**GOAL PROGRAMMING** → find the $x \in \mathbb{R}^n$ to get a particular value from $[f_1(x), \dots, f_m(x)]$

**MULTI-AGENT SYSTEMS** → each agent can
   1) however infos from environment
   2) act on environment
   3) receive reward from environment
}  Then they need to cooperate or compete to reach a goal

**OPTIMIZATION TYPES**
 ① Local vs Global
    - local: gradient methods
    - global: space filling curves (effective in low dimensional spaces)

 ② Deterministic vs Stochastic (eg. SGD or EA) or SIMULATED ANNEALING

 ③ Single vs. Multi Objective

 ④ Single vs Multi Agent

 ⑤ Heuristic vs Metaheuristic
    - Heuristic: problem specific methods like greedy
    - Metaheuristic: doesn't need domain knoledge, generates solutions typically with stochastic methods
      ( nature inspired algorithms → GA, Ant Colony, Swarm Intelligence )

  Swarm Intelligence (SI) → decentralized system of agents, interacting locally with one another and the environment
    - emerges intelligence that is unknown from single agent

**Multi Objective Optimization**
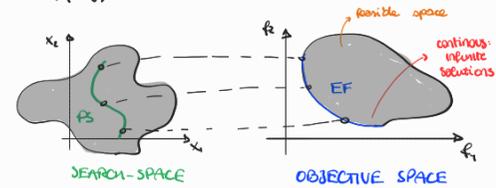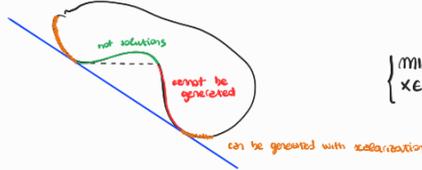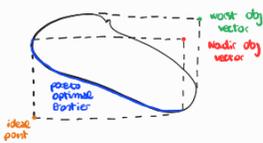$$\begin{cases} \min \ F(x) = (f_1(x), \dots, f_n(x)) \\ x \in D \ \curvearrowright \text{search space} \end{cases}$$
  - x può essere continua, discreto o una mixing variables

$\vec{x}$ dominates $\vec{y}$ ⟷
 1) $f_i(x) \leqslant f_i(y) \ \forall i$
 2) $\exists j : f_j(x) < f_j(y)$
  - strongly dominates if $\forall i \ f_i(x) < f_i(y)$

$\vec{x}$ is a pareto optimum if it is not dominated by any other solution
  - objectives are conflicting, otherwise we have an ideal point $\in D$
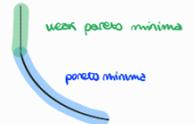
PS: pareto set
EF: efficient frontier


SEARCH-SPACE     OBJECTIVE SPACE

$$\begin{cases} \min \ w^T \cdot F(x) \\ x \in D \end{cases} \text{ with } \ \Sigma w_i = 1 \ \wedge \ w_i \geqslant 0$$

· can generate all the solutions ⟺ P is convex

$\bigcup_{w\geqslant 0} S_w = \{\text{Weak minimum}\}$   $(\nexists x \in S : f_i(x) < f_i(x^*) \ \forall i)$
$\bigcup_{w>0} S_w \subseteq \{\text{Minimum}\}$
    ↳ equal in the linear case

weak pareto minima
pareto minima

To obtain the final solution we need to weight the objectives
 1) a priori → single obj optimization
 2) a posteriori → get the pareto optimal net and the choose the preferred solution
 3) interactive → providing feed-backs and zone to focus (from humans or AI)

Objectives of the optimizers
 1) find solutions as close as possible to EF
 2) try to reach the extremes of EF
 3) solutions as uniform distributed as possible
} META-TARGETS

Two algorithms    1) scalarization    2) population based → more solutions per iteration (eg. GA)
  GA → avoid problems with non-convex P and generate with a single pass more points of EF
    - handles well non-differentiable objectives

**GENETIC ALGORITHMS** → meta-heuristic optimization algorithms, taking inspiration by evolution

1) A solution may be encoded as string, if binary:
   $\underbrace{\overset{2^4 \quad 2^6}{0100 \ 010010}}_{x_1 \quad x_2}$    we can map them with $\quad X_i = X_i^{min} + \dfrac{X_i^{max} - X_i^{min}}{2^{\ell_i} - 1} \ DV(s_i)$   $s_i$ string int

  · cross over operator → select two strings and recombine them in two off springs by mix together the genetic materials
    - single point cross-over: randomly select a crossing site and exchange all bits on its right
    - uniform cross-over → select a bit from one of the two parents with $p$ probability
   not all strings under goes the cross over, we use $P_c$ to model the P to select a string

  · mutation operator → randomly mutate some genes at random (probability of $P_m$)
  · Elitist → when the best solution (or all the non dominated ones) can be selected without mutation and crossover to the next generation
    (will be selected if not out-performed by the offspring)

2) Vector of real numbers $x \in \mathbb{R}^n$ → perturb with noise, combine $x_1 x_n$ with averaging, ...

*replace only a few individuals at each generation, more computationally efficient, but more generations needed to converge*

## 1) Steady State + Mutation
- elitist

```
create randomly initial population of pop_size
evaluate their fitness ( obj + penalization if not feasible)
while (not termination condition)
   - choose at random (or ∝ fitness) one individual from curr population
   - mutate and evaluate fitness
   - add to current population
   - remove worst individual
```

## 2) SS + Cross-Over + Mutation    → inside the while:

```
- choose at random two individuals
- cross-over to generate two offspring
- mutate both, evaluate
- remove the worst two
```

*potential full population replacement in one generation → more aggressive*

## 3) Generational GA + Cross-Over + Mutation
- still elitist

```
t=0 , create P_t randomly   (pop_size), evaluate P_t
while (not finish)
   - create Q_t of offspring from P_t of the same size (pop_size)
   - evaluate Q_t
   - R_t = Q_t ∪ P_t      (2 pop_size)
   - create P_{t+1} by removing the worst pop_size from R_t
```

To create $Q_t$ :

```
while ( size Q_t < pop_size)
   select first parent with BTS
   select second parent with BTS
   generate offsprings
   mutate and evaluate
   add to Q_t
```

→ Binary Tournament Selection (if P constrained)
Pick two individuals at random from P_t
- if both are feasible keep the best
- if only one is feasible keep that
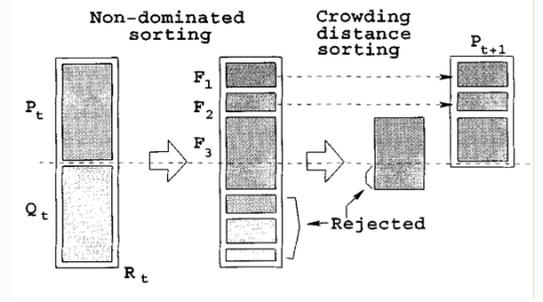- keep the one that violates less the constraints

## MOGA
- Family I → pareto dominance based algorithms (NSGA-II, PAES)
- Family II → indicator based algorithms (SMS-EMOA, HyPE, CHEA)
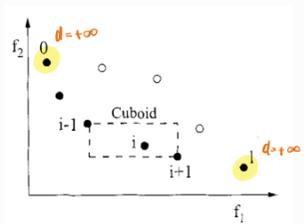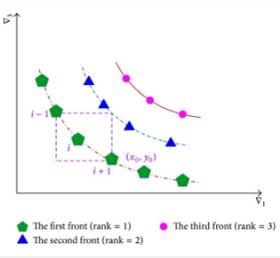- Family III → decomposition based algorithms (MOEA/D)

## NSGA-II (non-dominated sorted GA)
- params: pop_size, n-generation, p_cross, p_mut
- elitist     ↳ no more size as I

```
P_0 ← random initial population of size pop_size
Q_0 ← offsprings from P_0 of size pop_size
for i=1 to max_gen
   R_t = P_t ∪ Q_t
   F ← fast-non dominated sort of R_t      [ F = (F_1, ..., ) ]
   i ← 1 , P_{t+1} ← {∅}
   while |P_{t+1}| + |F_i| < pop_size
      P_{t+1} = P_{t+1} ∪ F_i
      i ← i+1
   compute crowding distance of F_i ⤳ last frontier
   K ← pop_size - |P_{t+1}|      # elements to be added from F_i
   P_{t+1} ← P_{t+1} ∪ best K elements in F_i
   Q_t ← generate offspring from P_{t+1}
```



