

Running Lecture Outline: MCS65C Multivariable Mathematics

Aaron Wang

Academic Year 2019-2020

Note: All addendums to a day made afterwards will be silently appended to that day

Contents

| | | |
|----------|---|----------|
| 1 | Fall 2019 | 3 |
| 1.1 | Sept 9: Foundations of \mathbb{R} | 3 |
| 1.2 | Sept 10: More axioms of \mathbb{R} ; sup and inf | 3 |
| 1.3 | Sept 11: Sequences and Convergence | 4 |
| 1.4 | Sept 12: Cauchy Sequences and the Bolzano-Weierstrass Theorem | 5 |
| 1.5 | Sept 13: Finishing up CCC | 5 |
| 1.6 | Sept 16: Rolles' Theorem and Convergence in R^d | 6 |
| 1.7 | Sept 17: EVT and Compact Sets | 6 |
| 1.8 | Sept 18: Open and Closed | 6 |
| 1.9 | Sept 19: Sequential Closure | 7 |
| 1.10 | Sept 20: Arbitrary Unions and Intersections | 7 |
| 1.11 | Sept 23: Open Close Open Close Open Close | 7 |
| 1.12 | Sept 24: Compactness and Covering | 8 |
| 1.13 | Sept 25: Heine-Borel and Limits and Continuity | 8 |
| 1.14 | Sept 26: Continuity Theorems | 9 |
| 1.15 | Sept 27: Identifying continuous f | 9 |
| 1.16 | Oct 3: Continuous Rigor | 10 |
| 1.17 | Oct 4: More abstraction | 10 |
| 1.18 | Oct 7: EVT for More Functions | 10 |
| 1.19 | Oct 8: IVT | 10 |
| 1.20 | Oct 10: Connected Sets | 11 |
| 1.21 | Oct 15: Connected Complements | 11 |
| 1.22 | Oct 21: The Short Short Proof of IVT | 11 |
| 1.23 | Oct 22: Uniform Continuity | 11 |
| 1.24 | Oct 23: Fire Drilled | 12 |
| 1.25 | Oct 24: UCT | 12 |
| 1.26 | Oct 25: Halved; UCT? | 12 |
| 1.27 | Oct 26: UCT!, and Differentiation | 12 |
| 1.28 | Oct 29: Differentiability | 12 |
| 1.29 | Oct 30: Derivatives | 12 |
| 1.30 | Oct 31: Partial Derivatives | 13 |
| 1.31 | Nov 1: Tangent Planes | 13 |
| 1.32 | Interlude I | 13 |
| 1.33 | Nov 5: Chain Rule | 14 |
| 1.34 | Nov 6: Proof of the Above | 14 |
| 1.35 | Nov 7: Matrix Algebra | 14 |
| 1.36 | Nov 8: Multiplication of Matrices | 15 |
| 1.37 | Nov 12: Utility of Linear Maps | 15 |
| 1.38 | Nov 13: Transposition and Inversion | 16 |
| 1.39 | Nov 14: Invertibility | 16 |
| 1.40 | Nov 15*: Squareness | 16 |

| | | |
|------|--|----|
| 2.20 | Mar 3: Still on Spectra | 37 |
| 2.21 | Mar 5*: More Numbers Time | 37 |
| 2.22 | Mar 6*: 0 | 38 |
| 2.23 | Mar 9: Numbers Review! | 38 |
| 2.24 | Mar 10: Fake Projectors | 38 |
| 2.25 | Mar 11: Megagram-Schmidt | 39 |
| 2.26 | Mar 12: Real Projectors | 39 |
| 2.27 | Mar 13: in Absentio | 39 |
| 2.28 | Mar 16: The Far Reaches | 39 |
| 2.29 | Mar 19: Spectral Thm Proved | 40 |
| 2.30 | Mar 20: Spectral Thm Proved II | 40 |
| 2.31 | Mar 23: Spectral Thm Proved III (for real) | 42 |
| 2.32 | Mar 24: Orthogonal Maps | 42 |
| 2.33 | Mar 25: Spectral Thm IV (we'll get them this time!) | 43 |
| 2.34 | Mar 30: Back to Sylvester | 44 |
| 2.35 | Apr 1: Sylvester Sylvester Sylvester Sylvester Sylvester | 45 |
| 2.36 | Apr 10: EVT and Friends | 46 |
| 2.37 | Apr 12: Critical Computation | 47 |
| 2.38 | Apr 13: Critical Proof | 47 |
| 2.39 | May 29: Lagrange Multipliers | 48 |
| 2.40 | Jun 1: et Ultra | 48 |
| 2.41 | Jun 5: Everybody Loves Proof | 49 |
| 2.42 | Jun 8: Everybody Kind of Loves Proof | 49 |
| 2.43 | That's a Wrap! | 50 |

1 Fall 2019

1.1 Sept 9: Foundations of \mathbb{R}

- \mathbb{R} can be defined with a set of axioms, and the additional elements $\{+, \cdot, 0, 1, \mathbb{R}^+\}$
- 0th order logic, or propositional logic, is not sufficient to define \mathbb{R} , but 1st order logic, which deals with predicates of individual elements, is.
 - The classical example of a 1st order conclusion is:

$$\frac{\begin{array}{l} \text{All humans are mortal} \\ \text{Socrates is a human} \end{array}}{\text{Socrates is mortal}}$$

- This can be quantified as:

$$\frac{\begin{array}{l} (\forall x)(H(x) \Rightarrow M(x)) \\ H(S) \end{array}}{M(S)}$$

where $H(x)$ represents " x is a human", $M(x)$ represents " x is mortal", and S is Socrates.

1.2 Sept 10: More axioms of \mathbb{R} ; sup and inf

- Comparison can be defined:
 - $x > y := x - y \in \mathbb{R}^+$
 - $x \geq y := x > y \vee x = y$
 - Less-than relations are defined similarly

- Transitivity of these relations can be proved:

$$\begin{aligned}x &< y \wedge y < z \\y - x, z - y &\in \mathbb{R}^+ \\y - x + z - y &= z - x \in \mathbb{R}^+ \\\therefore x &> z\end{aligned}$$

- The Completeness Axiom of \mathbb{R} states:

$$S \subseteq \mathbb{R} \wedge S \neq \emptyset \wedge \left(\exists b \in \mathbb{R} : S \subseteq (-\infty, b] \right) \Rightarrow \exists l : \left(S \subseteq (-\infty, l] \wedge (\forall \varepsilon > 0)(S \not\subseteq (-\infty, l - \varepsilon]) \right)$$

Then, $\sup S = l$, and $\inf S$ is defined similarly

1.3 Sept 11: Sequences and Convergence

- Armed with the Completeness Axiom, it is possible to define limits of sequences:

$$a_n \rightarrow l \Leftrightarrow (\forall \varepsilon > 0)(\exists N : (n \geq N \Rightarrow |l - a_n| < \varepsilon))$$

- It can be proved that elements of a set S can be arbitrarily close to $\sup S$:

$$(\forall \varepsilon > 0)(\exists x_\varepsilon \in S : |\sup S - x_\varepsilon| < \varepsilon)$$

- It can also be proved that a sequence of real numbers can approach at most one real limit; assume:

$$a_n \rightarrow l \wedge a_n \rightarrow l' \wedge l \neq l'$$

Then, with sufficiently large n ,

$$|l - a_n| < \varepsilon \wedge |l' - a_n| < \varepsilon$$

Let $\varepsilon = \frac{|l - l'|}{2}$, then add the inequalities:

$$|l - a_n| + |l' - a_n| < 2\varepsilon = |l - l'|$$

This violates the Triangle Inequality, so it must be true that $l = l'$

- Bounded Monotonic Convergence Theorem (BMCT):

$$(\forall n \geq 1)(a_n \leq a_{n+1} \wedge \exists b : b \geq a_n) \Rightarrow \exists l \in \mathbb{R} : a_n \rightarrow l$$

$$l = \sup\{a_n | n \geq 1\}$$

Let $l = \sup\{a_n | n \geq 1\}$, which exists since (a_n) is bounded above. The members of (a_n) can be arbitrarily close to l , so it can be proved that,

$$(\forall \varepsilon > 0)(\exists n : |l - a_n| < \varepsilon)$$

Because the sequence is increasing,

$$(\forall n \geq 1)(|l - a_{n+1}| \leq |l - a_n|)$$

Conclusion:

$$(\forall \varepsilon > 0)(\exists N : (n \geq N \Rightarrow |l - a_n| < \varepsilon))$$

$$a_n \rightarrow l$$

- A Cauchy sequence (c_n) converges, and satisfies:

$$(\forall \varepsilon > 0)(\exists N : (m, n \geq N \Rightarrow |c_m - c_n| < \varepsilon))$$

1.4 Sept 12: Cauchy Sequences and the Bolzano-Weierstrass Theorem

- If a sequence is convergent, it is cauchy
- Convergent sequences are cauchy; set $\varepsilon_1 = 2\varepsilon_2 > 0$ to be the tolerance of the sequence's limit:

$$a_n \rightarrow l \Rightarrow |a_n - a_m| = |a_n - l + l - a_m| \leq |a_n - l| + |l - a_m|$$

With sufficiently large n, m , the rightmost expression is less than $\varepsilon_1 + \varepsilon_1 = \varepsilon_2$

$$a_n \rightarrow l \Rightarrow (\forall \varepsilon_2 > 0)(\exists N : m, n > N \Rightarrow |a_n - a_m| < \varepsilon_2)$$

$\therefore a_n$ is cauchy

- The converse is provable too, but with more difficulty
 - The Sunrise Lemma states that every sequence contains a monotone subsequence
 - * Let n be *pretty* if $(\forall k \geq 1)(a_n > a_{n+k})$
 - * Assume there are infinitely many *pretty* n :
 - Let n_k be the k^{th} *pretty* n , (a_{n_k}) is a strictly increasing sequence
 - * Assume there are finitely many *pretty* n :
 - Let $m_1 = n_0 + 1$, where n_0 is the greatest *pretty* n
 - m_1 is not *pretty*, so $(\exists m_2 : a_{m_1} \leq a_{m_2})$
 - m_k , where $k \geq 1$ is not *pretty*, so $(\exists m_{k+1} : a_{m_k} \leq a_{m_{k+1}})$
 - By induction, (a_{m_k}) , where $k \geq 1$, is a decreasing sequence
 - A bounded sequence must also contain a monotone subsequence by the Sunrise Lemma, which is also bounded. Using the BMCT, the Bolzano-Weierstrass Theorem asserts that every bounded sequence has at least 1 convergent subsequence
 - Two additional lemmas:
 - * Every cauchy sequence is bounded:
 - Select any $\varepsilon > 0$, $\exists N_\varepsilon : (m \geq n \geq N \Rightarrow |a_m - a_n| < \varepsilon)$
 - All terms of the cauchy sequence past a_n are within ε of a_n
 - The remaining finite terms are obviously bounded
 - * A cauchy sequence (a_n) with a subsequence (a_{k_n}) converging to l converges to l :
 - Select any $\varepsilon_1, \varepsilon_2 > 0$, let N be the greater of the two cut-off points for the subsequence to be close enough to l , and the cauchy sequence's terms to be close enough to each other, and let $k \in \{k_n | n \geq 1\}$
 - $\exists N : (k, n > N \Rightarrow |a_n - a_k| < \varepsilon_1 \wedge |a_k - l| < \varepsilon_2)$
 - $|a_n - a_k| < \varepsilon_1 \wedge |a_k - l| < \varepsilon_2 \Rightarrow |a_n - a_k + a_k - l| = |a_n - l| < \varepsilon_1 + \varepsilon_2 = \varepsilon_3$
 - Consequently, $(\forall \varepsilon_3 > 0)(\exists N : (n \geq N \Rightarrow |a_n - l| < \varepsilon_3))$, so $a_n \rightarrow l$
 - By chaining the three results above, it is now easily proved that cauchy sequences are convergent

1.5 Sept 13: Finishing up CCC

- The two lemmas above were proved
- Reverse and Regular Triangle Inequality:

$$||x| - |y|| \leq |x + y| \leq |x| + |y|$$

1.6 Sept 16: Rolles' Theorem and Convergence in \mathbb{R}^d

- Rolles' Theorem:

$$f(x) \text{ is differentiable on } [a, b] \wedge f(a) = f(b) \Rightarrow \exists c \in [a, b] : f'(c) = 0$$

By the Extreme Value Theorem, $f(x)$ must have global extrema, which are also local extrema

Let $m \in (a, b)$ be a point where there is such a local maximum

If $n < m$ and n is close to m , then $f'(n) \geq 0$, and if $n > m$ and they are close, $f'(n) \leq 0$

Let n approach m from either side, the derivative at m must be both ≥ 0 and ≤ 0 , so it must be 0

If the global maximum is at the endpoints, but $f(x)$ has a global minimum in (a, b) , just invert the function, find that the derivative at the new global maximum is 0, and invert it again

If the global maximum and global minimum are both at the endpoints, then $f(x) = c$ and the proof is trivial

- Distance definitions: We use Euclidean Distance $dist(\mathbf{P}, \mathbf{Q}) = \sqrt{\sum_{k=1}^d (P_k - Q_k)^2}$ because arbitrary rotations of two points preserve the distance between them
The generalized Minkowski Distance is:

$$dist_m(\mathbf{P}, \mathbf{Q}) = \sqrt[m]{\sum_{k=1}^d |P_k - Q_k|^m}$$

For $m < 1$, the Triangle Inequality does not hold, so it is less useful

- There are useful bounds on the distance between two points:

$$\max_{1 \leq j \leq d} |P_j - Q_j| \leq dist(\mathbf{P}, \mathbf{Q}) \leq \sqrt{d} \max_{1 \leq j \leq d} |P_j - Q_j|$$

The first inequality is obvious and the second is derived from replacing each $|P_k - Q_k|$ in the sum with the maximum $|P_j - Q_j|$

- Convergence in \mathbb{R}^d is now possible to define:

$$\mathbf{P}_n \rightarrow \mathbf{q} \Leftrightarrow dist(\mathbf{q}, \mathbf{P}_n) \rightarrow 0$$

Combining this with the upper bound on the distance from earlier allows to use coordinatewise convergence:

$$\mathbf{P}_n \rightarrow \mathbf{q} \Leftrightarrow (\forall k \leq d)(P_{n_k} \rightarrow Q_k)$$

1.7 Sept 17: EVT and Compact Sets

- The norm of a vector, $|\mathbf{P}|$ is defined as $\sqrt{\sum_{j=1}^d P_j^2}$, and consequently, $dist(\mathbf{P}, \mathbf{Q}) = |\mathbf{P} - \mathbf{Q}|$
- A compact set is closed and bounded:
A bounded set is a subset of a sufficiently large n -sphere centered on the origin
A closed set includes its boundary;
A boundary point is one where an arbitrarily small n -sphere centered on it contains interior and exterior points
- The Extreme Value Theorem for real-valued functions:

$$f \text{ is continuous on a compact set } K \in \text{dom } f \Rightarrow \exists \mathbf{p}, \mathbf{q} \in K : (\forall \mathbf{x} \in K)(f(\mathbf{p}) \leq f(\mathbf{x}) \leq f(\mathbf{q}))$$

1.8 Sept 18: Open and Closed

- The existence of the boundary of a set K , $\text{bd } K$, can be defined precisely:

$$\mathbf{p} \in \mathbb{R}^d \wedge S \subseteq \mathbb{R}^d \wedge (\forall \varepsilon > 0)(\exists \mathbf{q}_\varepsilon \in S, \mathbf{r}_\varepsilon \in S^C : |\mathbf{q}_\varepsilon - \mathbf{p}| < \varepsilon \wedge |\mathbf{r}_\varepsilon - \mathbf{p}| < \varepsilon) \Rightarrow \mathbf{p} \in \text{bd } S$$

- Open set: $K \cap \text{bd } K = \emptyset$
- Closed set: $K \cap \text{bd } K = K$, so $\text{bd } K \subseteq K$
- A set S shares a boundary with its complement, so $\text{bd } S = \text{bd } S^C$
- Two sets in \mathbb{R}^d are both open and closed: \emptyset and \mathbb{R}^d , since their boundaries are both \emptyset
- If a set S in \mathbb{R}^d is closed, it is sequentially closed:
Any convergent sequence contained in S converges to a point in S

1.9 Sept 19: Sequential Closure

- Closed sets are sequentially closed:
Assume set S is not sequentially closed, $\exists \mathbf{p}_n \in S : \mathbf{p}_n \rightarrow \mathbf{q} \in S^C$, but it is closed, so $\text{bd } S$ is disjoint from S^C
Therefore, $\mathbf{q} \notin \text{bd } S \wedge \mathbf{q} \notin S$, so there are no elements of S arbitrarily close to \mathbf{q}
However, sequence (\mathbf{p}_n) , fully contained in S must have members arbitrarily close to \mathbf{q}
- Sequentially closed sets are closed:
Assume set S is open, $\text{bd } S \subsetneq S$, and let some sequence (\mathbf{p}_n) contained in S converge to some \mathbf{q} ;
Select an arbitrary $\mathbf{q} \in \text{bd } S$, so $\mathbf{q} \notin S$
This is possible because the elements of S can be arbitrarily close to a point on $\text{bd } S$
However, because (\mathbf{p}_n) is contained in S , and because the set S is sequentially closed, $\mathbf{q} \in S$
- The interior of set S , S^{int} is $S \setminus \text{bd } S$, and it is the largest open subset of S
Assume there was a set U , such that $S^{\text{int}} \subset U \subseteq S$
 U must have extra points in $\text{bd } S$ for this to be possible
 U is open, so every point in U is in $U^{\text{int}} \subseteq S^{\text{int}}$ However, $\text{bd } S \cap S^{\text{int}} = \emptyset$
- The closure of set S , \bar{S} is $S \cup \text{bd } S$, and it is the smallest closed subset of S
This can be proved by essentially negating the previous theorem

1.10 Sept 20: Arbitrary Unions and Intersections

- Four facts about open sets U_α and closed sets C_α :

$$\bigcup_{\alpha \in A} U_\alpha \text{ is open for any } A$$

$$\bigcap_{\alpha \in A} U_\alpha \text{ is open for finite } A$$

If A were infinite, the intersection might be closed, consider $\bigcap_{n \in \mathbb{N}} \left(-\frac{1}{n}, 1 + \frac{1}{n}\right) = [0, 1]$

$$\bigcap_{\alpha \in A} C_\alpha \text{ is closed for any } A$$

$$\bigcup_{\alpha \in A} C_\alpha \text{ is closed for finite } A$$

If A were infinite, the union might be open, consider $\bigcup_{n \in \mathbb{N}} \left[\frac{1}{n}, 1 - \frac{1}{n}\right] = (0, 1)$

1.11 Sept 23: Open Close Open Close Open Close

- Arbitrary unions of open sets are open because an arbitrary point in one of the sets will be fully contained in it, and so it must be fully contained in the union too
- Finite intersections of open sets are open; to fully contain a ball in an intersection of infinitely many, the radius must be a minimum of arbitrary reals; this minimum may not exist at all
- The similar theorems for closed sets are derived with De Morgan's Laws

1.12 Sept 24: Compactness and Covering

- Compact sets have already been defined as closed and bounded:

$$K \text{ is compact} \Leftrightarrow \text{bd } K \subseteq K \wedge \exists r : B_r(0) \supseteq K$$

- Arbitrary intersections of compact sets are compact, since to bound the intersection, it is sufficient to bound any one set, and intersections of closed sets have been established to be closed
- Finite unions of compact sets are compact, since to bound the union; to bound an infinite union, it is necessary to have a radius equal to the maximum of arbitrary reals, which may not exist
- For S^C to be compact, S must be open and a superset of $(B_r(0))^C$ for some r
- The Heine-Borel Covering Property:

$$\text{compact } K \subseteq \bigcup_{n=1}^{\infty} U_n \rightarrow \exists N : K \subseteq \bigcup_{n=1}^N U_n$$

This property even holds if the original set of sets is uncountable

1.13 Sept 25: Heine-Borel and Limits and Continuity

- Prove by contradiction:

$$\text{Assume } (\forall N) \left(\exists \mathbf{p}_n : \mathbf{p}_n \in \text{compact } K \wedge \mathbf{p}_n \notin \bigcup_{\alpha=1}^N U_\alpha \right)$$

K is bounded, so by the Bolzano-Weierstrass Theorem:

$$\exists (\mathbf{p}_{k_n}) : \mathbf{p}_{k_n} \rightarrow \mathbf{p}$$

K is closed, so it is sequentially closed:

$$\mathbf{p} \in K \subset \bigcup_{\alpha=1}^N U_\alpha$$

$$\mathbf{p} \in U_m \Rightarrow \exists r : B_r(\mathbf{p}) \subseteq U_m$$

\mathbf{p}_{k_n} must get arbitrarily close to \mathbf{p} , so for sufficiently large n and N ,

$$\mathbf{p}_{k_n} \in B_r(\mathbf{p}) \subseteq U_m \subset \bigcup_{\alpha=1}^N U_\alpha$$

- Let the function $\mathbf{f} : D \subseteq \mathbb{R}^d \rightarrow \mathbb{R}^e, D \neq \emptyset$ exist
- Limits of multivariate functions are defined like they are for real-valued functions:

$$(\forall \varepsilon > 0) \left(\exists \delta > 0 : (\forall \mathbf{x} \in D) (\mathbf{x} \in B_\delta(\mathbf{p}) \Rightarrow \mathbf{f}(\mathbf{x}) \in B_\varepsilon(\mathbf{l})) \right) \Rightarrow \mathbf{l} = \lim_{\mathbf{x} \rightarrow \mathbf{p}} \mathbf{f}(\mathbf{x})$$

Here, $\mathbf{p} \in \bar{D}$

- Continuity is familiar too:

$$\lim_{\mathbf{x} \rightarrow \mathbf{p}} \mathbf{f}(\mathbf{x}) = \mathbf{f}(\mathbf{p}) \Rightarrow \mathbf{f} \text{ is continuous at } \mathbf{p}$$

Here, $\mathbf{p} \in D$

1.14 Sept 26: Continuity Theorems

- Continuity can be defined in a more easily used form via limits:

$$(\forall \varepsilon > 0) \left(\exists \delta > 0 : (\forall \mathbf{x} \in D) (\mathbf{x} \in B_\delta(\mathbf{p}) \Rightarrow \mathbf{f}(\mathbf{x}) \in B_\varepsilon(\mathbf{f}(\mathbf{p}))) \right) \Rightarrow \mathbf{f} \text{ is continuous at } \mathbf{p}$$

- If \mathbf{f} is continuous, then for any \mathbf{q} , $S = \{\mathbf{x} \in D | \mathbf{f}(\mathbf{x}) = \mathbf{q}\}$ is closed (or empty)
Consider a sequence of points, (\mathbf{S}_n) all in S that converges to $\hat{\mathbf{S}}$
 \mathbf{S}_n approaches $\hat{\mathbf{S}}$, so $\mathbf{f}(\mathbf{S}_n) = \mathbf{q}$ must get arbitrarily close to $\mathbf{f}(\hat{\mathbf{S}})$, which must be \mathbf{q} , since \mathbf{f} is continuous
Then the set S is sequentially closed and closed
- If $\mathbf{f} : D_{\mathbf{f}} \in \mathbb{R}^j \rightarrow \mathbb{R}^k$ and $\mathbf{g} : D_{\mathbf{g}} \in \mathbb{R}^k \rightarrow \mathbb{R}^l$ are both continuous, then $\mathbf{g} \circ \mathbf{f}$ is too
For any ε radius for $(\mathbf{g} \circ \mathbf{f})(\mathbf{p})$ to be in, there exists some $\delta(\varepsilon)$ that $\mathbf{f}(\mathbf{p})$ can be in, since \mathbf{g} is continuous
Similarly, there exists some $\eta(\delta(\varepsilon))$ that \mathbf{p} can be in, since \mathbf{f} is continuous
- Therefore, it is possible to build up continuous functions from simpler continuous functions
Many univariate functions previously studied are continuous on some interval
Any differentiable univariate function f is continuous, since for $\lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$ to exist, it must be true that $\lim_{h \rightarrow 0} (f(x+h) - f(x)) = 0$

