↳ The rate of convergence is affected by the condition no. $\kappa(\Delta)$,

$$\frac{\|e_{k+1}\|_A}{\|e_k\|_A} \leq 2 \left(\frac{\sqrt{\kappa(\Delta)} - 1}{\sqrt{\kappa(\Delta)} + 1} \right)$$

Note when $\kappa(\Delta)$ is close to 1, e_k converges to 0 quickly.

- If the condition no. $\kappa(\Delta)$ is too large, the CG method may converge too slowly for practical use

→ we can pre-condition the problem, where we consider the transformed system

left preconditioning. → $P^{-1} \Delta x = P^{-1} b$

If $P^{-1} \approx \Delta^{-1}$, the condition no. of $P^{-1} \Delta$ is much better than Δ → CG method converges faster

→ In practice, it can be a balancing act - P^{-1} must be cheap to apply but close enough to Δ^{-1} to be effective.

For Personal Use Only -bkwk2

Singular value decomposition (SVD)

Matrix diagonalisation.

- We can diagonalise any Hermitian matrix $\underline{M} \in \mathbb{C}^{n \times n}$

$$\underline{M} = \underline{Q} \underline{\Delta} \underline{Q}^H$$

where the columns of \underline{Q} are the ℓ_2 -normalised eigenvectors of \underline{M} (which are orthogonal $\rightarrow \underline{Q}$ is unitary).

and $\underline{\Delta}$ is the diagonal matrix of the eigenvalues of \underline{M} (which are real).

- We can interpret the diagonalisation as a sum of rank 1 matrices.

$$\underline{M} = \sum_i \lambda_i \underline{u}_i \underline{u}_i^H$$

where $(\lambda_i, \underline{u}_i)$ is the i th eigenpair of \underline{M} .

- Note that the eigenvectors are not unique \rightarrow sign ambiguity $\underline{M} = \underline{Q} \underline{\Delta} \underline{Q}^H = (-\underline{Q}) \underline{\Delta} (\underline{Q})^H$

- Matrix diagonalisation is v. helpful as we see that in a particular basis, the action of a matrix on a vector is easy to interpret, and matrix-matrix operations become trivial.

- However, matrix diagonalisation has a few major drawbacks

\hookrightarrow only valid for square matrices

\hookrightarrow defective matrices cannot be diagonalised

\hookrightarrow The matrix \underline{Q} is guaranteed unitary only for Hermitian matrices \underline{M} .

Singular value decomposition (SVD)

- The SVD of a $m \times n$ matrix \underline{A} is

$$\underline{A} = \underline{U} \Sigma \underline{V}^H$$

where $\underline{U} \in \mathbb{C}^{m \times m}$ is a unitary matrix, $\underline{V} \in \mathbb{C}^{n \times n}$ is a unitary matrix

$\downarrow p = \min(m, n)$

$\Sigma \in \mathbb{R}^{m \times n}$ is a diagonal matrix, w/ diagonal entries or (singular values) sorted by $0 \geq \sigma_1 \geq \dots \geq \sigma_p > 0$

- Consider $\underline{S}_R = \underline{\Sigma}^H \underline{A}$ (which is Hermitian)

$$\underline{S}_R = \underline{\Sigma}^H \underline{A} = (\underline{U} \Sigma \underline{V}^H)^H (\underline{U} \Sigma \underline{V}^H) = \underline{U} \Sigma^H \underline{U} \Sigma \underline{V}^H = \underline{U} \Sigma^H \underline{V}^H$$

\rightarrow The columns of \underline{U} are the normalised eigenvectors of $\underline{A}^H \underline{A}$

The diagonal entries of $\Sigma^H \Sigma$ $\sigma_1^2, \sigma_2^2, \dots, \sigma_p^2$ are the eigenvalues of $\underline{A}^H \underline{A}$

- Consider $\underline{S}_L = \underline{A} \underline{\Sigma}^H$ (which is Hermitian)

$$\underline{S}_L = \underline{A} \underline{\Sigma}^H = (\underline{U} \Sigma \underline{V}^H) (\underline{U} \Sigma \underline{V}^H)^H = \underline{U} \Sigma \underline{V}^H \underline{U} \Sigma^H \underline{V}^H = \underline{U} \Sigma^H \underline{U}^H$$

\downarrow some eigenvalues

\rightarrow The columns of \underline{V} are the normalised eigenvectors of $\underline{A} \underline{A}^H$

The diagonal entries of $\Sigma^H \Sigma$ $\sigma_1^2, \sigma_2^2, \dots, \sigma_p^2$ are the eigenvalues of $\underline{A} \underline{A}^H$

- To determine the sign of the eigenvectors, postmultiply both sides of $\underline{A} = \underline{U} \Sigma \underline{V}^H$ by \underline{U} , so

$$\underline{A} \underline{U} = \underline{U} \Sigma$$

\hookrightarrow

$$\underline{A} \underline{u}_i = \sigma_i \underline{v}_i$$

where \underline{u}_i is the i th column of \underline{U} and \underline{v}_i is the i th column of \underline{V} .

* Every matrix has a SVD, which is unique up to the signs of the eigenvectors in \underline{U} and \underline{V} .

* The SVD of a Hermitian matrix is the usual matrix diagonalisation.

For Personal Use Only -bkwk2

The reduced SVD.

- For a $m \times n$ matrix w/ $m > n$ (skinny matrix), the SVD can be reduced by removing redundant entries

$$\begin{bmatrix} A \\ m \times n \end{bmatrix} = \begin{bmatrix} U \\ m \times n \end{bmatrix} \begin{bmatrix} \Sigma \\ m \times n \\ 0 \end{bmatrix} \begin{bmatrix} V^T \\ n \times n \end{bmatrix} \rightarrow \begin{bmatrix} A \\ m \times n \end{bmatrix} = \begin{bmatrix} U \\ m \times n \end{bmatrix} \begin{bmatrix} \Sigma \\ m \times n \end{bmatrix} \begin{bmatrix} V^T \\ n \times n \end{bmatrix}$$

No contribution

The terms below row n in Σ are always zero \rightarrow last $m-n$ columns of U make no contribution.

- For a $m \times n$ matrix w/ $m < n$ (fat matrix), the SVD can be reduced by removing redundant entries

$$\begin{bmatrix} A \\ m \times n \end{bmatrix} = \begin{bmatrix} U \\ m \times n \end{bmatrix} \begin{bmatrix} \Sigma \\ m \times n \\ 0 \end{bmatrix} \begin{bmatrix} V^T \\ n \times n \end{bmatrix} \rightarrow \begin{bmatrix} A \\ m \times n \end{bmatrix} = \begin{bmatrix} U \\ m \times n \end{bmatrix} \begin{bmatrix} \Sigma \\ m \times n \end{bmatrix} \begin{bmatrix} V^T \\ n \times n \end{bmatrix}$$

No contribution

The terms right of column n in Σ are always zero \rightarrow last $n-m$ rows of V^T make no contribution

Low rank approximations

- If we expand the SVD $\underline{A} = \underline{U} \Sigma \underline{V}^T$, we get a sum of rank-1 matrices.

$$\underline{A} = \sum_{i=1}^r \sigma_i \underline{u}_i \underline{v}_i^T$$

where r is the rank of \underline{A} (no. of non-zero singular values),

and \underline{u}_i is the i th column of \underline{U} , \underline{v}_i is the i th column of \underline{V} .

- Noting that σ_i are sorted in descending order, a rank- k approx of \underline{A} is given by

$$\underline{A}_k = \sum_{i=1}^k \sigma_i \underline{u}_i \underline{v}_i^T, \quad k < r$$

\underline{A}_k is a better approx of \underline{A} than any other matrix of rank k or less, \underline{B}_k .

$$(i) \|\underline{A} - \underline{A}_k\|_F = \sqrt{\sigma_{k+1}^2 + \dots + \sigma_r^2} \leq \|\underline{A} - \underline{B}_k\|_F \quad \text{or} \quad (ii) \|\underline{A} - \underline{A}_k\|_2 = \sigma_{k+1} \leq \|\underline{A} - \underline{B}_k\|_2$$

→ proof of (ii): $\|\underline{A} - \underline{A}_k\|_2^2 = \sum_{i=k+1}^r \sigma_i^2 \|\underline{u}_i \underline{v}_i^T\|_2^2 = \sum_{i=k+1}^r \sigma_i^2$

$$\therefore \|\underline{A} - \underline{A}_k\|_2 = \left\| \sum_{i=k+1}^r \sigma_i \underline{u}_i \underline{v}_i^T \right\|_2 = \sigma_{k+1} \quad \begin{array}{l} \text{l-2 norm of a matrix is} \\ \text{equal to the largest singular value} \end{array}$$

As \underline{B}_k has rank k , there exists a vector $w = \alpha_1 \underline{v}_1 + \dots + \alpha_k \underline{v}_k$ that lies in the null space of \underline{B}_k .

Denoting Σ_{k+1} as the first $k+1$ entries of Σ and w as $w = \begin{bmatrix} \uparrow & \uparrow \\ \underline{v}_1 & \dots & \underline{v}_{k+1} \\ \downarrow & \downarrow \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_{k+1} \end{bmatrix} = \underline{V}_{k+1} \alpha$,

for convenience, normalise w s.t. $\|w\|_2 = 1$ (so $|w|_1 < 1$)

$$\|\underline{A}\|_2^2 \|\underline{A} - \underline{B}_k\|_2^2 \geq \|\underline{A} - \underline{B}_k\|_F^2 = \left\| \underline{A} \underline{w} - \underline{B}_k \underline{w} \right\|_2^2 = \left\| \underline{U} \Sigma \underline{V}^T \underline{V}_{k+1} \alpha \right\|_2^2 = \left\| \sum_{i=k+1}^r \sigma_i \underline{u}_i \right\|_2^2$$

norm invariant under rotation
by construction of \underline{w}

$$= \sigma_1^2 \alpha_1^2 + \dots + \sigma_{k+1}^2 \alpha_{k+1}^2 \geq \sigma_{k+1}^2 = \|\underline{A} - \underline{A}_k\|_2^2 \rightarrow \|\underline{A} - \underline{B}_k\|_2 \geq \|\underline{A} - \underline{A}_k\|_2$$

Effective rank

- When taking measurements, noise is often unavoidable \rightarrow hard to detect linear dependencies.

→ e.g. $\underline{A} = \begin{bmatrix} 1 & 1 & 1 \\ 2 & 2 & 2 \\ 1 & 0 & 1 \end{bmatrix}, \quad \underline{A}_{\text{noisy}} = \begin{bmatrix} 1.0000022 & 1.0000006 & 1.0000014 \\ 2.0000018 & 2.0000015 & 2.0000011 \\ 1.0000008 & 2.1924 \times 10^{-6} & 1.0000012 \end{bmatrix} \quad (\text{noise } \in [0, 10^{-6}])$

rank 3

SVD on $\underline{A}_{\text{noisy}}$ gives $\sigma_1 = 4.048, \sigma_2 = 0.7811, \sigma_3 = 3.9 \times 10^{-7} \rightarrow$ effective rank = 2

- The effective rank is the no. of singular values that are greater than the noise level.

For Personal Use Only -bkwk2

least squares solution — full rank case.

- If \underline{A} is a $m \times n$ matrix w/ $m > n$, we can partition the SVD as

$$\underline{A} = [\underline{U}_1 \underline{U}_2] \begin{bmatrix} \Sigma \\ 0 \end{bmatrix} \underline{V}^H = \underline{U}_1 \Sigma \underline{V}^H$$

since \underline{A} is full rank, Σ will be non-zero, and \underline{U}_1 is $m \times n$.

