

Some problems where the region doesn't have to be a polyhedron, and the optimal solution doesn't have to be on the boundary. Assume the dual form in order to find the extreme point. The objective is simple and convex, so it's a convex feasible set. As a wild generalisation (depends on exact problem class and algorithm): **First-order methods** converge either sublinearly or linearly. **Second-order methods** converge superlinearly or quadratically.

**Information-Based Complexity** measures the number of iterations required to achieve a solution within a specified tolerance. This is a practical metric for understanding how efficiently an algorithm can approximate the optimal value.

The **Jacobian matrix** generalizes the gradient for vector-valued functions, containing partial derivatives of each output component with respect to each input variable. Every Row is the transpose of the derivative of that function  $f_1$  with respect to  $x_i$ .

The **Hessian matrix** provides a second-order derivative for scalar functions, showing curvature. Convex iff function lies on or above all of its tangents, explained below. Single variate – only second derivative, multi-variate – strictly convex, Hessian matrix is **positive definite** everywhere. Definiteness of Hessian at a stationary point indicates whether or not it is a minimum (pos definite), maximum (neg definite) or saddle point (indefinite).

**Local Minimum Condition:** In convex functions, if  $x^*$  is a local minimum, then it is also a global minimum. If the function is convex, first order condition becomes a **sufficient** condition for global minimum, if not then hessian being positive semidefinite is necessary.

In multi-variable constraints, visualizing sublevel ( $f(x) \leq c$ ) and super-level ( $g(x) \geq c$ ) sets helps in understanding the feasible region. For convex functions, sublevel sets are convex, which supports the optimization process by maintaining a convex feasible region.

Equality and inequality constraints together define what is known as the feasible set. Latter divide the region as represented by a hyperplane, with a positive/negative half space or null region i.e. the hyperplane itself. Large family of optimization problems, mainly convex optimization problems can be solved efficiently in a unified manner.

Polyhedral, which is intersection of half spaces, is a convex set. Region above the graph is called epigraph, and api of a convex function is convex (hypo for concave). Scalar multi or addition are convex. An optimization problem is convex if the objective function and constraints are convex or linear. (obj being lin and constraint being convex counts). If f is convex but g is concave, problem can still be convex.

**Principle of duality** – Take a hyperplane that supports my set in such a way that the set falls entirely on the positive side of this hyperplane. That means a set is convex when you recover the set in intersection of all the positive sides of all the hyperplanes.

Extra - While a twice-differentiable function with a positive definite Hessian everywhere is strictly convex (e) and will have at most one stationary point (a), which is a global minimum (b), it does not guarantee that Newton's method will converge in one step. Newton's method converges in one step only if the function is quadratic. For non-quadratic functions, it typically requires multiple iterations to converge.

**Interior Point Method** – Way to solve KKT conditions. Interior Point Methods are optimization techniques used to solve linear and nonlinear constrained problems. They work by keeping the solution iteratively "inside" the feasible region and gradually moving towards the boundary, as opposed to boundary methods like Simplex. Interior Point Methods avoid directly hitting constraint boundaries by using a "barrier" approach, which penalizes the solution as it nears the boundary, making it computationally efficient.

We are going to solve a modified perturbed version called KKT(t) (which are much easier to solve. where t is +ve parameter which controls the degree of perturbation. As this parameter goes to 0,  $x(t)$  converges to  $x^*$ .  $u = -t\mathbf{g}(x)$ .

Plugging into the gradient equation, we get  $t \cdot \mathbf{g}(x)$ , which acts as a barrier function at the objective, approaching infinity as the variables

$$\nabla(\mathbf{f}(x) - t \log(-\mathbf{g}(x))) = 0$$

and approach the boundary of a constraint, discouraging any further movement in that direction, smoothing out the sharp transition of a constraint. The **barrier method** (interior point method) solves a sequence of barrier problems with  $t$  converging to zero.

We can now perform newton's method to perform this minimisation. Initially, don't  $t$  to be too small otherwise it will result in a discontinuous function with the log, and newton would be too slow to converge, so start with a big  $t$ , meaning  $\mathbf{g}$  would be overshadowed, and only  $u$  would need to be minimised.  $x(t)$  would be the analytic center of the feasible region. Once we have a solution, try for a smaller  $t$ , i.e.  $t=0.9$ , apply newton again, but  $x(t)$  would be the starting point, so computing  $x(t)$  would be faster.

Working towards  $t=0$  formed by the path of  $x$ s of  $t$ , called the central path (trajectory that solutions follow in the interior point

method as they approach the optimal solution), falling in the interior of the feasible region.

Second order methods can often achieve higher rates of convergence over first order methods. However, this comes at the cost of more computation per iteration, e.g., in calculating Hessians and solving a large linear system of equations. First order methods might be necessary for applications where the problem is very large, where the Hessian may struggle to fit into memory.

For a function  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ , state a condition that  $f$  needs to satisfy in order to be a convex function if  $f$  is not differentiable. State a different condition for the case where  $f$  is differentiable.

For the first part you could use: the line segments between any two points on graph of  $f$ , must be above or equal to  $f$ . A second option would be to state that the graph of  $f$  must be a convex set.

When  $f$  is differentiable, different conditions from the above are: the Hessian of  $f$  must be positive semidefinite, or the graph of  $f$  must lie on or above all its tangent planes.