Non-dominated sorting   Crowding distance sorting

```
fast-non-dominated-sort(P)
for each p ∈ P
   S_p = ∅
   n_p = 0
   for each q ∈ P
      if (p ≺ q) then              If p dominates q
         S_p = S_p ∪ {q}           Add q to the set of solutions dominated by p
      else if (q ≺ p) then
         n_p = n_p + 1             Increment the domination counter of p
   if n_p = 0 then                 p belongs to the first front
      p_rank = 1
      F_1 = F_1 ∪ {p}              → non dominated solutions
i = 1                              Initialize the front counter
while F_i ≠ ∅
   Q = ∅                          Used to store the members of the next front
   for each p ∈ F_i   → last frontier
      for each q ∈ S_p
         n_q = n_q - 1
         if n_q = 0 then          q belongs to the next front
            q_rank = i + 1
            Q = Q ∪ {q}           → add to next frontier iif not dominated
   i = i + 1                         by any other, just F_1,...,F_i
   F_i = Q
```





```
crowding-distance-assignment(I)
l = |I|                                          number of solutions in I
for each i, set I[i]_distance = 0                initialize distance
for each objective m
   I = sort(I, m)                                sort using each objective value
   I[1]_distance = I[l]_distance = ∞            so that boundary points are always selected
   for i = 2 to (l-1)                            for all other points
      I[i]_distance = I[i]_distance + (I[i+1].m - I[i-1].m)/(f_m^max - f_m^min)  → normalization
```
↳ equal weight among objectives

- scelgo quelle con crowding distance maggiore
perchè zona meno affollata in quel fronte
→ soluzioni uniformi

- works also on non-convex ES

*pareto archived ES*

**$(\mu - \lambda)$ - PAES** → **mutation only** ( evolutionary strategies )

*initially only mutation*

$\mu$ : number of current solutions
$\lambda$ : number of mutants generated from current solutions     when $\lambda = 1$ : steady state algorithm
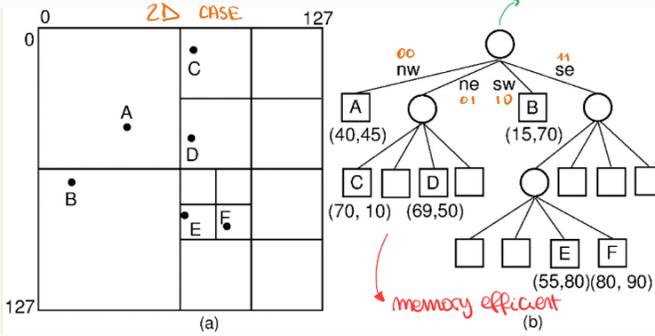
**(1-1) PAES**

*current solution*

*post ND solutions*
*approximation of current PS*

```
1    generate initial random solution c and add it to the archive
2    mutate c to produce m and evaluate m
3        if (c dominates m) discard m
4        else if (m dominates c)
5            replace c with m, and add m to the archive
6        else if (m is dominated by any member of the archive) discard m
7        else apply test(c,m,archive) to determine which becomes the new
             current solution and whether to add m to the archive
8    until a termination criterion has been reached, return to line 2
```

**test (c, m, archive)**

```
1    if the archive is not full
2        add m to the archive
3        if (m is in a less crowded region of the archive than c)
4            accept m as the new current solution
5        else maintain c as the current solution
6    else
7        if (m is in a less crowded region of the archive than x for
             some member x on the archive)
8            add m to the archive, and remove a member of the archive from
                 the most crowded region
9            if (m is in a less crowded region of the archive than c)
10               accept m as the new current solution
11           else maintain c as the current solution
12       else
13           if (m is in a less crowded region of the archive than c)
14               accept m as the new current solution
15           else maintain c as the current solution
```

→ *m may be the current solution but $\notin$ archive*

- Assuming the range is defined for each objective
  I can repeatedly bisection each objective range

  $2^{\ell \times d}$ hypercubes     $\ell$: # of space bisection
                                      $d$: space dimension

**Adaptive method**: specify $\ell$ but not the range ⟹ variable size hypercubes
  - the range for m objective is the one of current solutions
  - recompute the grid ⟺ the $\Delta$range > th



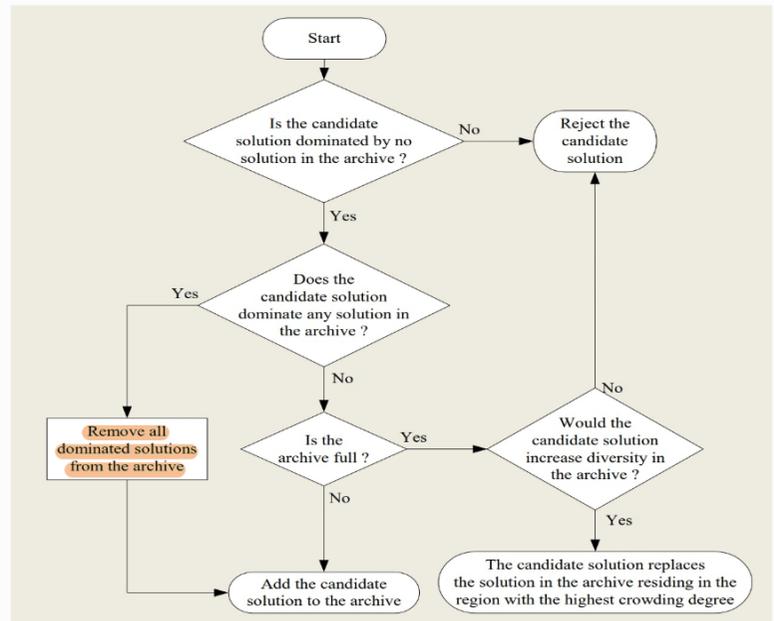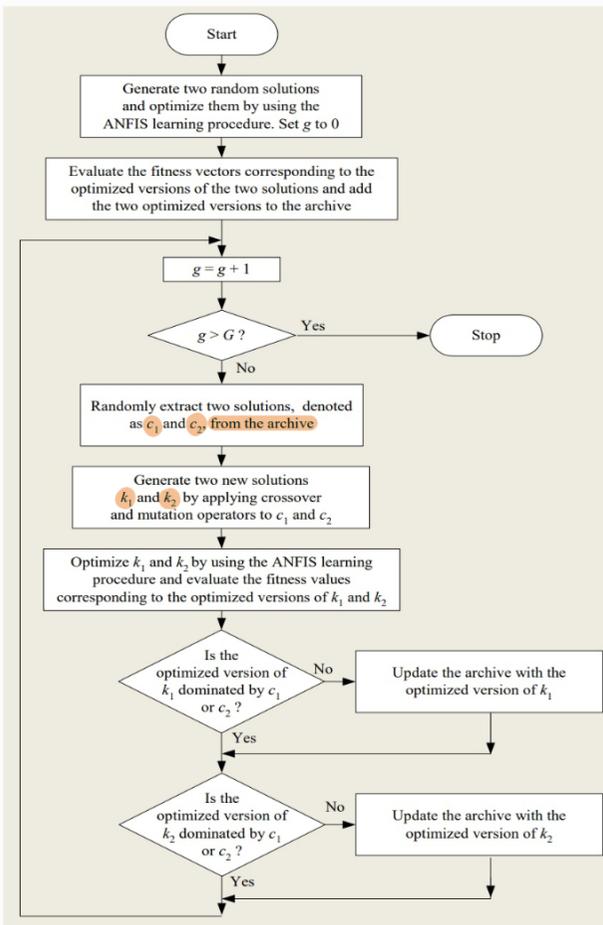(a)                                      (b)     **QUADTREE** → *memory efficient*