- * Note the columns of \underline{V} are orthonormal $\rightarrow \underline{V}^H \underline{V} = I_{nn}$ and $\underline{U}_2^H \underline{U}_2 = \Omega_{(m-n) \times n}$

- A least squares soln minimises $r = \|\underline{A}\underline{x} - \underline{b}\|_2^2$

$$r = \|\underline{A}\underline{x} - \underline{b}\|_2^2 = \left\| \begin{bmatrix} \underline{U}_1^H \\ \underline{U}_2^H \end{bmatrix} (\underline{A}\underline{x} - \underline{b}) \right\|_2^2 = \left\| \begin{bmatrix} \underline{U}_1^H \\ \underline{U}_2^H \end{bmatrix} (\underline{U}_1 \Sigma \underline{V}^H \underline{x} - \underline{b}) \right\|_2^2 = \left\| \begin{bmatrix} \Sigma \underline{V}^H \underline{x} - \underline{U}_1^H \underline{b} \\ -\underline{U}_2^H \underline{b} \end{bmatrix} \right\|_2^2 = \|\Sigma \underline{V}^H \underline{x} - \underline{U}_1^H \underline{b}\|_2^2 + \|\underline{U}_2^H \underline{b}\|_2^2$$

The residual r is minimised when $\Sigma \underline{V}^H \underline{x} = \underline{U}_1^H \underline{b}$, i.e.

$$\hat{\underline{x}} = \underline{U}_1 \Sigma^{-1} \underline{U}_1^H \underline{b}$$

Both Σ and \underline{V} are full rank, so the least-squares soln is unique. We could add "zero-pooling" to \underline{U}, Σ ,

$$\hat{\underline{x}} = \underline{U} \Sigma^{-1} \underline{U}^H \underline{b}, \quad \text{where } \underline{\Sigma}^{-1} = \begin{bmatrix} \Sigma_1 & \dots & \Sigma_n \\ \vdots & \ddots & \vdots \end{bmatrix}$$

- We can rederive the least-squares expression by expanding $\hat{\underline{x}} = \underline{A}^T \underline{b}$, $\underline{A}^T = (\underline{U}^H \Sigma) \underline{U}^H$, $\underline{A} = \underline{U} \Sigma \underline{V}^H$

$$\begin{aligned} \underline{A}^T (\underline{A} \hat{\underline{x}}) \underline{A}^T &= (\underline{U} \Sigma \underline{V}^H \underline{U}^H \Sigma \underline{V}^H) \underline{V} \Sigma^{-1} \underline{U}^H = (\underline{V} \Sigma^{-1} \underline{U}^H) \underline{V} \Sigma^{-1} \underline{U}^H \\ &= \underline{V} \Sigma^{-1} \underline{U}^H \underline{U} \Sigma \underline{V}^H \underline{U}^H = \underline{V} \Sigma^{-1} \underline{U}^H \Sigma \underline{U}^H \underline{V} = \underline{V} \Sigma^{-1} \underline{U}^H \underline{V} \rightarrow \hat{\underline{x}} = \underline{V} \Sigma^{-1} \underline{U}^H \underline{b} = \underline{U} \Sigma^{-1} \underline{U}^H \underline{b} \end{aligned}$$

- Consider the L-2 norm of the least-squares soln, squared, $\|\hat{\underline{x}}\|_2^2$

$$\|\hat{\underline{x}}\|_2^2 = (\underline{V} \Sigma^{-1} \underline{U}^H \underline{b})^H \underline{V} \Sigma^{-1} \underline{U}^H \underline{b} = \underline{b}^H \underline{U} \Sigma^{-1} \underline{U}^H \Sigma^{-1} \underline{U}^H \underline{b} = \|\Sigma^{-1} \underline{U}^H \underline{b}\|_2^2 \geq |(\Sigma^{-1} \underline{U}^H \underline{b})_1| = \|\underline{U}^H \underline{b}\|_1$$

i.e. if the smallest singular value σ_m is small, then the least-squares soln will be large and sensitive to changes in \underline{b} \rightarrow poor conditioning // if $\sigma_m = 0$, we cannot compute a unique least-squares soln.

least squares solution — rank deficient case.

- If \underline{A} is a $m \times n$ matrix w/ zero singular values, we can partition the SVD as

$$\underline{A} = [\underline{U}_1 \underline{U}_2] \begin{bmatrix} \Sigma & 0 \\ 0 & 0 \end{bmatrix} [\underline{V}_1 \underline{V}_2]^H = \underline{U}_1 \Sigma \underline{V}_1^H$$

where Σ has size $r \times r$, r is the rank of the matrix \underline{A} .

- * Note the columns of \underline{V} are orthonormal $\rightarrow \underline{V}_1^H \underline{V}_1 = I_{rr}$ and $\underline{V}_2^H \underline{V}_2 = \Omega_{(m-r) \times (m-r)}$

- A least-squares soln minimises $r = \|\underline{A}\underline{x} - \underline{b}\|_2^2$

$$r = \|\underline{A}\underline{x} - \underline{b}\|_2^2 = \left\| \begin{bmatrix} \underline{U}_1^H \\ \underline{U}_2^H \end{bmatrix} (\underline{A}\underline{x} - \underline{b}) \right\|_2^2 = \left\| \begin{bmatrix} \underline{U}_1^H \\ \underline{U}_2^H \end{bmatrix} (\underline{U}_1 \Sigma \underline{V}_1^H \underline{x} - \underline{b}) \right\|_2^2 = \left\| \begin{bmatrix} \Sigma \underline{V}_1^H \underline{x} - \underline{U}_1^H \underline{b} \\ -\underline{U}_2^H \underline{b} \end{bmatrix} \right\|_2^2 = \|\Sigma \underline{V}_1^H \underline{x} - \underline{U}_1^H \underline{b}\|_2^2 + \|\underline{U}_2^H \underline{b}\|_2^2$$

The residual r is minimised when $\Sigma \underline{V}_1^H \underline{x} = \underline{U}_1^H \underline{b}$, i.e.

$$\hat{\underline{x}} = \underline{U}_1 \Sigma^{-1} \underline{U}_1^H \underline{b}$$

- Note \underline{V}_1 is only a part of $\underline{V} \rightarrow$ has a non-trivial nullspace, specifically $\underline{V}_1^H (\underline{V}_2 \underline{z}) = 0$ for all $\underline{z} \in \mathbb{C}^{m-r}$

since the columns of \underline{V}_1 and \underline{V}_2 are orthogonal. We thus have infinitely many solns satisfying $\Sigma \underline{V}_1^H \underline{x} = \underline{U}_1^H \underline{b}$

$$\hat{\underline{x}} = \underline{U}_1 \Sigma^{-1} \underline{U}_1^H \underline{b} + \underline{V}_2 \underline{z}$$

Taking the norm and noting that \underline{V}_1 and \underline{V}_2 are orthogonal wrt each other,

$$\|\hat{\underline{x}}\|_2^2 = \|\underline{U}_1 \Sigma^{-1} \underline{U}_1^H \underline{b} + \underline{V}_2 \underline{z}\|_2^2 = \|\underline{U}_1 \Sigma^{-1} \underline{U}_1^H \underline{b}\|_2^2 + \|\underline{V}_2 \underline{z}\|_2^2$$

$\therefore \hat{\underline{x}} = \underline{U}_1 \Sigma^{-1} \underline{U}_1^H \underline{b}$ is the minimiser to the least-squares problem w/ minimal L-2 norm

For Personal Use Only -bkwk2

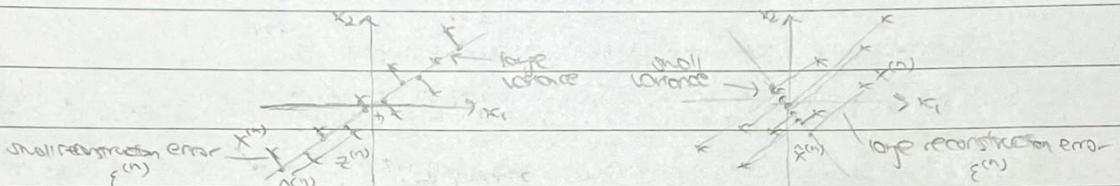
Fitting points and singular systems.

- To solve for the coefficients of the interpolating polynomial, we need to solve the matrieqn. $\Delta \mathbf{c} = \mathbf{f}$. If the data pts happen to lie on a line, then Δ will be singular.
- When Δ is singular, $\Delta^H \Delta$ is also singular, i.e. $(\Delta^H \Delta)^{-1}$ doesn't exist \rightarrow cannot find the least squares soln by solving $\Delta^H \Delta \mathbf{c} = \Delta^H \mathbf{f}$.
- To find the least squares soln, we perform SVD of Δ and compute $\mathbf{c} = \mathbf{U} \Sigma^{-1} \mathbf{V}^H \mathbf{f}$
- If the matrix Δ is near singular, we compute the SVD, ignore singular values less than the noise level, then compute $\mathbf{c} = \mathbf{U} \Sigma^{-1} \mathbf{V}^H \mathbf{f}$ as usual (unstable if we do not ignore those singular values).

Principal component analysis (PCA)

Principal component analysis (PCA).

- PCA is a dimension reduction technique. It finds new, uncorrelated variables from a dataset.
- The "new" variables are the dir. along which the variance of the sample is greatest / the reconstruction error is smallest. They can be defined as linear combination of the original variables, both interpretations are equivalent



- consider projecting from $\underline{x}^{(n)} \in \mathbb{R}^D$ to $\underline{z}^{(n)} \in \mathbb{R}^M$, w/ $M < D$. choose M projection vectors $w_1, \dots, w_M \in \mathbb{R}^D$, s.t. they are unit length $\|w_i\|=1$ and orthogonal, $w_i^T w_j = 0 \quad \forall i \neq j$. Also $\underline{x}^{(1)}, \dots, \underline{x}^{(N)}$ have zero mean.
- The projection of the nth item $\underline{x}^{(n)}$ onto the Mth dimension is $z_m^{(n)} = w_m^T \underline{x}^{(n)}$, so in general, $\underline{x}^{(n)}$ is mapped to $\underline{z}^{(n)} = \begin{bmatrix} z_1^{(n)} \\ \vdots \\ z_M^{(n)} \end{bmatrix} = \begin{bmatrix} w_1^T \underline{x}^{(n)} \\ \vdots \\ w_M^T \underline{x}^{(n)} \end{bmatrix} = \begin{bmatrix} \leftarrow w_1^T \rightarrow \\ \vdots \\ \leftarrow w_M^T \rightarrow \end{bmatrix} \begin{bmatrix} x_1^{(n)} \\ \vdots \\ x_D^{(n)} \end{bmatrix} = w^T \underline{x}^{(n)}$
- For construction, we have $\underline{x}^{(n)} = \underline{w}^T \underline{z}^{(n)}$
- The (j,k)th entry of the sample covariance matrix is given by $S_{jk} = \frac{1}{N-1} \sum_{n=1}^N x_j x_k$, given $\bar{x}_i = 0 \quad \forall i$
- In matrix form, the sample covariance matrix is given by $\underline{\Sigma} = \underline{X}^T \underline{X}$, for $\underline{X} = \begin{bmatrix} \leftarrow \underline{x}_1^T \rightarrow \\ \vdots \\ \leftarrow \underline{x}_N^T \rightarrow \end{bmatrix}$
- The sample variance of direction w_1 , $s_{z_1}^2$, is given by note $\underline{\Sigma} = \frac{1}{N-1} \underline{X}^T \underline{X} = \frac{1}{N-1} \sum_{n=1}^N \underline{x}^{(n)} \underline{x}^{(n)T}$

$$\hat{s}_{z_1}^2 = \frac{1}{N-1} \sum_{n=1}^N (z_1^{(n)} - \bar{z})^2 = \frac{1}{N-1} \sum_{n=1}^N (\underline{x}^{(n)T} w_1)^2 = \frac{1}{N-1} \sum_{n=1}^N w_1^T \underline{x}^{(n)} \underline{x}^{(n)T} w_1 = w_1^T \left(\frac{1}{N-1} \sum_{n=1}^N \underline{x}^{(n)} \underline{x}^{(n)T} \right) w_1 = w_1^T \underline{\Sigma} w_1$$

To max. this quantity, we choose w_1 to be the eigenvector of $\underline{\Sigma} = \frac{1}{N-1} \underline{X}^T \underline{X}$ that corresponds to the largest eigenvalue.

(we find w_2, w_3 by finding the eigenvectors of $\underline{\Sigma}$ that correspond to the second/third largest eigenvalue)

- Note that the eigenvalues and eigenvectors of $\underline{X}^T \underline{X}$ can be computed using the SVD of \underline{X} :

\hookrightarrow eigenvalues of $\underline{X}^T \underline{X}$ are the squared singular values of \underline{X} (diagonals of $\underline{\Sigma}$)

\hookrightarrow eigenvectors of $\underline{X}^T \underline{X}$ are the columns of the unitary matrix \underline{U}

Finite-space Markov chains

Markov chains

- A stochastic process X_0, X_1, \dots is a Markov chain iff for all times $i \geq 1$

$$P(X_{i+1} | X_0 = j_0, X_1 = j_1, \dots, X_i = j_i) = P(X_{i+1} | X_i = j_i)$$

where each time-indexed RV X_0, X_1, \dots takes a value from the state space S ($j \in S \forall i$)

- A Markov chain can be described by the transition matrix \underline{P} , where the (j,k) th element satisfies

$$P_{j,k} = P[\text{state } j \rightarrow \text{state } k]$$

$$\sum_k P_{j,k} = 1$$

- Denoting the distribution of states at time n by $\underline{x}^{(n)}$ and the initial distribution by $\underline{x}^{(0)}$, we have

$$\underline{x}^{(n+1)} = \underline{x}^{(n)} \underline{P}$$

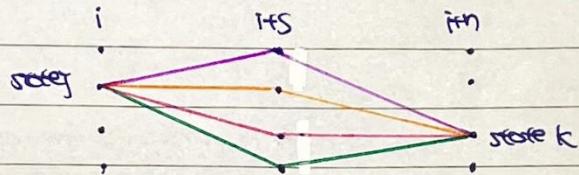
→

$$\underline{x}^{(n)} = \underline{P}^n \underline{x}^{(0)}$$

i.e. we can get the n -step transition matrix by raising the 1-step transition matrix to the power n

- Using the notation $(\underline{P}^n)_{j,k} = P_{j,k}^n$, we can derive the Chapman-Kolmogorov eqns as follows:

$$P[X_{t+n} = k | X_t = j] = (\underline{P}^n)_{j,k} = (\underline{P}^S \underline{P}^{nS})_{j,k} = \sum_l P_{j,l}^S P_{l,k}^{nS} = \sum_l P[X_{t+1} = l | X_t = j] P[X_{t+n} = k | X_{t+1} = l]$$



Stationary distributions and limiting distributions

- For a Markov process w/ transition matrix \underline{P} , a stationary distribution is any distribution s.t.

$$\underline{x}^{(\infty)} \underline{P} = \underline{x}^{(\infty)}$$

- A limiting distribution is one which for any initial distribution $\underline{x}^{(0)}$ satisfies

$$\underline{x}^{(0)} \underline{P}^n = \underline{x}^{(\infty)} \quad \text{as } n \rightarrow \infty$$

(A limiting distribution is the unique stationary position of a Markov process).

Transition matrix eigenvalues

- For a $n \times n$ transition matrix \underline{P} , (i) $\lambda=1$ is an eigenvalue, (ii) all eigenvalues satisfy $|\lambda| \leq 1$.