The constraints are:

$$\begin{aligned} & \min_{x \in \mathbb{R}^n} f(x) \\ & \text{s.t. } g(x) \leq 0 \end{aligned}$$

The Penalty Method turns a constrained optimisation problem into an unconstrained one, for a penalty parameter  $\eta_k$ , by applying a quadratic penalty to constraint violations:

$$\min_{x \in \mathbb{R}^n} f(x) + \eta_k \sum_i \max(0, g_i(x))^2$$

A sequence of the above penalty problems are solved for increasing  $\eta_k$ , until the constraint violations are within some tolerance.

Battery efficiency, denoted by  $\eta$  (eta), measures the fraction of energy effectively stored or discharged. For charging, the internal power  $\delta$  is less than or equal to the external power multiplied by the efficiency:  $\delta \leq b \cdot \eta$ . Roundtrip Efficiency: the efficiency over a full cycle of charging and discharging is calculated as the square of  $\eta$ : Roundtrip efficiency =  $\eta^2$ . For example, if  $\eta$  (efficiency) is 90%, then the roundtrip efficiency is  $0.9 \cdot 0.9 = 0.81$ , or 81%. Non-Const Efficiency: By adding constraints at multiple points along the efficiency curve, the model can approximate the true behaviour of the battery's efficiency more accurately, without resorting to complex nonlinear equations. The optimality gap is the difference between the upper and lower bounds of the solution range in a minimization problem. It indicates higher confidence i.e. how close the solution is to the true optimum. The optimality gap for a maximisation is the difference between the smallest upper bound and the largest found feasible solution, inverse for minimisation.

Dual: Scaling factor for each primal constraint. The dual problem aims to find the largest lower bound for the primal objective, constrained by inequalities that preserve the primal-dual relationship.

### Rules for Getting the Dual

Dual of the dual is the primal.

The dual of a minimization problem:

$$\min_{x \in \mathbb{R}^n} c^T x \quad \text{max}_{y \in \mathbb{R}^m} y^T b$$

The dual of a maximisation problem:

$$\max_{x \in \mathbb{R}^n} c^T x \quad \min_{y \in \mathbb{R}^m} y^T b$$

When a constraint is non-binding, its dual variable is zero. Constraints with higher priority are:

$$\begin{aligned} & \min_{x \in \mathbb{R}^n} x + 2y - 5z \\ & \text{s.t. } x + 2y \geq 2 \\ & \quad y \geq z \\ & \quad z \leq 6 \end{aligned}$$

Convert the above problem to an equality and non-negative variable form (only equality constraints and non-negative variables). Implement the transformation in PuLP and verify you get the same solution.

See transform.py:

$$\begin{aligned} & \min_{x \in \mathbb{R}^n} x + 2y - 2z - 5z + 5 \\ & \text{s.t. } x + 2y - 2z - 5z + 5 \geq 2 \\ & \quad x \geq 0 \\ & \quad y \geq 0 \\ & \quad z \geq 0 \end{aligned}$$

Or remember one side, and convert max to min or max to min before taking dual

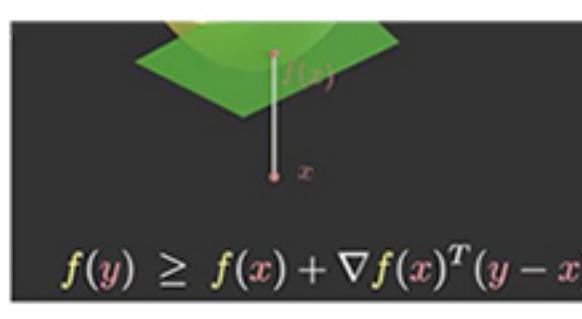
dual variable values have a stronger impact on the objective, indicating where resources are most valuable.

How a starting feasible solution is obtained for Simplex: Alternative simplex: Create a Feasibility Problem: This involves modifying the original problem by adding slack variables to relax the constraints. Objective: Minimize these slack variables to bring the system into feasibility. If the optimal solution to this feasibility problem has any non-zero slack variables, it indicates that the original problem is infeasible. Finding sv: The feasibility problem allows for a straightforward starting point by solving for a trivial starting vertex. Standard/Revised Simplex: Convert to Standard Form: Transform the problem into equality constraints with non-negative variables. Introduce Slack Variables: Add non-negative slack variables to handle equality constraints. The goal is to minimize the sum of the slack variables. Start with a BFS by setting all original variables to zero and using the slack variables in the basis. If all slack variables are zero at the optimal solution, the original problem is feasible, and this solution serves as the starting BFS.

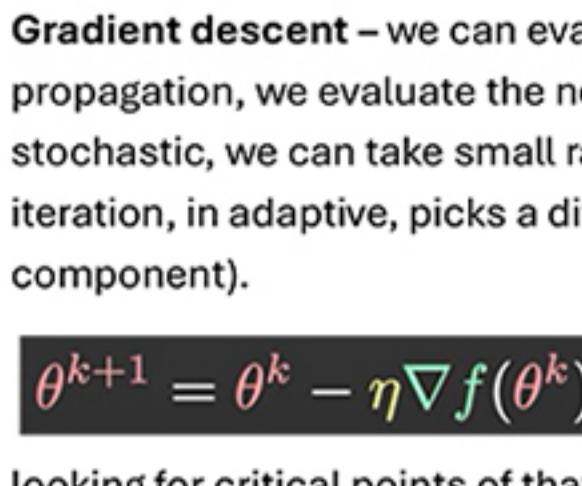
Pivoting in alternative simplex method, when a constraint is selected to leave the active set, the algorithm aims to move along the edge defined by the intersection of the remaining  $n-1$  active constraints. As we progress along this edge, the constraint that will enter the active set is the first one encountered. Determining this requires factoring in the alignment between the edge direction and the new constraint, as well as considering the slack available for the movement along this edge. There are two special cases to consider. First, if it's possible to move indefinitely along the edge without encountering a new constraint, the problem is deemed unbounded. Second, if multiple constraints are reached simultaneously, a pivot selection rule is employed to decide which constraint should enter the active set.