1.15 Sept 27: Identifying continuous \mathbf{f}

- Two functions where the range of one is a subset of the domain of the other are composable functions
- Define the Cartesian Product of two functions, $\mathbf{f} : D \subseteq \mathbb{R}^d \rightarrow \mathbb{R}^e$, $\mathbf{g} : K \subseteq \mathbb{R}^k \rightarrow \mathbb{R}^l$ as $\mathbf{f} \times \mathbf{g} : D \times E \subseteq \{(a, b) | a \in \mathbb{R}^d, b \in \mathbb{R}^k\}$ (isomorphic to \mathbb{R}^{d+k}) $\rightarrow \{(a, b) | a \in \mathbb{R}^e, b \in \mathbb{R}^l\}$ (isomorphic to \mathbb{R}^{e+l})
- In short, $(\mathbf{f} \times \mathbf{g})(\mathbf{x}, \mathbf{y}) = (\mathbf{f}(\mathbf{x}), \mathbf{g}(\mathbf{y}))$
- This function is continuous if both \mathbf{f} and \mathbf{g} are continuous:
Consider that $|(\mathbf{x}, \mathbf{y}) - (\mathbf{p}, \mathbf{q})| = \sqrt{|\mathbf{x} - \mathbf{p}|^2 + |\mathbf{y} - \mathbf{q}|^2}$
If (\mathbf{x}, \mathbf{y}) approaches (\mathbf{p}, \mathbf{q}) , then $\mathbf{f}(\mathbf{x})$ and $\mathbf{g}(\mathbf{y})$ must approach $\mathbf{f}(\mathbf{p})$ and $\mathbf{g}(\mathbf{q})$
- Define another function, $(\mathbf{f} \wedge \mathbf{g})(\mathbf{x}) = (\mathbf{f}(\mathbf{x}), \mathbf{g}(\mathbf{x}))$, the Cartesian Product of the two functions on the diagonal map $\Delta(\mathbf{x}) = (\mathbf{x}, \mathbf{x})$
- Since the diagonal map is obviously continuous, and compositions of continuous functions are continuous, the concatenation of continuous functions is continuous
- Define yet another function, the n^{th} projection map, $\pi_n(\mathbf{x})$ is the n^{th} component of \mathbf{x}
- With these rather basic definitions, it is possible to prove a myriad of functions continuous:

$$\mathbf{f}(\mathbf{x}, \mathbf{y}, \mathbf{z}) = (\sin(ye^{x+z}), x^2 + y \cos(yz))$$

$$g(x, y, z) = \sin(ye^{x+z}) \quad g = \sin \circ \cdot \circ (\pi_2 \wedge (\exp \circ + \circ (\pi_1 \wedge \pi_3)))$$

$$h(x, y, z) = x^2 + y \cos(yz) \quad h = + \circ (\cdot \circ (\pi_1 \wedge \pi_1) \wedge \cdot \circ (\pi_2 \wedge \cos \circ \cdot \circ (\pi_2 \wedge \pi_3)))$$

$$\mathbf{f} = g \wedge h$$

It's a bit ridiculous to write out functions this way all the time, but since all the root functions are continuous, then \mathbf{f} must be as well

1.16 Oct 3: Continuous Rigor

- $+$ is continuous;

$$\begin{aligned} |x - a| < \frac{\varepsilon}{2} \wedge |y - b| < \frac{\varepsilon}{2} \quad \delta(\varepsilon) = \varepsilon \\ \Rightarrow |x - a| + |y - b| \leq |(x, y) - (a, b)| < \varepsilon \end{aligned}$$

- \cdot is more tricky, work in reverse:

$$\begin{aligned} & | \cdot (x, y) - \cdot (a, b) | \\ &= |x(y - b) + b(x - a)| \leq |x||y - b| + |b||x - a| \end{aligned}$$

Here, let $|x - a|$ be at most 1, so $|x| \leq |a| + 1$

$$\begin{aligned} & \leq (|a| + 1)|y - b| + |b||x - a| \\ & < (|a| + 1)\delta(\varepsilon) + |b|\delta(\varepsilon) \leq \varepsilon \quad \delta(\varepsilon) = \min \left\{ 1, \frac{\varepsilon}{2 \max\{|a| + 1, |b|\}} \right\} \end{aligned}$$

- The distance between two outputs of a contractive function is less than the distance between the inputs

1.17 Oct 4: More abstraction

- The cartesian product is considered to be associative, even though it isn't really
- More practice:

$$\begin{aligned} \mathbf{f}(\mathbf{x}, \mathbf{y}, \mathbf{z}) &= (xye^{z \sin(x)}, z \ln(x + y)) \\ g &= \cdot \circ ((\cdot \circ (\pi_1 \hat{\circ} \pi_2)) \wedge (\exp \circ \cdot \circ (\pi_3 \hat{\circ} (\sin \circ \pi_1)))) \\ h &= \cdot \circ (\pi_3 \hat{\circ} (\ln \circ + \circ (\pi_1 \hat{\circ} \pi_3))) \end{aligned}$$

1.18 Oct 7: EVT for More Functions

- The EVT for real-valued functions is stated just like the one for real-to-real-valued functions, except the input set considered is some compact $K \subseteq D$
- The EVT for vector-valued functions asserts:

$$\mathbf{f} : D \subseteq \mathbb{R}^d \rightarrow \mathbb{R}^e, \text{ cont. on compact } K \subseteq D \Rightarrow \text{compact } \mathbf{f}[K]$$

First, show $K' = \mathbf{f}[K]$ is bounded:

Assume K' is unbounded, then $\exists (\mathbf{y}_n) = (\mathbf{f}(\mathbf{x}_n)), \mathbf{x}_n \in K : (\forall n) |\mathbf{y}_n| \geq n$

K is bounded, so \exists conv. subsequence $\mathbf{x}_{n_j} \rightarrow \mathbf{x}_0 \in K$, a sequentially closed set.

$\mathbf{x}_{j_n} \rightarrow \mathbf{x}_0 \Rightarrow \mathbf{y}_{j_n} \rightarrow \mathbf{f}(\mathbf{x}_0)$, by continuity

However, \mathbf{y}_{j_n} is unbounded by assumption, while $\mathbf{f}(\mathbf{x}_0)$ is not infinite

- Next, show K' is sequentially closed:

Let the same sequences as above exist, except now (\mathbf{y}_{j_n}) is not unbounded and approaches \mathbf{y}_0

Necessarily, $\mathbf{y}_0 = \mathbf{f}(\mathbf{x}_0)$

Since $\mathbf{x}_0 \in K$, then $\mathbf{y}_0 \in K'$

1.19 Oct 8: IVT

- Connected sets:

$$(\forall \mathbf{a}, \mathbf{b} \in C) (\exists \text{ continuous } \gamma(t) : \gamma(0) = \mathbf{a} \wedge \gamma(1) = \mathbf{b}) \wedge (m \in (0, 1) \Rightarrow \gamma(m) \in C) \Rightarrow C \text{ is connected}$$

Or, there is a continuous path from any point in the set to any other point entirely in the set

In convex sets, this path must be a straight line

- The IVT for real-valued functions is stated just like the one for real-to-real-valued functions, except the input set considered is some connected $C \subseteq D$, and an intermediate output is obtained from some input in C less the specified inputs
- The IVT for vector-valued functions asserts:

$$\mathbf{f} : D \subseteq \mathbb{R}^d \rightarrow \mathbb{R}^e, \text{ cont. on connected } C \subseteq D \Rightarrow \text{connected } \mathbf{f}[C]$$

1.20 Oct 10: Connected Sets

- $[\gamma]$ is the trace of the path γ
- $\mathbf{pq}(t)$ or $[\mathbf{p}, \mathbf{q}]$ is the line segment between \mathbf{p} and \mathbf{q}
- \emptyset and \mathbb{R}^d are connected
- Connected sets $C_1 \cap C_2 \neq \emptyset \Rightarrow C_1 \cup C_2$ is connected
Let $\gamma_1 : [a, b] \rightarrow \mathbb{R}^d, \gamma_2 : [c, d] \rightarrow \mathbb{R}^d$ be paths that connect two arbitrary points in C_1, C_2 such that $\gamma_1(b) = \gamma_2(c)$
Define, with domain $[a, b+d-c]$, $\gamma_3(t) = \begin{cases} \gamma_1(t), & t \in [a, b] \\ \gamma_2(t-b+c), & t \in [b, b+d-c] \end{cases}$ connecting $\gamma_1(a) \in C_1$ to $\gamma_2(d) \in C_2$
- Arbitrary intersections of convex sets are convex
- The only connected subsets of \mathbb{R} are finite or infinite intervals, single points, and \emptyset

1.21 Oct 15: Connected Complements

- Complements of connected sets may not be connected; consider an annulus in 2-space.
- More generally, If a path from \mathbf{p} to \mathbf{q} has a supremum of the time that the path remains in the set, then the point at that time must be on the boundary of the set, or else it is possible to advance in time.

1.22 Oct 21: The Short Short Proof of IVT

- Let $\mathbf{f} : C \rightarrow \mathbf{f}[C]$ Prove that there must exist a continuous path δ between any $\mathbf{u}, \mathbf{v} \in \mathbf{f}[C]$ entirely in $\mathbf{f}[C]$:
 $\exists \mathbf{p}, \mathbf{q}, \gamma : \mathbf{f}(\mathbf{p}) = \mathbf{u} \wedge \mathbf{f}(\mathbf{q}) = \mathbf{v} \wedge \gamma$ is a continuous path between \mathbf{p}, \mathbf{q} entirely in C
 \mathbf{f} is continuous $\Rightarrow \delta(t) = \mathbf{f}(\gamma(t))$ is continuous and entirely in $\mathbf{f}[C]$
- Connected sets in \mathbb{R} are convex:
Let $S \subseteq \mathbb{R}$ be connected, \exists continuous $\gamma : [0, 1] \rightarrow S : \gamma(0) = a \wedge \gamma(1) = b$ and choose any $x : a < x < b$
 $T = \{t \in [0, 1] \mid \gamma(t) \leq x\}$, $0 \in T$ so T is not empty, and T is bounded above by 1, so $t_0 = \sup T$ exists
 $0 \leq t_0 < 1$, since $\gamma(1) = b > x$
 $\gamma(t_0) < x \Rightarrow \gamma(t_0 + \varepsilon) < x$ since γ is continuous, then t_0 is not an upper bound of T
 $\gamma(t_0) > x \Rightarrow t_0 \notin T$
The only possibility is $\gamma(t_0) = x$, so any x is on the path from a to b , (a, b)

1.23 Oct 22: Uniform Continuity

- Continuity of $\mathbf{f} : D \rightarrow \mathbb{R}^e$ on $S \subseteq D$:

$$(\forall \mathbf{p} \in S)(\forall \varepsilon > 0)(\exists \delta > 0)(\forall \mathbf{x} \in D)(|\mathbf{x} - \mathbf{p}| < \delta \Rightarrow |\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{p})| < \varepsilon)$$

- Uniform Continuity of \mathbf{f} on $S \subseteq D$:

$$(\forall \varepsilon > 0)(\exists \delta > 0)(\forall \mathbf{x}, \mathbf{p} \in S)(|\mathbf{x} - \mathbf{p}| < \delta \Rightarrow |\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{p})| < \varepsilon)$$

- Logical scope shifts allow this; a universal quantifier may be expanded in scope to form a consequent of the original statement.
- If S is relatively open to D , $S = D \cap \text{open } U$, \mathbf{f} uniformly continuous on $S \Rightarrow \mathbf{f}$ continuous on S
- The additional condition is necessary because $\mathbf{x} \in D$ will approach $\mathbf{p} \in S$, if a path of \mathbf{x} exists that never enters S , then no conclusion is possible, then the boundary of S inside D cannot be in S
- The Uniform Continuity Theorem (UCT)

$$\mathbf{f} \text{ continuous on compact } K \Rightarrow \mathbf{f} \text{ uniformly continuous on } K$$

1.24 Oct 23: Fire Drilled

- The condition for uniform continuity to imply continuity has been updated; some progress was made towards a proof

1.25 Oct 24: UCT

- Assume that the UCT was false on compact K , with \mathbf{f} continuous on K :

$$\neg(\forall \varepsilon > 0)(\exists \delta > 0)(\forall \mathbf{p}, \mathbf{x} \in K)(|\mathbf{x} - \mathbf{p}| < \delta \Rightarrow |\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{p})| < \varepsilon)$$

$$\Leftrightarrow (\exists \varepsilon_{\mathbf{f}} > 0)(\forall \delta > 0)(\exists \mathbf{p}, \mathbf{x} \in K)(|\mathbf{x} - \mathbf{p}| < \delta \wedge |\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{p})| \geq \varepsilon_{\mathbf{f}})$$

Let $\delta_n = \frac{1}{n}$ and define the sequence terms $\mathbf{p}_n, \mathbf{x}_n \in K : |\mathbf{p}_n - \mathbf{x}_n| < \delta_n \wedge |\mathbf{f}(\mathbf{x}_n) - \mathbf{f}(\mathbf{p}_n)| \geq \varepsilon_{\mathbf{f}}$

By BW, in the bounded set, $\exists(\mathbf{p}_{n_k}) \rightarrow \mathbf{p}_e$, and because the set is closed, $\mathbf{p}_e \in K$

\mathbf{f} is continuous, so $|p_{n_k} - x_{n_k}| < \delta_{n_k} \rightarrow 0 \Rightarrow |\mathbf{f}(\mathbf{p}_{n_k}) - \mathbf{f}(\mathbf{x}_{n_k})| \rightarrow 0 < \varepsilon_{\mathbf{f}}$; this is a contradiction

1.26 Oct 25: Halved; UCT?

- The UCT was almost proved

1.27 Oct 26: UCT!, and Differentiation

- The UCT was finally proved
- $f : D \subseteq \mathbb{R} \rightarrow \mathbb{R}, p \in \text{int} D \Rightarrow \exists r > 0 : B_r(p) \subseteq D$
- If extant, $f'(p) = \lim_{h \rightarrow 0} \frac{f(p+h) - f(p)}{h}$
- This is extensible to higher dimensions, but not directly: $\frac{\mathbf{f}(\mathbf{p} + \mathbf{h}) - \mathbf{f}(\mathbf{p})}{\mathbf{h}}$ is rather meaningless
- $\mathbf{f} : D \subseteq \mathbb{R}^d \rightarrow \mathbb{R}^e$ is differentiable at \mathbf{p} if $\exists! \mathbf{A} : \lim_{\mathbf{h} \rightarrow 0} \frac{|\mathbf{f}(\mathbf{p} + \mathbf{h}) - \mathbf{f}(\mathbf{p}) - \mathbf{A}\mathbf{h}|}{|\mathbf{h}|} = 0$, $\mathbf{f}'(\mathbf{p}) = \mathbf{A}$, an $e \times d$ matrix

1.28 Oct 29: Differentiability

- For real functions $\exists! \mathbf{A}$; assume \mathbf{A} and \mathbf{B} both worked, and let $\Delta f = f(p+h) - f(p)$:

$$\frac{\Delta f - \mathbf{A}h}{|h|} \rightarrow 0 \wedge \frac{\Delta f - \mathbf{B}h}{|h|} \rightarrow 0 \Rightarrow \frac{\mathbf{A}h - \mathbf{B}h}{|h|} \rightarrow 0 - 0 \Rightarrow \left| \frac{(\mathbf{A} - \mathbf{B})h}{|h|} \right| \rightarrow 0 \Rightarrow |\mathbf{A} - \mathbf{B}| \rightarrow 0 \Rightarrow \mathbf{A} - \mathbf{B} = 0$$

- The \mathbf{A} requirement for differentiability and its value agree with the conventional definitions for real functions:

$$\lim_{h \rightarrow 0} \frac{|f(p+h) - f(p) - \mathbf{A}h|}{|h|} = 0 \Leftrightarrow \lim_{h \rightarrow 0} \frac{|f(p+h) - f(p)|}{|h|} = \mathbf{A}$$

- The uniqueness of \mathbf{A} is proved in almost the same way for multivariate functions:

$$\frac{\Delta \mathbf{f} - \mathbf{A}\mathbf{h}}{|\mathbf{h}|} \rightarrow 0 \wedge \frac{\Delta \mathbf{f} - \mathbf{B}\mathbf{h}}{|\mathbf{h}|} \rightarrow 0 \Rightarrow \frac{\mathbf{A}\mathbf{h} - \mathbf{B}\mathbf{h}}{|\mathbf{h}|} \rightarrow 0 - 0 \Rightarrow \left| \frac{(\mathbf{A} - \mathbf{B})\mathbf{h}}{|\mathbf{h}|} \right| \rightarrow 0 \Rightarrow |\mathbf{A} - \mathbf{B}|\mathbf{i} \rightarrow 0 \Rightarrow \mathbf{A} - \mathbf{B} = 0$$

where \mathbf{i} represents an arbitrary unit vector. Considering all basis vectors means that every column of $\mathbf{A} - \mathbf{B}$ is the zero vector. Matrices $\mathbf{A}\mathbf{h}, \mathbf{B}\mathbf{h}$ must be factored according to the order of factors

1.29 Oct 30: Derivatives

- The components of \mathbf{A} can be found by:

$$[\mathbf{f}'(\mathbf{p})]_{i,j} = \mathbf{e}_i \cdot \mathbf{f}'(\mathbf{p})\mathbf{e}_j = \frac{\partial f_i}{\partial x_j}(\mathbf{p})$$

where \mathbf{e}_i is the i^{th} basis vector, f_i is the i^{th} component of \mathbf{f} , and x_j is the j^{th} component of \mathbf{x}

1.30 Oct 31: Partial Derivatives

- The existence of partial derivatives at a point does not imply differentiability at that point
The partial derivatives must be continuous in the neighborhood for that to be true
Consider a continuous real mapping of intersecting lines in the basis directions of a plane
The partial derivatives of the point where the curves intersect exists, but this is true nowhere else
- The partial derivative of $z = f(\mathbf{p})$, $\mathbf{p} = (x, y)$ with respect to x is the slope of the tangent line to \mathbf{p} on the graph that is the intersection of the xz plane containing \mathbf{p} and the surface of f

1.31 Nov 1: Tangent Planes

- The span of vectors along the tangent lines to $z = f(x, y)$ at $\mathbf{p} = (a, b)$ is a tangent plane:

$$z - f(\mathbf{p}) = \frac{\partial f}{\partial x}(x - a) + \frac{\partial f}{\partial y}(y - b)$$

This is the best possible first-order approximation plane to the graph of f at \mathbf{p}

- A zero-order approximation plane $l(x, y)$ is any one where:

$$\lim_{(x,y) \rightarrow \mathbf{p}} |f(\mathbf{p}) - l(x, y)| = 0$$

The sole condition is that intersection at (a, b, c) occurs

- The unique or nonexistent first-order approximation plane is the one where:

$$\lim_{(x,y) \rightarrow \mathbf{p}} \frac{|f(\mathbf{p}) - l(x, y)|}{|\mathbf{p} - (x, y)|} = 0$$

The matrix of l is exactly the first derivative of $f(\mathbf{p})$

1.32 Interlude I

- A function that is Lipschitz continuous on S is uniformly continuous on S
- A neat function \mathbf{f} on S is one where $S = \mathbf{f}^{-1}[\mathbf{f}[S]]$, that is, no point outside S is mapped to the same point as a point in S
- Let \mathbf{f} be continuous and neat on S . The inverse image of any set relatively open in $\mathbf{f}[S]$ is relatively open in S
- If S is relatively open in D , and the inverse image of any set relatively open in $\mathbf{f}[S]$ is relatively open in S , then \mathbf{f} is continuous on S
- String the two above together, \mathbf{f} is continuous on D iff the inverse image of every set relatively open in R is relatively open in D
- If $\mathbf{f}|_K$ is one-to-one, where $K \subseteq D$ is compact, and \mathbf{f} is continuous on K , then $\mathbf{f}|_K^{-1}$ is continuous
- S is compact iff for every set of sets $\{U_a | a \in \mathbb{N}\}$ covering S , then there exists a covering of S with finitely many sets, i.e. the set of sets $\{U_a | a \in \mathbb{N} \wedge a < M\}$
- Let f be continuous on closed C . If there exists a point \mathbf{a} and radius $r > 0$ such that $f(\mathbf{x}) \geq f(\mathbf{a})$ for all $\mathbf{x} \in C \setminus B_r(\mathbf{a})$, then there exists a point $\mathbf{c} \in C \cap B_r(\mathbf{a})$ such that f has a global minimum at \mathbf{c}
- The cartesian products of two sets that are open, closed, bounded, or compact are open, closed, bounded, or compact respectively
- The distance between any two disjoint nonempty sets (i.e. the shortest distance between two points, one in one set and one in the other), one closed and the other compact, is positive
- The union of two nonempty open sets that do not overlap is not connected
- The union of two nonempty closed sets that do not overlap is not connected

1.33 Nov 5: Chain Rule

- Let the functions $\mathbf{f} : D \subseteq \mathbb{R}^c \rightarrow \mathbb{R}^d$ and $\mathbf{g} : E \subseteq \mathbb{R}^d \rightarrow \mathbb{R}^e$ exist
- The composition, $(\mathbf{g} \circ \mathbf{f})$ has domain $\{\mathbf{p} \in D : \mathbf{f}(\mathbf{p}) \in E\}$
If this set is empty, the two functions are impossible
- If additionally, \mathbf{f} is differentiable at \mathbf{p} and \mathbf{f} is differentiable at $\mathbf{f}(\mathbf{p})$, then the derivative $(\mathbf{g} \circ \mathbf{f})'(\mathbf{p})$ is defined as

$$\mathbf{g}'(\mathbf{f}(\mathbf{p}))\mathbf{f}'(\mathbf{p})$$

The two matrices have dimensions $e \times d$ and $d \times c$, so their product has dimensions $e \times c$, consistent with the dimensions of the composition

1.34 Nov 6: Proof of the Above

- \mathbf{f} and \mathbf{g} are differentiable at \mathbf{p} and $\mathbf{f}(\mathbf{p}) = \mathbf{q}$, respectively, is equivalent to the following statements:

$$(\exists! \mathbf{A})(\forall \varepsilon_1 > 0)(\exists \delta_1 > 0)(\forall \mathbf{h})(|\mathbf{h}| < \delta_1 \Rightarrow |\mathbf{f}(\mathbf{p} + \mathbf{h}) - \mathbf{f}(\mathbf{p}) - \mathbf{A}\mathbf{h}| < \varepsilon_1 |\mathbf{h}|)$$

$$(\exists! \mathbf{B})(\forall \varepsilon_2 > 0)(\exists \delta_2 > 0)(\forall \mathbf{k})(|\mathbf{k}| < \delta_2 \Rightarrow |\mathbf{g}(\mathbf{q} + \mathbf{k}) - \mathbf{g}(\mathbf{q}) - \mathbf{B}\mathbf{k}| < \varepsilon_2 |\mathbf{k}|)$$

The goal is to prove:

$$(\exists! \mathbf{C} = \mathbf{B}\mathbf{A})(\forall \varepsilon_3 > 0)(\exists \delta_3 > 0)(\forall \mathbf{h})(|\mathbf{h}| < \delta_3 \Rightarrow Q = |(\mathbf{g} \circ \mathbf{f})(\mathbf{p} + \mathbf{h}) - (\mathbf{g} \circ \mathbf{f})(\mathbf{p}) - \mathbf{C}\mathbf{h}| < \varepsilon_3 |\mathbf{h}|)$$

$$Q = |\mathbf{g}(\mathbf{f}(\mathbf{p} + \mathbf{h})) - \mathbf{g}(\mathbf{f}(\mathbf{p})) - \mathbf{B}(\mathbf{A}\mathbf{h})|$$

Fix $\mathbf{k} = \mathbf{f}(\mathbf{p} + \mathbf{h}) - \mathbf{f}(\mathbf{p})$, this will approach 0 as \mathbf{h} approaches 0, since \mathbf{f} is continuous at \mathbf{P}

$$= |\mathbf{g}(\mathbf{q} + \mathbf{k}) - \mathbf{g}(\mathbf{q}) - \mathbf{B}(\mathbf{A}\mathbf{h})|$$