↳ (i) $\lambda=1$ is an eigenvalue.

$$\underline{P} - \underline{I} = \begin{bmatrix} P_{1,1}-1 & P_{1,2} & \cdots & P_{1,N} \\ P_{2,1} & P_{2,2}-1 & \cdots & P_{2,N} \\ \vdots & \vdots & \ddots & \vdots \\ P_{N,1} & P_{N,2} & \cdots & P_{N,N}-1 \end{bmatrix}$$

Sum of each row is $1-1=0$
→ summing first $N-1$ col yields the last col

∴ $\underline{P} - \underline{I}$ is not full rank $\rightarrow \det(\underline{P} - \underline{I}) = 0 \rightarrow \lambda=1$ is an eigenvalue.

↳ (ii) All eigenvalues satisfy $|\lambda| \leq 1$

Select largest magnitude element of eigenvector \underline{v}_k , i.e. $|V_k| \geq |V_i| \forall i$

$$|\lambda| |V_k| = |\lambda V_k| = |(\lambda \underline{v})_k| = |(\underline{P} \underline{v})_k| = |P_{k,1} v_1 + \dots + P_{k,N} v_N| \leq P_{k,1} |V_1| + \dots + P_{k,N} |V_N|$$

$$\leq P_{k,1} |V_k| + \dots + P_{k,N} |V_k| = (P_{k,1} + \dots + P_{k,N}) |V_k| = |V_k| \rightarrow |\lambda| \leq 1$$

+ We can find the stationary distribution by finding the left eigenvector corr. to $\lambda=1$ $\underline{v} \underline{P} = \underline{v}$.

For Personal Use Only -bkwk2

Irreducibility, recurrence and transience.

- Two states j and k in a Markov chain communicate iff there exist integers $m, n \geq 0$

$$P_{jk}^m > 0 \quad \text{and} \quad P_{kj}^n > 0$$

i.e. there must be a route from state j/k to k/j in m/n steps

best determine by
drawing the state diagram

- If all states in a Markov chain communicate w/ each other, then the Markov chain is irreducible.

- Let S denote the set of all possible states of a Markov chain. A subset $\tilde{S} \subseteq S$ is a recurrent set if (i) all pairs of states in \tilde{S} communicate (ii) if $j \in \tilde{S}, k \notin \tilde{S}$, then $P_{jk}^i = 0$ for all i . i.e. it is possible to move between all states in a current set, but cannot move from a member of a current set to not a member.

- If a state belongs to a recurrent set, then it is recurrent, o/w it is transient.

↳ If state k is recurrent, then $P[\text{return to state } k | X_0 = k] = 1$

↳ If we start in a transient state, then we will leave that state and never return to it.

Periodicity

- The period of a state k of a Markov chain is the greatest common divisor of the set

$$\{t \geq 0, P_{kk}^{(t)} > 0\}$$

(e.g. $\{18, 24\} \rightarrow 18 = 2 \times 3^2, 24 = 2^3 \times 3 \text{ so } \text{GCD}(18, 24) = 2 \times 3 = 6$)

- A state k is aperiodic if it has period 1. This is true if for some time i

$$P_{kk}^i > 0 \quad \text{and} \quad P_{kk}^{i+1} > 0$$

- A Markov chain is aperiodic if all states are aperiodic.

- If the Markov chain is irreducible, all the states have the same period.

(All states in a recurrent set have the same period)

- For an irreducible Markov chain w/ period ℓ , there will be ℓ eigenvalues of \underline{P} w/ magnitude 1 (ℓ roots of unity)

Every regular MC is ergodic
but an ergodic MC may not be regular

$$\text{eg. } \underline{P} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

Regular ergodicity

- A Markov chain is regular if for some integer n , all entries of \underline{P}^n are positive.

An ergodic Markov chain is equivalent to an irreducible Markov chain.

- If a Markov chain is irreducible and aperiodic, then it is regular ergodic.

↳ A regular ergodic Markov chain has a limiting distribution which puts the mass on every element of the state space S

↳ If X_0, X_1, \dots is a regular ergodic Markov chain w/ limiting distribution $\underline{\pi}$, then for any function $f(x)$, as $N \rightarrow \infty$, $\frac{1}{N} \sum_{i=1}^N f(X_i) \rightarrow \sum_{x \in S} \pi_x f(x)$

For Personal Use Only -bkwk2

Detailed balance.

- A transition matrix \underline{P} and a distribution $\underline{\pi}$ are in detailed balance if

$$\pi_j p_{j,k} = \pi_k p_{k,j} \quad \forall j, k$$

Consider this distribution after one step, $\tilde{\underline{\pi}} = \underline{\pi}\underline{P}$.

$$\tilde{\pi}_k = \sum_j \pi_j p_{j,k} = \sum_j \pi_k p_{k,j} = \pi_k \left(\sum_j p_{k,j} \right) = \pi_k \rightarrow \tilde{\underline{\pi}} = \underline{\pi}$$

$\therefore \underline{\pi}$ is a stationary distribution of \underline{P} ,

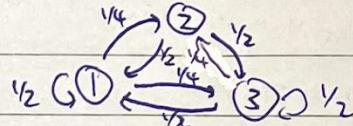
* If $\underline{\pi}$ satisfies detailed balance w/ \underline{P} , $\underline{\pi}$ is a stationary distribution for the transition matrix \underline{P} , but a stationary distribution of \underline{P} does not necessarily satisfy detailed balance.

Waiting time

- To calculate the expected waiting time to reach a certain state K from state j , q_j , we set up an expression q_j for each of the N states and solve.

- e.g. consider a MARKOV chain w/ transition matrix \underline{P} , ($N=3$)

$$\underline{P} = \begin{bmatrix} \frac{1}{2} & \frac{1}{4} & \frac{1}{4} \\ \frac{1}{2} & 0 & \frac{1}{2} \\ \frac{1}{4} & \frac{1}{4} & \frac{1}{2} \end{bmatrix}$$



Find the average wait time to transition to state 2, given we start at state 2.

$$\begin{aligned} q_1 &= \frac{1}{2}(1+q_1) + \frac{1}{4}(1) + \frac{1}{4}(1+q_2) \\ q_2 &= \frac{1}{2}(1+q_1) + \frac{1}{2}(1+q_3) \quad \rightarrow q_1 = 4, q_2 = 5, q_3 = 4 \\ q_3 &= \frac{1}{4}(1+q_1) + \frac{1}{4}(1) + \frac{1}{2}(1+q_2) \quad \therefore \text{average wait time req. is 5 steps.} \end{aligned}$$

continuous state space systems

Birth process

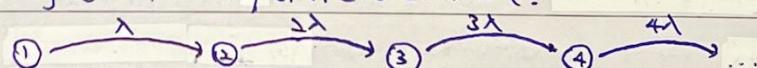
- Let $n(t)$ be the no. of cells at time t , $n(0)=n_0$ and λ be the birth rate for a cell per unit time.

$n(t)$ describes the state of the system \rightarrow Markovian. For a small time step Δt ,

$$n(t+\Delta t) = n(t) + n(t)\lambda \Delta t \xrightarrow{\Delta t \rightarrow 0} \frac{dn(t)}{dt} = \lambda n(t) \rightarrow n(t) = n_0 \exp(\lambda t)$$

* $n(t)$ gives the expected value of the no. of cells \rightarrow consider probabilistic approach for distribution of n .

- Let $P_n(t) = P[n(t)=n]$ be the probability of n cells at time t .



For a small time step Δt ,

$$P_n(t+\Delta t) = P_n(t)(1-n\lambda\Delta t) + P_{n-1}(t)((n-1)\lambda\Delta t)$$

As $\Delta t \rightarrow 0$,

$$\frac{dP_n(t)}{dt} = -n\lambda P_n(t) + (n-1)\lambda P_{n-1}(t)$$

- We can represent the system in terms of the transition rate matrix \underline{Q}

$$\frac{d\underline{x}(t)}{dt} = \underline{x}(t) \underline{Q} \quad \text{rows sum to zero}$$

$$\text{where } \underline{x}(t) = [P_1(t) \dots P_n(t) \dots], \sum_{n=1}^{\infty} P_n(t) = 1 \text{ and } \underline{Q} = \begin{bmatrix} -\lambda & \lambda & 0 & 0 & \dots \\ 0 & -2\lambda & 2\lambda & 0 & \dots \\ 0 & 0 & -3\lambda & 3\lambda & \dots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}$$

For Personal Use Only -bkwk2

- It can be shown that the solution to the birth process is

$$P_n(t) = C \sum_{n=1}^{\infty} \exp(-\lambda n t) (1 - \exp(-\lambda t))^{n-1}$$

- If the birth process did reach a steady state, then

$$\frac{d\pi(\infty)}{dt} = \pi(\infty) \underline{\Omega} = 0$$

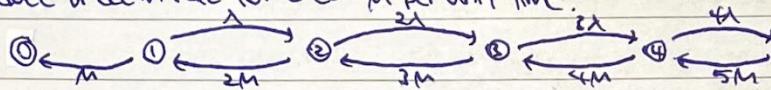
$$\rightarrow \begin{cases} 0 = -\lambda x_1(\infty) \\ 0 = \lambda x_1(\infty) - 2\lambda x_2(\infty) \\ \vdots \\ 0 = (n-1)x_{n-1}(\infty) - n\lambda x_n(\infty) \end{cases}$$

The set of eqns are only satisfied when $x_i(\infty) = 0$ for all $i \rightarrow$ not a valid PDF

\therefore there is no steady state distribution (the no. of cells keep growing)

Birth-death process

- Now we introduce a death rate for a cell M per unit time.



For a small time step Δt , $P_n(t+\Delta t) = P_n(t)(1 - (\lambda + M)\Delta t) + P_{n-1}(t)(\lambda\Delta t) + P_{n+1}(t)(M\Delta t)$

$$\text{As } \Delta t \rightarrow 0, \quad \frac{dP_n(t)}{dt} = -(\lambda + M)P_n(t) + (n-1)\lambda P_{n-1}(t) + (n+1)M P_{n+1}(t)$$

- We can represent the system in terms of the transition rate matrix $\underline{\Omega}$.

$$\frac{d\pi(t)}{dt} = \pi(t) \underline{\Omega}$$

$$\text{where } \pi(t) = [P_0(t) \dots P_n(t) \dots], \quad \sum_{n=0}^{\infty} P_n(t) = 1 \quad \text{and} \quad \underline{\Omega} = \begin{bmatrix} 0 & 0 & 0 & 0 & \dots \\ M & -(\lambda+M) & \lambda & 0 & \dots \\ 0 & 2M & -2(\lambda+M) & 2\lambda & \dots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}$$

- For steady state, we req.

$$\frac{d\pi(t)}{dt} = \pi(t) \underline{\Omega} = 0 \quad \rightarrow \quad -n(\lambda + M)P_n + (n-1)\lambda P_{n-1} + (n+1)M P_{n+1} = 0$$

If can be shown that $\pi(\infty) = [1 \ 0 \ 0 \ \dots]$ \rightarrow all states except 0 are transient.

- * If we λ_i, μ_i be the birth/death rate per unit time for state i respectively, then it can

$$\text{be shown using recursion that for } i \neq 0, \quad P_i(\infty) = \left(\frac{\lambda_i}{\lambda_i + \mu_i} \right) P_i(0)$$

random walk

- consider a random walk — at each time, we take step ± 1 , where the probability of dir.

at time K $S_K \in \{-1, 1\}$ is uniform, i.e. $P(S_K=1) = P(S_K=-1) = \frac{1}{2}$.

- After n steps, the position X_n is given by

$$X_n = \sum_{k=1}^n S_k$$

- At each instance, the average step is 0, and step variance is 1 \rightarrow by CLT, the distribution

of X_n as $N \rightarrow \infty$ is a Gaussian $N(0, N)$.

- If we now take a small step δ in dir. S_k every δ seconds, the distribution of the position W_t

at time $t=N\delta$ as $N \rightarrow \infty$ is a Gaussian $N(0, t) = N(0, N\delta)$

(AS $\delta \rightarrow 0$, we have Brownian motion)

For Personal Use Only -bkwk2

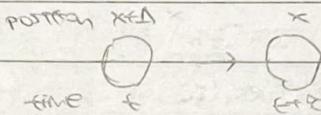
Brownian motion

- Brownian motion is the random motion of particles in a fluid resulting from collisions. since there are too many particles, we consider the particle density.
- Consider a simplified system where each dir. is independent, where the overall form is obtained by simply multiplying dimensions together. → only consider x dimension
- Let the particle density at time t and position x be $f(x,t)$. We can define the state of the system since probability of the next state only depends on the current state of time t → Markov process.
- Assuming smoothness in the particle density $f(x)$ over position x and time t, we can write

$$f(x+\Delta, t) = f(x, t) + \Delta \frac{\partial f(x, t)}{\partial x} + \frac{\Delta^2}{2} \frac{\partial^2 f(x, t)}{\partial x^2} + O(\Delta^3) \quad [1]$$

$$f(x, t+\tau) = f(x, t) + \tau \frac{\partial f(x, t)}{\partial t} + O(\tau^2) \quad [2]$$

- The change in density at time $t+\tau$ results from particles moving to position x from time t to $t+\tau$.