The Simplex algorithm typically requires a **basic feasible solution** as a starting point. A basic feasible solution is an extreme point of the feasible region. Since Alpha-Feas only guarantees that  $x_A$  is feasible (it satisfies  $Ax \leq b$ ) but doesn't guarantee it's a basic feasible solution,  $x_A$  may not be suitable as a starting point for the Simplex algorithm, and it's not guaranteed it would be an extreme point either.

To solve this optimization problem, we begin with the conditions that  $x-7 \geq 0$ ,  $x \geq 0$ , and  $y \geq 10$ , aiming to minimize  $x+5y$ . Dual objective function is  $\max 40y_1 + 10y_2 + 0y_3$  with constraints  $y_1 + y_3 \leq 1$  and  $y_2 - 7y_3 \leq 5$ . By simplex, substituting into the objective function, it simplifies to  $120 - 70s_1 - 10s_2 - 30y_1$ , so  $y_1=0$ , maximize the other two.



A function  $f$  is convex if, for any two points  $x$  and  $y$ , the line segment between these points lies above the function graph. Mathematically:  $f(\theta x + (1-\theta)y) \leq \theta f(x) + (1-\theta)f(y)$ , where  $\theta \in [0, 1]$ . This definition is central for proving that the solution is globally optimal in convex problems.



Tangent Hyperplane is a good approximation of the graph of  $f$  at point  $x$ . Dual def of convex fn,  $f$  is convex if and only if its graph is above its tangent hyperplane. So, for gradient 0, the hyperplane is horizontal – so global minimizer. So gradient 0 and solve for  $x$ . So non constant linear functions are unbounded. There can be more complex problems, to even solve the gradient, so there are other methods.

Gradient descent – we can evaluate the cost function at an arbitrary point, and with a process called back propagation, we evaluate the negative gradient and take a small step in that direction, doing it iteratively (in stochastic, we can take small random subset at each iteration, in adaptive, picks a different step size for each component).

Solve problem:  $\alpha_k := \arg \min_{\alpha} f(x_k + \nabla f(x_k)^T \alpha)$  where for gradient descent:  $\alpha_k = -\nabla f(x_k)$

Line Search

optimizing a single variable (moving in a straight line from our current iteration)

Either exactly which might take time, or inexactly using an approach such as backtracking line search.

We might not want to spend too much time doing this exactly because our current search direction might not be that great to start with!

1. Start with a large step estimate.

2. Iteratively shrink the step size until some conditions are met: i.e. the objective value has "adequately" improved.

3. Continue onto calculating a new search direction.

Newton's method-Hard to plot functions that have more than 2 variables(maybe even million), iterative optimisation algo, start with initial guess  $x_0$ (unknown, draw from some random distribution), and the closer the initial guess is to the true minimiser the better. We pick a direction  $d$ (descent direction) and we follow it to get  $x_1$ , track  $f(x_1)$  for progress in decline. How to pick good descent direction – gradient gives us direction of biggest increase.  $x_{k+1} = x_k + d_k$ .

Newton's method-Hard to plot functions that have more than 2 variables(maybe even million), iterative optimisation algo, start with initial guess  $x_0$ (unknown, draw from some random distribution), and the closer the initial guess is to the true minimiser the better. We pick a direction  $d$ (descent direction) and we follow it to get  $x_1$ , track  $f(x_1)$  for progress in decline. How to pick good descent direction – gradient gives us direction of biggest increase.  $x_{k+1} = x_k + d_k$ .

Newton's method-Hard to plot functions that have more than 2 variables(maybe even million), iterative optimisation algo, start with initial guess  $x_0$ (unknown, draw from some random distribution), and the closer the initial guess is to the true minimiser the better. We pick a direction  $d$ (descent direction) and we follow it to get  $x_1$ , track  $f(x_1)$  for progress in decline. How to pick good descent direction – gradient gives us direction of biggest increase.  $x_{k+1} = x_k + d_k$ .

Newton's method-Hard to plot functions that have more than 2 variables(maybe even million), iterative optimisation algo, start with initial guess  $x_0$ (unknown, draw from some random distribution), and the closer the initial guess is to the true minimiser the better. We pick a direction  $d$ (descent direction) and we follow it to get  $x_1$ , track  $f(x_1)$  for progress in decline. How to pick good descent direction – gradient gives us direction of biggest increase.  $x_{k+1} = x_k + d_k$ .