→ In n-D I would have need $2^n$ split in the tree

- Then concat the binary strings obtained and use it as index for a **map** containing the crowding number
  → In case of known obj. space (eg. ROC space) we could use static grids to reduce complexity
- This CD is more efficient than standard one

**(2-2) PAES**

**SMS-EMOA** → **S**-**M**etric **S**election-based **E**volutionary **M**ulti-**O**bjective **A**lgorithm
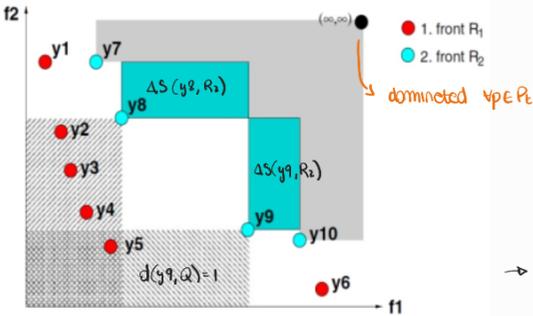
· Hyper Volume is called **S** measure → several definitions (exact or approximate)

· it is an *elitist* algorithm    $S(P_t) \leq S(P_{t+1})$

Lesbegue integral

$$S_2(P) = \Lambda\left(\left\{q \in \mathbb{R}^n : (\exists p \in P : p \prec q) \wedge (q \prec z)\right\}\right)$$

Maximize the Area between the pareto frontier and the Nadir Point

---

**Algorithm 1. SMS-EMOA**

1: $P_0 \leftarrow \text{init}()$                                          /* Initialise random population of $\mu$ individuals */
2: $t \leftarrow 0$
3: **repeat**
4:   $q_{t+1} \leftarrow \text{generate}(P_t)$    → **steady state**        /* generate offspring by variation */
5:   $P_{t+1} \leftarrow \text{Reduce}(P_t \cup \{q_{t+1}\})$              /* select $\mu$ best individuals */
6:   $t \leftarrow t+1$
7: **until** termination condition fulfilled

---

Point that if removed would affect S measure less

**1)** **Algorithm 2. Reduce(Q)**

1: $\{\mathcal{R}_1, \ldots, \mathcal{R}_v\} \leftarrow$ fast-nondominated-sort$(Q)$    /* all $v$ fronts of $Q$ */
2: $r \leftarrow \text{argmin}_{s \in \mathcal{R}_v}[\Delta_{\mathscr{S}}(s, \mathcal{R}_v)]$    $\Delta S = S(R_v) - S(R_v \setminus \{s\})$    /* $s \in \mathcal{R}_v$ with lowest $\Delta_{\mathscr{S}}(s, \mathcal{R}_v)$ */
3: **return** $(Q \setminus \{r\})$    ← last point    → $S(P_{t+1}) \geq S(P_t)$    /* eliminate detected element */

**2)** **Algorithm 3. Reduce(Q)**

1: $\{\mathcal{R}_1, \ldots, \mathcal{R}_v\} \leftarrow$ nondominated-sort$(Q)$    /* all $v$ fronts of $Q$ */
2: **if** $v > 1$ **then**
3:   $r \leftarrow \text{argmax}_{s \in \mathcal{R}_v}[d(s, Q)]$    /* $s \in \mathcal{R}_v$ with highest $d(s, Q)$ */
4: **else**
5:   $r \leftarrow \text{argmin}_{s \in \mathcal{R}_1}[\Delta_{\mathscr{S}}(s, \mathcal{R}_1)]$    /* $s \in \mathcal{R}_1$ with lowest $\Delta_{\mathscr{S}}(s, \mathcal{R}_1)$ */
6: **end if**
7: **return** $(Q \setminus \{r\})$    /* eliminate detected element */

1) more computationally efficient

2) keep points that cover the gaps of prev fronts

$$d(s, Q) = |\{y \in Q : y \prec s\}|$$    number of point in Q that dominates $s$
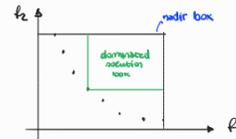- more computationally efficient    ($d(s,Q)$ not usable when $v=1$)

$$\Delta S \text{ function } \in O(\mu^3)$$
- maximize both the vicinity to the pareto point and also the sparseness of the solution



$d(y9, Q) = 1$
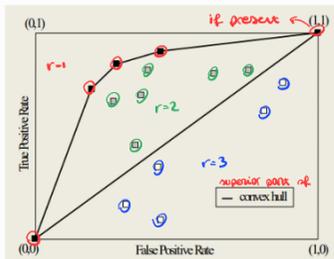
→ with first method I would prefer y8 with second y9



**HYPE** → Use a montecarlo approach to compute S
- generation of random points inside nadir point box
- compute the ratio between points inside the dominated solution box and the total
- it gives an idea of the importance of the solution

nadir box
dominated solution box

**FV-MOEA** → it is exact but faster then SMS-EMOA, since to compute $\Delta S(P, Q)$ we need just the neighbour of $P$, not the all set

---

**CHEA** (convex-hull based algorithm)



convex hull: smallest convex set containing all the binary classifiers

1) Generate random initial population of size $N_{pop}$, $P_0$ and evaluate TPR, FPR    $N_{pop}$ different classifiers
· create offsprings $Q_0$ from $P_0$ (random selection, crossover + mutation) of size $N_{pop}$

2) Evaluate FPR and TPR $\forall q \in Q_k$
· build  $R_k \leftarrow P_k \cup Q_k$    may not be present in $R_k$
· compute the convex hull of  $R_k \cup \{(0,0), (1,1)\}$

3) assign a rank (1,2,3) to each solution and call $N_1, N_2, N_3$ the # of points with i rank
rank 1 → points in the convex hull → they are also non-dominated
rank 2 → other candidate classifiers (above random line)
rank 3 → non-candidate classifiers

4) Generate $P_{k+1}$
- if $N_1 > N_{pop}$ → randomly select, otherwise add all
- if $N_1 < N_{pop}$  if $N_2 > N_{pop} - N_1$ → randomly select, otherwise merge all
- randomly extract $N_{pop} - N_1 - N_2$ from rank 3 and merge it, if still missing points

5) if $K = K_{max}$ → stop, otherwise create $Q_{k+1}$ offspring of size $N_{pop}$ from $P_{k+1}$ using binary tournament based on ranks and apply crossover + mutation
- $K \leftarrow K+1$ and goto step 2

---

**MOEA\D** → based on the decomposition into $N$ single-objective sub-problems

· problems are solved collaboratively → best solution found so far by one can affect the best solution found by another one in its neighborhood
· it uses both a population of size $N$ and an unbounded archive of non-dominated solutions found so far → elitist algorithm
  ↳ optional, get more than $N$ solutions

It uses $N$ lambda vectors (or weight vector) provided by the user

→ in this way I will get N solutions ∈ ES
→ WS (weighted sum scalarization) works well in convex problems



weighted sum scalarization

$n$ different sub-problems

$$\begin{cases} \min & g^{ws}(x, \lambda^i) = \lambda^i \cdot \vec{F}(x) = \lambda^i_1 F_1(x) + \ldots \\ \text{where} & \sum \lambda^i_i = 1 \quad \wedge \quad \lambda^i_i \geq 0 \quad \forall i \end{cases}$$

Chebyshev Scalarization (alternative to WS)

$$\begin{cases} \min_{x \in \Omega} g^c(x \mid \lambda, z) = \min_{x \in \Omega} \max_{1 \le i \le m} \{ \lambda_i \mid f_i(x) - z_i \mid \} = \min_{x \in \Omega} \| \lambda \circ (f(x) - z) \|_\infty \\ z_i = \min_{x \in \Omega} f_i(x) \\ \Sigma \lambda_i = 1 \wedge \lambda_i \ge 0 \end{cases}$$

this is not differentiable
each term is continuous wrt x (when divided ||)

- for any $x^*$ pareto optimal solution, $\exists \lambda : x^*$ is optimal for the above problem → · I can find any ideal point
  → $g^c(x \mid z, \lambda)$ is non-smooth wrt x and $\lambda$, so gradient method doesn't work
  · works also for non-convex problems

The neighbour of i-th subproblem are the ones with similar weight vectors (that identifies the search direction)
  - they have similar obj. functions and with high probability a similar optimal solution