Use the triangle inequality,

$$\leq |\mathbf{g}(\mathbf{q} + \mathbf{k}) - \mathbf{g}(\mathbf{q}) - \mathbf{B}\mathbf{k}| + |\mathbf{B}\mathbf{k} - \mathbf{B}(\mathbf{A}\mathbf{h})|$$

The left term can be upper-bounded by differentiability if $|\mathbf{k}| < \delta_2$, and the right can be simplified by C-S,

$$< \varepsilon_2 |\mathbf{k}| + \|\mathbf{B}\| |\mathbf{k} - \mathbf{A}\mathbf{h}|$$

$|\mathbf{k} - \mathbf{A}\mathbf{h}|$ can be upper-bounded by differentiability if $|\mathbf{h}| < \delta_1$,

$$< \varepsilon_2 |\mathbf{k}| + \|\mathbf{B}\| \varepsilon_1 |\mathbf{h}|$$

Triangle inequality once more,

$$\leq \varepsilon_2 (|\mathbf{k} - \mathbf{A}\mathbf{h}| + |\mathbf{A}\mathbf{h}|) + \varepsilon_2 \|\mathbf{B}\| |\mathbf{h}|$$

The same simplification as two above, and factoring $|\mathbf{h}|$

$$< (\varepsilon_2 \varepsilon_1 + \|\mathbf{A}\| \varepsilon_2 + \|\mathbf{B}\| \varepsilon_2) |\mathbf{h}|$$

The left factor of the result is arbitrarily small, so Q can be made less than $\varepsilon_3 |\mathbf{h}|$

1.35 Nov 7: Matrix Algebra

- Matrix \mathbf{A} of size $e \times d$ can be represented as the whole:

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} & \dots & a_{1d} \\ a_{21} & a_{22} & a_{23} & \dots & a_{2d} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ a_{e1} & a_{e2} & a_{e3} & \dots & a_{ed} \end{bmatrix}$$

or as compact representation $[a_{ij}]_{1 \leq i \leq e, 1 \leq j \leq d}$, or as a population of entries $[\mathbf{A}]_{ij} = a_{ij}$

- Matrix operations include scaling, $[r\mathbf{A}]_{ij} = r[\mathbf{A}]_{ij}$

$$0\mathbf{A} = \mathbf{O} \quad 1\mathbf{A} = \mathbf{A} \quad r(s\mathbf{A}) = (rs)\mathbf{A}$$

where \mathbf{O} represents the zero matrix of appropriate size

- Addition is entry-wise with matrices of the same size, $[\mathbf{A} + \mathbf{B}]_{ij} = [\mathbf{A}]_{ij} + [\mathbf{B}]_{ij}$

$$(r+s)\mathbf{A} = r\mathbf{A} + s\mathbf{A} \quad r(\mathbf{B} + \mathbf{A}) = r\mathbf{B} + r\mathbf{A} \quad \mathbf{A} + \mathbf{B} = \mathbf{B} + \mathbf{A} \quad \mathbf{A} + (\mathbf{B} + \mathbf{C}) = (\mathbf{A} + \mathbf{B}) + \mathbf{C} \quad \mathbf{A} + (-\mathbf{A}) = \mathbf{O} \quad \mathbf{A} + \mathbf{O} = \mathbf{A}$$

where $-\mathbf{A}$ is the matrix $-1\mathbf{A}$

- Multiplication is slightly more complicated, the matrices must have sizes in the form $e \times d$ and $d \times f$
Make the mapping $\mathbf{A} \in \mathbb{R}^{e \times d} \mapsto \mathbf{l}_{\mathbf{A}} : \mathbb{R}^d \rightarrow \mathbb{R}^e$

$$\mathbf{l}_{\mathbf{A}}(\mathbf{x}) = x_1 \text{col}_1(\mathbf{A}) + x_2 \text{col}_2(\mathbf{A}) + x_3 \text{col}_3(\mathbf{A}) + \cdots + x_d \text{col}_d(\mathbf{A}) = \sum_{j=1}^d x_j \mathbf{a}_j$$

1.36 Nov 8: Multiplication of Matrices

- It turns out that $\mathbf{l}_{\mathbf{BA}}(\mathbf{x}) = \mathbf{l}_{\mathbf{B}}(\mathbf{l}_{\mathbf{A}}(\mathbf{x}))$

$$\mathbf{l}_{\mathbf{B}}(\mathbf{l}_{\mathbf{A}}(\mathbf{x})) = \sum_{k=1}^e \left(\sum_{j=1}^d x_j \mathbf{a}_j \right) \mathbf{b}_k = \sum_{j=1}^d x_j \left(\sum_{k=1}^e [\mathbf{A}]_{kj} \mathbf{b}_k \right) = \mathbf{l}_{\mathbf{C}}(\mathbf{x})$$

Simply let $\mathbf{C} = \mathbf{BA}$, and $\mathbf{c}_j = \sum_{k=1}^e [\mathbf{A}]_{kj} \mathbf{b}_k$

$\mathbf{l}_{\mathbf{A}} = \mathbf{l}_{\mathbf{B}} \Rightarrow \mathbf{A} = \mathbf{B}$ by using the unit basis vectors \mathbf{e}_j , which establishes the equivalence of each column of \mathbf{A} and \mathbf{B}

To actually find \mathbf{C} ,

$$[\mathbf{C}]_{dj} = \sum_{k=1}^e [\mathbf{A}]_{kj} [\mathbf{B}]_{dk}$$

- Using this linear map, other matrix operations can be defined as well:

$$\mathbf{l}_{r\mathbf{A}} = r\mathbf{l}_{\mathbf{A}} \quad \mathbf{l}_{\mathbf{A}+\mathbf{B}} = \mathbf{l}_{\mathbf{A}} + \mathbf{l}_{\mathbf{B}} \quad \mathbf{l}_{\mathbf{BA}} = \mathbf{l}_{\mathbf{B}} \circ \mathbf{l}_{\mathbf{A}}$$

- \mathbf{l} is additive (distributive over $+$) and homogenous (distributive over scaling) so it is linear
- \mathbf{l} follows matrix distribution both ways ($\mathbf{l}_{\mathbf{C}(\mathbf{B}+\mathbf{A})} = \mathbf{l}_{\mathbf{CB}+\mathbf{CA}}$ and the reverse) and mixed scaling ($\mathbf{l}_{(r\mathbf{B})\mathbf{A}} = \mathbf{l}_{r(\mathbf{BA})} = \mathbf{l}_{\mathbf{B}(r\mathbf{A})}$)
- $(\mathbf{BA})^T = \mathbf{A}^T \mathbf{B}^T$, where T is the transpose

1.37 Nov 12: Utility of Linear Maps

- All linear maps are linear functions:

$$\mathbf{l}_{\mathbf{A}}(\mathbf{x}) = \sum_{k=1}^e x_k \mathbf{l}_{\mathbf{A}}(\mathbf{e}_k)$$

From the formula, this fact is obvious; simply realize that every term is first-degree in x_k

- All linear functions are linear maps:

Let an arbitrary $\mathbf{f}(\mathbf{x})$ be linear,

$$\begin{aligned} \mathbf{f}(\mathbf{x}) &= \mathbf{f}(x_1 \mathbf{e}_1 + x_2 \mathbf{e}_2 \dots) \\ &= \mathbf{f}(x_1 \mathbf{e}_1) + \mathbf{f}(x_2 \mathbf{e}_2) \dots \\ &= x_1 \mathbf{f}(\mathbf{e}_1) + x_2 \mathbf{f}(\mathbf{e}_2) \dots \end{aligned}$$

This is exactly the form of a linear map, just with the columns of \mathbf{A} equal to $\mathbf{f}(\mathbf{e}_j)$, and with the fact that $\mathbf{l}_{\mathbf{A}} = \mathbf{l}_{\mathbf{B}} \Rightarrow \mathbf{A} = \mathbf{B}$

- The equivalency of matrices and linear maps for them makes proving properties of matrices very easy:
For right-distributivity:

$$l_{(\mathbf{B}+\mathbf{C})\mathbf{A}} = l_{\mathbf{B}+\mathbf{C}} \circ l_{\mathbf{A}} = (l_{\mathbf{B}} + l_{\mathbf{C}}) \circ l_{\mathbf{A}} = l_{\mathbf{B}} \circ l_{\mathbf{A}} + l_{\mathbf{C}} \circ l_{\mathbf{A}} = l_{\mathbf{BA}} + l_{\mathbf{CA}} = l_{\mathbf{BA}+\mathbf{CA}}$$

Left-distributivity is basically the same proof

For multiplicative associativity:

$$l_{(\mathbf{CB})\mathbf{A}} = l_{\mathbf{CB}} \circ l_{\mathbf{A}} = l_{\mathbf{C}} \circ l_{\mathbf{B}} \circ l_{\mathbf{A}} = l_{\mathbf{C}} \circ l_{\mathbf{BA}} = l_{\mathbf{C}(\mathbf{BA})}$$

1.38 Nov 13: Transposition and Inversion

- Transposes are defined:

$$[\mathbf{A}^T]_{ij} = [\mathbf{A}]_{ji}$$

The transpose will have flipped size

- Transposition is linear, its own inverse, and has a special multiplicative law; the invertibility is a bit obvious:

$$[(r\mathbf{A})^T]_{ij} = [r\mathbf{A}]_{ji} = r[\mathbf{A}]_{ji} = r[\mathbf{A}^T]_{ij}$$

$$[(\mathbf{A} + \mathbf{B})^T]_{ij} = [\mathbf{A} + \mathbf{B}]_{ji} = [\mathbf{A}]_{ji} + [\mathbf{B}]_{ji} = [\mathbf{A}^T]_{ij} + [\mathbf{B}^T]_{ij} = [\mathbf{A}^T + \mathbf{B}^T]_{ij}$$

$$[(\mathbf{BA})^T]_{ij} = [\mathbf{BA}]_{ji} = \sum_d [\mathbf{B}]_{jd} [\mathbf{A}]_{di} = \sum_d [\mathbf{A}^T]_{id} [\mathbf{B}^T]_{dj} = [\mathbf{A}^T \mathbf{B}^T]_{ij}$$

- The identity matrix is simply all 1 on the main diagonal, or $[\mathbf{I}]_{ij} = \delta_{ij}$, the Kronecker delta, logically equal to $(i = j)$

$$\mathbf{AI} = \mathbf{A} \quad \mathbf{IA} = \mathbf{A}$$

$$[\mathbf{AI}]_{ij} = \sum_d a_{id} \delta_{dj} = a_{ij} \delta_{jj} = [\mathbf{A}]_{ij}$$

$$[\mathbf{IA}]_{ij} \dots$$

The appropriate \mathbf{I} must be used. The delta function sifts through the sum to return one term. The continuous analog is the Dirac delta, whose real version has a definite integral across the real line 1 and value 0 everywhere except $x = 0$

- \mathbf{A} is left/right invertible if $\exists \mathbf{L}/\mathbf{R} : (\mathbf{LA})/(\mathbf{AR}) = \mathbf{I}$ There may be one left/right inverse, or many, or none. \mathbf{O} has no left or right inverse. If a matrix has both a left and right inverse, they are unique and equal, and is denoted \mathbf{A}^{-1}

$$\mathbf{L} = \mathbf{L}\mathbf{I} = \mathbf{L}(\mathbf{AR}) = (\mathbf{LA})\mathbf{R} = \mathbf{I}\mathbf{R} = \mathbf{R}$$

$$\mathbf{L}' = \mathbf{L}'(\mathbf{AR}) = (\mathbf{L}'\mathbf{A})\mathbf{R} = \mathbf{R}$$

$$\mathbf{L} = \mathbf{L}(\mathbf{AR}) = (\mathbf{LA})\mathbf{R}' = \mathbf{R}'$$

1.39 Nov 14: Invertibility

- It is true that $\exists \mathbf{A}^{-1} \Rightarrow \mathbf{A} \in \mathbb{R}^{d \times d}$
- To prove this, define the notion of rank, the dimension of the span of the column vectors
- The rank of a matrix must lie in $[0, \min\{e, d\}]$; a matrix with every column linearly independent has rank $\min\{e, d\}$, since there are d vectors in \mathbb{R}^e
- The only matrices with rank 0 are \mathbf{O} . \mathbf{I}_n has rank n
- An essential fact is that $\text{rank } AB \leq \text{rank } A$

1.40 Nov 15*: Squareness

- Invertible matrices are square, and the proof is simple with the above properties assumed. Take $\mathbf{A} \in \mathbb{R}^{e \times d}$, $\mathbf{BA} = \mathbf{I}$, $\mathbf{AB} = \mathbf{I}$

$$e = \text{rank } \mathbf{I}_e = \text{rank } \mathbf{AB} \leq \text{rank } \mathbf{A} \leq d$$

$$d = \text{rank } \mathbf{I}_d = \text{rank } \mathbf{BA} \leq \text{rank } \mathbf{B} \leq e$$

1.41 Nov 18: Dimension

- The Steinitz Theorem:
Take two sets of vectors $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\}$ and $\{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_m\}$ in \mathbb{R}^d , such that their spans are equal, V . Then, $k = m \leq d$ and the dimension of V is k . k is also the number of vectors in the basis, or spanning set, of V .
- The proof of this involves a very fun algebraic fact; an underdetermined homogeneous system of linear equations has at least one nontrivial solution. Suppose that a system of n equations and $n + 1$ variables exists, so it is underdetermined by one equation. If $n = 1$, the system clearly has nontrivial solutions:

$$a_1x_1 + a_2x_2 = 0$$

Then use induction! Assuming any system with n equations and $n + 1$ variables has a nontrivial solution, prove that a system with $n + 1$ equations and $n + 2$ variables has one too:

$$a_{1,1}x_1 + a_{1,2}x_2 + \dots + a_{1,n+1}x_{n+1} + a_{1,n+2}x_{n+2} = 0$$

Solve for x_{n+2} in terms of the other variables. Then, remove this equation, and the system is now of n equations and $n + 1$ variables. If its coefficient is 0, pick a different x_n . If all of the coefficients are 0, this system is underdetermined by more than one equation.

Any system underdetermined by more than one equation will have a nontrivial solution from the same system with less variables and augmented 0's in every solution.

- Now, consider the sets of vectors. $\mathbf{w}_1 \in V$, so

$$\mathbf{w}_1 = a_{11}\mathbf{v}_1 + a_{21}\mathbf{v}_2 \dots a_{k1}\mathbf{v}_k$$

$$\mathbf{w}_2 = a_{12}\mathbf{v}_1 + a_{22}\mathbf{v}_2 \dots a_{k2}\mathbf{v}_k$$

$$\vdots$$

$$\mathbf{w}_m = a_{1m}\mathbf{v}_1 + a_{2m}\mathbf{v}_2 \dots a_{km}\mathbf{v}_k$$

Let \mathbf{A} be the matrix with the coefficients of this system, exactly as written, and \mathbf{B} the matrix with columns \mathbf{v}_n . Then $\mathbf{w}_n = \mathbf{B}(\text{col}_n \mathbf{A}) = \mathbf{B}\mathbf{a}_n$. Moreover, the matrix with columns \mathbf{w}_n , \mathbf{C} , is equivalent to $\mathbf{B}\mathbf{A}$.
To be continued in a while...

1.42 Nov 19: Subspaces and Spans

- A subspace V of \mathbb{R}^d is nonempty and is closed under addition and scaling, notated $V \subseteq \mathbb{R}^d$. $\dim V \leq d$
- $W \subseteq \mathbb{R}^d \wedge V \subseteq \mathbb{R}^d \Rightarrow (W \subseteq V \Rightarrow W \subseteq V)$
- A central theorem: $V \subseteq \mathbb{R}^d \Rightarrow \exists \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\} \subseteq V$. The set will be linearly independent and V is its span. Since 0 must be an element of V , either $V = \{0\}$ or $V \neq \{0\}$. If the first case is true, V is the span of \emptyset , which is linearly independent.
If the second is true, then $\mathbf{v}_1 \in V$. Either V is the span of the set of this one vector, which is linearly independent, or it is not.
If it is not, then $(\mathbf{v}_2 \neq r\mathbf{v}_1) \in V$. Either V is the span of the set of two vectors, which is linearly independent, or it is not.
This continues until the set grows to size $d + 1$, at which point it can no longer be linearly independent, from the below theorem, and V must be the span of the set, which can only be \mathbb{R}^d
- No set of more than d vectors in \mathbb{R}^d is linearly independent. Assume such a set $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$ existed.
- A related lemma solves this: $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$ is linearly independent iff $A\mathbf{x}$ has only the trivial solution, A is the matrix with columns \mathbf{v}_n .
Since the assumed system will have d equations and k variables, where it is specified that $k > d$, this system must have a nontrivial solution. Then, the set cannot be linearly independent.

1.43 Nov 20: Steinitz Resumed

- Recall that $\mathbf{C} = \mathbf{B}\mathbf{A}$. Then $\mathbf{C}\mathbf{x} = \mathbf{B}(\mathbf{A}\mathbf{x})$
- Set this equal to 0 to create a homogeneous system. Since \mathbf{C} is made of linearly independent vectors, $\mathbf{C}\mathbf{x} = 0$ has only the trivial solution. This means that $\mathbf{A}\mathbf{x} = 0$, a system with k variables and m equations must also have only the trivial solution, i.e. $k \leq m$
- By switching the roles of \mathbf{B} and \mathbf{C} throughout the proof, a system with m variables and k equations is formed which has only the trivial solution, i.e. $m \leq k$
- The only possibility is that $k = m$

1.44 Nov 21: Subspace Combinators

- Let $V, W \subseteq \mathbb{R}^d$, then $V \cap W \subseteq \mathbb{R}^d$ and $V + W \subseteq \mathbb{R}^d$, where space addition is defined as the set of all sums of a vector in one subspace and another in the other
- $V + W$ is also the smallest subspace containing $v \cup W$
- Reminiscent of the same counting paradigm found elsewhere, $\dim(V + W) = \dim V + \dim W - \dim(V \cap W)$. This can easily be extended by induction.

1.45 Nov 22: Complementary Subspaces

- V^\perp is the orthogonal complement to V , with every vector in V^\perp perpendicular to every vector in V
- W , a complement to V intersects it only at 0, and has $V + W = \mathbb{R}^d$, which together can be notated as $V \oplus W = \mathbb{R}^d$
- Inner products of two vectors must have certain special properties:

$$\begin{aligned}\langle \mathbf{x} | \mathbf{y} \rangle &\in \mathbb{R} \\ \langle \mathbf{x} | \mathbf{x} \rangle &\geq 0, \mathbf{x} = 0 \Rightarrow \langle \mathbf{x} | \mathbf{x} \rangle = 0 \\ \langle \mathbf{x} | \mathbf{y} \rangle &= \langle \mathbf{y} | \mathbf{x} \rangle \\ \langle c\mathbf{x} | \mathbf{y} \rangle &= c\langle \mathbf{x} | \mathbf{y} \rangle = \langle \mathbf{x} | c\mathbf{y} \rangle \\ \langle \mathbf{x} + \mathbf{y} | \mathbf{z} \rangle &= \langle \mathbf{x} | \mathbf{z} \rangle + \langle \mathbf{y} | \mathbf{z} \rangle\end{aligned}$$

Each inner product has a norm associated with it

$$|\mathbf{x}| = \sqrt{\langle \mathbf{x} | \mathbf{x} \rangle}$$

The simplified Cauchy-Schwarz inequality states:

$$|\langle \mathbf{x} | \mathbf{y} \rangle| \leq |\mathbf{x}| |\mathbf{y}|$$

- The special dot product is defined by $\mathbf{x} \cdot \mathbf{y} = \sum_j x_j y_j$, and its associated norm is the Euclidean one

1.46 Nov 25: Null Spaces

- The column space of \mathbf{A} , $\text{Col}\mathbf{A}$ is the span of the columns of \mathbf{A} , $\{\mathbf{A}\mathbf{x} | \mathbf{x} \in \mathbb{R}^d\}$, and is obviously a subspace of \mathbb{R}^e
- The row space, $\text{Row}\mathbf{A}$ is the same, but with the rows
- Essentially, $\dim \text{Row}\mathbf{A} = \dim \text{Col}\mathbf{A} = \text{rank}\mathbf{A}$
- The null space. $\text{Null}\mathbf{A}$, is $\{\mathbf{x} | \mathbf{A}\mathbf{x} = 0\}$, so it is the solutions to the homogeneous system $\mathbf{A}\mathbf{x} = 0$ with e equations and d variables. It is also a subspace due to closure under addition and scaling, with dimension $\text{null}\mathbf{A}$
- The count of free and determined variables also means that $\text{rank}\mathbf{A} + \text{null}\mathbf{A} = d$, so each dimension lost in a transformation must be lost to the origin

1.47 Nov 26: Ranking

- The rank of \mathbf{A} is equal to the number of leading 1's (or bound variables) in $\text{rref}(\mathbf{A})$
- The reduced row echelon form of \mathbf{A} is obtained purely by linearity-preserving matrix operations, so the rank of the matrix is not changed by the operation
- It is proved that $\text{rank} \mathbf{AB} \leq \text{rank} \mathbf{A}$:

\mathbf{AB} is the matrix with columns \mathbf{Ab}_k

These are all in the form $\{\mathbf{Ax} \mid x \in \mathbb{R}^e\}$, the column space of \mathbf{A}
Then all columns of \mathbf{AB} are in the column space of \mathbf{A}

$$\text{col}(\mathbf{AB}) \subseteq \text{col}(\mathbf{A})$$

$$\text{rank} \mathbf{AB} \leq \text{rank} \mathbf{A}$$

1.48 Nov 27*: Full Rank

- An invertible square matrix is of full rank

$$d = \text{rank} \mathbf{I} = \text{rank} \mathbf{AB} \leq \text{rank} \mathbf{A} \leq d$$

$$d \leq \text{rank} \mathbf{A} \leq d$$

The conclusion is apparent

1.49 Dec 2: Determinants

- The determinant is the only function $f : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}$ that is:

– Alternating:

$$f \left(\left[\begin{array}{c|c|c} \cdots & \mathbf{x} & \cdots \\ \hline & & \end{array} \right] \right) = -f \left(\left[\begin{array}{c|c|c} \cdots & \mathbf{y} & \cdots \\ \hline & & \end{array} \right] \right)$$

– Multilinear:

$$f \left(\left[\begin{array}{c|c|c} \cdots & a\mathbf{x} + b\mathbf{y} & \cdots \\ \hline & & \end{array} \right] \right) = a \left(\left[\begin{array}{c|c|c} \cdots & \mathbf{x} & \cdots \\ \hline & & \end{array} \right] \right) + b \left(\left[\begin{array}{c|c|c} \cdots & \mathbf{y} & \cdots \\ \hline & & \end{array} \right] \right)$$

– Normalized:

$$f(\mathbf{I}) = 1$$

- Uniqueness is easier to prove than existence:

$$f(A) = f \left(\left[\begin{array}{c|c|c} \mathbf{a}_1 & \mathbf{a}_2 & \cdots & \mathbf{a}_d \\ \hline & & & \end{array} \right] \right)$$

A given: $\mathbf{a} = \sum a_1 \mathbf{e}_1$

$$= f \left(\left[\begin{array}{c|c|c} \sum_{j_1=1}^d a_{j_1 1} \mathbf{e}_{j_1} & \sum_{j_2=1}^d a_{j_2 2} \mathbf{e}_{j_2} & \cdots & \sum_{j_d=1}^d a_{j_d d} \mathbf{e}_{j_d} \\ \hline & & & \end{array} \right] \right)$$

Use the multilinear property

$$= \sum_{j_1=1}^d \sum_{j_2=1}^d \cdots \sum_{j_d=1}^d a_{j_1 1} a_{j_2 2} \cdots a_{j_d d} f \left(\left[\begin{array}{c|c|c} \mathbf{e}_{j_1} & \mathbf{e}_{j_2} & \cdots & \mathbf{e}_{j_d} \\ \hline & & & \end{array} \right] \right)$$

Sift the sum, when any of the \mathbf{e}_{j_k} are the same, the overall term must be zero, because swapping the identical columns of the matrix cannot change the value of the determinant

$$= \sum_{(j_1, j_2, \dots, j_d) \in S_d} (-1)^{N(j_1, j_2, \dots, j_d)} a_{j_1 1} a_{j_2 2} \cdots a_{j_d d} f(I)$$