Newton's method-Hard to plot functions that have more than 2 variables(maybe even million), iterative optimisation algo, start with initial guess  $x_0$ (unknown, draw from some random distribution), and the closer the initial guess is to the true minimiser the better. We pick a direction  $d$ (descent direction) and we follow it to get  $x_1$ , track  $f(x_1)$  for progress in decline. How to pick good descent direction – gradient gives us direction of biggest increase.  $x_{k+1} = x_k + d_k$ .

Newton's method-Hard to plot functions that have more than 2 variables(maybe even million), iterative optimisation algo, start with initial guess  $x_0$ (unknown, draw from some random distribution), and the closer the initial guess is to the true minimiser the better. We pick a direction  $d$ (descent direction) and we follow it to get  $x_1$ , track  $f(x_1)$  for progress in decline. How to pick good descent direction – gradient gives us direction of biggest increase.  $x_{k+1} = x_k + d_k$ .

Newton's method-Hard to plot functions that have more than 2 variables(maybe even million), iterative optimisation algo, start with initial guess  $x_0$ (unknown, draw from some random distribution), and the closer the initial guess is to the true minimiser the better. We pick a direction  $d$ (descent direction) and we follow it to get  $x_1$ , track  $f(x_1)$  for progress in decline. How to pick good descent direction – gradient gives us direction of biggest increase.  $x_{k+1} = x_k + d_k$ .

Newton's method-Hard to plot functions that have more than 2 variables(maybe even million), iterative optimisation algo, start with initial guess  $x_0$ (unknown, draw from some random distribution), and the closer the initial guess is to the true minimiser the better. We pick a direction  $d$ (descent direction) and we follow it to get  $x_1$ , track  $f(x_1)$  for progress in decline. How to pick good descent direction – gradient gives us direction of biggest increase.  $x_{k+1} = x_k + d_k$ .

Newton's method-Hard to plot functions that have more than 2 variables(maybe even million), iterative optimisation algo, start with initial guess  $x_0$ (unknown, draw from some random distribution), and the closer the initial guess is to the true minimiser the better. We pick a direction  $d$ (descent direction) and we follow it to get  $x_1$ , track  $f(x_1)$  for progress in decline. How to pick good descent direction – gradient gives us direction of biggest increase.  $x_{k+1} = x_k + d_k$ .

Newton's method-Hard to plot functions that have more than 2 variables(maybe even million), iterative optimisation algo, start with initial guess  $x_0$ (unknown, draw from some random distribution), and the closer the initial guess is to the true minimiser the better. We pick a direction  $d$ (descent direction) and we follow it to get  $x_1$ , track  $f(x_1)$  for progress in decline. How to pick good descent direction – gradient gives us direction of biggest increase.  $x_{k+1} = x_k + d_k$ .

Newton's method-Hard to plot functions that have more than 2 variables(maybe even million), iterative optimisation algo, start with initial guess  $x_0$ (unknown, draw from some random distribution), and the closer the initial guess is to the true minimiser the better. We pick a direction  $d$ (descent direction) and we follow it to get  $x_1$ , track  $f(x_1)$  for progress in decline. How to pick good descent direction – gradient gives us direction of biggest increase.  $x_{k+1} = x_k + d_k$ .

Newton's method-Hard to plot functions that have more than 2 variables(maybe even million), iterative optimisation algo, start with initial guess  $x_0$ (unknown, draw from some random distribution), and the closer the initial guess is to the true minimiser the better. We pick a direction  $d$ (descent direction) and we follow it to get  $x_1$ , track  $f(x_1)$  for progress in decline. How to pick good descent direction – gradient gives us direction of biggest increase.  $x_{k+1} = x_k + d_k$ .

Newton's method-Hard to plot functions that have more than 2 variables(maybe even million), iterative optimisation algo, start with initial guess  $x_0$ (unknown, draw from some random distribution), and the closer the initial guess is to the true minimiser the better. We pick a direction  $d$ (descent direction) and we follow it to get  $x_1$ , track  $f(x_1)$  for progress in decline. How to pick good descent direction – gradient gives us direction of biggest increase.  $x_{k+1} = x_k + d_k$ .

Newton's method-Hard to plot functions that have more than 2 variables(maybe even million), iterative optimisation algo, start with initial guess  $x_0$ (unknown, draw from some random distribution), and the closer the initial guess is to the true minimiser the better. We pick a direction  $d$ (descent direction) and we follow it to get  $x_1$ , track <math

## DECOMP

**Column Generation:** This technique addresses complicating constraints by reformulating the problem into a "simpler" problem with many variables, called the Master Problem -> iteratively builds a (RMP) by selecting which variables or columns to add. Applications: **Constrained Shortest Path Problem:** Find the minimum-cost path from a start node to a goal node in a network, with each arc representing a cost, solvable by Dijkstra's. Addition of complexity like a time limit, A\* would require modifications to handle time constraints i.e. arrival times. **Flow Model:**

Formulate the problem to push one unit of flow from the start node and extract it at the goal node, each arc  $(i,j)$  has a binary variable  $x_{ij}$  to denote if it is used in the path, converted into ILP and solvable by branch and bound. Complicating constraint makes the problem challenging to solve directly, like the time constraint (not exceeding 14 units). **Reformulation Strategy:** Reformulate the problem by treating paths as individual variables. This approach shifts the focus to finding an optimal path without needing to consider all paths simultaneously. This separation enables us to introduce paths dynamically as needed and on demand generation. The constraint matrix is large, this approach leads to many variables (paths), especially in dense networks. However, we will use Column Generation to introduce only the paths that improve the solution, avoiding the need to consider all paths initially.