Considering all initial positions $x+\Delta$, the density at time $t+\tau$ is given by

$$f(x, t+\tau) = \int_{-\infty}^{\infty} f(x+\Delta, t) P(-\Delta) d\Delta \quad [3]$$

where $P(\Delta)$ is the probability of a particle moving Δ in time τ . continuous form of Chapman-Kolmogorov

- Sub [1] into [3], and by symmetry $P(-\Delta) = P(\Delta)$,

$$\begin{aligned} f(x, t+\tau) &\approx \int_{-\infty}^{\infty} \left[f(x, t) + \Delta \frac{\partial f(x, t)}{\partial x} + \frac{\Delta^2}{2} \frac{\partial^2 f(x, t)}{\partial x^2} \right] P(\Delta) d\Delta \\ &= f(x, t) \int_{-\infty}^{\infty} P(\Delta) d\Delta + \frac{\partial f(x, t)}{\partial x} \int_{-\infty}^{\infty} \Delta P(\Delta) d\Delta + \frac{1}{2} \frac{\partial^2 f(x, t)}{\partial x^2} \int_{-\infty}^{\infty} \Delta^2 P(\Delta) d\Delta. \\ &= f(x, t) + \frac{1}{2} \frac{\partial^2 f(x, t)}{\partial x^2} \int_{-\infty}^{\infty} \Delta^2 P(\Delta) d\Delta \end{aligned}$$

Comparing this w/ [2], we have

$$\begin{aligned} f(x, t+\tau) &= f(x, t) + \tau \frac{\partial f(x, t)}{\partial t} = f(x, t) + \frac{1}{2} \frac{\partial^2 f(x, t)}{\partial x^2} \int_{-\infty}^{\infty} \Delta^2 P(\Delta) d\Delta \\ \therefore \tau \frac{\partial f(x, t)}{\partial t} &= \frac{1}{2} \frac{\partial^2 f(x, t)}{\partial x^2} \int_{-\infty}^{\infty} \Delta^2 P(\Delta) d\Delta \rightarrow \frac{\partial f(x, t)}{\partial t} = D \frac{\partial^2 f(x, t)}{\partial x^2} \quad \text{where } D = \frac{1}{2\tau} \int_{-\infty}^{\infty} \Delta^2 P(\Delta) d\Delta. \end{aligned}$$

- The value of D can be determined by measurement of physical properties → no need to consider τ or $P(\Delta)$.

- The final soln depends on the initial conditions → particular soln.

For Personal Use Only -bkwk2

Solution for Brownian motion differential equation.

- Consider the Brownian motion DE. For a diff. X ,

$$\frac{\partial p(x,t)}{\partial t} = \alpha \frac{\partial^2 p(x,t)}{\partial x^2}$$

Assuming the soln can be written in the form, $p(x,t) = X(x)T(t)$,

$$X \frac{dT}{dt} = \alpha T \frac{d^2X}{dx^2} \rightarrow \frac{dT}{\alpha dt} = \frac{1}{X} \frac{d^2X}{dx^2} = -k^2 \quad \text{separation const.}$$

$$\therefore X(x) = A\sin(kx) + B\cos(kx), \quad T(t) = C\exp(-\alpha k t), \rightarrow p(x,t) = A(t)\exp(-\alpha k t)\exp(ikx)$$

The above soln is satisfied by only k , so the general soln is,

$$p(x,t) = \int_{-\infty}^{\infty} A(k) \exp(-\alpha k t) \exp(ikx) dk$$

- consider the initial condition as $p(x,0) = f(x)$,

$$f(x) = p(x,0) = \int_{-\infty}^{\infty} A(k) \exp(ikx) dk = \int_{-\infty}^{\infty} A(\tilde{k}) \exp(i\tilde{k}x) d\tilde{k}$$

Notice that this is the FT of $A(\tilde{k})$; i.e. $\text{FT}[A(\tilde{k})] = \delta(x)$, so

$$A(\tilde{k}) = \text{IFT}[\delta(x)] = \frac{1}{2\pi}$$

$$\therefore p(x,t) = \frac{1}{2\pi} \int_{-\infty}^{\infty} \exp(-\alpha k^2 t) \exp(ikx) dk = \frac{1}{2\pi} \int_{-\infty}^{\infty} \exp(-\alpha k^2 t + ikx) dk$$

The term in the exponential can be rewritten as $-\alpha k^2 t + ikx = -(\sqrt{\alpha t} k - \frac{ix}{\sqrt{\alpha t}})^2 - \frac{x^2}{4\alpha t}$

$$p(x,t) = \frac{1}{2\pi} \exp\left(-\frac{x^2}{4\alpha t}\right) \int_{-\infty}^{\infty} \exp\left[-\left(\sqrt{\alpha t} k - \frac{ix}{\sqrt{\alpha t}}\right)^2\right] dk = \frac{1}{2\pi} \exp\left(-\frac{x^2}{4\alpha t}\right) \int_{-\infty}^{\infty} dk$$

$$p(x,t) = \frac{1}{\sqrt{4\pi\alpha t}} \exp\left(-\frac{x^2}{4\alpha t}\right) = N(x|0, 2\alpha t)$$

complete the square in k .

Properties of a 1D wiener process.

- The simplified Brownian motion is a 1D wiener process. At time t , examine the wiener process w_t .

We can consider an instance of a "path" at time t : $w_t^{(i)}$ (the support of the RV x_t)

- The wiener process $\{w_t\}_{t \geq 0}$ has the following properties:

(i) Stationarity: $w_{t+s} - w_s$ is independent of s . e.g. $w_t = w_{t+k} - w_k$

(ii) Independence: $w_t - w_s$ is independent of w_u for $s < u$. e.g. $w_0 - w_1 \perp w_2$

(iii) Continuity: w_t is a continuous function w.r.t t

(iv) Gaussianity: w_t is a gaussian w/ $E[w_t] = 0, E[w_t w_s] = 2\alpha \min(t,s)$

general continuous state-space systems

- A range of assumptions have been made in the previous derivation - some may not be valid.

- A simple extension is to include a drift term: wiener process w/ drift.

$$\frac{\partial p(x,t)}{\partial t} = M \frac{\partial p(x,t)}{\partial x} + D \frac{\partial^2 p(x,t)}{\partial x^2}$$

- More complex models include the Ornstein-Uhlenbeck process.

$$\frac{\partial p(x,t)}{\partial t} = \frac{\partial (\beta p(x,t))}{\partial x} + \frac{\partial^2 (\gamma p(x,t))}{\partial x^2}$$

where β controls the drift and γ controls the diffusion

- When $\beta=0$, we have the wiener process. When $\beta \neq 0$, the diffusion term can be balanced by the drift term so the particle density tends to a limiting distribution (finite variance)

For Personal Use Only -bkwk2

Monte Carlo methods

Monte Carlo integration.

- consider a highly complicated, multi-dimensional ($\underline{x} \in \mathbb{R}^d$) integral: $\int h(\underline{x}) d\underline{x}$.

using a histogram approach (in d dimensions), we have

$$\int h(\underline{x}) d\underline{x} \approx \sum_{i=1}^N h(\underline{x}^{(i)}) \delta x_1 \delta x_2 \dots \delta x_d$$

where we have uniformly spread N histogram centres.

- the histogram approach is simple but wasteful for regions w/ very small/no contributions to the integral.

If also does not scale well w/ the no. of dimensions. ↵ to not lose accuracy in the approx, we req. n bins per dimension \rightarrow n^d bins (3rd dimension).

- For Monte Carlo integration, we sample from the normalised weighting function $w(\underline{x})$, $p(\underline{x})$

$$p(\underline{x}) = \frac{w(\underline{x})}{\int w(\underline{x}) d\underline{x}}$$

This is used to model the distribution of the pts.

- initially consider uniform distribution over a "volume" V , $p(\underline{x}) = \frac{1}{V}$. The integral becomes,

$$\boxed{\int h(\underline{x}) d\underline{x} = \int \frac{h(\underline{x})}{p(\underline{x})} p(\underline{x}) d\underline{x} = E_p[h(\underline{x})] = \frac{1}{N} \sum_{i=1}^N \frac{h(\underline{x}^{(i)})}{p(\underline{x}^{(i)})} = \frac{V}{N} \sum_{i=1}^N h(\underline{x}^{(i)})}$$

where $\underline{x}^{(i)}$ are drawn from $p(\underline{x})$.

- selecting an appropriate "volume" for $p(\underline{x})$ avoids wasted samples \rightarrow smaller variance for a fixed no. of samples.

Monte Carlo for expected values.

- consider computing the expected value of a function $f(\underline{x})$, $E[f(\underline{x})]$. The integral has the form

$$\boxed{\int f(\underline{x}) p(\underline{x}) d\underline{x} = E_p[f(\underline{x})] \approx \frac{1}{N} \sum_{i=1}^N f(\underline{x}^{(i)})}$$

where $\underline{x}^{(i)}$ are drawn from $p(\underline{x})$.

- denote the mean of the estimate using 1 sample w/ \hat{M}_1 , $\hat{M}_1 = \int f(\underline{x}) p(\underline{x}) d\underline{x}$.

The mean of the estimate using N samples, \hat{M}_N is given by

$$\hat{M}_N = \frac{1}{N} \sum_{i=1}^N \hat{M}_1 = \hat{M}_1 \leftarrow \text{unbiased}.$$

- denote the variance of the estimate using 1 sample w/ $\hat{\sigma}_1^2$, $\hat{\sigma}_1^2 = \int (f(\underline{x}) - M)^2 p(\underline{x}) d\underline{x}$.

$$\hat{\sigma}_N^2 = \frac{1}{N^2} \sum_{i=1}^N \hat{\sigma}_1^2 = \frac{1}{N} \hat{\sigma}_1^2$$

- In fact, for large N , the distribution of $\boxed{\frac{1}{N} \sum_{i=1}^N f(\underline{x}^{(i)})} \rightarrow N(\hat{M}_N, \hat{\sigma}_N^2) = N(\hat{M}_1, \frac{\hat{\sigma}_1^2}{N})$ by CLT.

\rightarrow the approach yield correct answers. (as $N \rightarrow \infty$).

For Personal Use Only -bkwk2

Importance sampling.

- If we want to sample from another pdf $q(x)$ rather than $p(x)$ (say it is not possible to sample from $p(x)$) or we want to sample from a nonuniform distribution, we use importance sampling.
- Importance sampling can be used for computing expected values / numerical integration:

① Computing expected value $E_p[f(x)]$

$$E_p[f(x)] = \int f(x) p(x) dx = \int f(x) \frac{p(x)}{q(x)} q(x) dx = E_q\left[f(x) \frac{p(x)}{q(x)}\right] \approx \frac{1}{N} \sum_{i=1}^N f(x^{(i)}) \frac{p(x^{(i)})}{q(x^{(i)})}$$

↳ the ratio $\frac{p(x^{(i)})}{q(x^{(i)})}$ is the importance of the drawn sample ↳ how much to "boost" a sample $f(x^{(i)})$.

$(q(x) > p(x))$: over-represented region ; $p(x) > q(x)$: under-represented region

↳ we can perform importance sampling w/o the normalised distributions. For $x^{(i)} \sim q(x)$,

$$\text{let } p(x) = \frac{1}{Z_p} p^*(x), \quad q(x) = \frac{1}{Z_q} q^*(x), \quad w^{(i)} = \frac{p^*(x^{(i)})}{q^*(x^{(i)})}, \quad \text{where } Z_p, Z_q \text{ cannot be computed.}$$

$$\text{The expected value can be shown to be } V = \int f(x) \frac{p(x)}{q(x)} q(x) dx \approx \frac{1}{N} \sum_{i=1}^N f(x^{(i)}) w^{(i)}$$

$$(\text{Numerator: } \frac{1}{N} \sum_{i=1}^N f(x^{(i)}) w^{(i)} \approx E_q[f(x) w^{(i)}] = \int f(x) \frac{p^*(x)}{q^*(x)} q(x) dx = \frac{Z_p}{Z_q} \int f(x) p(x) dx = \frac{Z_p}{Z_q} V)$$

$$(\text{Denominator: } \frac{1}{N} \sum_{i=1}^N w^{(i)} \approx E_q[w^{(i)}] = \int \frac{p^*(x)}{q^*(x)} q(x) dx = \frac{Z_p}{Z_q} \int p(x) dx = \frac{Z_p}{Z_q})$$

② Numerical integration $\int h(x) dx$

$$\int h(x) dx = \int \frac{h(x)}{q(x)} q(x) dx = E_q\left[\frac{h(x)}{q(x)}\right] \times \frac{1}{N} \sum_{i=1}^N \frac{h(x^{(i)})}{q(x^{(i)})}$$

↳ the ratio $\frac{h(x^{(i)})}{q(x^{(i)})}$ is similar to the "importance" of the drawn sample.

↳ we want $q(x)$ to be close to $h(x)$ → importance ≈ 1 everywhere

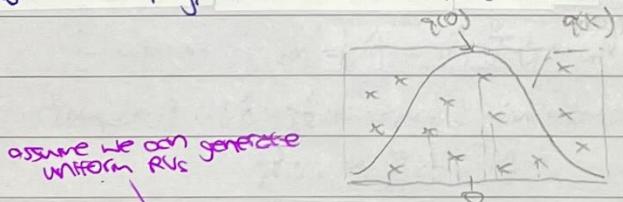
if $q(x) \propto h(x)$, then $\frac{h(x^{(i)})}{q(x^{(i)})}$ is a const. → we can get the answer in 1 sample.

↳ using importance sampling may give a smaller variance in the estimate for a fixed no. of samples.

Generating random samples – rejection sampling.

- We want to generate random samples from the distribution $q(x)$, $x \in \mathbb{R}^d$

- In rejection sampling, we enclose (more or) the distribution $q(x)$ in a 2D box



Drawing samples uniformly from the 2D box, we accept those under the curve (take x-values) and reject those above the line

- The fraction of samples accepted is given by

$$\text{fraction of samples accepted} = \frac{\text{area under curve}}{\text{area of box}}$$

slightly under 1.