S_d is the symmetric group on d objects, or every permutation of $(1, 2, \dots, d)$, and $N(j_1, j_2, \dots, j_d)$ is the number of inversions, or occurrences of $(j_a, j_b) : a < b \wedge j_a > j_b$

The definition of $f(\mathbf{A})$ does not depend on f at all, so this function must be unique

- Existence: f is multilinear, from its definition, since multiplications and summations are linear to each input vector
- Normality is easier too, when the permutation is the natural one, $N = 0$ and only $\prod_{k=1}^d a_{kk} = 1$ while any other elementary product is 0, so the sum is 1

1.50 Dec 3: Determinant Alternativity

- For simplicity, let $(j_1, j_2, \dots, j_d)^* = (j_1^*, j_2^*, \dots, j_d^*)$ be the swap of elements k and l of the permutation. The star function is a bijection on the permutations
- To be proven: $N(j_1^*, j_2^*, \dots, j_d^*) = N(j_1, j_2, \dots, j_d) + 2n + 1, n \in \mathbb{Z}$
- Then, $\det \mathbf{A}_{\text{swap}}$ will have every term negated, while the elementary products will cover every selection of entries in the matrix, so swapping the order of a_{jk} terms has no effect on the magnitude of each term. Alternativity is proven

1.51 Dec 4: Cauchy Schwarz and Triangles

- Let $f(t) = |\mathbf{w} - t\mathbf{v}|^2$
- The case where $\mathbf{z} = 0$ can be handled trivially. Otherwise,

$$\begin{aligned} 0 &\leq f(t) = (\mathbf{w} - t\mathbf{v}) \cdot (\mathbf{w} - t\mathbf{v}) \\ &= (\mathbf{w} \cdot \mathbf{w})^2 - 2(\mathbf{w} \cdot \mathbf{v})t + (\mathbf{v} \cdot \mathbf{v})t^2 \\ &= t^2|\mathbf{v}|^2 - 2t(\mathbf{w} \cdot \mathbf{v}) + |\mathbf{w}|^2 \end{aligned}$$

This quadratic in t is always nonnegative and has positive leading coefficient, so $\Delta \leq 0$

$$\begin{aligned} 4(\mathbf{w} \cdot \mathbf{v})^2 - 4|\mathbf{v}|^2|\mathbf{w}|^2 &\leq 0 \\ (\mathbf{w} \cdot \mathbf{v})^2 &\leq |\mathbf{w}|^2|\mathbf{v}|^2 \\ |\mathbf{w} \cdot \mathbf{v}| &\leq |\mathbf{w}||\mathbf{v}| \end{aligned}$$

- Equality is held under certain conditions: Either \mathbf{w} is codirectional to \mathbf{v} and $(\mathbf{w} \cdot \mathbf{v}) = |\mathbf{w}||\mathbf{v}|$, or they are antidualirectional, and $(\mathbf{w} \cdot \mathbf{v}) = -|\mathbf{w}||\mathbf{v}|$
- Defining a norm for matrices is a bit clunky, simply take the norm of the matrix as if it was a big vector
- Matrices satisfy these, assuming the matrices are of appropriate size:

$$\begin{aligned} \|r\mathbf{A}\| &= r\|\mathbf{A}\| \\ \|\mathbf{A}^T\| &= \|\mathbf{A}\| \\ \|\mathbf{AB}\| &\leq \|\mathbf{A}\|\|\mathbf{B}\| \\ \|\mathbf{A} + \mathbf{B}\| &\leq \|\mathbf{A}\| + \|\mathbf{B}\| \end{aligned}$$

The third property is essentially CS again

The fourth is an application of the triangle inequality, true for many vectors for every norm associated with an inner product and possibly so for those that are not, and whose proof is below:

$$\begin{aligned} |\mathbf{x} + \mathbf{y}|^2 &= (\mathbf{x} + \mathbf{y}) \cdot (\mathbf{x} + \mathbf{y}) = |\mathbf{x}|^2 + |\mathbf{y}|^2 + 2(\mathbf{x} \cdot \mathbf{y}) \\ &\leq |\mathbf{x}|^2 + |\mathbf{y}|^2 + 2|\mathbf{x} \cdot \mathbf{y}| \\ &\leq |\mathbf{x}|^2 + |\mathbf{y}|^2 + 2|\mathbf{x}||\mathbf{y}| = |\mathbf{x}| + |\mathbf{y}|^2 \end{aligned}$$

Equality is held if all vectors are codirectional

1.52 Dec 5: Matrix Breakdown

- CS for matrices, which works because an entry in the product of two matrices is basically a dot product:

$$\begin{aligned}
 \|\mathbf{AB}\|^2 &= \sum_i \sum_j [\mathbf{AB}]_{ij}^2 = \sum_i \sum_j \left(\sum_k (a_{ik} b_{kj}) \right)^2 \\
 &= \sum_i \sum_j (\text{row}_i \mathbf{a} \cdot \mathbf{b}_j) \\
 &\leq \sum_i \sum_j |\text{row}_i \mathbf{a}|^2 |\mathbf{b}_j|^2 \\
 &= \sum_i |\text{row}_i \mathbf{a}|^2 \sum_j |\mathbf{b}_j|^2 \\
 &= \|\mathbf{A}\|^2 + \|\mathbf{B}\|^2
 \end{aligned}$$

Equality holds when every single row of \mathbf{A} is parallel to every single column of \mathbf{B} , not the most exciting case

- There is a more versatile matrix norm, the operator norm, defined as:

$$\|\mathbf{A}\| = \sup\{|\mathbf{Ax}| \mid |\mathbf{x}| = 1\}$$

Simply put, how far does the matrix transform the unit sphere?

A nice property: $\|\mathbf{I}\| = 1$; this was not true for the Hilbert-Schmidt norm

1.53 Dec 6: Chain Rule Examined

- The chain rule can already be applied to functions, and the result is a derivative matrix.
- The individual entries of the derivative matrix are sometimes more useful by themselves:

$$[\mathbf{h}']_{ij} = \frac{\partial h_i}{\partial x_j} = \sum_{k=1}^e \frac{\partial g_i}{\partial y_k}(\mathbf{f}(\mathbf{p})) \cdot \frac{\partial f_k}{\partial x_j}(\mathbf{p})$$

e is the intermediate dimension

- The chain rule with multiple variables and multiple dependencies makes a distinction between the total and partial derivative, which is exactly as observed

Suppose f depends on x, y, t and x, y depend on t , then simply call the higher t by an identical u :

$$\frac{\partial f}{\partial u} = \frac{\partial x}{\partial t} \frac{\partial y}{\partial t}$$

1.54 Dec 9: Gradient and MVT

- The MVT for vector-to-scalar functions is exactly the same as it is for scalar-to-scalar functions:

$$f : \mathbb{R}^d \rightarrow \mathbb{R}$$

$$g(t) := f((1-t)\mathbf{p} + t\mathbf{q})$$

$$g(1) - g(0) = f(\mathbf{q}) - f(\mathbf{p})$$

Let the domain of g be $(-\varepsilon, 1 + \varepsilon)$, so it is differentiable on $[0, 1]$

$$g'(t) = f'((1-t)\mathbf{p} + t\mathbf{q})(\mathbf{q} - \mathbf{p})$$

Because g is real-to-real, the MVT applies, so

$$\exists c : g'(c) = \frac{g(1) - g(0)}{1 - 0} = f(\mathbf{q}) - f(\mathbf{p})$$

$$g'(c) = f'((1-c)\mathbf{p} + c\mathbf{q})(\mathbf{q} - \mathbf{p})$$

Therefore,

$$f(\mathbf{q}) - f(\mathbf{p}) = f'((1-c)\mathbf{p} + c\mathbf{q})(\mathbf{q} - \mathbf{p})$$

$f'((1-c)\mathbf{p} + c\mathbf{q})$ is now notated as $f'(r)$, where r is obviously on the line segment

- $f'(r)^T$ is the vector defined as $\nabla f(\mathbf{r})$, so $\nabla f(\mathbf{r})_j = \frac{\partial f}{\partial x_j}(\mathbf{r})$

1.55 Dec 10: Gradient Properties

- The MVT can be expressed in terms of the gradient, $f(\mathbf{q}) - f(\mathbf{p}) = \nabla f(\mathbf{r}) \cdot (\mathbf{q} - \mathbf{p})$
- Naturally then, the all-powerful Cauchy-Schwarz inequality is applicable:

$$|f(\mathbf{q}) - f(\mathbf{p})| \leq |\nabla f(\mathbf{r})| |\mathbf{q} - \mathbf{p}|$$

- If f is additionally continuously differentiable, the gradient will have continuous components on the compact set, so the EVT applies, and $|\nabla f(\mathbf{r})| \leq M$, so a C^1 function from \mathbb{R}^d to \mathbb{R} is necessarily Lipschitz continuous on the line segment
- If $f : D \supseteq U \rightarrow \mathbb{R}$ has gradient 0 at every point in region U , then f must be constant

1.56 Dec 11: MVI

- The MVT does not extend to vector-valued functions:
- Consider $\mathbf{f} : \mathbb{R}^d \rightarrow \mathbb{R}^2$, where $\mathbf{f} = g \wedge h$
- A simple lemma: \mathbf{f} diff $\Rightarrow f_1, f_2, \dots$ diff, since the error term of each component must be smaller or equal to the overall error term. Then, assuming \mathbf{f} is differentiable, g and h are too
- By the older MVT:

$$g(\mathbf{q}) - g(\mathbf{p}) = g'(\mathbf{r})(\mathbf{q} - \mathbf{p})$$

$$h(\mathbf{q}) - h(\mathbf{p}) = h'(\mathbf{s})(\mathbf{q} - \mathbf{p})$$

\mathbf{r} and \mathbf{s} , while on the same line segment, are not held to be the same point, which is the strongest argument that can be made with this theorem

- A new MVT would assert:

$$\mathbf{f}(\mathbf{q}) - \mathbf{f}(\mathbf{p}) = \mathbf{f}'(\mathbf{u})(\mathbf{q} - \mathbf{p})$$

Therefore,

$$g(\mathbf{q}) - g(\mathbf{p}) = g'(\mathbf{u})(\mathbf{q} - \mathbf{p})$$

$$h(\mathbf{q}) - h(\mathbf{p}) = h'(\mathbf{u})(\mathbf{q} - \mathbf{p})$$

However, this assertion is too strong, since g and h are completely independent, so it is unreasonable to assume that the same point u will work for both

- The MVI is a weakening:

$$|\mathbf{f}(\mathbf{q} - \mathbf{p})| = |\mathbf{f}'(\mathbf{r})(\mathbf{q} - \mathbf{p})| \leq |\mathbf{f}'(\mathbf{r})| |\mathbf{q} - \mathbf{p}| \leq \max_{\mathbf{u} \in [\mathbf{p}, \mathbf{q}]} |\mathbf{f}'(\mathbf{u})| |\mathbf{q} - \mathbf{p}|$$

Now, the right side no longer references one specific point. A more practical way to use the MVI is to upper bound the maximum norm with an easy value, which must exist because \mathbf{f}' is continuous and $[\mathbf{p}, \mathbf{q}]$ is compact

$$|\mathbf{f}(\mathbf{q} - \mathbf{p})| \leq M |\mathbf{q} - \mathbf{p}|$$

This theorem can be extended to any convex compact set K

1.57 Dec 12: MVT too

- The MVT was reviewed

1.58 Dec 13: Gradients too

- Fun with gradients;
- The only time $\text{size} f(\mathbf{p}) = \text{size} f'(\mathbf{p})$ is when $\dim f(\mathbf{p}) = 1$

1.59 Dec 16: MVI too

- The MVI was reviewed

1.60 Dec 17: MVI fwee

- A "short" proof of MVI on K by continuous induction:

Fix $\varepsilon > 0$, which can be pushed to 0 later

$$\mathbf{q}_t = (1-t)\mathbf{p} + t\mathbf{q} \quad T = \{t \in [0, 1] \mid |\mathbf{f}(\mathbf{q}_t) - \mathbf{f}(\mathbf{p})| \leq (M + \varepsilon)|\mathbf{q}_t - \mathbf{p}|\}$$

$$\exists \sup T = t_0 \in [0, 1]$$

$t_0 \in T$ is the next step. Take some sequence of $t_n \in T$ increasing to t_0

$$\forall n |\mathbf{f}(\mathbf{q}_{t_n}) - \mathbf{f}(\mathbf{p})| \leq (M + \varepsilon)|\mathbf{q}_{t_n} - \mathbf{p}|$$

Limits and continuity ensure the truth of this at $n \rightarrow \infty$

Now, to prove $t_0 = 1$, by assuming it is not, then $\exists \delta : t_0 + \delta \leq 1$

If it is the case that $|\mathbf{f}(\mathbf{q}_{t_0+\delta}) - \mathbf{f}(\mathbf{p})| \leq (M + \varepsilon)|\mathbf{q}_{t_0+\delta} - \mathbf{p}|$, then $t_0 + \delta \in T$, so $\sup T \neq t_0$

$$\begin{aligned} & |\mathbf{f}(\mathbf{q}_{t_0+\delta}) - \mathbf{f}(\mathbf{p})| \\ & \leq |\mathbf{f}(\mathbf{q}_{t_0+\delta}) - \mathbf{f}(\mathbf{q}_{t_0})| + |\mathbf{f}(\mathbf{q}_{t_0}) - \mathbf{f}(\mathbf{p})| \\ & \leq |\mathbf{f}(\mathbf{q}_{t_0}) + \delta(\mathbf{q} - \mathbf{p}) - \mathbf{f}(\mathbf{q}_{t_0})| + M|\mathbf{q}_{t_0} - \mathbf{p}| \\ & \leq |\mathbf{f}(\mathbf{q}_{t_0} + \delta(\mathbf{q} - \mathbf{p})) - \mathbf{f}(\mathbf{q}_{t_0}) - \delta \mathbf{f}'(\mathbf{q}_{t_0})(\mathbf{q} - \mathbf{p})| + \delta M|\mathbf{q} - \mathbf{p}| \\ & \leq \varepsilon \delta |\mathbf{q} - \mathbf{p}| + \delta M|\mathbf{q} - \mathbf{p}| \\ & |\mathbf{q} - \mathbf{p}|(\varepsilon \delta + M\delta + (M + \varepsilon)t_0) \\ & |\mathbf{q} - \mathbf{p}|(\delta + t_0)(M + \varepsilon) \end{aligned} \quad \begin{aligned} & (M + \varepsilon)|\mathbf{q}_{t_0+\delta} - \mathbf{p}| \\ & (M + \varepsilon)|\mathbf{p} - \delta \mathbf{p} - t_0 \mathbf{p} + t_0 \mathbf{q} + \delta \mathbf{q} - \mathbf{p}| \\ & (M + \varepsilon)(t_0 + \delta)|\mathbf{q} - \mathbf{p}| \end{aligned}$$

$$\leq (M + \varepsilon)(t_0 + \delta)|\mathbf{q} - \mathbf{p}|$$

That was terrible.

1.61 Dec 18: ImpFT

- Let some $F(x, y, z)$ exist, then the relation $F(x, y, z) = 0$ can sometimes be reworked as $f(x, y) = z$, even if f will never be defined in terms of named functions
- Affine singlevariate real functions provide the best analogy: $F(x, y) = ax + by + c = 0$ can be expressed as $y = dx + e$ when $\frac{\partial F}{\partial y} \neq 0$, or else the line is vertical
- Similarly, the Implicit Function Theorem states that $F(x_1, x_2, \dots, x_d)$ can be redone as $x_k = f(x_1, x_2, \dots, x_{k-1}, x_{k+1}, \dots, x_d)$ in any domain (commonly a box) where $\frac{\partial F}{\partial x_k} \neq 0$
- f is at least as smooth as F where it is defined

1.62 Dec 19: Differentiability Revisited

- Let $f : D \subseteq \mathbb{R}^d \rightarrow \mathbb{R}$, such that the partial derivatives are continuous at \mathbf{p} and exist around \mathbf{p}

The candidate for A , or the derivative, is $\left[\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_d}\right]$, of which all entries exist. Now, an interesting sum

$$f(\mathbf{x}) - f(\mathbf{p}) = \sum_{j=1}^d (f(\mathbf{p}_j) - f(\mathbf{p}_{j-1}))$$

$$\mathbf{p}_0 = \mathbf{p} \quad \mathbf{p}_j = \mathbf{p}_{j-1} + (x_j - p_j)\mathbf{e}_j$$

Each increment of j switches just one coordinate of \mathbf{p} to the coordinate of \mathbf{x} . By the MVT, each term in the sum is equal to some partial derivative

$$\begin{aligned} & = \sum_{j=1}^d \frac{\partial f}{\partial x_j}(\mathbf{q}_j)(x_j - p_j) = \sum_{j=1}^d \frac{\partial f}{\partial x_j}(\mathbf{p})(x_j - p_j) + \sum_{j=1}^d \left(\frac{\partial f}{\partial x_j}(\mathbf{q}_j) - \frac{\partial f}{\partial x_j}(\mathbf{p}) \right) (x_j - p_j) \\ & = A(\mathbf{x} - \mathbf{p}) + \varepsilon(\mathbf{x}) \end{aligned}$$

$$\frac{|\varepsilon(\mathbf{x})|}{|\mathbf{x} - \mathbf{p}|} \leq \sum_{j=1}^d \left| \frac{\partial f}{\partial x_j}(\mathbf{q}_j) - \frac{\partial f}{\partial x_j}(\mathbf{p}) \right| \frac{|x_j - p_j|}{|\mathbf{x} - \mathbf{p}|} \rightarrow 0 \cdot k, k \leq 1$$

Then f has satisfied the test for differentiability

1.63 Dec 20: ImpFT again

- A partial proof, for one equation in three variables, called x, y, z for simplicity, WLOG $\frac{\partial F}{\partial z}(a, b, c) > 0$:
Choose $h > 0$, so the box $C_h^3(a, b, c) \subseteq \text{dom } F$ and $\frac{\partial F}{\partial z} > 0$ everywhere inside, so F is strictly increasing in z
Choose $r > 0$, so that $F(a, b, c-h) < 0$ for all $a, b \in C_r^2(a, b, c-h)$ and $F(a, b, c+h) > 0$ for all $a, b \in C_r^2(a, b, c+h)$
For every a, b in the square, there is one c_0 where $F(a, b, c_0) = 0$, by the IVT and strict increasing,
then write $c_0 = f(a, b)$ on that square

1.64 Jan 2: ImpFT again again

- Even if $\frac{\partial F}{\partial y} = 0$, there may still exist a function f such that $f(\mathbf{x}) = y$, but the ImpFT cannot help
- To prove the smoothness preserving,
 - f is continuous on the relevant box
 - Partial derivatives of f exist in the box
 - $\frac{\partial f}{\partial x_j}(\mathbf{x}) = -\frac{\frac{\partial F}{\partial x_j}(\mathbf{x}, y)}{\frac{\partial F}{\partial y}(\mathbf{x}, y)}$
 - f is C^1
 - f is C^k by induction

1.65 Jan 3: ImpFT proooofs

- To prove point 1:

$$0 = F(\mathbf{x}, f(\mathbf{x})) - F(\mathbf{x}_0, f(\mathbf{x}_0))$$

\mathbf{x} is fixed, and \mathbf{x}_0 is fixed. Since F is by hypothesis C^1 , the MVT applies, but one term of the dot product is separated:

$$0 = \nabla_{\mathbf{x}} F(\mathbf{x}^*, y^*) \cdot (\mathbf{x} - \mathbf{x}_0) + \frac{\partial F}{\partial y}(\mathbf{x}^*, y^*)(f(\mathbf{x}) - f(\mathbf{x}_0))$$

Let $\mathbf{x} \rightarrow \mathbf{x}_0$

$$0 = [\text{bounded by EVT}] \cdot 0 + [\text{not 0 in the box}] \cdot (f(\mathbf{x}) - f(\mathbf{x}_0))$$

$$f(\mathbf{x}) - f(\mathbf{x}_0) \rightarrow 0$$

So, f is continuous

- To prove points 2, 3, 4, make $\mathbf{x} = \mathbf{x}_0 + \mathbf{e}_j h$:

$$\frac{f(\mathbf{x}) - f(\mathbf{x}_0)}{h} = \frac{-\nabla_{\mathbf{x}} F(\mathbf{x}^*, y^*) \cdot (\mathbf{e}_j) h}{h \frac{\partial F}{\partial y}(\mathbf{x}^*, y^*)} = \frac{-\frac{\partial F}{\partial x_j}(\mathbf{x}^*, y^*)}{\frac{\partial F}{\partial y}(\mathbf{x}^*, y^*)}$$

Let $h \rightarrow 0, \mathbf{x} \rightarrow \mathbf{x}_0, y \rightarrow y_0$

$$\frac{\partial f}{\partial x_j}(\mathbf{x}) = -\frac{\frac{\partial F}{\partial x_j}(\mathbf{x}, y)}{\frac{\partial F}{\partial y}(\mathbf{x}, y)}$$

Since both partial derivatives on the right are continuous, and the denominator is not 0, their quotient is continuous

- To prove point 5:
If f is C^{k-1} and F is C^k ,

$$\frac{\partial^{k-1} f}{\partial x_j^{k-1}}(\mathbf{x}) = \frac{\frac{\partial^{k-1} F}{\partial x_j^{k-1}}(\mathbf{x}, y) \dots}{\frac{\partial^{k-1} F}{\partial y^{k-1}}(\mathbf{x}, y) \dots}$$

The right is C^1 so the left is too, then F is C^k

- What about when there are multiple functions $F, G, H \dots$?