At each generation and sub-problem:
  1) Obtain the current solutions from some neighbours (T-1)
  2) Generate new solution by applying crossover + mutation from (its own solution and) the borrowed ones
  3) - Replace its own solution with the newer one if improve its objective
     - Pass the new solution to some neighbours that will replace if it improves their objective
       ↳ T-1

1) Initialization
  - compute $L_2$ distances between any pair of weight vectors and determine the T closest to each $\lambda_i$. $B(i) = \{ i_1, \ldots, i_T \}$
     neighbourhood size                                                                                      contains i
  - randomly generate the initial population and evaluate $F(x_i)$ $\forall i \in 1, \ldots, N$ (for each individuals)
       ↳ N = pop.size
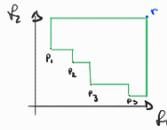  - initialize z as $z_j = \min_{1 \le i \le N} f_j(x_i)$

2) UPDATE
  For i=1 to N do:     N agents
    - randomly select e and k in $B(i)$ and generate $X_{cried}$ from $X_e$ and $X_k$ with crossover + mutation
    - repair with some heuristic if $X_{cried}$ violates the constraints
    - evaluate $F(X_{cried})$
    - $\forall j \in 1, \ldots, m$ if $z_j > f_j(X_{cried})$ update it
    - $\forall j \in B(i)$ if $g^c(X_{cried} \mid \lambda_j, z) < g^c(X_j \mid \lambda_j, z)$ update it (Neighbour/self update)
    - here we could maintain an archive of non-dominated solutions

3) STOPPING CRITERION → if stop criterion satisfied return $(\{x_1, \ldots, x_N\}, \{F(x_1), \ldots, F(x_N)\})$ else step 2

MOEA performance comparison → can be done in different ways

1) Hyper-Volume → exactly or by using montecarlo methods (we need a reference nadir point)

2) GD - generational distance    $GD(A) = \frac{1}{n} \left( \sum_{i=1}^{n} d_i^p \right)^{\frac{1}{p}}$    where $A = \{a_1, \ldots, a_n\}$ non-dominated objective vector set
                                        avg distance of A to Z                    $Z = \{z_1, \ldots, z_m\}$ the reference set
                                        ↳ $L_2$ distance between                  $p \in \mathbb{N}$
                                          $a_i$ and nearest $z_j \in Z$

3) IGD - inverted generational distance    $IGD(A) = \frac{1}{m} \left( \sum_{i=1}^{m} d_i^p \right)^{\frac{1}{p}}$    · how well-spreaded are the solutions wrt reference points
                                        ↳ $L_2$ distance from $z_i$ and NN in A

4) Delta metric    $\Delta(A) = \max \{ GD(A), IGD(A) \}$

low GD
high IGD
↳ non covered zone

- low IGD
- high GD

# MANY OBJECTIVE OPTIMIZATION  (from 4 to 30, if >30 then massive-objective)

Problems: 1) most of the solutions are non-dominated → no solutions to eliminate

2) difficult to use crowding distance or hypervolume because becomes computationally heavy (and curse of dimensionality)

3) difficult to visualize the ES and needs exponentially more points to be covered

4) Ineffective recombination operations since two distant solutions tend to generate an offspring distant from both, not inheriting beneficial characteristics
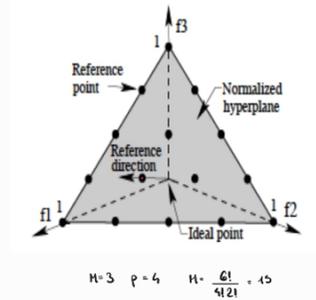
## NSGA-III   · based on a set of reference points provided by the user



Reference point
Normalized hyperplane
Reference direction
Ideal point

$M=3$   $p=4$   $H=\dfrac{6!}{4!2!}=15$

0) Generation of reference points (optional)  → just at start  If not provided we can generate them on the normalized hyper-plane

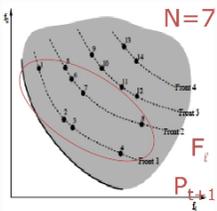#reference points $=\binom{M+p-1}{p}$   $M \to$ # of obj functions
typically = pop.size   $p \to$ # of divisions per edge of the simplex
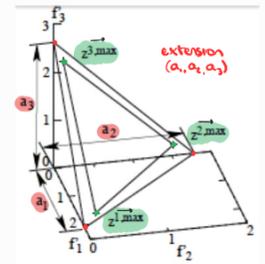
Das-Dennis Algorithm

· each point is associated with a reference (search) direction

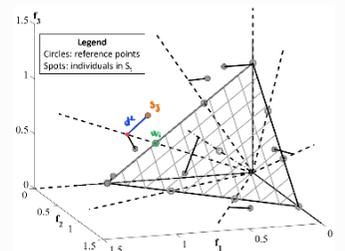## 1) Identifying the non-dominated frontiers   → generation of $P_{t+1}$



N=7

· Parent $(P_t)$ and offsprings $(Q_t)$ are combined into $R_t$   ← crossover + mutation
· non-dominated sorting of $R_t$
· collect and save fronts from 1 to $\ell$ to $S_t$, where $\ell$ is the first integer that $\sum_{i=1}^{\ell} |F_i| \geq N$   (elitism)
· if $|S_t| = N$   → $P_{t+1} = S_t$ (return $P_{t+1}$)
  else $P_{t+1} = \bigcup_{i=1}^{\ell-1} F_i$ then remaining points will be added to $P_{t+1}$ from $F_\ell$ but not using crowding distance (step 2 and 3)

## 2) Normalization of population members   $(|S_t| \neq N)$



extension $(a_1, a_2, a_3)$

· identify ideal point of $S_t$  $(\overrightarrow{z^{min}})$   $z_i^{min} = \text{argmin}_{s \in S_t} f_i(s)$
· translate everything, making the ideal point the origin   $f'_j(s) = f_j(s) - z_j^{min}$
· identify extreme points for each objective (nadir points could be calculated)   $\overrightarrow{z^{i,max}} = f(s^{i,max}),$   $s^{i,max} = \text{argmax}_{s \in S_t} f_i(s)$  ← M extreme vectors
· generate the M-dimensional hyper-plane with this M extreme vectors, compute $a_j \in \mathbb{R}$ as the intercepts of this plane with the axis
· $f''_j(s) = \dfrac{f_j(s) - z_j^{min}}{a_j} = \dfrac{f'_j(s)}{a_j}$   $\forall x \in S_t, \forall j = 1, \dots, M$
  ↳ can handle different scaled objectives

· the reference points if generated in step 0 are already in this hyperplane, otherwise (provided by user) are projected in $\Pi$

## 3) Association of population members   $(|S_t| \neq N)$



Legend
Circles: reference points
Spots: individuals in $S_t$

Assign each $s_i \in S_t$ to a reference point $\in Z$ using the minimum perpendicular distance

$d^\perp(s, w) = \left\| f''(s) - w^T f''(s) \dfrac{w}{\|w\|} \right\|_2$   ← w versor   assign s to $\pi(s) = \text{argmin}_{w \in Z^r} d^\perp(s, w)$  ← reference points set

· create the vector $\overrightarrow{\pi}$ in this way:   $\overrightarrow{\pi} = \begin{bmatrix} \pi_1 \\ \vdots \\ \pi_{|P_{t+1}|} \\ \vdots \\ \pi_{|S_t|} \end{bmatrix} \begin{matrix} \} s \in P_{t+1} \\ \\ \} s \in F_\ell \end{matrix} \Big\} S_t$

## 4) Compute Niche Count   $\forall$ reference point   $w_i \in Z^r$.   $\rho_i = \sum_{s \in S_t \setminus F_\ell} \mathbb{I}[\pi(s) = w_i]$   ← current $P_{t+1}$

## 5) The Niching Procedure

· NSGA II/III no parameters needed, while MOEA\D requires $T$ and $\overrightarrow{\lambda}$
· NSGA-III uses niche measure while NSGA-II uses crowding distance to select points from $F_\ell$
  ↳ favors better speed since uses information of all $P_{t+1}$ points selected
  - with higher M exponentially more points are needed to cover correctly the space