**Reduced Master Problem (RMP):** current version of problem with only a subset( $P$ ) of paths (columns). We repeatedly solve the RMP as more paths are added. In Column Generation, the **Pricing Problem** identifies paths that can enhance the current RMP solution by calculating their reduced cost. **Reduced Cost:** represents the marginal increase in the objective when each non-basic variable is marginally increased. Here  $\text{root}(A)$  is the basic part of  $A$  in  $\text{root}(A)x = b$  for basic variables, and  $A^*$  is the non-basic part of  $A$ .  $c^*$  is the obj coeff of the NB vars, and  $\lambda$  are the dual variables for the basic variables. **Lambda 0** won't affect our reduced cost, because we already have chosen one path, no need to incorporate, already been satisfied, and it's 1 if no path is chosen. Since the pricing problem does not have a comp cstr, it's easy to solve, and only paths with negative reduced cost are added to the RMP, as these are paths that could improve the objective function. Tldr:

Start with a set of initial columns and solve the RMP using these columns. Then, use the dual variables from this solution to define a new subproblem known as the Pricing Problem. Solve the Pricing Problem to find the column with the minimum reduced cost. If this new column has a negative reduced cost, add it to the RMP and repeat the process. If not negative, then the optimal solution to the linear relaxation of the problem has been found. **Branch & Price:** If a path variable  $x_{ij}$  has a fractional value in the RMP solution  $x_{12}=0.6$ , we create branches to enforce integer solutions. **Branch 1:**  $x_{12}=0$ , **Branch 2:**  $x_{12}=1$ , for 0, remove all paths that use 1->2, inverse for 1. Very similar process for cutting stock, that identifies new cutting patterns that minimize waste while

changes can be made to the solution, allowing the algorithm to move from one solution to another in search of an improved outcome. In the n-queens problem, The Random Walk algo for n-queens selects random moves from the neighbourhood and executes them continuously, asymptotically complete – visits every state. Hill Climbing where the best move—or one of the best moves—is chosen at each step to reduce conflicts. This means evaluating all possible moves for each queen and selecting the one that most improves the solution (i.e., minimizes conflicts). This method requires evaluating the entire neighborhood to identify the optimal move. Unlike greedy search, which evaluates the entire neighborhood to select the best move, first-found randomizes the neighborhood evaluation. This approach picks the first improving move it encounters rather than the absolute best. **Randomized greedy** combines elements of randomness with greedy selection. With a probability  $p$ , the algorithm will perform a greedy or first-found move; otherwise, it will make a random move. This probabilistic element ensures that the algorithm is asymptotically complete. Biased, moves are selected with a probability inversely proportional to their score. For example, a move that reduces conflicts more significantly is more likely to be chosen.

**VRP:** The 1-Move operator relocates a single customer from one route to another. Swap 1-1 swaps two customers between different routes. Swap 2-1 operator exchanges two customers in one route with a single customer from another. Swap Tails involves swapping the end segments, or "tails," of two routes. The 2-Opt operator removes two arcs within a route and reconnects them in a way that shortens the path, while 3-Opt removes three arcs and reconnects them in various possible configurations. These operators provide a way to systematically reduce travel distances, enhancing the efficiency of the overall routing solution. **Neighbourhoods - Knapsack problem:** Swap 2 items, Swap 1 item with multiple items of equal size, Scheduling - swapping jobs between machines or the order of jobs, Map Coloring (k-coloring), adjusting neighborhoods to allow single-color changes, local search can more easily explore feasible solutions and satisfy colouring constraints. **Soft constraints** allow some degree of flexibility. Instead of completely blocking certain solutions, they impose penalties for violations but allow the algorithm to explore those solutions, like allow adjacent regions to have the same colour but apply a penalty to the obj fn. Effectively defining a neighbourhood structure provides a balance between exploring new areas and intensifying the search around promising solutions.

**Local minima** are points where the solution is better than all neighboring solutions, making further improvement difficult. **Plateaus** are flat areas where many neighboring solutions have similar values, causing the search to stagnate. To address these issues, **random restarts** allow the search to start fresh from a new random solution whenever it gets stuck, helping explore different regions. **Sideways moves** in 8-queens permit non-improving moves across plateaus, with a set limit of 100 to avoid getting stuck indefinitely.

**Simulated Annealing:** Temperature influences the acceptance of new solutions. At higher temperatures, the algorithm allows greater diversification by accepting worse solutions, which broadens the search for new areas. As the temperature decreases, the algorithm shifts towards intensification, focusing on refining and exploiting the 'good' solutions already found. In Simulated Annealing, maintaining a constant low temperature does not guarantee that only improving solutions are accepted because there's always a small probability of accepting worse solutions (considering  $P(\text{accept increase } \Delta) = e^{-\Delta/T}$ ). Additionally, using a constant low temperature throughout the search limits the algorithm's ability to explore the solution space effectively, increasing the risk of becoming trapped in local minima. The **cooling schedule** is crucial here; common approaches use a geometric cooling rate (e.g.,  $T_k \leftarrow q \cdot T_k$  with  $q$  usually in  $[0.9, 0.999]$ ) to gradually reduce the temperature. Don't want too many temp changes. Steps involve updating T, re-boil(reinitialize t), handle randomly, and adaptive parameters like maxChangesThisT.