For higher dimensions, $x \in \mathbb{R}^d$, the fraction of samples accepted would be reduced to the power of d → highly wasteful (curse of dimensionality).

For Personal Use Only -bkwk2

Generating random samples — Metropolis-Hastings algorithm.

- We want to generate random samples from the distribution $p(x)$, i.e. $x \sim p(x)$.

- The Metropolis-Hastings algorithm is as follows:

↳ Given the current sample $x^{(i)}$, generate another sample $x^{(i+1)}$ from $p(\underline{x}|x^{(i)})$

$$\underline{x}^{(i+1)} = x^{(i)} + z$$

where \underline{z} is the proposal distribution (e.g. $\underline{z} \sim N(0, I)$)

↳ Accept the sample, i.e. $x^{(i+1)} = \underline{x}^{(i+1)}$ w/ probability $\alpha = \min\left\{\frac{p(x^{(i+1)}) p(x^{(i)}|\underline{x}^{(i)})}{p(x^{(i)}) p(\underline{x}^{(i)}|x^{(i)})}, 1\right\}$

otherwise reject the sample, i.e. $x^{(i+1)} = x^{(i)}$

- It can be shown that the distribution of $\{x^{(n)}\}$ as $n \rightarrow \infty$ has the same distribution as $p(x)$.

- Overall samples are drawn but

↳ initial samples depend on $x^{(0)}$ → ignore the initial samples (burn-in phase).

↳ samples correlated w/ neighbouring samples → take every nth sample (thinning).

- We need to decide on the proposal distribution \underline{z} . If \underline{z} is symmetric, i.e. $p(x^{(i)}|\underline{x}^{(i)}) = p(\underline{x}^{(i)}|x^{(i)})$,

then the probability of accepting the sample is $\alpha = \min\left\{\frac{p(x^{(i+1)})}{p(x^{(i)})}, 1\right\}$, true for $\underline{z} \sim N(0, I)$.

Markov chain Monte Carlo (MCMC).

- A MCMC approach is based on the Metropolis-Hastings algorithm to generate samples from a process w/ a limiting distribution Π .

- Consider a proposal function — a transition matrix \underline{R} , where $r_{j,k} \geq 0$ and $r_{j,k} > 0$ ($j \neq k$).
(The true transition matrix is not known).

- The MCMC algorithm is as follows:

↳ Given the current state x_i , select \hat{x}_{i+1} by sampling from the x_i -th row of \underline{R} .

↳ Accept the sample, i.e. $x_{i+1} = \hat{x}_{i+1}$ w/ probability $\alpha = \min\left\{\frac{\Pi_{k|x_i} r_{k|i}}{\Pi_{k|x_i} r_{k|i}}, 1\right\}$

otherwise reject the sample, i.e. $x_{i+1} = x_i$

- The process is equivalent to a transition matrix \underline{R}' where,

$$r'_{j,k} = r_{j,k} \min\left\{\frac{\Pi_k r_{k|i}}{\Pi_j r_{j|i}}, 1\right\} \quad \text{if } j \neq k, \quad r'_{j,j} = 1 - \sum_{k \neq j} r'_{j,k} \min\left\{\frac{\Pi_k r_{k|i}}{\Pi_j r_{j|i}}, 1\right\}$$

- We can show that \underline{R}' and Π are in detailed balance.

$$\Pi_j^T r'_{j,k} = \Pi_j r_{j,k} \min\left\{\frac{\Pi_k r_{k|i}}{\Pi_j r_{j|i}}, 1\right\} = \min\left\{\Pi_k r_{k|i}, \Pi_j r_{j,k}\right\} = \Pi_k r_{k|i} \min\left\{1, \frac{\Pi_k r_{k|i}}{\Pi_j r_{j|i}}\right\} = \Pi_k r_{k|i}$$

\underline{R}' and Π are in detailed balance → Π is a stationary distribution for the transition matrix \underline{R}' .

By construction, \underline{R} is irreducible and aperiodic → regular ergodic → Π is the limiting distribution.

Optimisation

The optimisation problem

- The goal of optimisation is to find the independent variables that min./max. a given quantity, possibly subject to some restrictions on the allowed range of parameters.
- The quantity to be max/min is the objective / cost / utility / loss function, denoted by $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$.
- The parameters that can be changed are the control / decision variables

The restrictions on the allowed parameter values are the constraints.

- Mathematically, the optimisation problem is

Inequality constraints on single variables are bounds, e.g. $x_{i,\min} \leq x_i \leq x_{i,\max}$.

$$\text{minimise } f(\mathbf{x}), \quad \mathbf{x} = [x_1, x_2, \dots, x_N]^T \quad \text{subject to } h_i(\mathbf{x}) = 0, i=1,2,\dots,m \text{ and } g_j(\mathbf{x}) \leq 0, j=1,2,\dots,n$$

where $f(\mathbf{x})$ is the objective function, \mathbf{x} is the col. vector of N control variables,

$\{h_i(\mathbf{x}), g_j(\mathbf{x})\}$ is the set of constraint functions.

- + A min. of the function $-f(\mathbf{x})$ is a max. of the function $f(\mathbf{x}) \rightarrow$ max. problems can be cast as min. problems

- Common forms of the objective function $f(\mathbf{x})$ include:

↳ Linear :

$$f(\mathbf{x}) = \mathbf{b}^T \mathbf{x} + c$$

↳ Quadratic :

more difficult to solve

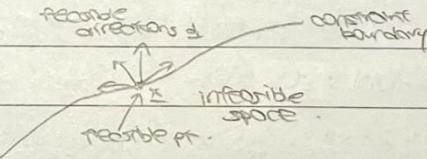
$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{x} + c$$

↳ Non-linear (not linear/quadratic) : (eg) $f(\mathbf{x}) = \exp(x_1) + \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{x} + c$.

- When min. $f(\mathbf{x})$ subject to constraints, S is the feasible region, and only $\mathbf{x} \in S$ is a feasible sol'n.

For an unconstrained problem, S is infinitely large.

- At a feasible pt. \mathbf{x} , a direction d is a feasible direction if an arbitrarily small move from \mathbf{x} in direction d remains feasible.



Stationary points and types of minima.

- If $f(\mathbf{x})$ is smooth, s.t. $\nabla f(\mathbf{x})$ exists, then \mathbf{x}^* is a stationary pt. of f if

$$\boxed{\nabla f(\mathbf{x}^*) = 0}$$

↳ minima, maxima and saddle pts are stationary pts.

- Types of minima includes the following:

↳ Global min.

$$f(\mathbf{x}^*) \leq f(\mathbf{x}) \quad \forall \mathbf{x} \in S$$

↳ Weak local min.

$$f(\mathbf{x}^*) \leq f(\mathbf{x}) \quad \forall \mathbf{x} = \mathbf{x}^* + \mathbf{e}_d \in S, \mathbf{y} \neq \mathbf{x}^*$$

↳ strong global min.

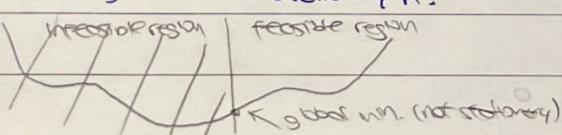
$$f(\mathbf{x}^*) < f(\mathbf{x}) \quad \forall \mathbf{x} \in S, \mathbf{y} \neq \mathbf{x}^*$$

↳ Strong local min.

$$f(\mathbf{x}^*) < f(\mathbf{x}) \quad \forall \mathbf{x} = \mathbf{x}^* + \mathbf{e}_d \in S$$

- + Local max/min. are local extrema, and refer to interior local max/min., unless o/w specified.

- When subject to constraints, a global min. might not be a stationary pt.



For Personal Use Only -bkwk2

Necessary and sufficient conditions for a local minimum

- A necessary condition for \underline{x}^* to be a local min. of $f(\underline{x})$ in S is

$$\nabla f(\underline{x}^*) \cdot \underline{d} \geq 0 \quad \text{for all feasible directions } \underline{d}$$

- If \underline{x}^* is an interior pt., i.e. if it is not at the boundary of the feasible region, then all dir. are feasible $\rightarrow \underline{x}^*$ must be a stationary pt. when $\nabla f(\underline{x}^*) = 0$

(the condition $\nabla f(\underline{x}^*) = 0$ is necessary but not sufficient for \underline{x}^* to be a min. as the above condition holds for max/ saddle pts)

- For the univariate case ($x \in \mathbb{R}$), the Taylor expansion about x^* is

$$f(x) = f(x^*) + (x - x^*)f'(x^*) + \frac{1}{2}(x - x^*)^2 f''(x^*) + \text{H.O.T.}$$

$$f(x) - f(x^*) = (x - x^*)f'(x^*) + \frac{1}{2}(x - x^*)^2 f''(x^*) + \text{H.O.T.}$$

If x^* is an interior stationary pt, $f'(x^*) = 0$, then we have a strong local min. if

$$f(x) - f(x^*) \approx \frac{1}{2}(x - x^*)^2 f''(x^*) > 0 \rightarrow f''(x^*) > 0$$

- At a stationary pt x^* , $f''(x^*) > 0$ is a sufficient condition of strong local min.

$$f''(x^*) \geq 0 \text{ is a necessary condition of strong local min.}$$

(If $f''(x^*) = 0$, we need to analyse the higher order terms)

- For the multivariate case ($\underline{x} \in \mathbb{R}^N$), the Taylor expansion about \underline{x}^* is

$$f(\underline{x}) = f(\underline{x}^*) + (\underline{x} - \underline{x}^*)^T \nabla f(\underline{x}^*) + \frac{1}{2}(\underline{x} - \underline{x}^*)^T \underline{\underline{H}}(\underline{x}^*)(\underline{x} - \underline{x}^*) + \text{H.O.T.}$$

$$f(\underline{x}) - f(\underline{x}^*) = (\underline{x} - \underline{x}^*)^T \nabla f(\underline{x}^*) + \frac{1}{2}(\underline{x} - \underline{x}^*)^T \underline{\underline{H}}(\underline{x}^*)(\underline{x} - \underline{x}^*) + \text{H.O.T.} \quad \text{the Hessian } \underline{\underline{H}}(\underline{x}) \text{ is symmetric by defn.}$$

where $\underline{\underline{H}}(\underline{x}) = \nabla(\nabla f(\underline{x})) = \nabla([\frac{\partial^2 f}{\partial x_1^2}, \dots, \frac{\partial^2 f}{\partial x_N^2}]^T) = [\frac{\partial^2 f}{\partial x_1^2} \dots \frac{\partial^2 f}{\partial x_N^2} \dots \frac{\partial^2 f}{\partial x_1 \partial x_N} \dots \frac{\partial^2 f}{\partial x_N \partial x_1}]$ is the Hessian ↴

If \underline{x}^* is an interior stationary pt, $\nabla f(\underline{x}^*) = 0$ then we have a strong local min. if

$$f(\underline{x}) - f(\underline{x}^*) \approx \frac{1}{2} \underline{d}^T \underline{\underline{H}}(\underline{x}^*) \underline{d} > 0 \rightarrow \underline{d}^T \underline{\underline{H}}(\underline{x}^*) \underline{d} > 0, \text{ where } \underline{d} = \underline{x} - \underline{x}^*$$

- At a stationary pt \underline{x}^* , $\underline{d}^T \underline{\underline{H}}(\underline{x}^*) \underline{d} > 0$ (PD Hessian) is a sufficient condition of strong local min.

$$\underline{d}^T \underline{\underline{H}}(\underline{x}^*) \underline{d} \geq 0 \text{ (SPP Hessian) is a necessary condition of strong local min.}$$

- If $\underline{\underline{H}}(\underline{x})$ is PD everywhere, then $f(\underline{x})$ is strictly convex \rightarrow the min. is unique and a global min.

- A matrix is PD iff all its eigenvalues are pos. For a 2×2 matrix, if it is PD if

$$H_{11} > 0$$

and

$$\det(\underline{\underline{H}}) > 0$$

(If $\underline{\underline{H}}(\underline{x}^*) = 0$, we need to analyse the higher order terms)

For Personal Use Only -bkwk2

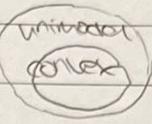
Unimodal and convex functions

- A function is unimodal if it has a single extremum. It is strictly unimodal if the function strictly increases/decreases.

- A function is convex if its graph at any pt y is never below the tangent of any other pt x .

$$f(y) \geq f(x) + \nabla f(x)^T (y-x)$$

(If the function is strictly convex, we replace \geq with $>$)



- convex functions are unimodal but not all unimodal functions are convex

- the local min. of a unimodal/convex function is unique \rightarrow local min. is the global min.

Unconstrained optimization — search methods for univariate functions

General search algorithm

- The general search algorithm is as follows:

↳ 1) Start w/ an initial guess x_0 for the min. of $f(x)$. (random or informed)

↳ 2) Propose a search dir. d_k

↳ 3) Propose a step size α_k along d_k (typically by an inner line search loop)

↳ 4) Update the estimate of the min. location $x_{k+1} = x_k + \alpha_k d_k$

↳ 5) Repeat steps 2-5 until convergence

Convergence criteria and rate of convergence

- convergence can be assessed in diff. ways. For a user-defined tolerance ε , we stop the search when

↳ Norm of the residual

$$\|f(x_{k+1}) - f(x_k)\| < \varepsilon_f$$

↳ Norm of the error

$$\|x_{k+1} - x_k\| < \varepsilon_x$$

↳ Norm of the gradient

$$\|\nabla f(x_k)\| < \varepsilon_g$$

← used when the curvature is small, and combined w/ test on norm of the error,

- The rate of convergence of an algorithm quantifies how fast the approximate soln approaches the

true min. x^* , close to the min., iteration after iteration. Mathematically,

$$\lim_{K \rightarrow \infty} \frac{\|x_{k+1} - x^*\|}{\|x_k - x^*\|^p} = \beta$$

where β is the convergence ratio (smaller is better)

and p is the order of convergence (larger is better)

$p=1$: linear convergence
 $p=2$: quadratic convergence

* The rate of convergence is an asymptotic quantity — it is the speed at which a sequence converges to its limit, if it converges.