- If there are constrains there is a binary tournament selection (using $CV(P_i)$ if both unfeasible or random if both feasible)
  for generating the mating pool for $Q_t$   ← constrain violation

N - $|P_{t+1}|$ (# of points to add)

```
Algorithm 4 Niching (K, ρⱼ, π, d, Zʳ, Fₗ)   procedure
Input: K, ρⱼ, π(s ∈ Sₜ), d(s ∈ Sₜ), Zʳ, Fₗ
Output: P_{t+1}  → niche count   → reference points
1:  k = 1
2:  while k ≤ K do
3:    J_min = {j : argmin_{j∈Zʳ} ρⱼ}   → may be more than one
4:    j̄ = random(J_min)
5:    I_j̄ = {s : π(s) = j̄, s ∈ Fₗ}
6:    if I_j̄ ≠ ∅ then
7:      if ρ_j̄ = 0 then   → no p ∈ P_{t+1} associated with reference j̄
8:        P_{t+1} = P_{t+1} ∪ (s : argmin_{s∈I_j̄} d(s))  → because j̄ is under-represented
9:      else                                              ↳ d(s) = d⊥(s, π(s))
10:       P_{t+1} = P_{t+1} ∪ random(I_j̄)
11:     end if
12:     ρ_j̄ = ρ_j̄ + 1, Fₗ = Fₗ\s
13:     k = k + 1
14:   else
15:     Zʳ = Zʳ /{j̄}   → exclude reference point from further calculations
16:   end if
17: end while
```

# PURE LEXICOGRAPHIC MULTI-OBJECTIVE OPTIMIZATION

$\begin{cases} \text{lexmin} \quad f_1(x), \dots, f_n(x) \\ x \in \Omega \end{cases}$ where the importance of minimizing $f_1$ is infinitely more important than $f_2$ and so on

→ the lexicographic order allows to compare vectors $[3,4] \leq_{lex} [3,7]$ → defines a total order

vectors $\qquad$ euclidean scalars

$[3,7] \leq_{lex} [2,8] \iff 3+7\eta \leq_{euc} 2+8\eta$

## PRIORITY CHAINS (PC) → there maybe some priority chains among objectives



$$\min \begin{bmatrix} \text{lexmin}(f_1^{(1)}(x), f_1^{(2)}(x), \dots, f_1^{(p_1)}(x)) \\ \text{lexmin}(f_2^{(1)}(x), f_2^{(2)}(x), \dots, f_2^{(p_2)}(x)) \\ \vdots \\ \text{lexmin}(f_m^{(1)}(x), f_m^{(2)}(x), \dots, f_m^{(p_m)}(x)) \end{bmatrix}$$

$$\min \begin{bmatrix} f_1^{(1)}(x) + \textcircled{1}^{-1}f_1^{(2)}(x) + \dots + \textcircled{1}^{1-p_1}f_1^{(p_1)}(x) \\ f_2^{(1)}(x) + \textcircled{1}^{-1}f_2^{(2)}(x) + \dots + \textcircled{1}^{1-p_2}f_2^{(p_2)}(x) \\ \vdots \\ f_m^{(1)}(x) + \textcircled{1}^{-1}f_m^{(2)}(x) + \dots + \textcircled{1}^{1-p_m}f_m^{(p_m)}(x) \end{bmatrix} \triangleq \min \begin{bmatrix} \underline{f_1(x)} \\ \underline{f_2(x)} \\ \vdots \\ \underline{f_m(x)} \end{bmatrix}$$

Gross scalar $\qquad$ PAN (polynomial algorithmic number) $\quad \Sigma \cdot S(\alpha) \eta^n$

$\underline{f_i(x)}$ is a chain in the directed graph

$f_1^1$ may be the profit

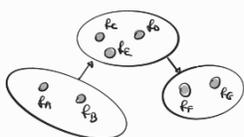$f_1^2$ may be the profit from sustainable investments → secondary objective

$\textcircled{1}$ → grossone number : cardinality of $\mathbb{N}$

$\textcircled{1}^{-1}$ → reciprocal of grossone number : smallest number greater than zero

**Pareto lexicographic Dominance** (PC) : A dominates B : $A \prec B \iff \begin{cases} f_i(x_A) \leq f_i(x_B) \quad \forall i=1,\dots,m \\ \exists j : f_j(x_A) < f_j(x_B) \end{cases}$



· I could adapt both NGSA-II and MOEA\D Algorithms

## PRIORITY LEVELS (PL)



level 1 $\qquad$ level P

$$\text{lexmin} \begin{bmatrix} \min\begin{pmatrix} f_1^{\,1}(x) \\ \vdots \\ f_{m_1}^{\,1}(x) \end{pmatrix}, \dots, \min\begin{pmatrix} f_1^{\,P}(x) \\ \vdots \\ f_{m_p}^{\,P}(x) \end{pmatrix} \end{bmatrix} = \min \begin{bmatrix} \begin{pmatrix} f_1^{\,1}(x) \\ \vdots \\ f_{m_1}^{\,1}(x) \end{pmatrix} + \dots + \textcircled{1}^{1-P}\begin{pmatrix} f_1^{\,P}(x) \\ \vdots \\ f_{m_p}^{\,P}(x) \end{pmatrix} \end{bmatrix}$$

- find pareto non-dominated solutions from $L_i$ and those will be used to search on $L_{i+1}$

**PL Dominance** Given A and B belonging to the subfronts $F_i$ and $F_j$ : $A <^{pl} B \iff i < j$

· **SUBFRONTS**: we take the fronts set of $L_1$, for each of them we further split based on the $L_2$ and recursively until $L_p$

$F_i$ front is defined by a gross scalar eg : $i = 2\textcircled{1}^0 + 7\textcircled{1}^{-1} + 3\textcircled{1}^{-2}$ all the solutions belonging to the second front of $L_1$ objectives,

global rank $\qquad$ lev.2 rank $\qquad$ 7th front to $L_2$ objectives ---

**problem**: a posteriori nature of the definition. To compute $F_i$ and $F_j$ I need to find all other individuals

## PL-NSGA-II

determines final rank, called with lev = 0

**Algorithm 1** Priority Levels fast non-dominated sort.

1: /* This function is recursive */
2: /* P is the population, lvl is the current level to consider */
3: **procedure** PL_FAST_NONDOMINATED_SORT(P, lvl)
4: $\quad$ /* Base case of recursion */ $\quad$ eg = -2 if 3 levels are present
5: $\quad$ **if** $lvl < min\_lvl$ **then return** P
6: $\quad$ /* The first iteration also initializes all ranks to 0 */
7: $\quad$ **if** $lvl == 0$ **then**
8: $\quad\quad$ **for all** $p \in P$ **do**
9: $\quad\quad\quad$ $p_{rank} = 0$
10: $\quad$ /* Ranking within the priority level to determine subfronts */
11: $\quad$ $F^{(lvl)} = fast\_nondom\_sort\_in\_level(P, lvl) = [F_A, \dots]$
12: $\quad$ /* Repeat for every subfront found */
13: $\quad$ **for all** $F_i \in F^{(lvl)}$ **do**
14: $\quad\quad$ /* In the next priority level */
15: $\quad\quad$ $F_i = PL\_fast\_nondominated\_sort(F_i, lvl-1)$