**Meta-heuristics** are versatile optimization strategies that guide the search process in both discrete and continuous domains to find good feasible solutions without guaranteeing exact optimality. They differ from exact methods, as they are not based on algebraic models and often use **randomization** to explore the solution space. Typically non-deterministic, meta-heuristics may yield different results with each run and often stop based on time or iteration limits rather than achieving a guaranteed optimum. They are generally **problem-independent**, meaning they can be applied to a variety of optimization tasks, though they can incorporate **problem-specific heuristics** when beneficial. Many meta-heuristics also utilize **memory**, tracking past solutions to better direct the search process.

## ADVANCED OPTIMISATION TOPICS

You have seen several heuristic search algorithms in the AI class to solve unconstrained combinatorial problems. You have seen that A\* using an admissible heuristic is guaranteed to find the optimal solution to the problem.

Recall that A\* works by minimizing the cost of reaching the goal and an admissible heuristic for A\* is a function  $h(s)$  such that  $h(s) \leq h^*(s)$  for all states  $s$  of the problem where  $h^*(s)$  represents the optimal cost to reach the goal from  $s$ . In other words, an admissible heuristic  $h(s)$  is a lower bound on the optimal solution for all  $s$ .

In the multi-objective version of A\*, a multi-objective cost function  $\vec{c}(x) = [c_1(x), c_2(x), \dots, c_m(x)]$  is used for  $m \geq 2$  and  $h^*(s)$  now represents the Pareto Front for the state  $s$ .

Explain how to generalize the concept of an admissible heuristic for the multi-objective case and describe the property an admissible multi-objective heuristic should satisfy.

Furthermore, use your definition to provide an example of an admissible heuristic value for the Pareto front  $\left\{ \begin{pmatrix} 6.2 \\ 7.2 \\ 34 \end{pmatrix}, \begin{pmatrix} 4.9 \\ 7.0 \\ 38 \end{pmatrix}, \begin{pmatrix} 1.4 \\ 7.5 \\ 45 \end{pmatrix}, \begin{pmatrix} 1.52 \\ 8.2 \\ 29 \end{pmatrix} \right\}$

## Multi-Objective Optimisation

The formal definition of an admissible heuristic for multi-objective A\* is:

A multi-objective heuristic  $h$  is a set of vectors and it is admissible if and only if, for all states  $s$ , for all vector  $\vec{u}^*$  in the Pareto Front of  $s$  there exists a vector  $\vec{v} \in h(s)$  such that  $\vec{v} \preceq \vec{u}^*$ .

The key differences from the single objective case are:

- replace  $\leq$  by  $\preceq$
- make sure that every value in the Pareto Front has at least 1 vector in the heuristic that is "better than or equal to" it, i.e.,  $\vec{v} \preceq \vec{u}^*$ .

Also, keep in mind that the best admissible heuristic possible is the optimal solution, i.e., a lower bound that is equal to the function being bounded. For multi-objective, this is the same thing, i.e., the Pareto-Front is the best possible admissible heuristic. The ideal-point or using an individual admissible for each direction is an **example** of admissible multi-objective heuristic and it can be greatly improved in practice.

## Stochastic Programming

In this question, let's consider that our random variable  $\epsilon$  represents all the uncertainty in our stochastic problem and is "discrete", meaning  $P(\epsilon = i) > 0$  for  $i \in \{o_1, \dots, o_n\}$ . This is exactly the case we discussed in the lecture.

An important property of the expected value operator is its "linearity": given a random variable  $\epsilon$ ,  $E[\epsilon a + b] = aE[\epsilon] + b$  for  $a \in \mathbb{R}$  and  $b \in \mathbb{R}$ . It is easy to see this when you expand the expected value operator, i.e.,

$$E[\epsilon a + b] = \sum_i P(\epsilon = i)(a + b) = a(\sum_i P(\epsilon = i)) + b = aE[\epsilon] + b$$

In class, we learned that the optimal solution of the recourse-problem is less than or equal to the optimal solution of the expected-value problem. Formally, as we saw in class, let  $g(x, \epsilon)$  represent the stochastic problem for a fixed value  $\epsilon \in \{o_1, \dots, o_n\}$  for the random variable  $\epsilon$ . Moreover, for this question, assume that this stochastic problem takes the form of a **Linear Program**, that is, the objective function is a linear function of  $x$  and  $\epsilon$ , and the feasible region is defined by a set of linear constraints. Then we have  $\min_{x \in \Omega} E[g(x, \epsilon)] \leq \min_{x \in \Omega} g(x, E[\epsilon])$ .

Explain why the recourse-problem  $\min_{x \in \Omega} E[g(x, \epsilon)]$  can potentially be lower than  $\min_{x \in \Omega} g(x, E[\epsilon])$  for a **Linear Program**, despite the **linearity** of the expected value operator.

Using the same notation as the L-Shape Method slides, we have that  $\min_{x \in \Omega} g(x, E[\epsilon])$  is:

$$\begin{aligned} \min_{x, \epsilon} f^\top x + (\sum_i P(\epsilon = i)c_i)^\top z \\ \text{s.t. } (\sum_i P(\epsilon = i)B_i)x + Dz = (\sum_i P(t = i)d_i) \end{aligned}$$

where all the occurrences of " $(\sum_i P(\epsilon = i) \cdot \cdot \cdot)$ " are computing the expectation regarding to the uncertainty. In other words, they will become the average vector or matrix depending on what the target of the expectation is.