For Personal Use Only -bkwk2

Line search

- Line search is the search of a min. for a univariate function. Interval reduction is the best we can do if the gradient $f'(x)$ is not available.

- The line search algorithm is as follows:

↳ 1) Start w/ three pts x_1, x_2, x_3 s.t. $f(x_1) > f(x_3)$ and $f(x_3) < f(x_2)$

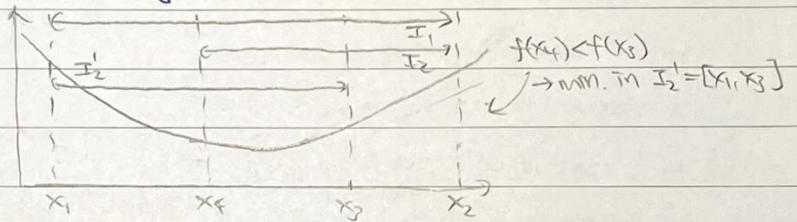
x_1 , and x_2 are the bounds and x_3 is the interior pt.

↳ 2) By Bolzano thm, there exists a min. in the interval $I_1 = [x_1, x_2]$ (bracketing)

↳ 3) Compute $f(x_4)$ at a fourth pt $x_4 \in I_1$ for the case $x_4 < x_3$

↳ 4) If $f(x_4) > f(x_3)$, the min. must be in $I_2 = [x_4, x_2]$, otherwise, the min. is in $I'_2 = [x_1, x_3]$

↳ 5) Repeat steps 2-5 until convergence (interval smaller than tolerance)



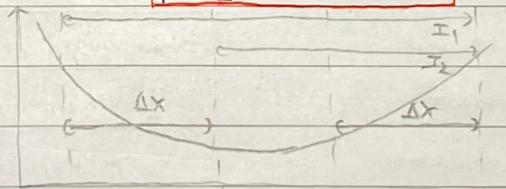
Golden section search

- To find the optimal location of the interior pt x_4 , we impose a const. reduction factor β for the interval length. w/ linear convergence,

$$\frac{I_2}{I_1} = \frac{I_1 - I_2}{I_2} = \frac{\Delta x}{I_2}$$

Setting $\beta = \frac{I_2}{I_1}$ and solving the quadratic yields

$$\beta = \frac{\sqrt{5}-1}{2} \approx 0.618 = \frac{1}{\phi} \quad \begin{matrix} \text{inverse of the} \\ \text{golden ratio} \end{matrix}$$



Note that $\Delta x = I_1 - I_2 = I_1(1-\beta)$. Find $\Delta x \rightarrow x_4 = x_1 + \Delta x, x_3 = x_2 - \Delta x$

- ∵: No need to compute the first or second derivatives

↓: Faster convergence than other "non-gradient" methods

✗: Slower convergence than gradient methods

✗: Not robust — goes wrong if the initial interval does not contain a min.

For Personal Use Only -bkwk2

Fitting with a quadratic polynomial of three points

- we can approx. the objective function $f(x)$ w/ a quadratic function $q(x)$ by matching them at three pts.

$$q(x) = a_0 + a_1 x + a_2 x^2$$

We can solve the min. of $q(x)$ analytically $x^* = -\frac{a_1}{2a_2}$, which approx. the min. of $f(x)$.

- To obtain the coefficients, we impose.

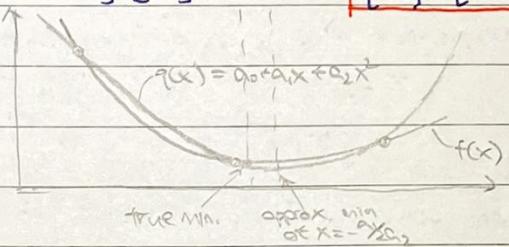
$$f(x_1) = a_0 + a_1 x_1 + a_2 x_1^2$$

$$f(x_2) = a_0 + a_1 x_2 + a_2 x_2^2$$

$$f(x_3) = a_0 + a_1 x_3 + a_2 x_3^2$$

using matrix form and solving, we have

$$\begin{bmatrix} f(x_1) \\ f(x_2) \\ f(x_3) \end{bmatrix} = \begin{bmatrix} 1 & x_1 & x_1^2 \\ 1 & x_2 & x_2^2 \\ 1 & x_3 & x_3^2 \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \\ a_2 \end{bmatrix} \rightarrow \boxed{\begin{bmatrix} a_0 \\ a_1 \\ a_2 \end{bmatrix} = \begin{bmatrix} 1 & x_1 & x_1^2 \\ 1 & x_2 & x_2^2 \\ 1 & x_3 & x_3^2 \end{bmatrix}^{-1} \begin{bmatrix} f(x_1) \\ f(x_2) \\ f(x_3) \end{bmatrix}}$$



Fitting with a quadratic polynomial of one point - Newton's method

- Instead of fitting the function of three pts, we can fit the function value, first and second derivatives evaluated at a single pt.

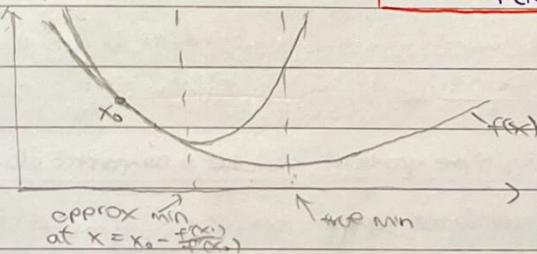
$$f(x) \approx f(x_0) + (x-x_0)f'(x_0) + \frac{1}{2}(x-x_0)^2 f''(x_0) + \text{H.O.T.}$$

We get a quadratic approx. $q(x)$ of $f(x)$ by neglecting the higher order terms,

$$q(x) = f(x_0) + (x-x_0)f'(x_0) + \frac{1}{2}(x-x_0)^2 f''(x_0)$$

- To calculate the min. of $q(x)$, we set its derivative to zero,

$$q'(x) = f'(x_0) + (x-x_0)f''(x_0) = 0 \rightarrow x = x_0 - \frac{f'(x_0)}{f''(x_0)}$$



- ✓ : When it works, has quadratic convergence ($\rho=2$)

✓ : could be generalised to multidimensional searches (but could be computationally expensive)

✗ : sensitive to the choice of the initial pt $x_0 \rightarrow$ could make the estimate worse or even converge to the wrong extremum (usually when $f''(x)$ is zero or negative)

For Personal Use Only -bkwk2

Quasi-Newton methods

- The computation of the second derivatives can be expensive (especially for higher dim.)

A large no. of methods are designed to approx. the second derivative.

- The simplest approx is the secant method, which takes the current and previous pts. to estimate the second derivative at the current pt.

$$f''(x_i) \approx \frac{f'(x_i) - f'(x_0)}{x_i - x_0}$$

- The Taylor expansion abt. x_i is approx. as

$$S(x) = f(x_i) + (x-x_i) f'(x_i) + \frac{1}{2} (x-x_i)^2 \frac{f'(x_i) - f'(x_0)}{x_i - x_0}$$

Setting $S'(x)=0$ yields the estimate for the min.

$$S'(x) = f'(x_i) + (x-x_i) \frac{f'(x_i) - f'(x_0)}{x_i - x_0} = 0 \rightarrow x = x_i - f(x_i) \frac{x_i - x_0}{f(x_i) - f(x_0)}$$

- convergence for the secant method is superlinear ($p = \frac{\ln \frac{f(x_0)}{f(x_1)}}{\ln \frac{x_0 - x_1}{x_1 - x_2}} \approx 1.618 > 1$)

- Quasi-Newton methods have higher dim. generalisations: Broyden's method, Broyden-Fletcher-Goldfarb-Shanno method.

Unconstrained optimisation — search methods for multivariate functions

steepest-descent method

- Access to the gradient enables gradient search — more efficient than direct (gradient-free) search

- we can exploit the fact that the gradient is the dir. of max first-order change of a multivariate function $f(x)$. The steepest-descent algorithm is as follows:

↳ 1) start w/ an initial guess for the min. x_0 .

↳ 2) Set the search dir. of the -ve gradient,

$$\underline{d}_k = -\nabla f(x_k)$$

main loop change
the direction

↳ 3) determine by line search the step size α_k that min. $f(x)$ along \underline{d}_k .

↳ 4) Update the sol'n to $\underline{x}_{k+1} = \underline{x}_k + \alpha_k \underline{d}_k$

inner loop is
optimisation in 1D

↳ 5) repeat steps 2-4 until convergence.

- Ideally, the line search in step 3 should satisfy $\frac{df(\underline{x}_k + \alpha_k \underline{d}_k)}{d\alpha_k} = 0$, which implies that the

gradient at $\underline{x}_{k+1} = \underline{x}_k + \alpha_k \underline{d}_k$ is orthogonal to \underline{d}_k , i.e. successive search dir. are orthogonal $\underline{d}_k^T \underline{d}_{k+1} = 0$

$$\hookrightarrow 0 = \frac{df(\underline{x}_k + \alpha_k \underline{d}_k)}{d\alpha_k} = \frac{df(\underline{x}_{k+1})}{d\alpha_k} = \frac{df(\underline{x}_{k+1})}{d\underline{x}_{k+1}} \cdot \frac{d\underline{x}_{k+1}}{d\alpha_k} = \nabla f(\underline{x}_{k+1}) \cdot \frac{d(\underline{x}_k + \alpha_k \underline{d}_k)}{d\alpha_k} = (-\underline{d}_k) \cdot \underline{d}_k$$

- Practically, an approximate line search is performed, so $\frac{df(\underline{x}_k + \alpha_k \underline{d}_k)}{d\alpha_k} \approx 0 \rightarrow \underline{d}_k^T \underline{d}_{k+1} \approx 0$

→ approx too poor
→ may not converge
approx too good
→ waste resources.

- If we approximate the objective function $f(x)$ w/ a 2nd order Taylor expansion, we can get a

closed form optimum step size α_k ,

$$\alpha_k = -\frac{\nabla f(x_k)^T \underline{d}_k}{\underline{d}_k^T H(x_k) \underline{d}_k}$$

$$\hookrightarrow f(\underline{x}_k + \alpha_k \underline{d}_k) = f(\underline{x}_k) + \nabla f(\underline{x}_k)^T (\alpha_k \underline{d}_k) + \frac{1}{2} (\alpha_k \underline{d}_k)^T H(\underline{x}_k) (\alpha_k \underline{d}_k) + \text{H.O.T.}$$

$$\frac{df(\underline{x}_k + \alpha_k \underline{d}_k)}{d\alpha_k} = \nabla f(\underline{x}_k)^T \underline{d}_k + \alpha_k \underline{d}_k^T H(\underline{x}_k) \underline{d}_k = 0 \rightarrow \alpha_k = -\frac{\nabla f(\underline{x}_k)^T \underline{d}_k}{\underline{d}_k^T H(\underline{x}_k) \underline{d}_k}$$

- However, if we have the Hessian $H(\underline{x}_k)$, better search methods (e.g. Newton-Raphson) are available.

For Personal Use Only -bkwk2

convergence of steepest-descent method

- the steepest-descent algorithm converges linearly, i.e. $\rho = 1$

$$\|\underline{x}_{k+1} - \underline{x}^*\| = \beta \|\underline{x}_k - \underline{x}^*\|$$

β is close to 1 \rightarrow steepest-descent is slow.

- It can be shown using Kantorovich inequality that

$$\beta \leq \left(\frac{k-1}{k+1} \right)^2 \quad \begin{cases} k \geq 1 \rightarrow \beta \rightarrow 0 & (\text{good}) \\ k \rightarrow \infty \rightarrow \beta \rightarrow 1 & (\text{bad}) \end{cases}$$

where k is the l_2 condition no. of $H(\underline{x}^*)$, $K = \frac{\max(H(\underline{x}^*))}{\lambda_{\min}(H(\underline{x}^*))}$ note that the Hessian $H(\underline{x}^*)$ is hermitian (symmetric).

- In typical engineering problems, $H(\underline{x}^*)$ is quite ill-conditioned $\rightarrow \beta$ is close to 1.

- Ill-conditioned quadratic functions have elongated valleys or isocontours \rightarrow long valley of small gradients and a large distance to move \rightarrow slow to converge.

Newton-Raphson method.

- Assume we can evaluate the gradient $\nabla f(\underline{x})$ and the Hessian $H(\underline{x})$.