$F(x, y, z, w) = 0 \quad G(x, y, z, w) = 0$ is solved by (a, b, c, d)

F, G are C^1

$$\text{At } (a, b, c, d), \quad \frac{\partial F}{\partial z} \frac{\partial G}{\partial w} - \frac{\partial F}{\partial w} \frac{\partial G}{\partial z} \neq 0$$

Because of the non-degeneracy condition, one of $\frac{\partial F}{\partial z}$ or $\frac{\partial F}{\partial w}$ is not 0, WLOG assume the former is not 0
By ImpFT,

$$z = f(x, y, w)$$

Non-degeneracy for G and w is also necessarily true now so,

$$w = g(x, y, z)$$

Now define $H(x, y, z) = G(x, y, z, g(x, y, z))$

$$\frac{\partial H}{\partial z} = \frac{\partial G}{\partial z} + \frac{\partial G}{\partial w} \frac{\partial f}{\partial z} = \frac{\partial G}{\partial w} - \frac{\partial G}{\partial w} \frac{\partial F}{\partial z} = \frac{\frac{\partial F}{\partial z} \frac{\partial G}{\partial w} - \frac{\partial F}{\partial w} \frac{\partial G}{\partial z}}{\frac{\partial F}{\partial w}} \neq 0$$

Because $w = g$ is also C^1 , H is now C^1 .

$$z = \phi(x, y)$$

$$w = g(x, y, \phi(x, y)) = \psi(x, y)$$

1.66 Jan 6: ImpFT prooooooooooooofs

- The proof of the ImpFT for more than one equation was reviewed
- The generalization for solving for n of $n + d$ variables with n equations requires a consistent non-degeneracy condition:

$$\det \frac{\partial(F_1, F_2, \dots, F_n)}{\partial(x_1, x_2, \dots, x_n)} = \det \begin{bmatrix} \frac{\partial F_1}{\partial x_1} & \frac{\partial F_1}{\partial x_2} & \dots & \frac{\partial F_1}{\partial x_n} \\ \frac{\partial F_2}{\partial x_1} & \frac{\partial F_2}{\partial x_2} & \dots & \frac{\partial F_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial F_n}{\partial x_1} & \frac{\partial F_n}{\partial x_2} & \dots & \frac{\partial F_n}{\partial x_n} \end{bmatrix} \neq 0$$

This represents a partial Jacobian that is constrained in variables to a square, which is contrasted to the regular Jacobian that contains all the variables

1.67 Jan 7: ImpFT to the Max

- General ImpFT:

$$\mathbf{F}(\mathbf{x}, \mathbf{y}) = 0 \text{ is } C^1$$

$$\mathbf{F}(\mathbf{a}, \mathbf{b}) = 0$$

With the non-degeneracy condition above, in some box around (\mathbf{a}, \mathbf{b}) ,

$$\exists \mathbf{f} : \mathbf{y} = \mathbf{f}(\mathbf{x})$$

$$\frac{\partial \mathbf{f}}{\partial \mathbf{x}} = \mathbf{f}' = - \left(\frac{\partial \mathbf{F}}{\partial \mathbf{y}}(\mathbf{x}, \mathbf{f}(\mathbf{x})) \right)^{-1} \left(\frac{\partial \mathbf{F}}{\partial \mathbf{x}}(\mathbf{x}, \mathbf{f}(\mathbf{x})) \right)$$

1.68 Jan 8: ImpFT; You thought it was over!

- The determinant of a small matrix is commonly found via cofactor expansion; the proof of this is truly marvelous, which this page is too narrow to contain; a hat index indicates omission of a row or column.

$$\det \mathbf{A} = \sum_{j=1}^n (-1)^{j-1} a_{1j} \det \mathbf{A}_{\hat{1}\hat{j}}$$

1.69 Jan 9: ImpFT ImpFT ImpFT ImpFT ImpFT ImpFT

- The induction step that will finally be resolved: A system of $m+1$ equations in $m+1+d$ variables that fulfills non-degeneracy can be turned into a system of m equations in $m+d$ variables

$$\begin{cases} F(\mathbf{x}, y, \mathbf{z}) = 0 \\ G(\mathbf{x}, y, \mathbf{z}) = 0 \end{cases}$$

Or,

$$H(\mathbf{x}, y, \mathbf{z}) = 0$$

By assumption,

$$\left| \frac{\partial H}{\partial(y, \mathbf{z})} \right| = \begin{vmatrix} \frac{\partial F}{\partial y} & \frac{\partial F}{\partial \mathbf{z}} \\ \frac{\partial G}{\partial y} & \frac{\partial G}{\partial \mathbf{z}} \end{vmatrix} = \frac{\partial F}{\partial y} \left| \frac{\partial G}{\partial \mathbf{z}} \right| - \dots \neq 0$$

At least one term of the determinant is not 0, so WLOG make the nonzero one y ;

$$\frac{\partial F}{\partial y} \left| \frac{\partial G}{\partial \mathbf{z}} \right| \neq 0$$

Now, solve F for y , $0 = G(\mathbf{x}, f(\mathbf{x}, \mathbf{z}), \mathbf{z}) = K(\mathbf{x}, \mathbf{z})$ A second non-degeneracy condition is necessary to solve K for \mathbf{z}

$$\frac{\partial K}{\partial \mathbf{z}} = \begin{bmatrix} \frac{\partial G_1}{\partial z_1} & \frac{\partial G_1}{\partial z_2} & \cdots \\ \frac{\partial G_2}{\partial z_1} & \frac{\partial G_2}{\partial z_2} & \cdots \\ \vdots & \vdots & \ddots \end{bmatrix} = \begin{bmatrix} \frac{\partial G_1}{\partial y} \frac{\partial f}{\partial z_1} & \frac{\partial G_1}{\partial y} \frac{\partial f}{\partial z_2} & \cdots \\ \frac{\partial G_2}{\partial y} \frac{\partial f}{\partial z_1} & \frac{\partial G_2}{\partial y} \frac{\partial f}{\partial z_2} & \cdots \\ \vdots & \vdots & \ddots \end{bmatrix} + \begin{bmatrix} \frac{\partial G_1}{\partial z_1} & \frac{\partial G_1}{\partial z_2} & \cdots \\ \frac{\partial G_2}{\partial z_1} & \frac{\partial G_2}{\partial z_2} & \cdots \\ \vdots & \vdots & \ddots \end{bmatrix} = \frac{\partial G}{\partial y} \frac{\partial F}{\partial \mathbf{z}} + \frac{\partial G}{\partial \mathbf{z}}$$

1.70 Jan 10: ImpFT

- To continue:

$$= -\frac{\partial G}{\partial y} \frac{\partial F}{\partial \mathbf{z}} + \frac{\partial G}{\partial \mathbf{z}} = \frac{\frac{\partial G}{\partial \mathbf{z}} \frac{\partial F}{\partial y} - \frac{\partial G}{\partial y} \frac{\partial F}{\partial \mathbf{z}}}{\frac{\partial F}{\partial y}}$$

By analogy to determinants and induction with matrix blocks,

$$\left| \frac{\partial K}{\partial \mathbf{z}} \right| = \frac{\left| \frac{\partial H}{\partial(y, \mathbf{z})} \right|}{\frac{\partial F}{\partial y}} \neq 0$$

Now, $\mathbf{z} = \phi(\mathbf{x})$, $y = f(\mathbf{x}, \phi(\mathbf{x})) = \psi(\mathbf{x})$

1.71 Jan 13: ImpFT

- Induction with matrix blocks:

$$\begin{vmatrix} a & \mathbf{c}^T \\ \mathbf{b} & \mathbf{D} \end{vmatrix} = a \begin{vmatrix} 1 & \frac{\mathbf{c}^T}{a} \\ \mathbf{b} & \mathbf{D} \end{vmatrix}$$

By a multiple row shear, where each row from 2 to m is increased by some multiple of the previous row so that the left entry is 0,

$$= \begin{vmatrix} 1 & \frac{\mathbf{c}^T}{a} \\ 0 & \mathbf{D} - \mathbf{b} \frac{\mathbf{c}^T}{a} \end{vmatrix}$$

Then cofactor expansion on the first row, following block multiplication laws, and extraction of a constant from every entry,

$$= a \begin{vmatrix} \mathbf{A} - \mathbf{b} \frac{\mathbf{c}^T}{a} \end{vmatrix} = \begin{vmatrix} a \frac{\mathbf{A}}{a} - \mathbf{b} \frac{\mathbf{c}^T}{a} \end{vmatrix} = \frac{1}{a^{m-1}} |a\mathbf{A} - \mathbf{b}\mathbf{c}^T|$$

- The Inverse Function Theorem: given

$$\mathbf{f} : D \subseteq \mathbb{R}^d \xrightarrow{C^1} \mathbb{R}^d$$

and

$$\forall \mathbf{p} \in D, |\mathbf{f}'(\mathbf{p})| \neq 0$$

then

$$\exists U \ni \mathbf{p}, V \ni \mathbf{p}, (\mathbf{g} : V \xrightarrow{C^1} U) : \forall \mathbf{x} \in U, \mathbf{y} \in V,$$

$$\mathbf{f}(\mathbf{g}(\mathbf{y})) = \mathbf{y}$$

$$\mathbf{g}(\mathbf{f}(\mathbf{x})) = \mathbf{x}$$

or, \mathbf{g} is a local inverse of \mathbf{f} near \mathbf{p}

- The justification for this arises from ImpFT:

$$\mathbf{F}(\mathbf{x}, \mathbf{y}) := \mathbf{y} - \mathbf{f}(\mathbf{x})$$

with $\dim \mathbf{x} = \dim \mathbf{y}$. By its definition,

$$\mathbf{F}(\mathbf{x}, \mathbf{y}) = 0 \Rightarrow \mathbf{y} = \mathbf{f}(\mathbf{x})$$

By ImpFT, given non-degeneracy,

$$\mathbf{F}(\mathbf{x}, \mathbf{y}) = 0 \Rightarrow \mathbf{x} = \mathbf{g}(\mathbf{y})$$

1.72 Jan 15: InvFT

- Let $F(\mathbf{y}, \mathbf{x}) = \mathbf{y} - \mathbf{f}(\mathbf{x})$. With input (\mathbf{y}, \mathbf{x}) , $F = 0$, so there is a solution point.

- Non-degeneracy of F : $0 \neq \left| \frac{\partial \mathbf{f}}{\partial \mathbf{x}} \right|$

$$= \left| \frac{\partial \mathbf{f}}{\partial x_1} \cdots \right| = \left| \frac{\partial \mathbf{F}}{\partial x_1} \cdots \right| = \frac{\partial \mathbf{F}}{\partial \mathbf{x}}$$

- Continuity of the partials:

$$\frac{\partial \mathbf{F}}{\partial \mathbf{x}} = -\frac{\partial \mathbf{f}}{\partial \mathbf{x}}, C^1$$

$$\frac{\partial \mathbf{F}}{\partial \mathbf{y}} = I, C^\infty$$

- Now, the ImpFT applies!

1.73 Jan 16: InvFT sets

- It remains to be seen that the open sets U, V actually exist from the InvFT.
- \mathbf{g} is injective, and \mathbf{f} is surjective.
- Let $U = \mathbf{g}[V]$, for simplicity, then what is V ? The C^1 image of an open set may not be open, but its preimage will definitely be. $V = \mathbf{f}^{-1}[\mathbf{f}[C_r(\mathbf{q})]] \cap C(\mathbf{q})$, an open set.

1.74 Jan 17: Optimization

- Constrained and unconstrained extrema of a real function can be:

- Local: on a neighborhood around each point
- Relative: on some set
- Global: on the domain

- A local strict maximum of $f : D \subseteq \mathbb{R}^d \rightarrow \mathbb{R}$, \mathbf{p} , occurs when

$$(\exists r > 0)(\forall \mathbf{x} \in B_r(\mathbf{p}) \cap D)(f(\mathbf{x}) < f(\mathbf{p}))$$

This can be true for any function f , but to actually find p , it is necessary for f to be C^1

- Interval analysis is impossible due to the simple fact that intervals are not in \mathbb{R}^d . The second derivative test is the preferred method, which has a specific analogue, dealing with all $\binom{d}{2} + d$ second partial derivatives.
- If \mathbf{p} is a local extremum of f , then \mathbf{p} is a critical point, or $f'(\mathbf{p})$ is not full rank, or $\nabla f(\mathbf{p}) = 0$.

$$\frac{\partial f}{\partial x_j}(\mathbf{p}) = \lim_{h \rightarrow 0} \frac{f(\mathbf{p} + h\mathbf{e}_j) - f(\mathbf{p})}{h} = \lim_{h \rightarrow 0} \frac{K}{h}$$

This limit exists because of C^1 -ness. WLOG let f be a local minimum, so $K \geq 0$

$$\lim_{h \rightarrow 0^-} \frac{K}{h} \leq 0$$

$$\lim_{h \rightarrow 0^+} \frac{K}{h} \geq 0$$

The conclusion is apparent.

2 Spring 2020

2.1 Jan 28: Critical Yeehaw

- The analogue of the second derivative test simply uses the eigenvalues of the matrix formed by all the partial derivatives (even repeated mixed ones). For 2×2 matrices this is not too bad.
- If the second derivative test is inconclusive, move on to higher derivatives. If an even derivative is not zero, there is a local extremum; if an odd derivative is not zero, there is a saddle point; if the required derivative does not exist or every derivative is 0, the test fails catastrophically or does not even halt.
- A practical example: $f(x, y) = x^2 - y^2$, the simplest hyperbolic paraboloid (which is named for its traces along the axial planes)

$$f'(x, y)|_{(0,0)} = \begin{bmatrix} 0 & 0 \end{bmatrix}$$

However, the origin is clearly not a local extremum, for

$$f(0, \varepsilon) = -\varepsilon^2 \quad f(\varepsilon, 0) = \varepsilon^2$$

Calculate the determinant of the 2×2 Hessian, the product of its eigenvalues:

$$\det \mathbf{H} = \begin{vmatrix} \frac{\partial^2 f}{\partial x^2} & \frac{\partial^2 f}{\partial x \partial y} \\ \frac{\partial^2 f}{\partial y \partial x} & \frac{\partial^2 f}{\partial y^2} \end{vmatrix} = \begin{vmatrix} 2 & 0 \\ 0 & -2 \end{vmatrix} = -4$$

Then the eigenvalues must have opposite sign, and a saddle point is confirmed.

2.2 Jan 29: Taylor's Theorem

- Taylor's Theorem:

$$f : I \subseteq \mathbb{R} \xrightarrow{C^{k+1}} \mathbb{R}, (a-r, a+r) \in I$$

$$\Rightarrow (\forall h : |h| < r) \left(f(a+h) = \sum_{j=0}^k \frac{f^{(j)}(a)}{j!} h^j + \frac{h^{k+1}}{k!} \int_0^1 (1-t)^k f^{(k+1)}(a+th) dt \right) = (T_k^a f)(h) + (R_k^a f)(h)$$

- The error term must be $o(h^j)$, or else it subsumes terms of the polynomial
- To prove the theorem, establish the base case for degree 0:

$$f(a) + h \int_0^1 f'(a+th) dt = f(a) + \int_a^{a+h} f'(u) du = f(a+h)$$

- Then, induction: assume the theorem holds for degree m , then establish it for degree $m + 1$ by integration by parts:

$$\begin{aligned}
f(a+h) &= (T_m^a f)(h) + \frac{h^{m+1}}{m!} \int_0^1 (1-t)^m f^{(m+1)}(a+th) dt \\
&= (T_m^a f)(h) + \frac{h^{m+1}}{m!} \left(f^{(m+1)}(a+th) \cdot \frac{-(1-t)^{m+1}}{m+1} \right]_0^1 + \int_0^1 \frac{(1-t)^{m+1}}{m+1} \cdot h f^{(m+2)}(a+th) dt \\
&= (T_m^a f)(h) + \frac{h^{m+1}}{m!} \left(0(\dots) + \frac{f^{(m+1)}(a+th)}{m+1} \right) + \frac{h^{m+2}}{(m+1)!} \int_0^1 (1-t)^{m+1} f^{(m+2)}(a+th) dt \\
&= (T_{m+1}^a f)(h) + \frac{h^{m+2}}{(m+1)!} \int_0^1 (1-t)^{m+1} f^{(m+2)}(a+th) dt
\end{aligned}$$

This only applies for $m \leq k-1$, for the $m+2^{th}$ derivative to be integrable

2.3 Jan 30: One Index was Never Enough!

- Multi-indexes!

Let the following exist:

$$\mathbf{x} \quad \alpha$$

Both are k -dimensional vectors, with α composed of strictly non-negative integers.

- The degree:

$$|\alpha| = \sum \alpha_j$$

- The factorial:

$$\alpha! = \prod (\alpha_j!)$$

- The power of a vector:

$$\mathbf{x}^\alpha = \prod (x_j^{\alpha_j})$$

- The multi-differential, in the appropriate order:

$$\partial^\alpha f = \frac{\partial^{|\alpha|} f}{\prod (\partial_{x_j}^{\alpha_j})}$$

A practical example:

$$\partial^{(2,0,5)} f = \frac{\partial^7 f}{\partial_{x_3}^5 \partial_{x_1}^2}$$

- To apply, any polynomial of degree at most k in \mathbf{x} can be written as:

$$P(\mathbf{x}) = \sum_{|\alpha| \leq k} c_\alpha \mathbf{x}^\alpha$$

- The multinomial theorem mirrors the binomial theorem:

$$(x_1 + x_2)^n = \sum_{|\alpha| \leq n} \frac{n!}{\alpha!} \mathbf{x}^\alpha$$

$$(x_1 + x_2 + \dots + x_k)^n = \sum_{\substack{|\alpha| \leq n \\ \dim(\alpha) = k}} \frac{n!}{\alpha!} \mathbf{x}^\alpha$$

This can be derived by induction on k

- The dot gradient:

$$\mathbf{h} \cdot \nabla = \sum_j h_j \frac{\partial}{\partial x_j}$$

Miraculously,

$$(\mathbf{h} \cdot \nabla)^{(j)} = \sum_{|\alpha|=j} \mathbf{h}^\alpha \partial^\alpha$$

2.4 Jan 31: protobowl.com

- Try it, it's fun

2.5 Feb 1: Cross Partial Derivative Equality

- Use the wonderful square on C^2 f: Let $\Delta(h) = f(x+h, y+h) + f(x, y) - f(x+h, y) - f(x, y+h)$, where $h \neq 0$, and h is sufficiently small so that each point of the square is in the domain,

$$f(x+h, y+h) - f(x, y+h) - (f(x+h, y) - f(x, y)) = f(x+h, y+h) - f(x+h, y) - (f(x, y+h) - f(x, y))$$

Each pair can be redefined as a C^1 function of just the changing variable, so MVT applies, with $\theta_j \in (0, 1)$,

$$hf'_x(x + \theta_1 h, y + h) - hf'_x(x + \theta_2 h, y) = hf'_y(x + h, y + \theta_3 h) - hf'_y(x, y + \theta_4 h)$$

Divide by h , take limits, do the same,

$$hf''_{yx}(x, y + \theta_5 h) = hf''_{xy}(x + \theta_6 h, y)$$

Arrive at:

$$f''_{yx}(x, y) = f''_{xy}(x, y)$$

- The same general property is true for higher-order partials, since swaps of consecutive elements allows any permutation of a sequence to be reached, and enough smoothness ensures that every swap happens to a C^2 function

$$f \in C^k \wedge i, j \text{ are permutations} \Rightarrow f^{(k)}_{x_{i_k}, x_{i_{k-1}}, \dots, x_{i_2}, x_{i_1}}(\mathbf{x}) = f^{(k)}_{x_{j_k}, x_{j_{k-1}}, \dots, x_{j_2}, x_{j_1}}(\mathbf{x})$$

2.6 Feb 5: Taylor's Theorem

- Armed with the mighty multi-index, attempt to state Taylor's Theorem for multivariate scalar functions $f : D \subset \mathbb{R}^d \rightarrow \mathbb{R}$:

$$f(\mathbf{a} + \mathbf{h}) = \sum_{\substack{|\alpha| \leq k \\ \dim(\alpha) = d}} \frac{(\partial^\alpha f)(\mathbf{a})}{\alpha!} \mathbf{h}^{(\alpha)} + (R_k^{\mathbf{a}} f)(\mathbf{h})$$

2.7 Feb 6: Directional Derivatives

- Derivatives with a direction!

$$(D_{\mathbf{u}} f)(\mathbf{p}) = f'_{\mathbf{u}}(\mathbf{p}) = \lim_{h \rightarrow 0} \frac{f(\mathbf{p} + h\mathbf{u}) - f(\mathbf{p})}{h}$$

The intuitive definition of the direction vector dotted with the gradient only works with $|\mathbf{u}| = 1$, otherwise scale it to represent velocity

- A useful relation immediately derived from the intuitive definition:

$$\nabla f(\mathbf{p}) \cdot \mathbf{u} = \sum_{j=1}^d u_j \frac{\partial f}{\partial p_j}(\mathbf{p}) = f'_{\mathbf{u}}(\mathbf{p})$$

- From CS, the machine-learning concept of steepest descent can be understood through gradients; as the cosine of the angle between \mathbf{u} and ∇ increases, the directional derivative gets closer to $|\nabla|$:

$$\max_{|\mathbf{u}|=1} f'_{\mathbf{u}}(\mathbf{p}) = |\nabla f(\mathbf{p})| \text{ attained iff } \mathbf{u} = \frac{\nabla f(\mathbf{p})}{|\nabla f(\mathbf{p})|}$$

A gradient is now easily seen to be normal to the level sets, or isovalues, where it is computed.

2.8 Feb 7: Directional Directional Derivatives Derivatives

- What about $f''_{\mathbf{u}\mathbf{u}}$?

$$\begin{aligned} f'_{\mathbf{u}} &= (\nabla f)(\mathbf{p}) \cdot \mathbf{u} = \sum_{j=1}^d u_j \cdot \frac{\partial f}{\partial x_j}(\mathbf{p}) \\ f''_{\mathbf{u}\mathbf{u}} &= (\nabla((\nabla f) \cdot \mathbf{u}))(\mathbf{p}) \cdot \mathbf{u} = \sum_{i=1}^d u_j \frac{\partial}{\partial x_i} \left(\sum_{j=1}^d u_j \frac{\partial f}{\partial x_j}(\mathbf{x}) \right) (\mathbf{p}) \\ &= \sum_{i=1}^d \sum_{j=1}^d u_i u_j \frac{\partial^2 f}{\partial x_i \partial x_j} \end{aligned}$$

This definition is easily extensible to higher-order directional derivatives all in one direction! It is also equivalent to the more compact, but more cryptic,

$$\begin{aligned} &= (\mathbf{u} \cdot \nabla)(\mathbf{u} \cdot \nabla)f \\ &= \mathbf{u}^T \mathbf{H}_f(\mathbf{p}) \end{aligned}$$

2.9 Feb 10: Multinomial Proofs

- Define $\sigma(\mathbf{x}) = \sum x_k$
- By induction on the dimension of the vector \mathbf{x} :
- Base Case:

$d = 2$, equivalent to the Binomial Theorem, which is considered trivial

- Inductive Step:
By the Binomial Theorem,

$$(\sigma(\mathbf{x}) + y)^n = \sum_{j=0}^n \frac{n!}{j!(n-j)!} (\sigma(\mathbf{x}))^j y^{n-j}$$

Assume the multinomial theorem for $d = n$

$$= \sum_{j=0}^n \frac{n!}{j!(n-j)!} \sum_{|\alpha|=j} \frac{j!}{\alpha!} \mathbf{x}^\alpha y^{n-j} = \sum_{j=0}^n \sum_{|\alpha|=j} \frac{n!}{(n-j)!\alpha!} \mathbf{x}^\alpha y^{n-j}$$

Let $\beta = (\alpha, n-j)$, then there is a bijection from α, n to β

$$= \sum_{|\beta|=n} \frac{n!}{\beta!} (\mathbf{x}, y)^\beta$$

- Functional multinomials are interpreted differently but obey a similar relationship, as from earlier:

$$(\mathbf{h} \cdot \nabla)^n = \sum_{|\alpha|=n} \frac{n!}{\alpha!} \mathbf{h}^\alpha \partial^\alpha$$

2.10 Feb 11: A Tense Tactical Retreat; One Index was Actually Enough!

- Multi-indexes are hard, retreat to (obviously easier to understand) tensors:

$$f(\mathbf{a} + \mathbf{h}) = \sum_{j=0}^k \frac{f^{(j)}(\mathbf{a})}{j!} \mathbf{h}^{\otimes j} + (R_k^{\mathbf{a}} f)(\mathbf{h})$$

The same limit on the remainder must apply

- $\mathbf{h}^{\otimes n}$ is the n-th tensor power of \mathbf{h} , or the n-tensor with entries $h_{i_1} h_{i_2} \dots h_{i_n}$

- $f^{(n)}(\mathbf{a})$ is harder to materialize; it is the n -tensor with entries $\frac{\partial^j f}{\partial x_{i_1} \dots \partial x_{i_j} \partial x_{i_1}}$
- The tensor contraction is equivalent to sifting out the terms of a sum where corresponding indices are equal:

$$\begin{aligned} \frac{f^{(j)}(\mathbf{a})}{j!} \mathbf{h}^{\otimes j} &= \sum_{i_1, i_2, \dots, i_j, l_1, l_2, \dots, l_j=1}^d \left(\prod_{m=1}^j h_{i_m} \right) \left(\frac{\partial f^{(j)}}{\partial x_{l_j} \dots \partial x_{l_2} \partial x_{l_1}} \right) \left(\prod_{m=1}^j \delta_{i_m l_m} \right) \\ &= \sum_{i_1, i_2, \dots, i_j}^d \prod_{m=1}^j h_{i_m} \frac{\partial f}{\partial x_{i_m}} \end{aligned}$$

2.11 Feb 12: Sift Sift Sift

- A alternate, but non-extensible way to represent the second Taylor term is:

$$\begin{aligned} &\mathbf{h}^T(\mathbf{H}_f)(\mathbf{a})\mathbf{h} \\ &= \begin{bmatrix} h_1 & h_2 & \dots & h_d \end{bmatrix} \begin{bmatrix} f_{x_1 x_1} & f_{x_1 x_2} & \dots & f_{x_1 x_d} \\ f_{x_2 x_1} & f_{x_2 x_2} & \dots & f_{x_2 x_d} \\ \vdots & \vdots & \ddots & \vdots \\ f_{x_d x_1} & f_{x_d x_2} & \dots & f_{x_d x_d} \end{bmatrix} \begin{bmatrix} h_1 \\ h_2 \\ \vdots \\ h_d \end{bmatrix} \\ &= \left[\sum_{i=1}^d f_{x_i x_1} h_i \quad \sum_{i=1}^d f_{x_i x_2} h_i \quad \dots \quad \sum_{i=1}^d f_{x_i x_d} h_i \right] \begin{bmatrix} h_1 \\ h_2 \\ \vdots \\ h_d \end{bmatrix} = \sum_{i,j=1}^d f_{x_i x_j} h_i h_j = f''(\mathbf{a}) \mathbf{h}^{\otimes 2} \end{aligned}$$

2.12 Feb 13: Quadratizing

- Let $f : D \in \mathbb{R}^d \xrightarrow{C^3} \mathbb{R}$, with a critical point at \mathbf{a} . The Taylor series is now

$$f(\mathbf{a} + \mathbf{h}) = f(\mathbf{a}) + 0 + \mathbf{h}^T(\mathbf{H}_f)(\mathbf{a})\mathbf{h} + \varepsilon$$

- Every square matrix has an associated quadratic form, or the sum of all quadratic terms that is like the tensor contraction:

$$Q_{\mathbf{A}}(\mathbf{h}) = \mathbf{a}^T \mathbf{A} \mathbf{a} = \sum_{i,j} \mathbf{A}_{ij} h_i h_j$$

This can be divided into the pure and cross terms, which when applied to the Hessian results in the two types of second partials:

$$= \sum_i \mathbf{A}_{ii} h_i^2 + 2 \sum_{i < j} \mathbf{A}_{ij} h_i h_j$$

- Scaling either part of a quadratic form proceeds simply:

$$Q_{c\mathbf{A}} = cQ_{\mathbf{A}} \quad Q_{\mathbf{A}}(c\mathbf{a}) = c^2 Q_{\mathbf{A}}(\mathbf{a})$$

- A symmetric non-zero square matrix can be classified on its definiteness:
 - A positive or negative definite matrix has a quadratic form that is positive or negative for all \mathbf{a} , and has all positive or negative eigenvalues
 - A positive or negative semidefinite matrix has a quadratic form that is not positive or not negative for all \mathbf{a} , and is zero for some \mathbf{a} ; zero is some of its eigenvalues, and the rest are positive or negative
 - An indefinite matrix has a quadratic form that takes both signs, and its eigenvalues do too
- The second derivative test can also now be stated in terms of definiteness. If \mathbf{H}_f is positive or negative definite, there is a local minimum or maximum; an indefinite \mathbf{H}_f indicates a saddle point. A semidefinite matrix is inconclusive.