return $F^{(lvl)}$

**procedure** FAST_NONDOM_SORT_IN_LEVEL(P, lvl)

1: **procedure** FAST_NONDOM_SORT_IN_LEVEL(P, lvl)
2: $\quad$ **for all** $p \in P$ **do**
3: $\quad\quad$ $S_p = \emptyset$ → points that p dominates
4: $\quad\quad$ $n_p = 0$ → # of points that dominate p
5: $\quad\quad$ **for all** $q \in P$ **do**
6: $\quad\quad\quad$ /* non-dominance at specified priority level */
7: $\quad\quad\quad$ **if** $p \prec^{lvl} q$ **then**
8: $\quad\quad\quad\quad$ $S_p = S_p \cup \{q\}$
9: $\quad\quad\quad$ **else if** $q \prec^{lvl} p$ **then**
10: $\quad\quad\quad\quad$ $n_p = n_p + 1$
11: $\quad\quad$ **if** $n_p == 0$ **then**
12: $\quad\quad\quad$ /* First subfront within P */
13: $\quad\quad\quad$ $p_{rank} = p_{rank} + \textcircled{1}^{lvl}$ $\quad$ first sub-front set
14: $\quad\quad\quad$ $F_1 = F_1 \cup \{p\}$
15: $\quad$ /* Start from the first subfront */
16: $\quad$ $i = \textcircled{1}^{lvl}$
17: $\quad$ **while** $F_i \neq \emptyset$ **do** → i-th subfront set accumulator
18: $\quad\quad$ $Q = \emptyset$
19: $\quad\quad$ **for all** $p \in F_i$ **do**
20: $\quad\quad\quad$ **for all** $q \in S_p$ **do**
21: $\quad\quad\quad\quad$ $n_q = n_q - 1$ → because we removed p
22: $\quad\quad\quad\quad$ **if** $n_q == 0$ **then**
23: $\quad\quad\quad\quad\quad$ /* q belongs to the i-th subfront */
24: $\quad\quad\quad\quad\quad$ $q_{rank} = q_{rank} + i + \textcircled{1}^{lvl}$
25: $\quad\quad\quad\quad\quad$ $Q = Q \cup \{q\}$ → next sub front
26: $\quad\quad$ /* Index of the next subfront */
27: $\quad\quad$ $i = i + \textcircled{1}^{lvl}$
28: $\quad\quad$ $F_i = Q$

return $[F_{\textcircled{1}^{lvl}}, \dots, F_i]$

$cd = cd_1 \textcircled{1}^0 + \dots + cd_p \textcircled{1}^{1-p}$

**Algorithm 2** Priority Levels crowding distance assignment.

1: /* F is a leaf-front in the hierarchy of population fronts partitioning */
2: **procedure** PL_CROWDING_DIST_ASSIGNMENT(F)
3: $\quad$ $n = |F|$
4: $\quad$ **for all** $i \in F$ **do**
5: $\quad\quad$ $F[i]_{dist} = 0$
6: $\quad$ /* For each level of priority q; p is the index of the last level */
7: $\quad$ **for** $q = 1 \dots p$ **do**
8: $\quad\quad$ **for** $j = 1 \dots m_q$ **do**
9: $\quad\quad\quad$ $F = \text{sort}\left(F, f_j^{(q)}\right)$
10: $\quad\quad\quad$ /* +Inf means "full-scale" (IEEE 754 standard) */
11: $\quad\quad\quad$ $F[1]_{dist} += +\text{Inf } \textcircled{1}^{1-q}$
12: $\quad\quad\quad$ $F[n]_{dist} += +\text{Inf } \textcircled{1}^{1-q}$
13: $\quad\quad\quad$ **for** $i = 2 \dots (n-1)$ **do**
14: $\quad\quad\quad\quad$ /* PL_crowd_dist_ass. has infinitesimal parts */
15: $\quad\quad\quad\quad$ $F[i]_{dist} += \textcircled{1}^{1-q} \frac{f_j^{(q)}(F[i+1]) - f_j^{(q)}(F[i-1])}{f_j^{(q)max} - f_j^{(q)min}}$

→ level 1 rank will have $i \in \mathbb{N}$

**Algorithm 3** PL-NSGA-II algorithm.

1: /* $P_0$ is the starting population, $T$ is the total number of iterations */
2: /* The function returns the approximated Pareto front $P_T$ */
3: **procedure** PL_NSGA_II $(P_0, T)$
4:    **for** $t = 0 \ldots T - 1$ **do**
5:       $Q_t = make\_new\_pop(P_t)$
6:       $R_t = P_t \cup Q_t$
7:       /* $F$ is the set of all the non-dominated subfronts */
8:       $F = PL\_fast\_nondom\_sort(R_t, 0)$ = $[F_a, \cdots]$
9:       $P_{t+1} = \emptyset$
10:      /* $i$ is a gross-index, a.k.a. gross-scalar */
11:     ~~$i = 1$~~    $i \leftarrow first\_index(F)$    $(1 + 0^{-1} + \cdots + 10^{-9})$
12:      **while** $|P_{t+1}| + |F_i| \leq N$ **do**
13:         /* Crowding distance is computed ==within subfront== $F_i$ */
14:         PL_crowding_dist_assignment$(F_i)$
15:         $P_{t+1} = P_{t+1} \cup F_i$
16:         /* Move to the next subfront */     I may increase in order:
17:         $i = next\_index(F, i)$   → eg.   $i = 20^0 + 30^{-1} + 70^{-2}$
18:      Sort$(F_i, ~~All~~$ by CD$)$
19:      $P_{t+1} = P_{t+1} \cup F_i[1 : (N - [P_{t+1}])]$
   **return** $P_T$

*(handwritten, green, top right)* Real word Application

- Vehicle crash worthiness    $lexmin \left[ min \begin{pmatrix} mass(x) \\ toe(x) \end{pmatrix}, min \begin{pmatrix} toe(x) \\ -acc(x) \end{pmatrix} \right]$

*(note)* I caed pt $-acc(x)$ if a company wants
   ↳ intrusion resistance

- Aircraft design : 10 objectives that caued be danded by importance
      (eg. noise, mass, ... )

**PC-NSGA-II** → It is completely identical to NSGA-II but substitute $f_i$ with $\underline{f_i}$ in CD
      $F[1]_{dist} = ①$

**PC-MOEA-D**

1: **procedure** PC_MOEA/D
2:   $EP = \emptyset$
3:   **for** $i = 1 \ldots N$ **do**
4:      $\lambda^{i_1}, \ldots, \lambda^{i_T} = find\_closest\_weights(\lambda_i, T)$
5:      $B_i = \{i_1, \ldots, i_T\}$
6:   $x^1, \ldots, x^N = initialize\_population(N)$
7:   ==**for** $i = 1 \ldots N$ **do**==
8:      $FV^i = F(x^i)$
9:   $\underline{z}_1, \ldots, \underline{z}_m = initialize\_ref\_point(m)$
10:   **while** $stop\_criteria() = False$ **do**
11:      **for** $i = 1 \ldots N$ **do**
12:         $k, l = rand(B_i, 2)$
13:         $y = make\_new\_sol(x^k, x^l)$
14:         $y' = mutate(y)$
15:         **for** $j = 1 \ldots m$ **do**
16:            **if** $\underline{z}_j < \underline{f}_j(y')$ **then**
17:              $\underline{z}_j = \underline{f}_j(y')$
18:         **for all** $j \in B_i$ **do**
19:            **if** $g^{te}(y'|\lambda^j, z) \leq g^{te}(x^j|\lambda^j, z)$ **then**
20:              $x^j = y'$
21: *(archive →)*              $FV^j = F(y')$
22:         $DD = PC\_find\_dominated\_set(EP, F(y'))$ *dominated by sol*
23:         $DG = PC\_find\_dominating\_set(EP, F(y'))$ *if no one dominate $y'$*
24:         $EP = EP \setminus DD$
25:         **if** $DG = \emptyset$ **then** *add to archive*
26:            $EP = EP \cup F(y')$

# PARALLELIZATION of MOEA → typically is really simple to // MOEA


coordinator
③ compare results and compute $P_{t+1}$ and $Q_{t+1}$
CPU 1   CPU 2
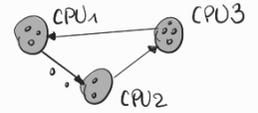
## COORDINATOR WORKERS APPROACH
- a coordinator generates the offspring population and distributes the computation to the workers
- the workers compute the solution of the current generation that are collected by the coordinator that will generate the next generation

} needs coordination

## THE ISLAND MODEL
- decentralized (if a CPU fails we can continue)   - population is spreaded
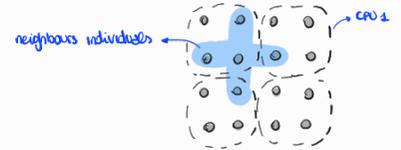- On each CPU (island) evolve a sub-population independently
- After k generations shares the best solutions with other islands → random or determined migration