Linearity of expectation still holds, e.g., we can move  $z$  in the objective function inside the summation computing the expectation and value is still the same:  $(\sum_i P(\epsilon = i)c_i)^\top z = \sum_j (\sum_i P(\epsilon = i)c_i)_j * z_j = \sum_i (\sum_j P(\epsilon = i)c_{ij} * z_j)$ . Using this, we can create an specific value of  $z$  (the recourse action) for each possible value of the uncertainty, i.e., change  $z$  to  $z_i$ . But this alone is not enough. We need to specialize the constraints, i.e., change  $(\sum_i P(\epsilon = i)B_i)x + Dz = (\sum_i P(t = i)d_i)$  to  $B_i x + Dz_i = d_i \quad \forall i$  and this change in the feasible space is what allows us to obtain a potential smaller objective function.

satisfying demand. Here  $a$  is the no of items of length  $i$  for pattern  $j$ . The max gives us a knapsack problem which is np easy. **Dantzig-Wolfe Decomposition** is a technique used to solve large-scale optimization problems with a **block structure** by breaking them into smaller, more manageable subproblems. Each subproblem, or "block," is optimized independently through **Column Generation**, where only the most promising solutions (columns) are generated. These subproblem solutions are then coordinated in a **Master Problem** that enforces global constraints, ensuring all blocks align to meet the overall objective. The process iterates, with new constraints added to the subproblems based on the Master Problem's solution, making it an efficient method for complex problems with strong interdependencies. **Complicating Constraints:** These constraints link variables in ways that make the feasible region difficult to define or compute efficiently, often adding complexity to the problem. Examples include the maximum time constraint in the constrained shortest path problem and the decision on cutting lengths in the cutting stock problem. **Complicating Variables:** These variables make solving the problem challenging because they are not continuous and/or link multiple constraints. Fixing their values simplifies the optimization problem. For instance, deciding whether to build facility  $k$  ( $y_{k0} \in \{0,1\}$ ) in facility location problems.

## METAHEURISTICS

Construction algorithms are essential in combinatorial optimization, where they create a foundational solution for problems like the Traveling Salesman Problem (TSP) or Knapsack. These algorithms construct solutions by adding one element at a time, setting up a structure that more advanced methods can later refine. A common approach in construction is the "**Nearest Neighbor**" method, often applied to TSP. Starting at a "home" city, it iteratively adds the nearest unvisited city to build a route. Greedy, performs locally optimal moves, often make "structural" errors, really good to ok for knapsack and scheduling. **"Minimum Insert Cost"** algorithm considers multiple potential positions for each addition and chooses the insertion point that minimally impacts the current solution's cost. By carefully balancing the cost of each insertion, this method helps avoid the structural flaws of simpler greedy approaches and performs well in more complex scenarios like the Vehicle Routing Problem (VRP), where routes and costs must be optimized across many stops. It allows insert between any pair of customers-not just at the end. **Randomization** in construction algorithms introduces variability by selecting a random element, such as a customer, and inserting it in the least costly position within the solution. This approach, unlike strict greedy methods, allows any element to be placed optimally rather than only at the end, which can prevent common structural errors and avoid overly deterministic solutions that lead to local minima. It's also to introduce unconventional solution paths that might not emerge in a deterministic approach and to mitigate deterministic biases. For Randomized greedy, choose the  $r$ th best of  $n$  choices i.e.  $r = \text{rand.uniform}(0,1)^y$ .

**Vehicle Routing Problem (VRP):** The VRP involves finding optimal routes for a fleet of vehicles to serve a set of customers at minimum cost. Solutions are built by adding customers in ways that minimize added cost to the route.

This problem is foundational in logistics and transport optimization. **Regret Heuristic:** The regret heuristic adds foresight by considering the difference (or "regret") between the cost of the best route for a customer and the second-best alternative. The customer with the highest regret value is prioritized, aiming to minimize potential losses from future choices. The concept extends to  $k$ -regret, where choices are based on differences among multiple top routes, ensuring better overall optimization by accounting for future limitations. **Regret with Seeds:** To improve the initial structure in VRP, "seeds" are used to select initial customers based on distance criteria. For example, Seed 1 is the customer farthest from the depot, Seed 2 is the farthest from Seed 1, and so on. This method helps in laying out an initial skeleton for routes, enabling more balanced and structured solutions as regret-based choices build on this framework.

**Bespoke Method:** Involve ordering customers in a systematic way, often based on geographic or other logical sequences, to simplify route construction. Like "Sweep": Sweeping through an area, clustering those customers until max capacity is reached, route is formed. Not good for constraints involving perishable goods.

**Problem Definition:** 3D packing addresses the challenge of fitting a collection of items into a set of trucks or containers efficiently. Each item has specific dimensions (width, height, depth), weight, and a due date, while each truck has compartments with set dimensions, weight limits, and departure times. The goal is to maximize the number of items packed while respecting all constraints. **Problem Statement:** Given a fleet of trucks and a set of items with various attributes, determine how to fit as many items as possible into the trucks. Constraints include observing weight limits, fitting items within truck dimensions without overhangs (bigger item on top of the smaller item), respecting load and departure dates, and adhering to placement rules (e.g., "top-load only" items must not be stacked under other items). **Solution:** Using a metaheuristic to solve this i.e. LNS, with 2 related subroutines – distribution of items among the trucks, fitting the items inside a truck. The solution involves packing items starting from corners to utilize boundaries effectively. Corners are easier to track and provide stability. Tops of items act as additional "containers" for further stacking, and items with the same top height can be joined to create flat surfaces for optimized stacking. This layered approach facilitates stable, efficient packing while following the constraints, reducing wasted space and maximizing load efficiency.