- There is a simple justification, for a positive definite Hessian,

$$f(\mathbf{a} + \mathbf{h}) = f(\mathbf{a}) + 0 + D + \varepsilon$$

D is always positive, and ε is subsumed by it, so any displacement results in a higher value. A negative definite is essentially the same, but D is always negative; an indefinite matrix necessarily means changes in both directions with any displacement. A semidefinite matrix means that along some directions, ε prevails, meaning that the sign is hard to determine.

2.13 Feb 14: Definitely Indefinite

- Notation is man's best friend! Definiteness can be symbolized nicely as:

$$\mathbf{A} \succ 0 \quad \mathbf{A} \succcurlyeq 0 \quad \mathbf{A} \prec 0 \quad \mathbf{A} \preccurlyeq 0 \quad \mathbf{A} \asymp 0$$

The meanings are quite obvious, really.

- How is the definiteness of a matrix actually determined? Larger $n \times n$ matrices are much better determined by Sylvester's Criterion: take the determinant of each upper-left square submatrix. The n determinants determine the definiteness just like the quadratic form in a weird way, except only finite values need to be checked!

2.14 Feb 24: Shackles

- The essential point of the argument linking definiteness to extrema:

$$(\forall \mathbf{A})(\exists \alpha, \beta)(\forall \mathbf{h}) (\alpha |\mathbf{h}|^2 \leq |Q_{\mathbf{A}}(\mathbf{h})| \leq \beta |\mathbf{h}|^2)$$

Factor out $|\mathbf{h}|^2$ from all expressions:

$$\alpha \leq \left| \frac{\mathbf{h}^T}{|\mathbf{h}|} \mathbf{A} \frac{\mathbf{h}}{|\mathbf{h}|} \right| \leq \beta$$

Rename the unit vector in the direction of \mathbf{h} , \mathbf{u}

$$\alpha \leq |Q_{\mathbf{A}}(\mathbf{u})| \leq \beta$$

The quadratic form of a matrix is continuous to the vector, and \mathbf{u} resides in a compact set. By EVT, the output is bounded!

$$\alpha = \min_{\mathbf{u} \in S_1^d} |Q_{\mathbf{A}}(\mathbf{u})| \quad \beta = \max_{\mathbf{u} \in S_1^d} |Q_{\mathbf{A}}(\mathbf{u})|$$

If the matrix is definite, then $\alpha > 0$, since the absolute value of a non-zero is positive. Otherwise, $\alpha = 0$

2.15 Feb 25: Sylvester

- Analyzing the quadratic form need only be done for symmetric matrices. Take \mathbf{A} to be $\begin{bmatrix} a & b \\ c & d \end{bmatrix}$

$$Q_{\mathbf{A}}(x, y) = ax^2 + 2 \left(\frac{b+c}{2} \right) xy + dy^2 = Q_{\mathbf{B}}(x, y)$$

$$\text{where } \mathbf{B} = \begin{bmatrix} a & \frac{b+c}{2} \\ \frac{b+c}{2} & d \end{bmatrix}$$

This proof obviously will not work out for higher matrices. Let $\mathbf{B} = \frac{\mathbf{A} + \mathbf{A}^T}{2}$,

$$Q_{\mathbf{B}}(\mathbf{h}) = \mathbf{h}^T \mathbf{B} \mathbf{h} = \mathbf{h}^T \frac{\mathbf{A} + \mathbf{A}^T}{2} \mathbf{h} = \frac{\mathbf{h}^T \mathbf{A} \mathbf{h}}{2} + \frac{\mathbf{h}^T \mathbf{A}^T \mathbf{h}}{2}$$

Establish the simple equalities:

$$(\mathbf{h}^T \mathbf{A}^T \mathbf{h})^T = \mathbf{h}^T \mathbf{A}^T \mathbf{h}$$

Because a scalar is its own transpose

$$(\mathbf{h}^T \mathbf{A}^T \mathbf{h})^T = \mathbf{h}^T \mathbf{A} \mathbf{h}$$

By the properties of the transpose of a matrix product

$$\frac{\mathbf{h}^T \mathbf{A} \mathbf{h}}{2} + \frac{\mathbf{h}^T \mathbf{A}^T \mathbf{h}}{2} = \mathbf{h}^T \mathbf{A} \mathbf{h} = Q_{\mathbf{A}}(\mathbf{h})$$

- The Sylvester criterion for a 2×2 matrix:

$$B \succ 0 \Rightarrow \det a > 0 \wedge \det B > 0$$

$$Q_{\mathbf{A}}(x, y) = ax^2 + 2bxy + cy^2$$

If $y = 0$, then the expression collapses to ax^2 , which is easy to analyze

$$= y^2 \left(a \left(\frac{x}{y} \right)^2 + 2b \frac{x}{y} + c \right)$$

If $a = 0$, then the expression is linear and analyzed separately

$$= ay^2 \left(\left(\frac{x}{y} \right)^2 + \frac{2b}{a} \frac{x}{y} \right) + y^2(c)$$

Complete the square,

$$= ay^2 \left(\left(\frac{x}{y} \right)^2 + \frac{2b}{a} \frac{x}{y} + \left(\frac{b}{a} \right)^2 \right) + y^2 \left(c - \frac{b^2}{a} \right) = y^2 \left(a \left(\frac{x}{y} + \frac{b}{a} \right)^2 + \frac{\det B}{a} \right)$$

For this result to always be positive, with a bit of casework, the criterion is shown to be true. For the negative definite case, simply negate $\det a$

- Extended to higher matrices, the n^{th} principal submatrix determinant must be always positive for a positive definite matrix, and negative for odd n and positive for even n for a negative definite matrix

2.16 Feb 26: Eigenthings

- Eigenvectors are characteristics of a linear transformation. For these vectors, the transformation will be equivalent to scaling by a complex number; the number is the associated eigenvalue
- Interestingly, a symmetric real matrix is guaranteed to have only real eigenvalues, making this character even better when visualized!
- These eigenvalues λ_n are elements of the spectrum of the matrix, $\sigma(\mathbf{A})$, and the eigenvectors of each eigenvalue are elements of the specific eigenspace $E_{\lambda}(\mathbf{A})$
- The Spectral Theorem asserts that a symmetric real matrix is diagonalizable in a different basis, with the only non-zero entries being the eigenvalues on the principal diagonal. This is fantastic for finding the quadratic form, since the principal diagonal terms are all squares!

$$Q_{\mathbf{A}}(\mathbf{h}) = \lambda_1 u_1^2 + \lambda_2 u_2^2 + \dots + \lambda_d u_d^2$$

$$\mathbf{u} = \mathbf{C}\mathbf{h}$$

\mathbf{C} is the transformation that represents the change in basis

- Another method for finding definiteness, given the eigenvalues is possible. If all are of one sign, then the matrix is definite; if all are of one sign except for a zero somewhere, then the matrix is semidefinite; if some are negative and some are positive, then the matrix is indefinite.

2.17 Feb 27: Eigenspaces

- The eigenspace can be restricted to real vectors only, so only the eigenvectors that are also real count: $E_{\lambda}^{\mathbb{R}}$
- If there is a real λ , even with a complex vector, then $E_{\lambda}^{\mathbb{R}}$ is not the trivial space:

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$$

$$\mathbf{A}(\mathbf{x} + i\mathbf{y}) = \lambda(\mathbf{x} + i\mathbf{y})$$

$$\mathbf{A}\mathbf{x} + i\mathbf{A}\mathbf{y} = \lambda\mathbf{x} + i\lambda\mathbf{y}$$

$$\mathbf{A}\mathbf{x} = \lambda\mathbf{x} \wedge \mathbf{A}\mathbf{y} = \lambda\mathbf{y} \wedge (\mathbf{x} \neq 0 \vee \mathbf{y} \neq 0)$$

- The condition of an eigenvalue is equivalent to $(\lambda \mathbf{I} - \mathbf{A})\mathbf{v} = 0$, and this turns out to imply that $(\lambda \mathbf{I} - \mathbf{A})$ is noninvertible,

$$(\lambda \mathbf{I} - \mathbf{A})^{-1}(\lambda \mathbf{I} - \mathbf{A})\mathbf{v} = (\lambda \mathbf{I} - \mathbf{A})^{-1}0$$

$$\mathbf{v} = 0$$

This is contradictory to the condition!

- The characteristic polynomial of \mathbf{A} is $\det(t\mathbf{I} - \mathbf{A})$ for each t , and this has at most degree d

$$P_{\mathbf{A}}(t) = \det(t\mathbf{I} - \mathbf{A}) = \begin{vmatrix} t - a_{11} & -a_{12} & \dots & -a_{1d} \\ -a_{21} & t - a_{22} & \dots & -a_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ -a_{d1} & -a_{d2} & \dots & t - a_{dd} \end{vmatrix}$$

The determinant is just the sum of all elementary products, and the maximum degree attained is when the principal diagonal is selected. There are no products with degree $d - 1$, since swaps can't perturb only one element. Substitute 0 for t to find the constant term.

$$P_{\mathbf{A}}(t) = t^d - (\text{tr } \mathbf{A})t^{d-1} + \dots + (-1)^d \det \mathbf{A}$$

The other terms are, unfortunately, very complicated. The trace of \mathbf{A} is the sum of the principal entries

- From its definition, the roots of the characteristic polynomial are the eigenvalues, which leads to the fact that the determinant is the product of the eigenvalues, and that the trace is the sum of the eigenvalues

2.18 Feb 28: Numbers Time!

- It is finally time again!

$$\mathbf{A} = \begin{bmatrix} -1 & 3 \\ 2 & 0 \end{bmatrix}$$

Find the eigenvalues!

$$P_{\mathbf{A}}(t) = \det(t\mathbf{I} - \mathbf{A}) = \begin{vmatrix} t+1 & -3 \\ -2 & t \end{vmatrix} = t^2 + t + 6 = (t+3)(t-2)$$

$$\lambda_1 = -3, \lambda_2 = 2$$

Find the eigenspaces!

$$[2\mathbf{I} - \mathbf{A} | 0] = \begin{bmatrix} 3 & -3 & 0 \\ -2 & 2 & 0 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & -1 & 0 \\ -2 & 2 & 0 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & -1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

The 2-eigenspace has basis $(1, 1)$

$$[-3\mathbf{I} - \mathbf{A} | 0] = \begin{bmatrix} -2 & -3 & 0 \\ -2 & -3 & 0 \end{bmatrix} \rightarrow \begin{bmatrix} 2 & 3 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

The -3-eigenspace has basis $(3, -2)$

- To graphically transform any vector, simply change the basis into the eigenspaces. The transformation's effect on the eigenvectors is much easier to see, then at the end change the basis back
- How about a matrix that graphically rotates every single line? Use the counterclockwise 90 degree matrix:

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

$$P_{\mathbf{A}}(t) = t^2 + 1 = (t+i)(t-i)$$

$$\lambda_1 = i, \text{ basis } E_i = (i, 1), \lambda_2 = -i, \text{ basis } E_{-i} = (1, i)$$

The eigenspaces are complex, so clearly drawing the transformation on the real plane would not reveal the eigenspaces!

- What about 3×3 matrices? Any one must have an axis of rotation with eigenvalue 1, intuitively. Because the degree of the characteristic polynomial is odd, one real root must exist, and it will be 1.

2.19 Mar 2: Back to Spectra

- Lofty goals with a symmetric real square matrix:
 - All eigenvectors are real
 - All real eigenspaces are orthogonal
 - The real eigenspaces together can form an ON, or orthonormal, basis of \mathbb{R}^d
- Expanding an arbitrary vector in an ON basis is very easy; a few establishing steps are required!
 - An ON basis is linearly independent:

$$\mathbf{A}\mathbf{t} = \mathbf{0}$$

$$\mathbf{a}_1 \cdot (\mathbf{a}_1 t_1 + \mathbf{a}_2 t_2 + \dots) = 0$$

$$\mathbf{a}_1 \cdot \mathbf{a}_1 t_1 + 0 + 0 + \dots = 0$$

$$|\mathbf{a}_1|^2 t_1 = 0$$

$$t_1 = 0$$

- Then \mathbf{a}_n will span \mathbb{R}^d .

$$\mathbf{v} = \sum t_i \mathbf{a}_i$$

$$\mathbf{a}_1 \cdot \mathbf{v} = |\mathbf{a}_1|^2 t_1 + 0 + 0 + \dots$$

$$t_i = \mathbf{v} \cdot \mathbf{a}_i$$

- The components are now easily found. Now, substitute them into the expansion:

$$\mathbf{v} = (\mathbf{v} \cdot \mathbf{a}_1) \mathbf{a}_1 + (\mathbf{v} \cdot \mathbf{a}_2) \mathbf{a}_2 + \dots$$

$$= \mathbf{a}_1 (\mathbf{a}_1^T \mathbf{v}) + \mathbf{a}_2 (\mathbf{a}_2^T \mathbf{v}) + \dots$$

$$= (\mathbf{a}_1 \mathbf{a}_1^T + \mathbf{a}_2 \mathbf{a}_2^T + \dots) \mathbf{v}$$

$$\mathbf{I} = \mathbf{a}_1 \mathbf{a}_1^T + \mathbf{a}_2 \mathbf{a}_2^T + \dots$$

The reason for this is surprisingly simple. Consider just one transformation represented as $\mathbf{a}\mathbf{a}^T$:

$$\mathbf{v} \mapsto \mathbf{a}\mathbf{a}^T \mathbf{v} = (\mathbf{v} \cdot \mathbf{a}) \mathbf{a} = |\mathbf{v}| |\mathbf{a}| \cos \angle(\mathbf{v}, \mathbf{a}) \mathbf{a}$$

This is just the projection of \mathbf{v} on \mathbf{a} !

- As more $\mathbf{a}\mathbf{a}^T$ terms are summed together in the transformation, the projection becomes onto the plane spanned by the two vectors, then the 3-space, ..., eventually the projection of \mathbf{v} on the space it resides in is reached, and that is just \mathbf{v} , so the transformation is the identity
- Expanding in a non-ON basis is harder:

$$\mathbf{v} = \mathbf{A}\mathbf{t}$$

$$\mathbf{A}^{-1} \mathbf{v} = \mathbf{A}^{-1} \mathbf{A} \mathbf{t}$$

The requirement of an inverse matrix makes the task very complicated

2.20 Mar 3: Still on Spectra

- An equivalent statement of the Spectral Theorem for real symmetric \mathbf{A} is:

There are d real eigenvalues, possibly repeating

$$\lambda_i \neq \lambda_j \Rightarrow E_{\lambda_i}^{\mathbb{R}} \perp E_{\lambda_j}^{\mathbb{R}}$$

$\dim(E_{\lambda_i}^{\mathbb{R}}) = m(\lambda_i)$, where m represents multiplicity

- The Gram-Schmidt Process is an algorithm that takes d linearly independent vectors spanning \mathbb{R}^d , rotates them to be orthogonal, and then rescales to normality, with the ending ON vectors still spanning \mathbb{R}^d , so they are an ON basis
- From this, each eigenspace has an ON basis. Take the union of all the bases; since all the eigenspaces are orthogonal, all the vectors will be too. They are also obviously all normal. Each multiple of an eigenvalue has one vector associated, and there are necessarily d multiples, so there will be d vectors, all ON, which now must be a basis of \mathbb{R}^d . Equivalently, $\mathbf{v}_i \cdot \mathbf{v}_j = \delta_{ij}$
- A great theorem concerning square matrices whose columns are all ON:

$$[\mathbf{S}^T \mathbf{S}]_{ij} = \sum_l [\mathbf{v}_i]_l [\mathbf{v}_j]_l = \mathbf{v}_i \cdot \mathbf{v}_j = \delta_{ij}$$

$$\mathbf{S}^T \mathbf{S} = \mathbf{I}$$

The inverse of \mathbf{S} must exist, since taking determinants on both sides gives $\det \mathbf{S}^2 = 1$, and there is only one matrix it can be!

$$\mathbf{S}^T = \mathbf{S}^{-1}$$

2.21 Mar 5*: More Numbers Time

- Let $\mathbf{A} = \begin{bmatrix} 4 & 2 & 2 \\ 2 & 4 & 2 \\ 2 & 2 & 4 \end{bmatrix}$

- The characteristic polynomial:

$$\begin{aligned} P_{\mathbf{A}}(t) &= t - 4(t^2 - 8t + 16 - 4) + 2(-2t + 8 - 4) - 2(4 + 2t - 8) \\ &= t^3 - 12t^2 + 36t - 32 \\ &= (t - 2)(t - 2)(t - 8) \end{aligned}$$

Thus the eigenvalues are $\{2, 8\}$

- Bases of the eigenspaces:

$$E_2^{\mathbb{R}} : \begin{bmatrix} -2 & -2 & -2 & 0 \\ -2 & -2 & -2 & 0 \\ -2 & -2 & -2 & 0 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 1 & 1 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

One basis is $\{(1, -1, 0), (-1, 1, 0)\}$

$$E_8^{\mathbb{R}} : \begin{bmatrix} 4 & -2 & -2 & 0 \\ -2 & 4 & -2 & 0 \\ -2 & -2 & 4 & 0 \end{bmatrix} \rightarrow \begin{bmatrix} 4 & -2 & -2 & 0 \\ -2 & 4 & -2 & 0 \\ 0 & -6 & 6 & 0 \end{bmatrix} \rightarrow \begin{bmatrix} 2 & -1 & -1 & 0 \\ 0 & 3 & -3 & 0 \\ 0 & -2 & 2 & 0 \end{bmatrix} \rightarrow \begin{bmatrix} 2 & -1 & -1 & 0 \\ 0 & 1 & -1 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

One basis is $\{(1, 1, 1)\}$

- Not coincidentally, the multiplicity of the eigenvalue determines the dimension of its eigenspace
- Gram-Schmidt Processing:

$$E_2^{\mathbb{R}} : (1, -1, 0) \cdot (x, y, -x - y) = 0 = x - y \Rightarrow (1, 1, -2)$$

$$\{(1, -1, 0), (1, 1, -2)\} \xrightarrow{|\mathbf{v}|=1} \left\{ \left(\frac{1}{\sqrt{2}}, \frac{-1}{\sqrt{2}}, 0 \right), \left(\frac{1}{\sqrt{6}}, \frac{1}{\sqrt{6}}, \frac{-2}{\sqrt{6}} \right) \right\}$$

$$E_8^{\mathbb{R}} : \{(1, 1, 1)\} \xrightarrow{|\mathbf{v}|=1} \left\{ \left(\frac{1}{\sqrt{3}}, \frac{1}{\sqrt{3}}, \frac{1}{\sqrt{3}} \right) \right\}$$

- These three vectors serve as an ON basis of \mathbb{R}^d !

2.22 Mar 6*: 0

- Idleness is simply the opposite of labor, not the opposite of work!

2.23 Mar 9: Numbers Review!

- The problem was reviewed

2.24 Mar 10: Fake Projectors

- Projection on a 1-space l is simple enough:

$$\text{proj}(\mathbf{v}|\hat{l}) = \mathbf{p}$$

There are two properties of a projection:

$$\mathbf{p} = t\hat{l}$$

$$(\mathbf{v} - \mathbf{p}) \perp \hat{l}$$

Derive from here:

$$(\mathbf{v} - \mathbf{p}) \cdot \hat{l} = 0$$

$$\mathbf{v} \cdot \hat{l} - \mathbf{p} \cdot \hat{l} = 0$$

$$\mathbf{v} \cdot \hat{l} = t\hat{l} \cdot \hat{l}$$

$$t = \frac{\mathbf{v} \cdot \hat{l}}{|\hat{l}|^2}$$

$$\mathbf{p} = \frac{\mathbf{v} \cdot \hat{l}}{|\hat{l}|^2} \hat{l}$$

- This definition is independent from \hat{l} , so long as it remains in the same direction

$$\frac{\mathbf{v} \cdot \hat{l}}{|\hat{l}|^2} \hat{l} = \frac{k^2}{k^2} \frac{\mathbf{v} \cdot \hat{l}}{|\hat{l}|^2} \hat{l} = \frac{\mathbf{v} \cdot k\hat{l}}{|k\hat{l}|^2} k\hat{l}$$

- What about projection onto higher spaces S with dimension $k \leq d$? It is now convenient to define the altitude vector along with the projection vector:

$$\mathbf{p} = \text{proj}(\mathbf{v}|S)$$

$$\mathbf{a} = \text{alt}(\mathbf{v}|S)$$

Each pair has the same properties:

$$\mathbf{p} + \mathbf{a} = \mathbf{v}$$

$$\mathbf{p} \in S, \mathbf{a} \perp S$$

- There exists a unique pair for every vector and every space.
- Uniqueness: Assume two pairs existed, $(\mathbf{p}, \mathbf{a}), (\mathbf{q}, \mathbf{b})$:

$$\mathbf{p} + \mathbf{a} = \mathbf{v} = \mathbf{q} + \mathbf{b}$$

$$S \ni \mathbf{p} - \mathbf{q} = \mathbf{b} - \mathbf{a} \perp S$$

$$\mathbf{p} - \mathbf{q} \perp \mathbf{p} - \mathbf{q}$$

$$\mathbf{p} - \mathbf{q} = 0$$

- Existence: The space S must have an ON basis $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k\}$. To find the projection on a higher space, just sum each projection onto the basis vectors:

$$\mathbf{p} = \sum_n \text{proj}(\mathbf{v}|\mathbf{u}_n)$$

2.25 Mar 11: Megagram-Schmidt

- Every nontrivial $S \subseteq \mathbb{R}^d$ has an ON basis
- Induction on $k = \dim S$, start with $k = 1$, where the space is the span of a single vector \mathbf{w} . The ON basis is simply $\{\frac{\mathbf{w}}{|\mathbf{w}|}\}$
- Assume every space of dimension $k - 1$, T , has an ON basis. Consider S , dimension k , with some basis:

$$\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\}$$

Let T be the span of the first $k - 1$, then those are also a basis of T . By assumption, an ON basis of T exists,

$$\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_{k-1}\}$$

- The Gram-Schmidt Formula gives a \mathbf{u}_k ,

$$\mathbf{u}_k = \frac{\mathbf{v}_k - \sum_{i=1}^{k-1} (\mathbf{v}_k \cdot \mathbf{u}_i) \mathbf{u}_i}{|\mathbf{v}_k - \sum_{i=1}^{k-1} (\mathbf{v}_k \cdot \mathbf{u}_i) \mathbf{u}_i|}$$

- This new \mathbf{u}_k is clearly unitary, but is it orthogonal to the other ON \mathbf{u} ?

$$\mathbf{u}_k \cdot \mathbf{u}_j = \mathbf{v}_k \cdot \mathbf{u}_j - \sum_{i=1}^{k-1} (\mathbf{v}_k \cdot \mathbf{u}_i) (\mathbf{u}_i \cdot \mathbf{u}_j)$$

Sift! The dot product of ON vectors is collapsible:

$$= \mathbf{v}_k \cdot \mathbf{u}_j - (\mathbf{v}_k \cdot \mathbf{u}_j) = 0$$

Success! The \mathbf{u} vectors are ON!