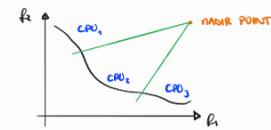CPU1   CPU3
CPU2

## THE CELLULAR MODEL
- decentralized   - population is spreaded
- The individuals are spatially arranged → solution of one migrates only to the neighbours (defined by a pre-det. schema)


neighbours individuals   CPU 1

## CONE SEPARATION MEA (cspMEA)   → is not general purpose as the others, but specific (redefined the NSGA-II algorithm)
- Divide the objective space in n cones   n= #CPU
- usable for 2D/3D problems
- At each population Pt the nadir point change (dynamic)


$f_2$   CPU₁   NADIR POINT   CPU₂   CPU₃   $f_1$

Typically a master-slave approach (or just a decentralized communication)

● → requires inter-process communication

---

## Algorithm 1 Cone-separated NSGA-II

Initialize the different sub-populations
● *Normalize fitness values* → ALL $F_i$ within unit square (or unit hypercube) → NADIR POINT
● *Determine region constraints* → divide 90° in equal parts
Non-dominated sorting → needed for TS in the offspring generation
REPEAT
    Generate Offspring
    IF (migration)
●       *Normalize fitness values*
●       *Determine region constraints*
●       *Migrate individuals violating constraints*
    Non-dominated sorting
    Prune population to original size → via CD
UNTIL stop-condition

CROWDED TS:
1) compare the two front ranks
2) if equal use crowding distance

---

- It is better then random individuals assignments since it exploit the search space

- Local pareto front may be dominated by other front (∉ global front)
- The solutions of the global pareto front may be unevenly distributed
- The migration rate can be really high


dominated   denser   denser

MICRO-CONE SEPARATION PARALLEL MOEA  (mspMEA) → Do not only divide the unit square into n cones, but each cone is divided in Ngroup microcones → each of them a CPU

—when $N_g = 1$ → cspMEA ≈ mspMEA

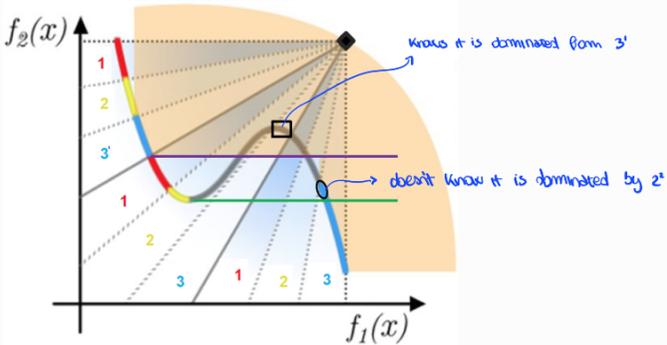· drawback → the offspring generated is more likely to fall outside the microcones of the CPU  (parents may come from different microcones)
- thus the constraint is considered as soft instead of hard

- at each generation a fixed number of random individuals violating the soft constraints are migrated to the right CPU
→ allow to control the migration rate
→ Added the binary tournament to favour recombination of the individuals falling inside the microcones assigned to the CPU
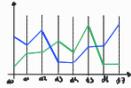


→ Each CPU has an idea on the global front  (glocal front)

1) get Rt and compute global extremes

2) Perform non dominated sort between all microcones of a CPU

3) For last front select iteratively one point for microcone until reached pop-size or no more point

  To select the point for each microcone:
   a) If ∃ points inside ⇒ use crowding distance
   b) If ∄ points inside ⇒ use microcone distance (angular)

4) Once pop-size reached perform the random migration
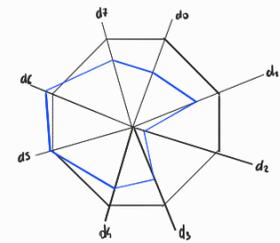
---

VISUALIZE DATA in HIGH DIMENSION

1) Parallel coordinate



→ Two non-dominated solutions

·If they intersect they are non-dominated

2) Radar Chart



3) Scatter plot matrix → triangular matrix of plots  (dim$_i$ vs dim$_j$)

4) Nested Axes Plot for 4 and 5D → works particularly well if some are are discrete



5) Heatmap → similar to // coordinates but instead of lines the points are represented with different coloured segment
- no dominance can be seen

6) Bubble chart → 4D = 2D + bubble size + colour

7) Self Organizing Maps (SOM) → NN that lower the dimension from M to 2
  - count how many points fell inside each cell → visualize them in a 2D grid
  - their problems is that close cells could map very distant zone of a manifold if it overlap
    → iSOM has been introduced for efficient sets → near cells maps for near points in the original manifold

8) RadViz → normalize the obj between (0,1), map each obj to the vertex of a polygon
  - the single point position can be understand imagine it connected by all the obj vertexes by a spring



3D RadViz → adds a dimension as the distance of the solution to the hyperplane
  $(1,0,\dots), (0,1,\dots), \dots$

  -in this way we understand the dominated/non-dominated solutions

# QUANTUM COMPUTING

A qubit is a physical entity (electron, atom, ...) that can be at the same time in two distinct states

→ It is in a superposition of two states
- Decay when measured → it destroys superposition

$\Psi$ function measures the probability the electron is in a specific point in space at a given time
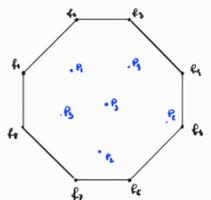→ follows a PDE called Shrödinger Equation → wave equation

state     pure state
$|\Psi\rangle = \alpha|0\rangle + \beta|1\rangle$
- $\alpha, \beta \in \mathbb{C}$, $\alpha^2 + \beta^2 = 1$
- $\alpha^2$ is the probability of being in state $|0\rangle$
- super position can be represented using a sphere

- Quantum Physics can be axiomatized

When we have 2 qubits (4 pure states) they can also entangle
→ the two are not independent ↝ some superstate is preferred

> two can prepared H-T, if I see one I immediately know the second, the outcome was decided at start

Different from classical correlation (when two states have common causes and knowing one, reveal the other)
1) Violates local realism (Bell Inequalities) → object influenced only by it's surroundings
2) Instantaneous but infos max at speed of light
3) Measurement has an effect

Since Shrödinger is a wave equation then exists interference between paths

→ eg. Mac-Zender Interferometer or double-slit
→ used in quantum computing to increase the probability
of finding the qubits in the correct final state for the problem

Tunneling → when a particle cross on energy barrier even with not sufficient energy (thanks to its wave nature)
1) In optimization can be used to escape local minima (by a thermal jump)
2) In sensing can be used to denoise the signal and reach the true one

## ADIABATIC QUANTUM COMPUTERS
→ some NP-hard problems could be solved in polynomial time

Problems are written in terms of the Ising Model (spin ↑ or ↓)
- each possible arrangement is a candidate solution
- total energy of the arrangement → quality of the solution

Condition: the energy needs to change slowly and smoothly → ADIABATIC
- if rush the system could jump in the wrong state

- Quite tolerant to noise wrt other implementations

Problem is solved via Quantum Annealing
1) start via a simple land-scape with obvious global minimum
2) slowly reshape the landscape until the one that encode the problem → quantum tunnelling helps to not get stuck

Applications     1) Quantum NN     2) Quantum Kernel Methods     3) QRL     4) Q Generative Models

**MARKOV PROCESS** → Stocastic Process that satisfies the markov property



- markovian property: $P(S_t = \bar{s} \mid S_{t-1}, \ldots, S_1) = P(S_t = \bar{s} \mid \bar{s}_{t-1})$
- can be defined a transition matrix $P \in \mathbb{R}^{|S| \times |S|}$ where $P_{ij} = P(S_i = i \mid S_t = j) \geq 0$, $\sum_{i=1}^{|S|} P_{ij} = 1 \quad \forall j$
- Starting distribution $s^{(0)} \in \mathbb{R}^{|S|}$ (eg $s^{(0)} = [1, 0, 0]^T$)