**Local search** is an optimization technique that improves an initial solution by exploring its "neighborhood"—a set of solutions reachable through small adjustments. A neighborhood is defined by an **operator** that specifies how

IILS improves solutions by performing repeated local searches combined with periodic perturbations to escape local minima. Once a local minimum is found, ILS applies a "kick" or "shake" with 'strength' of how many elements to change. This perturbation strength must be balanced—to weak, and the search risks returning to the same minimum; too strong, and it becomes a random restart. A common perturbation for the Traveling Salesman Problem (TSP) is the **Double Bridge** move, which replaces four arcs to generate a new, slightly altered solution.

**GRASP (Greedy Randomised Ascent Procedure)** - constructs solutions in a randomized, greedy fashion. It starts with an empty solution and gradually adds elements using a **Reduced Candidate List (RCL)** based on the cost of each addition (e.g., minimum insert cost in TSP). Each new element is chosen randomly from the RCL, allowing controlled randomness. The parameter alpha controls the balance: low  $\alpha$  yields a greedy solution, while high  $\alpha$  introduces more randomness. Like ILS, GRASP can partially destroy and reconstruct solutions, combining systematic construction with controlled randomness to balance exploration and exploitation. A key characteristic of the construction phase in GRASP is that it constructs a solution by selecting components from a Restricted Candidate List (RCL) based on certain strategies. The RCL is created using a greedy function to evaluate and rank potential components. Instead of always choosing the top-ranked component (as in a purely greedy algorithm), GRASP introduces randomness by selecting from the RCL, which includes several high-quality candidates.

**VNS - VNS** dynamically changes neighborhoods to escape local minima. Starting with a list of neighborhoods, the algorithm performs a local search in one neighborhood. If no improvement is found, it switches to a larger neighborhood. When an improvement is achieved, it returns to the smallest neighborhood to refine the solution. For VRP, neighborhood moves might include 1-move, 1-1 swap, 2-2 swap, Or-opt, and k-opt moves. Define multiple neighborhoods, create arcs that jump things, visit some other states or solutions that our other neighborhood didn't allow us to do, changing how you get to the optimal min, but the local minima might change, you want the neighborhoods to get more complicated (from smaller to larger branch of factors, computationally cheap to expensive since it resets, computationally efficient) as  $k$  gets larger (increment count of neighbors), helps to escape the stuck local minima, exhaust the search space of the neighbor before moving onto the next one. No requirement there should be no intersection, i.e. neighbors can intersect, the conjunction doesn't have to be empty.

**LNS:** optimizes solutions by partially destroying and reconstructing them. This "destroy-and-recreate" approach removes parts of the solution (such as customers in VRP) and reinserts them using a preferred insertion method. The destruction phase can vary, for example, removing the longest arc, choosing a sequence of related customers, or dumping all customers from certain routes. The extent of destruction is typically random within set limits to ensure flexibility. For VRP, remove some visits, move them to assigned lists.

During reconstruction, LNS may use methods like **Minimum Insert Cost** or **k-Regret** to rebuild the solution systematically. This balance between structured removal and strategic reinsertion makes LNS effective in complex optimization problems, allowing significant changes without losing the overall solution quality.

**Adaptive Large Neighborhood Search (ALNS)** builds on LNS by adapting both the removal and insertion methods based on past performance, increasing the chances of finding optimal or near-optimal solutions.

To destroy, rand (min, max) where min is a number that is meaningful to destroy and max is not too large, otherwise it might be harder to reconstruct. For milp reconstruction, the original solution is a feasible solution to the milp, so you are either guaranteed to have either or better solution, which is a nice property to have. Validating against a milp is very effective, as the number of free variables is going to be very small and is guaranteed to get the global minimum given enough time, since the variables are fixed on the milp, so the model is going to be 1000 or millions bigger, but it's effective due to the low number of free variables.

**Adaptive insert method** - Method keeps accumulating that reward, starting with random uniform. Start with a small reward, and after 100 iterations, we use the rewards as a basis to rerate the probabilities, so the method that got more chosen is going to be chosen more often, keep doing that every 100 iterations. Same with **adaptive destruction method**, but you don't want any methods to have zero probability, so guard against zero should be there.

$\lambda$ : a parameter controlling the penalty strength,  $p_i$ : Penalty for feature  $i$ , which increases as the feature is penalized.  $l_i(s)$ : Indicator function that equals 1 if feature  $i$  appears in solution  $s$ , and 0 otherwise.

**Initialize:** Begin with an initial solution  $s$  and set penalties  $p_i=0$  for all features. **Local Search:** Perform a local search to minimize the original objective  $g(s)$ . **Select Features to Penalize:** Calculate the utility of each feature. Features with high utility are prioritized for penalization. **Update Penalties:** Increase the penalty  $p_i$  by 1 for all selected features. Perform local search on the augmented objective  $h(s)$  and continue the process.