- It remains to be seen whether the \mathbf{u} vectors are a basis for S . The following compound statement suffices:

$$S \subseteq \text{span}(\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k) \subseteq S$$

The first statement; show that each \mathbf{u} vector is in S . The first $k - 1$ are in $T \subseteq S$, and the last, from the formula, is a linear combination of \mathbf{v}_k and \mathbf{u} vectors

The second statement:

$$\mathbf{x} \in S = \sum a_i \mathbf{v}_i = \mathbf{y} + a_k \mathbf{v}_k = \sum a_i \mathbf{u}_i + \ell \mathbf{u}_k + \sum b_i \mathbf{u}_i \in \text{span}(\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k)$$

2.26 Mar 12: Real Projectors

- Review of subspaces! Two special spaces are defined, the trivial space $\mathbf{0}$ and the entire space $\mathbf{1}$
- The operators inclusion, $S \subseteq T$, intersection $S \cap T$, and addition $S + T$, are from before, and the unary orthocomplement is

$$S^\perp = \{\mathbf{v} | \mathbf{v} \perp S\}$$

The fundamental inclusion relations:

$$\mathbf{0} \subseteq S \cap T \subseteq S, T \subseteq S + T \subseteq \mathbf{1}$$

The intersection and sum of the two spaces is the nearest space to the two in this inclusion; no spaces are in between

2.27 Mar 13: in Absentio

- Order has begun to break down!

2.28 Mar 16: The Far Reaches

- The start of the illustrious ~~YouTube Show Me~~ YouTube career of Joseph Stern!

2.29 Mar 19: Spectral Thm Proved

- Orthogonal Matrices are the name for those matrices whose columns are ON
- The Spectral Thm guarantees too that $\exists \mathbf{S}$ orthogonal, where $\mathbf{S}^{-1}\mathbf{A}\mathbf{S}$ is the diagonal matrix \mathbf{D} with the entries the eigenvalues of \mathbf{A}
- Additionally, \exists projector matrices $\mathbf{P}_1, \mathbf{P}_2, \dots, \mathbf{P}_d : \mathbf{P}_i^2 = \mathbf{P}_i = \mathbf{P}_i^T, \mathbf{P}_i\mathbf{P}_j = \mathbf{P}_j\mathbf{P}_i = \mathbf{O}$ when $i \neq j$, and finally, $\mathbf{A} = \sum_j \lambda_j \mathbf{P}_j$, the spectral decomposition. The projectors project onto the columns of \mathbf{S} from above
- To start, tackle $\sigma(\mathbf{A}) \subseteq \mathbb{R}$ with a complex inner product, like dot products but one factor is conjugated (because of this the associated norm is essentially the same):

$$\langle \mathbf{v} | \mathbf{A} \mathbf{v} \rangle$$

This inner product is also the dot product of the adjoint of one vector and the other $\mathbf{v}^\dagger \mathbf{w}$, which relies on the definition of the matrix adjoint: $\mathbf{M}^\dagger = \overline{\mathbf{M}^T}$

$$= \mathbf{v}^\dagger \mathbf{A} \mathbf{v} = \mathbf{v}^\dagger \lambda \mathbf{v} = \lambda \mathbf{v}^\dagger \mathbf{v} = \lambda |\mathbf{v}|^2$$

Use the associative property and adjoints of products on the expression, though, and get a different expression! Remember \mathbf{A} is real symmetric

$$= \mathbf{v}^\dagger \mathbf{A} \mathbf{v} = (\mathbf{A}^\dagger \mathbf{v})^\dagger \mathbf{v} = (\mathbf{A} \mathbf{v})^\dagger \mathbf{v} = (\lambda \mathbf{v})^\dagger \mathbf{v} = \bar{\lambda} |\mathbf{v}|^2$$

The eigenvector is never 0, then $\lambda = \bar{\lambda}$, and the eigenvalue is real.

- Next target: Orthogonality of distinct real eigenspaces. Let \mathbf{v} and \mathbf{w} be elements of the two eigenspaces, and λ and μ the eigenvalues

$$\langle \mathbf{w} | \mathbf{A} \mathbf{v} \rangle = \mathbf{w}^\dagger \mathbf{A} \mathbf{v} = \mathbf{w}^\dagger \lambda \mathbf{v} = \lambda (\mathbf{w}^\dagger \mathbf{v})$$

But also

$$\mathbf{w}^\dagger \mathbf{A} \mathbf{v} = (\mathbf{A}^\dagger \mathbf{w})^\dagger \mathbf{v} = (\mathbf{A} \mathbf{w})^\dagger \mathbf{v} = (\mu \mathbf{w})^\dagger \mathbf{v} = \bar{\mu} (\mathbf{w}^\dagger \mathbf{v})$$

Now, $\lambda \langle \mathbf{w} | \mathbf{v} \rangle = \bar{\mu} \langle \mathbf{w} | \mathbf{v} \rangle = \mu \langle \mathbf{w} | \mathbf{v} \rangle$, but the eigenvalues are assumed to be different! Then the inner product is 0, the dot product is too (since the vectors are real anyway) and \mathbf{w} and \mathbf{v} are orthogonal

- The multiplicity of an eigenvalue in the characteristic polynomial is equal to the dimension of the eigenspace: r will stand for $\dim E_\lambda^{\mathbb{R}}(\mathbf{A})$, between 1 and d , and from Gram-Schmidt construct a ON basis $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r\}$ for the eigenspace.
Extend the basis to the entire space, to $\{\mathbf{u}_1, \dots, \mathbf{u}_r, \mathbf{u}_{r+1}, \dots, \mathbf{u}_d\}$, also ON by Gram-Schmidt. Recall for the early vectors

$$\mathbf{A} \mathbf{u}_i = \lambda \mathbf{u}_i$$

This is not true for the later \mathbf{u} vectors anymore, resort to the larger basis;

$$\mathbf{A} \mathbf{u}_i = \sum_{j=1}^r b_j \mathbf{u}_j + \sum_{j=r+1}^d c_j \mathbf{u}_j$$

However, ON-ality still gives

$$\mathbf{u}_i \cdot \mathbf{u}_j$$

Two parts now: if a is the multiplicity of λ as a root, $a \geq r$ and $a \neq r$

2.30 Mar 20: Spectral Thm Proved II

- Let \mathbf{M} be the matrix of ON columns \mathbf{u}_i , it is orthogonal ($\in O(d)$) i.e. $\mathbf{M}^T = \mathbf{M}^{-1}$, and proof:

$$\mathbf{M}^T \mathbf{M} = \begin{bmatrix} - & \mathbf{u}_1^T & - \\ - & \mathbf{u}_2^T & - \\ & \vdots & \\ - & \mathbf{u}_d^T & - \end{bmatrix} \begin{bmatrix} | & | & & | \\ \mathbf{u}_1 & \mathbf{u}_2 & \dots & \mathbf{u}_d \\ | & | & & | \end{bmatrix} = \begin{bmatrix} \mathbf{u}_1 \cdot \mathbf{u}_1 & \mathbf{u}_1 \cdot \mathbf{u}_2 & \dots & \mathbf{u}_1 \cdot \mathbf{u}_d \\ \mathbf{u}_2 \cdot \mathbf{u}_1 & \mathbf{u}_2 \cdot \mathbf{u}_2 & \dots & \mathbf{u}_2 \cdot \mathbf{u}_d \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{u}_d \cdot \mathbf{u}_1 & \mathbf{u}_d \cdot \mathbf{u}_2 & \dots & \mathbf{u}_d \cdot \mathbf{u}_d \end{bmatrix} = \mathbf{I}$$

The other side inverse is only a little more complicated:

- If $\mathbf{AB} = \mathbf{I}$ and \mathbf{A} and \mathbf{B} are square, then $\mathbf{BA} = \mathbf{I}$.

- Let S be the space of $d \times d$ matrices, then $S \supseteq \mathbf{B}S \supseteq \mathbf{B}^2S \supseteq \dots \supseteq \mathbf{B}^kS \supseteq \dots$. This is because $\mathbf{B}^{k+1}\mathbf{C} = \mathbf{B}^k\mathbf{BC}$, and $\mathbf{BC} \in S$.
- S is also isomorphic to the d^2 -vector space, since the operations are essentially the same, so it is possible to say that \mathbf{B}^kS is a space too.
- Spaces must have a nonnegative integer dimension, and since there are only d^2 in the biggest space, there can only be that many distinct spaces in the progression; $\exists l : \mathbf{B}^lS = \mathbf{B}^{l+1}S$

$$\mathbf{B}^l \in \mathbf{B}^lS = \mathbf{B}^{l+1}S$$

$$\mathbf{B}^l = \mathbf{B}^{l+1}\mathbf{C}$$

$$\mathbf{A}^k\mathbf{B}^k = \mathbf{A}^k\mathbf{B}^k\mathbf{BC}$$

$$\mathbf{I} = \mathbf{IBC} = \mathbf{BC}$$

- \mathbf{B} is square and has one left-inverse and one right-inverse, then both inverses are \mathbf{A}

- Now that \mathbf{M} is orthogonal,

$$\mathbf{AM} = \left[\begin{array}{c|c|c|c|c|c} | & | & | & | & | & | \\ \mathbf{Au}_1 & \dots & \mathbf{Au}_r & \mathbf{Au}_{r+1} & \dots & \mathbf{Au}_d \\ | & | & | & | & | & | \end{array} \right] = \left[\begin{array}{c|c|c|c|c|c} | & | & | & | & | & | \\ \lambda\mathbf{u}_1 & \dots & \lambda\mathbf{u}_r & \mathbf{Au}_{r+1} & \dots & \mathbf{Au}_d \\ | & | & | & | & | & | \end{array} \right]$$

Recall that the column space of \mathbf{M} now represents the entire space,

$$\begin{aligned} &= \left[\begin{array}{c|c|c|c|c|c} | & | & | & | & | & | \\ \lambda\mathbf{Me}_1 & \dots & \lambda\mathbf{Me}_r & \mathbf{Mc}_{r+1} & \dots & \mathbf{Mc}_d \\ | & | & | & | & | & | \end{array} \right] = \mathbf{M} \left[\begin{array}{c|c|c|c|c|c} | & | & | & | & | & | \\ \lambda\mathbf{e}_1 & \dots & \lambda\mathbf{e}_r & \mathbf{c}_{r+1} & \dots & \mathbf{c}_d \\ | & | & | & | & | & | \end{array} \right] \\ &= \mathbf{M} \left[\begin{array}{c|c} \lambda\mathbf{I}_r & | \\ \hline \mathbf{O} & \mathbf{C} \end{array} \right] \end{aligned}$$

Now investigate:

$$\mathbf{M}^T\mathbf{AM} = \mathbf{M}^{-1}\mathbf{AM} = \left[\begin{array}{c|c} \lambda\mathbf{I}_r & | \\ \hline \mathbf{O} & \mathbf{C} \end{array} \right]$$

But also this is a symmetric matrix,

$$(\mathbf{M}^T\mathbf{AM})^T = \mathbf{M}^T\mathbf{A}^T\mathbf{M} = \mathbf{M}^T\mathbf{AM}$$

$$\left[\begin{array}{c|c} \lambda\mathbf{I}_r & | \\ \hline \mathbf{O} & \mathbf{C} \end{array} \right] = \left[\begin{array}{c|c} \lambda\mathbf{I}_r & | \\ \hline \mathbf{O} & \mathbf{D} \end{array} \right]$$

- Let $\mathbf{B} = \mathbf{M}^T\mathbf{AM}$,

$$\begin{aligned} P_{\mathbf{B}}(t) &= \det(t\mathbf{I} - \mathbf{B}) \\ &= \det(t\mathbf{M}^T\mathbf{M} - \mathbf{M}^T\mathbf{AM}) \\ &= \det(\mathbf{M}^T(t\mathbf{I} - \mathbf{A})\mathbf{M}) \\ &= \det \mathbf{M}^T \det(t\mathbf{I} - \mathbf{A}) \det \mathbf{M} \end{aligned}$$

(Determinants distribute over multiplication)

$$= P_{\mathbf{A}}(t)$$

Anytime $\mathbf{B} = \mathbf{S}^{-1}\mathbf{AS}$, then the two matrices are similar and the characteristic polynomials are exactly the same

2.31 Mar 23: Spectral Thm Proved III (for real)

- Return to $P_{\mathbf{B}}(t)$

$$\begin{aligned} P_{\mathbf{A}}(t) &= P_{\mathbf{B}}(t) = \det(t\mathbf{I} - \mathbf{B}) \\ &= \left| \begin{array}{c|c} (t-\lambda)\mathbf{I}_r & \mathbf{O} \\ \hline \mathbf{O} & t\mathbf{I} - \mathbf{D} \end{array} \right| = (t-\lambda)^r \det(t\mathbf{I} - \mathbf{D}) \\ P_{\mathbf{A}}(t) &= (t-\lambda)^r P_{\mathbf{D}}(t) \end{aligned}$$

Then the multiplicity of λ is at least r since that many can be factored out, and $a \geq r$

- Now suppose $a > r$, then λ is also a root of $P_{\mathbf{D}}(t)$, or $P_{\mathbf{D}}(\lambda) = 0$, $\lambda \in \sigma(\mathbf{D})$.

$$\exists \mathbb{R}^{d-r} \ni \mathbf{v} \neq 0 : \mathbf{D}\mathbf{v} = \lambda\mathbf{v}$$

Look at $\mathbf{w} = (0, \mathbf{v}) \in \mathbb{R}^d$, made by concatenating a zero vector

$$\mathbf{B}\mathbf{w} = (0, \mathbf{D}\mathbf{v}) = (\lambda 0, \lambda \mathbf{v}) = \lambda \mathbf{w}$$

\mathbf{w} is another eigenvector of B ! i.e. $\mathbf{M}^T \mathbf{A} \mathbf{M} \mathbf{w} = \lambda \mathbf{w}$

Let $\mathbf{z} = \mathbf{M}\mathbf{w}$, this is also not 0,

$$\begin{aligned} \mathbf{M}^T \mathbf{A} \mathbf{M} \mathbf{w} &= \lambda \mathbf{w} \\ \mathbf{M} \mathbf{M}^T \mathbf{A} \mathbf{M} \mathbf{w} &= \mathbf{M} \lambda \mathbf{w} = \lambda \mathbf{M} \mathbf{w} \end{aligned}$$

\mathbf{z} is also an eigenvector of \mathbf{A} , then it is expressible in the basis of the eigenspace,

$$\begin{aligned} \mathbf{z} &= \sum_{j=1}^r a_j \mathbf{u}_j \\ \mathbf{w} = \mathbf{M}^{-1} \mathbf{z} &= \mathbf{M}^T \mathbf{z} = \sum_{j=1}^r a_j (\mathbf{M}^T \mathbf{u}_j) \end{aligned}$$

The summed terms are no strangers, however:

$$[\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_d] = \mathbf{I} = \mathbf{M}^T \mathbf{M} = \mathbf{M}^T [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_d] = [\mathbf{M}^T \mathbf{u}_1, \mathbf{M}^T \mathbf{u}_2, \dots, \mathbf{M}^T \mathbf{u}_d]$$

$$\mathbf{w} = \sum_{j=1}^r a_j \mathbf{e}_j = (a_1, a_2, \dots, a_r, 0, 0, \dots, 0)$$

This does not at all satisfy the definition of \mathbf{w} as $(0, 0, \dots, 0, \mathbf{v})$! Then $a \not\geq r$, or $a = r$

2.32 Mar 24: Orthogonal Maps

- What is the linear transformation of an orthogonal matrix \mathbf{M} ? All of the standard basis vectors stay ON, so it is just a rotation and possibly reflection of space, intuitively; now prove it
- To prove centrality: $\mathbf{M}\mathbf{0} = \mathbf{0}$; all linear maps are central
- To prove dot-preservation:

$$(\mathbf{M}\mathbf{v}) \cdot (\mathbf{M}\mathbf{w}) = (\mathbf{M}\mathbf{v})^T (\mathbf{M}\mathbf{w}) = \mathbf{v}^T \mathbf{M}^T \mathbf{M} \mathbf{w} = \mathbf{v}^T \mathbf{w} = \mathbf{v} \cdot \mathbf{w}$$

- Using that, to prove isometry:

$$\begin{aligned} |\mathbf{M}\mathbf{v}|^2 &= (\mathbf{M}\mathbf{v}) \cdot (\mathbf{M}\mathbf{v}) = \mathbf{v} \cdot \mathbf{v} = |\mathbf{v}|^2 \\ |\mathbf{M}\mathbf{p} - \mathbf{M}\mathbf{q}| &= |\mathbf{M}(\mathbf{p} - \mathbf{q})| = |\mathbf{p} - \mathbf{q}| \end{aligned}$$

- To prove a determinant of unity (or its negative):

$$\det \mathbf{I} = \det(\mathbf{M}^T \mathbf{M}) = \det \mathbf{M}^T \det \mathbf{M} = (\det \mathbf{M})^2$$

$(\det \mathbf{A} = \det \mathbf{A}^T)$ \mathbf{M} can flip the space, or it might not; in this way it is either a rotation or a rotation and a reflection

2.33 Mar 25: Spectral Thm IV (we'll get them this time!)

- The eigenvectors of \mathbf{A} can make an ON basis of \mathbb{R}^d . Start with the characteristic polynomial

$$P_{\mathbf{A}}(t) = (t - \lambda_1)^{m_1}(t - \lambda_2)^{m_2} \dots (t - \lambda_k)^{m_k}$$

By an earlier part of the Spectral Thm, m_n is also the dimension of each eigenspace

- Now take a list of d eigenvectors, such that each set of m_n eigenvectors forms a ON basis of an eigenspace. Each eigenspace is orthogonal so the whole list is ON and linearly independent. Of course, d linearly independent vectors that all reside in \mathbb{R}^d must also span it!
- \mathbf{A} is orthogonally diagonalizable. Let \mathbf{S} have columns the ON basis discovered above. \mathbf{S} is now orthogonal, and is a convenient change of basis matrix!

$$\begin{aligned}\mathbf{AS} &= [\mathbf{A}\mathbf{v}_1 \quad \mathbf{A}\mathbf{v}_2 \quad \dots \quad \mathbf{A}\mathbf{v}_d] \\ &= [\lambda_1\mathbf{v}_1 \quad \lambda_2\mathbf{v}_2 \quad \dots \quad \lambda_d\mathbf{v}_d] \\ &= [\lambda_1\mathbf{S}\mathbf{e}_1 \quad \lambda_2\mathbf{S}\mathbf{e}_2 \quad \dots \quad \lambda_d\mathbf{S}\mathbf{e}_d] \\ &= \mathbf{S} [\lambda_1\mathbf{e}_1 \quad \lambda_2\mathbf{e}_2 \quad \dots \quad \lambda_d\mathbf{e}_d]\end{aligned}$$

$$\mathbf{AS} = \mathbf{SD}$$

$$\mathbf{S}^{-1}\mathbf{AS} = \mathbf{D}$$

\mathbf{D} is the fabled diagonal matrix with the elements the eigenvalues of \mathbf{A} - finally! Since \mathbf{S} is orthogonal, \mathbf{A} is further orthogonally diagonalizable.

- \mathbf{A} is the sum of each projector matrix onto an eigenvector multiplied by its eigenvalue

$$\begin{aligned}\mathbf{Ax} &= \mathbf{SDS}^T\mathbf{x} = \mathbf{S} \begin{bmatrix} \lambda_1\mathbf{v}_1 \cdot \mathbf{x} \\ \lambda_2\mathbf{v}_2 \cdot \mathbf{x} \\ \vdots \\ \lambda_d\mathbf{v}_d \cdot \mathbf{x} \end{bmatrix} \\ &= \mathbf{S} (\lambda_1(\mathbf{v}_1 \cdot \mathbf{x})\mathbf{e}_1 + \lambda_2(\mathbf{v}_2 \cdot \mathbf{x})\mathbf{e}_2 + \dots + \lambda_d(\mathbf{v}_d \cdot \mathbf{x})\mathbf{e}_d) \\ &= \lambda_1(\mathbf{v}_1 \cdot \mathbf{x})\mathbf{v}_1 + \lambda_2(\mathbf{v}_2 \cdot \mathbf{x})\mathbf{v}_2 + \dots + \lambda_d(\mathbf{v}_d \cdot \mathbf{x})\mathbf{v}_d \\ &\quad (\lambda_1\mathbf{v}_1\mathbf{v}_1^T + \lambda_2\mathbf{v}_2\mathbf{v}_2^T + \dots + \lambda_d\mathbf{v}_d\mathbf{v}_d^T)\mathbf{x} \\ &= (\lambda_1P(\mathbf{v}_1) + \lambda_2P(\mathbf{v}_2) + \dots + \lambda_dP(\mathbf{v}_d))\mathbf{x}\end{aligned}$$

- It is done!
- Now it all comes together, The Principal Axes Thm: Given familiar \mathbf{A} , there exists a rotation \mathbf{R} such that

$$Q_{\mathbf{A}}(\mathbf{x}) = \mathbf{x}^T\mathbf{Ax} = \sum_{j=1}^d \lambda_j[\mathbf{Rx}]_j^2 = Q_{\mathbf{D}}(\mathbf{Rx})$$

$Q_{\mathbf{D}}(\mathbf{x})$ is much much much easier to compute than the other quadratic form!

$$\mathbf{x}^T\mathbf{Ax} = \mathbf{x}^T\mathbf{SDS}^T\mathbf{x} = (\mathbf{S}^T\mathbf{x})^T\mathbf{D}(\mathbf{S}^T\mathbf{x}) = Q_{\mathbf{D}}(\mathbf{S}^T\mathbf{x})$$

If \mathbf{S} is positively directed, then it is the rotation we are looking for.

$$\mathbf{R}^T = (\mathbf{S}^T)^T = (\mathbf{S}^T)^{-1} = \mathbf{R}^{-1}$$

Otherwise, just modify \mathbf{D} slightly so that one column swap is performed, then \mathbf{S} has a column swapped too, and the determinant will be negated. The resultant \mathbf{R} just has one basis vector swapped with another.