- $f(\underline{x})$ is approx. by a quadratic function using the function value, gradient and Hessian at the current iteration.

$$f(\underline{x}) = f(\underline{x}_k) + \nabla f(\underline{x}_k)^T (\underline{x} - \underline{x}_k) + \frac{1}{2} (\underline{x} - \underline{x}_k)^T H(\underline{x}_k) (\underline{x} - \underline{x}_k) + \text{H.O.T.}$$

We get a quadratic approx $q(\underline{x})$ of $f(\underline{x})$ by neglecting the higher order terms,

$$q(\underline{x}) = f(\underline{x}_k) + \nabla f(\underline{x}_k)^T (\underline{x} - \underline{x}_k) + \frac{1}{2} (\underline{x} - \underline{x}_k)^T H(\underline{x}_k) (\underline{x} - \underline{x}_k)$$

- To calculate the min. of $q(\underline{x})$, we set its gradient to zero,

$$\nabla q(\underline{x}) = \nabla f(\underline{x}_k) + H(\underline{x}_k)(\underline{x} - \underline{x}_k) = 0 \quad \rightarrow \quad \underline{x} = \underline{x}_k - [H(\underline{x}_k)]^{-1} \nabla f(\underline{x}_k)$$

- We use the min. of $q(\underline{x})$ as our next guess for the min. of $f(\underline{x})$.

$$\underline{x}_{k+1} = \underline{x}_k + \alpha_k \underline{d}_k, \quad \alpha_k = 1 \quad \text{and} \quad \underline{d}_k = -[H(\underline{x}_k)]^{-1} \nabla f(\underline{x}_k)$$

- We can extend the Newton-Raphson method by finding the step size α_k using line search at each iteration \rightarrow speeds up convergence, will find a min. (not max/saddle).

convergence of Newton-Raphson method.

- The Newton-Raphson algorithm converges quadratically, i.e. $\rho = 2$

$$\|\underline{x}_{k+1} - \underline{x}^*\| \leq \beta \|\underline{x}_k - \underline{x}^*\|^2$$

It converges to the min. in 1 iteration in a quadratic function $f(\underline{x})$.

- Newton-Raphson is not very robust — it may take the iteration away from a min. / not converge if the initial guess is not appropriate.
- w/o line search for step size α_k , if $H(\underline{x})$ is not PD, the method may converge to max/saddle.
- Newton-Raphson requires computing + inverting the Hessian, computation $\sim O(N^2)$, inversion $\sim O(N^3)$

For Personal Use Only -bkwk2

Bazilei-Borwein method

- We can generalise the secant method for higher dimensions. The Hessian $\nabla^2 f(x_k)$ in Newton-Raphson is approx as follows in the Bazilei-Borwein method.

Quasi-Newton method
$$x = x_1 - \nabla f(x_1) \frac{(x_1 - x_0)^T (\nabla f(x_1) - \nabla f(x_0))}{(\nabla f(x_1) - \nabla f(x_0))^2}$$

- Borwein-Borwein method only requires computing the gradient $\nabla f(x_k)$, and is usually more stable than Newton's method or other Quasi-Newton methods.

- Not much is known about its convergence.

Conjugate gradient method.

- If we could step along the eigendirections, the convergence would be faster than steepest descent.
(even for ill-conditioned quadratic function)
(Eigendirections are good search dir. as they represent the principal axes of ellipsoids)
- However, computing the eigenvectors $N(O(N^2))$ - as bad as Hessian-based method \rightarrow don't use eigenvectors.
- Consider a set of Δ -conjugate vectors $\{d_1, \dots, d_N\}$, i.e.

true for symmetric A . $d_i^T A d_j = 0 \quad \text{for } i \neq j.$

(Orthogonal eigenvectors fulfill the conjugacy condition, but there are other vectors that satisfy this)

- Consider transforming the coords $x \rightarrow y$, where $x = \Sigma y$, $\Sigma = [d_1^T \ d_2^T \ \dots \ d_N^T]$,

For a quadratic function $f(x) = \frac{1}{2} x^T A x - b^T x$, we have

$$f(x) = f(\Sigma y) = \frac{1}{2} (\Sigma y)^T A (\Sigma y) - b^T (\Sigma y) = \frac{1}{2} y^T (\Sigma^T A \Sigma) y - (\Sigma^T b)^T y.$$

By construction, $\Sigma^T A \Sigma = I$.

$$\Sigma^T A \Sigma = \begin{bmatrix} \uparrow & \uparrow & \uparrow \\ d_1 & d_2 & \dots & d_N \end{bmatrix}^T \begin{bmatrix} \uparrow & \uparrow & \uparrow \\ d_1 & d_2 & \dots & d_N \end{bmatrix} = \begin{bmatrix} d_1^T d_1 & d_1^T d_2 & \dots & d_1^T d_N \\ d_2^T d_1 & d_2^T d_2 & \dots & d_2^T d_N \\ \vdots & \vdots & \ddots & \vdots \\ d_N^T d_1 & d_N^T d_2 & \dots & d_N^T d_N \end{bmatrix} = \begin{bmatrix} d_1^T d_1 & 0 & \dots & 0 \\ 0 & d_2^T d_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & d_N^T d_N \end{bmatrix} = I$$

\rightarrow The problem in the y space is a quadratic form along the axes (Search dir. $\{d_1, \dots, d_N\}$).

- At iteration k , we have the gradient $\nabla f(x_k)$ and the previous search dir. $\{d_1, \dots, d_{k-1}\}$.

We construct the new search dir. as

$$d_k = -\nabla f(x_k) + \beta_k d_{k-1}$$

β_k has memory of the previous dir. d_{k-1} .

Imposing the conjugacy condition $d_k^T A d_k = 0$,

$$d_k^T A (-\nabla f(x_k) + \beta_k d_{k-1}) = -d_{k-1}^T A \nabla f(x_k) + \beta_k d_{k-1}^T A d_{k-1} \rightarrow \beta_k = \frac{d_{k-1}^T A \nabla f(x_k)}{d_{k-1}^T A d_{k-1}}$$

It can be shown that w/ this choice of β_k , the search dir. d_k is Δ -conjugate to all previous search dir.

- The conjugate gradient algorithm is as follows:

4 1) Initialize w/ steepest descent $d_0 = -\nabla f(x_0)$

4 2) Update location and determine α_k w/ line search, $x_{k+1} = x_k + \alpha_k d_k$

4 3) Find a new search dir d_{k+1} , $d_{k+1} = -\nabla f(x_{k+1}) + \beta_{k+1} d_k = -\nabla f(x_{k+1}) + \left[\frac{\|\nabla f(x_{k+1})\|^2}{\|\nabla f(x_k)\|^2} \right] d_k$

4 4) Repeat steps 2-3 until convergence

For Personal Use Only -bkwk2

Convergence of Conjugate gradient method.

- The conjugate gradient algorithm converges linearly, i.e. $P=1$

$$\|\underline{x}_{k+1} - \underline{x}^*\| = \beta \|\underline{x}_k - \underline{x}^*\|$$

β is close to 0 \rightarrow Conjugate Gradient is superlinear.

- For large ill-conditioned problems, it can be shown that

$$\boxed{\beta \geq 1 - \frac{2}{JF}}.$$

where J is the l-2 condition no.

- The conjugate gradient method does not req. the Hessian $H(\underline{x}^*)$ and its inversion $(H(\underline{x}^*))^{-1}$.

- The conjugate gradient method is a Krylov subspace method — for a N -dimensional objective function $f(\underline{x})$, it converges after N iterations.

- Numerical errors can lead to loss of conjugacy \rightarrow periodically reset by setting $\underline{s}_k = -\nabla f(\underline{x}_k)$.

- The conjugate gradient method is originally used to solve linear systems $\underline{A}\underline{x} = \underline{b}$.

For optimisation of quadratic function, $f(\underline{x}) = \frac{1}{2} \underline{x}^T \underline{A} \underline{x} - \underline{b}^T \underline{x} \rightarrow \nabla f(\underline{x}) = \underline{A} \underline{x} - \underline{b} = 0$

Summary of multi-dimensional search-direction methods.

Method.	Advantages	Disadvantages
steepest-descent	Only needs gradient	Linear convergence Slow for ill-conditioned problems
Newton-Raphson	Quadratic convergence (more reliable w/ line search)	Needs Hessian Convergence
Barzilai-Borwein	Very stable. Faster than steepest descent	Convergence analysis difficult
conjugate gradient	Faster than steepest descent	Linear convergence

Least squares fitting of a model to data.

- Let $y^{(j)}$ be the measured data, and $\phi_j(\underline{x})$ be the model's prediction for $y^{(j)}$, $j=1, 2, \dots, m$

$\underline{x} = [x_1, \dots, x_n]^T$ is the vector of independent parameters, which can be tuned to fit the model to data.

- The elements of the residual vector $r(\underline{x})$ are the errors in the model's prediction for each data pt.

$$\boxed{r_j(\underline{x}) = \phi_j(\underline{x}) - y^{(j)}}$$

The cost function we want to min. is the sum of squared errors,

$$\boxed{f(\underline{x}) = \sum_{j=1}^m r_j(\underline{x})^2 = r(\underline{x})^T r(\underline{x})}$$

The gradient $\nabla f(\underline{x})$ is given by

$$\frac{\partial f}{\partial x_i} = \frac{\partial}{\partial x_i} \left(\sum_{j=1}^m r_j(\underline{x})^2 \right) = 2 \sum_{j=1}^m \frac{\partial \phi_j(\underline{x})}{\partial x_i} r_j(\underline{x}) \rightarrow \boxed{\nabla f(\underline{x}) = 2 J(\underline{x})^T r(\underline{x})}$$

where $\boxed{J(\underline{x}) = \begin{bmatrix} \frac{\partial \phi_1}{\partial x_1}, \dots, \frac{\partial \phi_m}{\partial x_1} \\ \vdots \\ \frac{\partial \phi_1}{\partial x_n}, \dots, \frac{\partial \phi_m}{\partial x_n} \end{bmatrix}}$ is the Jacobian matrix of $r(\underline{x})$.

The Hessian $H(\underline{x})$ is given by

$$\frac{\partial^2 f}{\partial x_i^2} = \frac{\partial}{\partial x_i} \left(2 \sum_{j=1}^m \frac{\partial \phi_j(\underline{x})}{\partial x_i} r_j(\underline{x}) \right) = 2 \sum_{j=1}^m \frac{\partial^2 \phi_j(\underline{x})}{\partial x_i^2} r_j(\underline{x}) + 2 \sum_{j=1}^m \frac{\partial \phi_j(\underline{x})}{\partial x_i} \frac{\partial r_j(\underline{x})}{\partial x_i} \rightarrow \boxed{H(\underline{x}) = 2 \sum_{j=1}^m r_j(\underline{x}) R(\underline{x}) + 2 J(\underline{x})^T J(\underline{x})}$$

where $\boxed{R(\underline{x}) = \nabla(\nabla f(\underline{x}))}$ is the Hessian of the residual $r(\underline{x})$

For Personal Use Only -bkwk2

Gauss-Newton method.

- The Gauss-Newton method assumes that the residuals at the optimum $r_j(x^*)$ are small, therefore $r_j(x)$ can be neglected, so the Hessian $H(x)$ can be approximated as,

$$H(x) \approx \tilde{H}(x) = 2J(x)^T J(x).$$

for x near the minimum x^* , and $\tilde{H}(x)$ is PD.

- Applying the approx. for the Hessian $H(x)$ to the Newton-Raphson method, the search dir. g_k is

$$g_k = -[\tilde{H}(x_k)]^{-1} \nabla f(x_k) = -[2J(x_k)^T J(x_k)]^{-1} 2J(x_k)^T r(x_k) = -J(x_k)^+ r(x_k)$$

cannot distribute the -1 as $J(x_k)$ may not be square.
where $J(x_k)^+ = [J(x_k)^T J(x_k)]^{-1} J(x_k)^T$ is the pseudo inverse matrix

- The estimate is updated as for the Newton-Raphson method, $x_{k+1} = x_k + g_k$.
- The Gauss-Newton method only requires the Jacobian (first derivatives) and avoids computing or inverting the Hessian $H(x)$.
- The Gauss-Newton method converges nearly quadratically if the initial residuals are small (otherwise, it has poor performance)

Constrained optimisation — Linear programming.

Linear programming (LP)

- LP is an optimisation problem in which both the objective and constraints are linear
- The standard / canonical form of a linear program is as follows:

Minimise	$f(x) = c^T x$	$x \in \mathbb{R}^n$
subject to m equality constraints	$Ax = b$	$A \in \mathbb{R}^{m \times n}$
and bounds	$x_i \geq 0$	$i = 1, \dots, n$

* We assume that $\text{rank}(A) = m$

- If there are inequality constraints of the form $a_j^T x \geq b_j$, we can convert them into equality constraints by introducing a new slack variable s_j , $s_j = a_j^T x - b_j$, where $s_j \geq 0$.
- The m constraints $Ax = b$ define a subspace of \mathbb{R}^n (usually $m < n$ o/w the feasible set is empty / $A^T b$)
- The bounds $x_i \geq 0$ cut out a portion of this subspace to define the feasible region (convex polyhedron)
- The feasible region has a no. of corners - extremal pts, which are the basic solns. The basic solns have m non-zero basic variables and $n-m$ zero free (non-basic) variables
- The isocontours of f are a series of parallel $(n-1)$ -dim hyperplanes. As f increasing, the first isocontour to intersect the feasible region will do so at the minimum.
- * We assume the problem is non-degenerate \rightarrow there are m (non-zero) basic variables (o/w, cycling may occur)
- * For low dimensional problems, LP can be solved graphically.

For Personal Use Only -bkwk2

Fundamental theorem of linear programming.

- Graph a LP in standard form where $A \in \mathbb{R}^{m \times n}$ and $\text{rank}(A) = m$, if there is a feasible soln, there is a basic feasible soln; if there is an optimal feasible soln, there's an optimal basic feasible soln.
- We can solve a linear program by inspecting basic feasible solns until we find the optimum.

The simplex algorithm.

- Large LPs cannot be easily solved by hand → simplex algorithm systematically solves LPs.
- The simplex algorithm involves two phases:
 - (i) Phase 1: start the algorithm.
 - ↳ Find an initial basic feasible soln if there is not an obvious one
 - ↳ Start from the initial basic feasible soln (vertex of feasible region w/ m (non-zero) basic variables)
 - (ii) Phase 2: move from one vertex to another until the min. is found
 - ↳ Move along an edge of the feasible region to another vertex where f is smaller
 - ↳ Repeat until you find a vertex where edge leading away would increase f → this is the min.