- Lastly, to prove the definiteness arising from eigenvalues; start by converting $Q_{\mathbf{A}}(\mathbf{x})$ to $Q_{\mathbf{D}}(\mathbf{R}\mathbf{x})$ to $\sum_{j=1}^d \lambda_j [\mathbf{R}\mathbf{x}]_j^2$

If $\mathbf{x} \neq 0$ then $\mathbf{R}\mathbf{x} \neq 0$ by rotation;

Assuming all the eigenvalues are positive,

$$\sum_{j=1}^d \lambda_j [\mathbf{R}\mathbf{x}]_j^2 > 0$$

All negative,

$$\sum_{j=1}^d \lambda_j [\mathbf{R}\mathbf{x}]_j^2 < 0$$

Some negative and positive,

$$\sum_{j=1}^d \lambda_j [\mathbf{R}\mathbf{x}]_j^2 \in (-\infty, \infty)$$

Not less than 0 is a bit trickier, let $\lambda_j = 0$

$$\mathbf{x} = \mathbf{R}^T \mathbf{e}_j = \mathbf{R}^{-1} \mathbf{e}_j$$

$$Q_{\mathbf{D}}(\mathbf{R}\mathbf{x}) = Q_{\mathbf{D}}(\mathbf{e}_j) = \mathbf{e}_j^T \mathbf{D} \mathbf{e}_j = \mathbf{e}_j^T \lambda_j \mathbf{e}_j = \lambda_j |\mathbf{e}_j|^2 = 0$$

$$\sum_{j=1}^d \lambda_j [\mathbf{R}\mathbf{x}]_j^2 \geq 0$$

Not greater than 0,

$$\sum_{j=1}^d \lambda_j [\mathbf{R}\mathbf{x}]_j^2 \leq 0$$

2.34 Mar 30: Back to Sylvester

- To reiterate, positive definite matrices have all subdeterminants positive, negative definite matrices have negative-positive alternating subdeterminants, indefinite matrices have subdeterminants that are almost like definite matrices, but the first subdeterminant that messes up the pattern has the wrong sign, and semidefinite matrices are the rest, i.e. when the first deviant subdeterminant is zero.
- In the 2×2 case, a negative determinant always implies indefiniteness, and a zero determinant implies semidefiniteness
- Positive definiteness:

$$\begin{aligned} P_{\mathbf{A}}(t) &= t^2 - (a+c)t + ac - b^2 \\ \lambda &= \frac{a+c \pm \sqrt{a^2 + 2ac + c^2 - 4ac + 4b^2}}{2} \\ &= \frac{a+c \pm \sqrt{(a-c)^2 + 4b^2}}{2} \end{aligned}$$

Symmetry grants the discriminant positiveness. For both eigenvalues to be the same, $a = c$ and $b = 0$, making the whole matrix $a\mathbf{I}$, or a uniform scaling

Consider eigenvalue signs: Positive definiteness requires the lesser eigenvalue be positive:

$$\begin{aligned} \sqrt{(a-c)^2 + 4b^2} &< a+c \Rightarrow \\ \Leftrightarrow (a-c)^2 + 4b^2 &< (a+c)^2 \\ 4b^2 &< 4ac \\ ac - b^2 &> 0 \end{aligned}$$

The complete reverse requires an additional condition

$$(a-c)^2 + 4b^2 < (a+c)^2 \Rightarrow \sqrt{(a-c)^2 + 4b^2} < |a+c| \stackrel{a \geq 0}{\underset{(ac \geq 0)}} \sqrt{(a-c)^2 + 4b^2} < |a+c| < a+c$$

Negative definiteness is the same

- Indefiniteness requires that the eigenvalues take opposite signs

$$-\sqrt{(a-c)^2 + 4b^2} < a+c < \sqrt{(a-c)^2 + 4b^2}$$

$$\sqrt{(a-c)^2 + 4b^2} > |a+c|$$

Simply follow from the above case, except with switched inequality

- It's Numbers time again!

$$f(x, y) = 3x^2 + y^2 - x + 4y - 5$$

$$\nabla f = 0 = (6x - 1, 2y + 4)$$

$$\text{Critical point at } \left(\frac{1}{6}, -2\right)$$

$$\mathbf{H} = \begin{bmatrix} 6 & 0 \\ 0 & 2 \end{bmatrix}$$

$$\det \mathbf{H} > 0$$

The critical point represents a minimum

$$g(x, y) = 2x^2 + y^2 - 3xy + x + 2y + 4$$

$$\nabla g = 0 = (4x - 3y + 1, 2y - 3x + 2)$$

$$\begin{bmatrix} 4 & -3 & -1 \\ -3 & 2 & -2 \end{bmatrix} \rightarrow \begin{bmatrix} 12 & -9 & -3 \\ -12 & 8 & -8 \end{bmatrix} \rightarrow \begin{bmatrix} 0 & -1 & -11 \\ -12 & 8 & -8 \end{bmatrix}$$

$$\text{Critical point at } (8, 11)$$

$$\mathbf{H} = \begin{bmatrix} 4 & -3 \\ -3 & 2 \end{bmatrix}$$

$$\det \mathbf{H} < 0$$

The critical point represents a saddle

2.35 Apr 1: Sylvester Sylvester Sylvester Sylvester Sylvester

- Unsurprisingly, negation also reverses definiteness, since the eigenvalues are negated. Alternatively, odd sub-determinants are negated too.
- Sylvester's positive definite criterion for symmetric $d \times d$ matrices! Let $\Delta_n(\mathbf{A})$ be the n^{th} principal subdeterminant

$$\mathbf{A} = \begin{bmatrix} \mathbf{B} & \mathbf{c} \\ \mathbf{c}^T & r \end{bmatrix} = \mathbf{A}^T = \begin{bmatrix} \mathbf{B}^T & \mathbf{c} \\ \mathbf{c}^T & r \end{bmatrix}$$

$$\mathbf{z} = (\mathbf{x}, y)$$

$$Q_{\mathbf{A}}(\mathbf{z}) = \begin{bmatrix} \mathbf{x}^T & y \end{bmatrix} \begin{bmatrix} \mathbf{B} & \mathbf{c} \\ \mathbf{c}^T & r \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ y \end{bmatrix} = \begin{bmatrix} \mathbf{x}^T & y \end{bmatrix} \begin{bmatrix} \mathbf{B}\mathbf{x} + \mathbf{c}y \\ \mathbf{c}^T\mathbf{x} + ry \end{bmatrix} = \mathbf{x}^T\mathbf{B}\mathbf{x} + \mathbf{x}^T\mathbf{c}y + y\mathbf{c}^T\mathbf{x} + yry = Q_{\mathbf{B}}(\mathbf{x}) + 2y\mathbf{c} \cdot \mathbf{x} + ry^2$$

Importantly, $\mathbf{A} \succ 0 \Rightarrow \mathbf{B} \succ 0$:

$$\mathbf{z} = (\mathbf{x}, 0) \Rightarrow Q_{\mathbf{A}}(\mathbf{z}) = Q_{\mathbf{B}}(\mathbf{x}) + 0 + 0$$

For any \mathbf{x} this is true

- Start with positive definiteness and prove positive subdeterminants; base case $d = 1$ is biconditional: a positive number is positive definite. Assume the $d - 1$ case,

$$\mathbf{B} \succ 0 \Rightarrow \Delta_{1\dots d-1}(\mathbf{B}) = \Delta_{1\dots d-1}(\mathbf{A}) > 0$$

$$\mathbf{B} \succ 0 \Leftarrow \mathbf{A} \succ 0 \Rightarrow (\forall j) \lambda_j > 0$$

$$\det \mathbf{A} = \Delta_d(\mathbf{A}) = \prod \lambda > 0$$

Induction is complete: $\mathbf{A} \succ 0 \Rightarrow \Delta_{1\dots d} > 0$

- Now, start with positive subdeterminants and prove positive definiteness. The base case is the same; assume the $d - 1$ case

$$\Delta_{1\dots d}(\mathbf{A}) > 0$$

$$\Delta_{1\dots d-1}(\mathbf{B}) > 0 \Rightarrow \mathbf{B} \succ 0$$

\mathbf{A} cannot have one negative eigenvalue, or else not every subdeterminant would be positive: $\det \mathbf{A} \neq 0$. The same prohibits \mathbf{A} from having a 0 eigenvalue

- Suppose \mathbf{A} had two or more negative eigenvalues,

$$\lambda_1, \lambda_2 < 0$$

$$\mathbf{A}\mathbf{v} = \lambda_1\mathbf{v} \quad \mathbf{A}\mathbf{w} = \lambda_2\mathbf{w}$$

Either the eigenvectors are of two different orthogonal eigenspaces or they reside in at least a planar eigenspace, so it is possible to pick \mathbf{v}, \mathbf{w} to have

$$\mathbf{v} \perp \mathbf{w}$$

Contrive:

$$\mathbf{z} = w_d\mathbf{v} - v_d\mathbf{w} = (\mathbf{x}, 0)$$

$$Q_{\mathbf{A}}(\mathbf{z}) = Q_{\mathbf{B}}(\mathbf{x})$$

- Consider \mathbf{B} , then $v_d \neq 0 \wedge w_d \neq 0$. Assume that $v_d = 0$

$$\mathbf{A}\mathbf{v} = \begin{bmatrix} \mathbf{B} & \mathbf{c} \\ \mathbf{c}^T & r \end{bmatrix} \begin{bmatrix} \mathbf{u} \\ v_d \end{bmatrix} = \begin{bmatrix} \mathbf{B}\mathbf{u} + \mathbf{c}v_d \\ \mathbf{c}^T\mathbf{u} + v_d \end{bmatrix} = \begin{bmatrix} \mathbf{B}\mathbf{u} \\ \mathbf{c}^T\mathbf{u} \end{bmatrix}$$

$$\lambda_1\mathbf{v} = \begin{bmatrix} \lambda_1\mathbf{u} \\ 0 \end{bmatrix}$$

$$\mathbf{B}\mathbf{u} = \lambda_1\mathbf{u}$$

But $\mathbf{B} \succ 0$, so it cannot have the negative eigenvalue λ_1 . The same argument holds for w_d . Now, $\mathbf{z} \neq 0$ since it is a nonzero linear combination of linearly independent nonzeros, and \mathbf{x} is also $\neq 0$, so that is good for computing the quadratic forms

- Go all the way back to $Q_{\mathbf{A}}(\mathbf{z}) = Q_{\mathbf{B}}(\mathbf{x}) > 0$

$$\begin{aligned} &= (w_d\mathbf{v}^T - v_d\mathbf{w}^T)\mathbf{A}(w_d\mathbf{v} - v_d\mathbf{w}) = (w_d\mathbf{v}^T - v_d\mathbf{w}^T)(w_d\mathbf{A}\mathbf{v} - v_d\mathbf{A}\mathbf{w}) \\ &= (w_d\mathbf{v}^T - v_d\mathbf{w}^T)(w_d\lambda_1\mathbf{v} - v_d\lambda_2\mathbf{w}) = w_d\mathbf{v}^T w_d\lambda_1\mathbf{v} - v_d\mathbf{w}^T w_d\lambda_1\mathbf{v} - w_d\mathbf{v}^T v_d\lambda_2\mathbf{w} + v_d\mathbf{w}^T v_d\lambda_2\mathbf{w} \\ &= w_d^2\lambda_1|\mathbf{v}|^2 + v_d^2\lambda_2|\mathbf{w}|^2 - v_d w_d(\mathbf{v} \cdot \mathbf{w})(\lambda_1 + \lambda_2) = w_d^2\lambda_1|\mathbf{v}|^2 + v_d^2\lambda_2|\mathbf{w}|^2 < 0 \end{aligned}$$

Contradiction! So \mathbf{A} cannot have two or more negative eigenvalues, it also cannot have one negative eigenvalue, and it cannot have any zero eigenvalues. Therefore, it must be positive definite

2.36 Apr 10: EVT and Friends

- Let $f : D^\circ \rightarrow \mathbb{R}$ be C^0 . If on the boundary of the restriction on S , f does not have an infinite discontinuity, then it must have extrema.
- EVT guarantees that S being compact is good enough; there are other weaker theorems!
- MinVT: the restriction of f on set S has a minimum if $\exists \bar{B}(\mathbf{c}) : \forall \mathbf{x} \in S \setminus \bar{B}, f(\mathbf{x}) \geq f(\mathbf{c})$; the minimum will be some $\mathbf{q} \in \bar{B}$
This condition is most useful when f has end behavior approaching ∞
- MaxVT is essentially the same thing but negated
- Both are proved by EVT because \bar{B} is compact, so f has an extremum in \bar{B} and every other point in S is less extreme
- Of course, \bar{B} can be replaced with any other compact set

2.37 Apr 12: Critical Computation

- To find critical points, $f \in C^1$, and let S be nonempty closed. Assume $\operatorname{argmax}(f|_S) \neq \emptyset$, argmax being the points where the maximum is achieved, and

$$\operatorname{bd} S \subseteq \bigcup_{\alpha \in A} P_\alpha \subseteq S$$

$$U_\alpha \overset{\circ}{\subseteq} \mathbb{R}^{j(\alpha)}, j(\alpha) \geq 0$$

$$g_\alpha : U_\alpha \xrightarrow{C^1} S$$

$$P_\alpha = \operatorname{ran}(g_\alpha)$$

- The complicated condition is, in short, that $\operatorname{bd} S$ can be in the union of the range of a bunch of g -functions, or in the union of many n -dimensional "open" patches P_α - the patch may only be open in its preimage dimension, with each patch also in S
- Then, $\operatorname{argmax}(f|_S) = \operatorname{argmax}(f|_C)$, where C is the union of the critical points in the interior of S and the images under g of the critical points in each patch ($g_\alpha[\operatorname{crit}(f \circ g_\alpha, U_\alpha)]$) i.e. the maximum must occur at critical points
- It is necessary to consider patches because they transform the closed boundary of S into sets that are open in lower dimension, and critical points only make sense in open sets
- Additionally, an easy consideration is that the argmax on a larger set intersected with a smaller set inside it is a subset of the argmax on the smaller set; if it:

$$S \subseteq T \Rightarrow S \cap \operatorname{argmax}(f|_T) \subseteq \operatorname{argmax}(f|_S)$$

2.38 Apr 13: Critical Proof

- An even more marvelous proof of critical point extrema, which this page is just wide enough to contain.
- Let $\mathbf{p} \in \operatorname{amx}(f|_S)$, then $\forall \mathbf{x} \in S, f(\mathbf{p}) \geq f(\mathbf{x})$. Either $\mathbf{p} \in S^\circ$ or $\mathbf{p} \in \operatorname{bd} S$
- If $\mathbf{p} \in S^\circ$, let B be a small ball around it in S .

$$\mathbf{p} \in B \cap \operatorname{amx}(f|_S)$$

$$\mathbf{p} \in \operatorname{amx}(f|_B)$$

\mathbf{p} is a local maximum in S

$$\mathbf{p} \in \operatorname{crit}(f, S^\circ)$$

$$\mathbf{p} \in C \cap \operatorname{amx}(f|_S)$$

$$\mathbf{p} \in \operatorname{amx}(f|_C)$$

- If $\mathbf{p} \in \operatorname{bd} S$, then \mathbf{p} must be in one patch, $\mathbf{p} = g_\beta(\mathbf{u}_0)$. Because of its maximality, \mathbf{u}_0 must be an element of $\operatorname{amx}(f \circ g_\beta)$. Let B_2 be the ball around \mathbf{u}_0 which exists.

$$\mathbf{u}_0 \in B_2 \cap \operatorname{amx}(f \circ g_\beta)$$

$$\mathbf{u}_0 \in \operatorname{amx}((f \circ g_\beta)|_{B_2})$$

\mathbf{u}_0 is a local maximum in U_β , and the proof is now the same

- Let $\mathbf{q} \in \operatorname{amx}(f|_C) \subseteq C \subseteq S$, then $\forall \mathbf{x} \in C, f(\mathbf{x}) \leq f(\mathbf{q})$. Assume $\mathbf{q} \notin \operatorname{amx}(f|_S)$, by hypothesis, $f(\mathbf{p}) > f(\mathbf{q}) \Leftarrow \mathbf{p} \in \operatorname{amx}(f|_S) \Rightarrow \mathbf{p} \in \operatorname{amx}(f|_C)$; this is impossible, so $\mathbf{p} \in \operatorname{amx}(f|_S)$

2.39 May 29: Lagrange Multipliers

- Previous exploration into optimization restricted \mathbf{x} in an explicitly specified domain, but what if that domain S is just the solution set of another equation(s)?
- Find the maximum of the objective $6 - x - 2y$ satisfying the constraint $x^2 + y^2 = 1$

$$g(\mathbf{x}) = c \rightarrow f(\mathbf{x})?$$

First use EVT... to find existence! In the specific case, S is just a circle inside the domain of f , which is compact and so continuous f will have extrema

Next, if $\mathbf{p} \in \text{aex}(f|_S)$, then either

$$\nabla f(\mathbf{p}) = \lambda \nabla g(\mathbf{p})$$

where λ is the eponymous multiplier, or

$$\nabla g(\mathbf{p}) = \mathbf{0}$$

- Execution:

$$\nabla g(\mathbf{p}) = (2x, 2y)$$

$$\nabla f(\mathbf{p}) = (-1, -2)$$

Degeneracy only occurs at $(0,0)$, which does not satisfy g , however parallelism is possible at $(\frac{1}{\sqrt{5}}, \frac{2}{\sqrt{5}})$ and $(\frac{-1}{\sqrt{5}}, \frac{-2}{\sqrt{5}})$, which if unequal through f (they are) must represent a point of either extrema!

- Generalization to many constraints:

$$\nabla f(\mathbf{p}) = \sum \lambda_k \nabla g_k(\mathbf{p})$$

so the gradient of f is in the smaller subspace spanned by the gradients of g_k

Or degeneracy,

$$\det(\mathbf{g}')(\mathbf{g}')^T = 0$$

where the derivative matrix is invertible. This is equivalent to the gradients of the components being linearly dependent

2.40 Jun 1: et Ultra

- The problem from before can be reformulated with three variables:

$$\text{Maximize } z \text{ if } \begin{bmatrix} 6 - x - 2y - z \\ x^2 + y^2 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

Existence is first, the constraint set S is the intersection of closed (and one bounded) sets, so EVT does apply; the solution of a continuous equation is closed, which will be proven later

$$6 - x - 2y - z = 0 \quad x^2 + y^2 = 1$$

$$\nabla g_1 = (-1, -2, -1)$$

$$\nabla g_2 = (2x, 2y, 0)$$

$$\nabla f = (0, 0, 1)$$

$$(0, 0, 1) = \lambda_1(-1, -2, -1) + \lambda_2(2x, 2y, 0)$$

$$(-1, -2, 0) = \lambda_2(2x, 2y, 0)$$

$$(-1, -2, 0) = \lambda_2(2x, 2y, 0)$$

$$\left(-\frac{1}{\lambda_2}, -\frac{1}{\lambda_2}, 0\right) = (2x, 2y, 0)$$

The solution is the same from here; the degenerate case also has no solution, since the only time

$$(-1, -2, -1) \parallel (2x, 2y, 0)$$

is when the second is $(0, 0, 0)$, which fails in the same way as before

- Time to formalize:

$$f : U \subseteq \overset{\circ}{\mathbb{R}^d} \xrightarrow{C^1} \mathbb{R}, g : V \subseteq \overset{\circ}{\mathbb{R}^d} \xrightarrow{C^1} \mathbb{R}, U \cap V = \overset{\circ}{W} \neq \emptyset, 1 \leq \dim g \leq d$$

$$S = \{\mathbf{x} \in W | g(\mathbf{x}) = \mathbf{c}\} \neq \emptyset$$

$$\exists \mathbf{p} \in \text{aex } f|_S$$

\Downarrow

$$\begin{cases} \nabla f(\mathbf{p}) = \nabla f(\mathbf{p}) = \sum \lambda_k \nabla g_k(\mathbf{p}) \\ g(\mathbf{p}) = \mathbf{c} \end{cases} \vee \begin{cases} \det(g'(\mathbf{p}))(g'(\mathbf{p}))^T = 0 \\ g(\mathbf{p}) = \mathbf{c} \end{cases}$$

2.41 Jun 5: Everybody Loves Proof

- For this proof, let $f|_S \leq f(\mathbf{p})$
- By way of ImpFT:

$$\mathbf{F}(\mathbf{x}) = \mathbf{F}(\mathbf{y}, \mathbf{z}) = g(\mathbf{x}) - \mathbf{c} = \mathbf{0}$$

Solve for k variables (\mathbf{y}) in terms of $d - k$ variables (\mathbf{z}) near the solution $\mathbf{p} = (\mathbf{a}, \mathbf{b})$, then non-degeneracy is required

$$\det \frac{\partial \mathbf{F}}{\partial \mathbf{y}}(\mathbf{p}) = \det \frac{\partial g}{\partial \mathbf{y}}(\mathbf{p}) \neq 0$$

$$\Rightarrow \exists \psi : (\text{nbhd } \mathbf{a} \ni \mathbf{y}) \xrightarrow{C^1} (\text{nbhd } \mathbf{b} \ni \mathbf{z})$$

- To prove $LS(\mathbf{p}) \vee DS(\mathbf{p})$, it is often easier to rephrase it as $\neg LS(\mathbf{p}) \Rightarrow DS(\mathbf{p})$, or better yet, $LS(\mathbf{p}) \Leftarrow \neg DS(\mathbf{p})$ so do that

$$\det(g'(\mathbf{p}))(g'(\mathbf{p}))^T = \det(\mathbf{F}'(\mathbf{p}))(\mathbf{F}'(\mathbf{p}))^T \neq 0$$

- A fun useful theorem, $\det \text{gram } \mathbf{A} \neq 0 \Leftrightarrow \text{rank } \mathbf{A} = \text{full}$
- Also $\det \text{gram } \mathbf{A} = 0 \Leftrightarrow \det \text{gram } \mathbf{A}^T = 0$
- Therefore, the failure of DS is equivalent to $\text{rank } \mathbf{F}'(\mathbf{p}) = \text{full} = k$, now assume that!

2.42 Jun 8: Everybody Kind of Loves Proof

- Since $\mathbf{F}'(\mathbf{p})$ has full rank, it has a set of k linearly independent columns out of d ; without too much LOG just pick the first k and modify the proof if not.
- Each row vector \mathbf{x} will have the structure $(y_1, y_2, \dots, y_k, z_1, \dots, z_{d-k})$. Also, the left of the derivative matrix is just full rank $\frac{\partial \mathbf{F}}{\partial \mathbf{y}}(\mathbf{p})$, and its determinant is not 0 by FTLA.
- Now ImpFT takes effect, so $\mathbf{y} = \psi(\mathbf{z})$ near \mathbf{p} . Also, $S \cap (\text{nbhd } a \times \text{nbhd } b)$, which is not empty (\mathbf{p}), is just the range of $\phi(\mathbf{z}) = (\psi(\mathbf{z}), \mathbf{z})$, a parameterization with the z variables, since that is well defined exactly on the neighborhoods. ϕ is also C^1 since ψ is too and z_n is C^A
- Now force \mathbf{x} in the range of ϕ , so $f|_S(\mathbf{x}) = f(\mathbf{x}) = (f \circ \phi)(\mathbf{z})$. Additionally, since \mathbf{p} is a maximum, $(f \circ \phi)(\mathbf{z}) \leq (f \circ \phi)(\mathbf{b})$, so \mathbf{b} is also a maximum of its set. $(f \circ \phi)'(\mathbf{b})$ can now be taken to be $\mathbf{0}$, $= f'(\phi(\mathbf{b}))\phi'(\mathbf{b}) = f'(\mathbf{p})\mathbf{M}$, $\mathbf{0} = \mathbf{M}^T \nabla f(\mathbf{p})$. Since

$$\mathbf{M}^T = \begin{bmatrix} \psi'_1(\mathbf{z}) \\ \psi'_2(\mathbf{z}) \\ \vdots \\ \psi'_k(\mathbf{z}) \\ \mathbf{I} \end{bmatrix}$$

must have rank $d - k$ by the last rows, and nullity k , adding to d . Now, $\nabla f(\mathbf{p})$ is in the k -space, but how is that space made? $g(\phi(\mathbf{z})) = \mathbf{c}$ from before, so the derivatives are also equal. $g'(\phi(\mathbf{z}))\phi'(\mathbf{z}) = \mathbf{0}$. Taking transposes eventually leads to

$$[\mathbf{M}^T \nabla g_1(\mathbf{p}) \quad \mathbf{M}^T \nabla g_2(\mathbf{p}) \quad \dots \quad \mathbf{M}^T \nabla g_k(\mathbf{p})] = \mathbf{0}$$

All k of those vectors are all in the null space; are they linearly independent? It turns out so by assumption of full rank, so these can make a basis!

$$\nabla f(\mathbf{p}) = \sum \lambda_k \nabla g_k(\mathbf{p})$$

LS is proven

2.43 That's a Wrap!