Phase 2: Move along one vertex to another (Tabular method)

- Step 0: Formulate the problem

$$\text{minimise } c^T x = 200x_0 + 400x_1 + 680x_2 + 800x_3$$

$$\text{subject to } Ax = b \Leftrightarrow \begin{bmatrix} 3 & 7 & 5 & 3 \\ 0 & 3 & 5 & 7 \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 80 \\ 40 \end{bmatrix} \Leftrightarrow \begin{array}{l} 3x_0 + 7x_1 + 5x_2 + 3x_3 = 80 \\ 3x_1 + 5x_2 + 7x_3 = 40 \end{array}$$

$$\text{and } x_i \geq 0$$

↳ We can collect the constraint eqns and the objective function coeff into a single tableau ((m+1)x(n+1) matrix)

$$\left[\begin{array}{cc|c} A & b \\ \hline c^T & 0 \end{array} \right] = \left[\begin{array}{cccc|c} 3 & 7 & 5 & 3 & 80 \\ 0 & 3 & 5 & 7 & 40 \\ \hline 200 & 400 & 680 & 800 & 0 \end{array} \right] \begin{array}{l} \text{constraints} \\ \text{variables} \end{array}$$

- Step 1: start from a basic soln.

↳ From those 1, we get an initial basic feasible soln $\begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 10 \\ 2 \\ 0 \end{bmatrix} \rightarrow x_1, x_2 \text{ are the basic variable.}$

↳ Manipulate the tableau to get canonical form. (under basic variables, we get $\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$)

Eliminate the coeff. of the basic variables in the last line (get reduced cost, -ve of objective function)

$$\begin{aligned} [4]: & \frac{1}{4}x([1]-[2]) & \begin{bmatrix} x_0 & x_1 & x_2 & x_3 & \text{RHS} \\ \hline 3/4 & 1 & 0 & -1 & 10 \end{bmatrix} & [4] \\ [5]: & \frac{1}{2}x([2]-3x[4]) & \begin{bmatrix} -9/20 & 0 & 1 & 2 & 2 \\ 20/6 & 0 & 0 & -160 & -520 \end{bmatrix} & [5] \\ [6]: & [3]-400[4]-400[5] & \begin{bmatrix} 0 & 0 & 0 & 0 & 0 \end{bmatrix} & [6] \end{aligned}$$

- Step 2: Identify the edge to follow

↳ choose the entering variable to be the one w/ the most -ve reduced cost $\rightarrow x_3$ is EV

↳ To choose the leaving variable, consider the ratios $\frac{\text{RHS}}{x_3}$ for each row, i.e. $\frac{10}{-1} \text{ vs } \frac{2}{2}$

We choose the variable corresponding to the smallest, -ve ratio $\rightarrow x_2 \text{ is LV}$

- Step 3: calculate the next vertex

↳ Repeat steps 1,2 until all reduced costs are true. $\begin{bmatrix} x_0 & x_1 & x_2 & x_3 \\ \hline 0.525 & 0 & 0.5 & 1 \\ -0.225 & 0 & 0.5 & 1 \\ 170 & 0 & 80 & 0 \end{bmatrix} \rightarrow \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$

For Personal Use Only -bkwk2

Phase 1: start the algorithm

- To find an initial basic feasible soln, we construct an auxiliary optimisation problem, where the objective function is the "error" that measures by how much the constraints are violated.
- The optimal sol'n of this problem has zero error \rightarrow basic feasible sol'n.
- We create m slack variables, add one to each of the constraint eqns. The objective function is the sum of the slack variables. (the slack variables are the initial basic variables).
- Consider the previous example. The auxiliary problem is

$$\begin{aligned} 3x_0 + 7x_1 + 5x_2 + 3x_3 + s_1 + 0 &= 80 \\ 0 + 3x_1 + 5x_2 + 7x_3 + 0 + s_2 &= 40 \end{aligned}$$

The auxiliary objective function is $s_1 + s_2$, so the auxiliary tableau is

$$\left[\begin{array}{cccccc|c} x_0 & x_1 & x_2 & x_3 & s_1 & s_2 & RHS \\ \hline 3 & 7 & 5 & 3 & 1 & 0 & 80 \\ 0 & 3 & 5 & 7 & 0 & 1 & 40 \\ 0 & 0 & 0 & 0 & 1 & 1 & 0 \end{array} \right] \begin{matrix} [1] \\ [2] \\ [3] \end{matrix}$$

Eliminating the coeff of the basic variables in the last line,

$$\left[\begin{array}{cccccc|c} x_0 & x_1 & x_2 & x_3 & s_1 & s_2 & RHS \\ \hline 3 & 7 & 5 & 3 & 1 & 0 & 80 \\ 0 & 3 & 5 & 7 & 0 & 1 & 40 \\ 0 & -10 & -10 & -10 & 0 & 0 & -120 \end{array} \right] \begin{matrix} [1] \\ [2] \\ [4]: [3]-[1]-[2] \end{matrix}$$

Then solve the auxiliary tableau as in phase 2.

- After solving the auxiliary tableau, if $s_1 = s_2 = 0$ (slack variables are non-basic variables), and the cost function equals 0, then we have found the initial basic feasible soln.
- If $s_1 > 0$ and $s_2 > 0$, there is no feasible soln to the original problem.

Constrained optimisation — nonlinear constrained optimisation

Nonlinear constrained optimisation

- consider the optimisation problem

$$\begin{aligned} \text{minimise} \quad f(\underline{x}) &\quad \text{generally nonlinear} \\ \text{subject to} \quad h_i(\underline{x}) &= 0 \quad i=1,2,\dots,m \quad [\text{equality constraints}] \\ g_j(\underline{x}) &\leq 0 \quad j=1,2,\dots,n \quad [\text{inequality constraints}] \end{aligned}$$

The set of equality and inequality constraints defines the feasible region.

- At a location \underline{x} , constraints are

equations, active inequalities

\hookrightarrow active if they do not allow \underline{x} to infinitesimally change in at least one dir.

\hookrightarrow inactive if they allow \underline{x} to infinitesimally change in any dir. \hookrightarrow inactive inequalities.

- In N dimensions, it is not possible to have a sol'n w/ more than N independent constraints active.

For Personal Use Only -bkwk2

Optimality condition on active constraints

- A necessary condition for \underline{x}^* to be a local min. is

$$\nabla f(\underline{x}^*) \cdot \underline{d} \geq 0 \quad \text{for all feasible directions } \underline{d}.$$

In constrained optimization, some directions \underline{d} may not be feasible, so in general,

$$\nabla f(\underline{x}^*) \cdot \underline{d} \geq 0 \Rightarrow \nabla f(\underline{x}^*) \cdot \underline{d} = 0$$

- consider that \underline{x} is subject to one equality constraint $h(\underline{x}) = 0$.

Let \underline{d} be the tangent dir. to the constraint. Suppose $\exists \underline{d}$ is a feasible dir.,

↳ If $\nabla f(\underline{x}^*) \cdot \underline{d} < 0 \rightarrow \underline{x}$ is not a min.

↳ If $\nabla f(\underline{x}^*) \cdot \underline{d} > 0 \rightarrow$ in opposite dir. $\underline{d}' = -\underline{d}$, $\nabla f(\underline{x}^*) \cdot \underline{d}' < 0 \rightarrow \underline{x}$ is not a min.

$\Rightarrow \underline{x}$ is a min. if $\nabla f(\underline{x}^*) \cdot \underline{d} = 0$.

Geometrically, $\nabla f(\underline{x}^*) + \underline{d} \rightarrow \nabla f(\underline{x}^*) \parallel \nabla h(\underline{x}^*)$, i.e. $\nabla f(\underline{x}^*) = -\lambda \nabla h(\underline{x}^*)$ where λ is a scalar

- w/ multiple independent active constraints, at the min., ∇f must be perpendicular to the tangent plane of the constraint intersection, which has dimensionality N_m . The optimality condition is thus

$$\nabla f(\underline{x}^*) = - \sum_{i=1}^m \lambda_i \nabla h_i(\underline{x}^*) = - \nabla h(\underline{x}^*)^T \underline{\lambda}, \quad \text{where } \lambda_i, i=1,2,\dots,m \text{ are scalars}$$

Lagrangian and optimality condition for equality constraints

- we can recover the optimality condition by defining an unconstrained min. problem w/ a new cost function

$$L(\underline{x}, \lambda_i) = f(\underline{x}) + \sum_{i=1}^m \lambda_i h_i(\underline{x}) \quad \text{note } h_i(\underline{x}) = 0$$

where $L(\underline{x}, \lambda_i)$ is the Lagrangian and λ_i are the Lagrangian multipliers (scalars)

- \underline{x}^* is a stationary pt. of the Lagrangian $L(\underline{x}, \lambda_i)$ if

$$\nabla L(\underline{x}^*, \lambda) = 0 \Rightarrow \nabla f(\underline{x}^*) + [\nabla h(\underline{x}^*)]^T \underline{\lambda} = 0 \quad \text{and} \quad h_i(\underline{x}) = 0 \quad \text{for } i=1,2,\dots,m$$

A sufficient condition for \underline{x}^* to be a minimum is

$$\underline{d}^T \nabla^2 L(\underline{x}^*, \lambda) \underline{d} > 0 \quad \text{where } \underline{d} \in \text{space tangent to all constraints.}$$

i.e. the Hessian $\nabla^2 L(\underline{x}^*, \lambda)$ is PD in the tangent space.

- + IF the Lagrangian multiplier is zero, $\lambda_i = 0$, then $\nabla f(\underline{x}^*) = 0$ at that pt.

For Personal Use Only -bkwk2

Lagrangian and optimality conditions for inequality constraints

- For inequality constraints, we define the Lagrangian

$$L(\underline{x}, M_j) = f(\underline{x}) + \sum_{j=1}^m M_j g_j(\underline{x}) \quad \text{Note } g_j(\underline{x}) \leq 0$$

\underline{x}^* is a stationary pt. of the Lagrangian $L(\underline{x}, M_j)$ if

$$\nabla L(\underline{x}^*, M_j) = 0 \Rightarrow \nabla f(\underline{x}^*) + [\nabla g(\underline{x}^*)]^T M = 0 \quad \text{and} \quad M_j g_j(\underline{x}^*) = 0, \quad M_j \geq 0 \quad \text{for } j=1, 2, \dots, n$$

These are the Karush-Kuhn-Tucker (KKT) conditions, M_j are the KKT multipliers.

- For active constraints, $g_j(\underline{x}) = 0$ [constraint becomes an equality constraint]

For inactive constraints, $g_j(\underline{x}) < 0$ and $M_j = 0$

- In the general case w/ both equalities and inequality constraints,

$$L(\underline{x}, \lambda_i, M_j) = f(\underline{x}) + \sum_{i=1}^m \lambda_i h_i(\underline{x}) + \sum_{j=1}^m M_j g_j(\underline{x})$$

Sensitivity.

- Consider perturbing the sol'n \underline{x}^* by an infinitesimally small amount, i.e. $\underline{x}^* \rightarrow \underline{x}^* + \delta \underline{x}$.

$$\nabla f(\underline{x}^*) + [\nabla h(\underline{x}^*)]^T \lambda + [\nabla g(\underline{x}^*)]^T M = 0 \rightarrow \delta \underline{x}^* (\nabla f(\underline{x}^*) + [\nabla h(\underline{x}^*)]^T \lambda + [\nabla g(\underline{x}^*)]^T M)$$

$$\therefore \delta f(\underline{x}^*) = -[\nabla h(\underline{x}^*)]^T \lambda - [\nabla g(\underline{x}^*)]^T M$$

→ change in cost function due to a change in the sol'n is prop. to the -ve Lagrangian/KKT multipliers.

i.e. $-\lambda_i$ and M_j are the sensitivities of f to an infinitesimal change in the constraint h_i and g_j respectively.

Penalty functions and barrier functions

- w/ penalty functions, we do not enforce the constraints exactly, but only approx and controllably.

(i.e. soft constraints), we define a new unconstrained problem

$$q(\underline{x}, K) = f(\underline{x}) + K \sum_{i=1}^m h_i(\underline{x})^2$$

where K is the penalty parameter and $K \sum_{i=1}^m h_i(\underline{x})^2$ is the penalty function.

- It can handle both equalities and inequalities (for inequalities, we use $K \sum_{j=1}^m \max[0, g_j(\underline{x})]^2$)

- For large K , $h(\underline{x})$ needs to be very small → problem close to constrained problem. For small K ,

$h(\underline{x})$ does not need to be small → constraint violated more, usually start w/ small K and increase it.

- w/ barrier functions, the feasibility is enforced exactly by using a penalty that diverges as we

approach the boundary of the constraint. We define a new unconstrained problem

inverse barrier function

$$q(\underline{x}, K) = f(\underline{x}) - \frac{1}{K} \sum_{j=1}^m \frac{1}{g_j(\underline{x})}$$

logarithmic barrier function

$$q(\underline{x}, K) = f(\underline{x}) - \frac{1}{K} \sum_{j=1}^m \ln[g_j(\underline{x})]$$

- It can only handle inequalities and always yields a feasible answer.

- If req. a feasible starting pt. and cannot iterate through infeasible space.

- It basically restricts the search to the feasible region by penalizing sol'n's infeasible edge.

- The larger the K , the sharper the divergence of the singularity near the constraint boundary.