$s^{(t+1)} = P s^{(t)} \quad \rightarrow \quad s^{(t+1)} = P^t s^{(0)}$

- The stationary distribution $s^*$ satisfy: $s^* = P s^* \quad \Rightarrow \quad s^*$ is an eigenvector of $P$ associated with $\lambda = 1$
  - $s^*$ represents the probability of being in each state in the long run

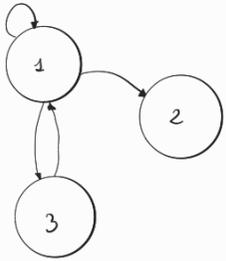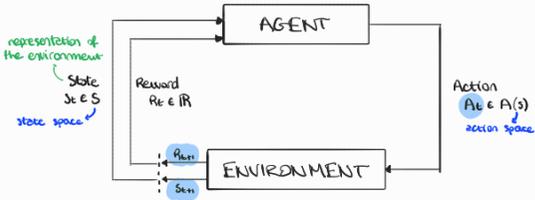Ergodic Theorem: If a chain is irreducible (can reach any state from any other) and a-periodic (not locked in cycles)

State space $S = \{1, 2, 3\}$

$\Rightarrow$ 1) $\exists! \, s^*$: $s^* = P s^*$
2) $s^{(t)} \to s^*$ regardless $s^{(0)}$

Implications
1) Sampling long run episodes estimates $\mathbb{E}_{s^*}$
2) Policy gradient exists uniquely $\left[ \nabla_\theta J \propto \sum_s d_{\pi_\theta}(s) \sum_a \nabla_\theta \pi_\theta(a|s) Q(a,s) \right]$ and independently from $s^{(0)}$

---

**RL** → branch of ML where an agent learns to take actions within an environment, by maximizing the rewards obtained



representation of the environment
State $S_t \in S$ ← state space
Reward $R_t \in \mathbb{R}$
Action $A_t \in A(s)$ ← action space

$p(s', r \mid s, a) \doteq P(S_t = s', R_t = r \mid S_{t-1} = s, A_{t-1} = a)$ ← depends only on last state and action

dynamics function $p: S \times R \times S \times A \to [0,1]$
$\sum_{s' \in S} \sum_{r \in R} p(s', r \mid s, a) = 1 \quad \forall s, a \in S \times A(s)$

- To satisfy the markov property $S$ must carry all the important historical infos

$p(s' \mid s, a) \doteq P(S_t = s' \mid S_{t-1} = s, A_{t-1} = a) = \sum_{r \in R} p(s', r \mid s, a)$
state transition function $p: S \times S \times A \to [0,1]$

$r(s, a) \doteq \mathbb{E}[R_t \mid S_{t-1} = s, A_{t-1} = a] = \sum_{r \in R} r \underbrace{\sum_{s' \in S} p(s', r \mid s, a)}_{p(r \mid s, a)}$
reward function $r: S \times A \to \mathbb{R}$

---

The objective is to maximize the cumulative reward it receives in the long run

simplest one $\left\{ \begin{array}{l} G_t \doteq R_{t+1} + R_{t+2} + \cdots + R_T \\ \text{final time step (could be finite eg. TicTacToe or } T = +\infty) \\ \text{return (must be a function of the reward sequence)} \end{array} \right.$

$G_t \doteq R_{t+1} + \gamma R_{t+2} + \cdots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$ ← consider absorbing final state $t > T$ with $R_t = 0$
discount rate $\gamma \in [0,1]$

1) $R_{t+k}$ is worth only $\gamma^{k-1} R_{t+k}$ what it would if received immediately
2) If $\gamma < 1 \wedge \{R_t\}$ bounded $\Rightarrow$ $G_t$ has a finite value
   - If $\gamma = 0$ the agent is "myopic", trying to maximize only immediate rewards
3) $G_t = R_{t+1} + \gamma(R_{t+2} + \gamma R_{t+3} + \cdots) = R_{t+1} + \gamma G_{t+1}$

---

**Policy**: $\pi(a|s) = P(A_t = a \mid S_t = s)$     If agent is following policy $\pi \Rightarrow p(a|s) = \pi(a|s)$

$r(s, \pi) \doteq r_s^\pi = \sum_{a \in A} \pi(a|s) \, r_s^a = \mathbb{E}_\pi[r_s^a]$     $P_{ss'}^\pi = \sum_{a \in A} \pi(a|s) \, P_{ss'}^a = \mathbb{E}_\pi[P_{ss'}^a]$
$\hookrightarrow P_{ss'}^a = p(s'|s, a)$

**Value function**: $v_\pi(s) = \mathbb{E}_\pi[G_t \mid S_t = s] = \mathbb{E}_\pi\left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right]$

$p(a_t, r_{t+1}, s_{t+1}, \ldots, s_T, r_T \mid s_t) = \prod_{k=0}^{\infty} \pi(a_{t+k} \mid s_{t+k}) \cdot p(s_{t+k+1}, r_{t+k+1} \mid s_{t+k}, a_{t+k})$
distribution over trajectories from $s_t$     $p(a_{t+1}, s_{t+k+1}, r_{t+k+1} \mid s_{t+k})$

Expected return when in $s$ and follow the $\pi$ policy
- how good is to be in that state, it doesn't inform on what action is better
- $v_\pi$ and $q_\pi$ can be estimated from experience → Monte Carlo Approximation
- $q_\pi$ says also what is the best action to take

**Action-Value function** $q_\pi(a, s) \doteq \mathbb{E}_\pi[G_t \mid S_t = s, A_t = a]$     after $a$, follow $\pi$ policy

$v_\pi(s) = \mathbb{E}_\pi[q_\pi(s, a)] = \sum_{a \in A} \pi(a|s) \, q_\pi(s, a)$

**Markov Decision Process** $\text{MDP} = \langle S, A, P, r, \gamma \rangle$     ← state space
$\hookrightarrow$ can be seen as an extension of Markov chains with $a_t$ and $r_t$

- model that formalizes decision making in scenarios where outcomes are potentially random and partly under control of a decision maker → MDP is the environment
- follows the markov property     - solving a MDP means to find an optimal policy $\pi^*$

**Agent** $\langle \pi, v_\pi \rangle$ or $\langle \pi, q_\pi \rangle$     - evaluate and improve $\pi$ under a MDP

---

**POLICY Evaluation** → how worth a policy $\pi$ is     - since $v_\pi$ is unknown we can keep track of rewards achieved through time following $\pi$, in this way we can estimate $v_\pi$ by taking the mean of $G_t$ following a visit to $s$

total expectation $\mathbb{E}[X] = \mathbb{E}[\mathbb{E}(X|A)]$

$v_\pi(s) \doteq \mathbb{E}_\pi[G_t \mid S_t = s] = \mathbb{E}_\pi[R_{t+1} + \gamma G_{t+1} \mid S_t = s] = \mathbb{E}_\pi[R_{t+1} + \gamma \mathbb{E}_\pi[G_{t+1} \mid S_{t+1}] \mid S_t = s] = \mathbb{E}_\pi[R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s]$

$= \sum_a \pi(a|s) \sum_{s', r} p(s', r \mid s, a) \left[ r + \gamma v_\pi(s') \right] = \sum_a \pi(a|s) \left[ r(s, a) + \gamma \sum_{s'} p(s'|s, a) v_\pi(s') \right]$

**BELLMAN EQUATIONS**

$= \sum_{a \in A} \pi(a|s) \left[ R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a v_\pi(s') \right] = R_s^\pi + \sum_{s'} P_{ss'}^\pi v_\pi(s')$

$Q_\pi(s_t, a_t) = \mathbb{E}_\pi[R_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) \mid s_t, a_t] = R_s^a + \gamma \sum_{s \in S} P_{ss'}^a \sum_{a' \in A} \pi(a'|s') Q_\pi(s', a')$
                                                                                                    $v_\pi(s')$

# Closed form solution of Bellman Equation

Can be rewritten in matrix form $V_\pi = R^\pi + \gamma P^\pi V_\pi$ $\Rightarrow$ $V_\pi = (I - \gamma P^\pi)^{-1} R_\pi$

$R^\pi \in \mathbb{R}^{|S|}$, $P^\pi \in \mathbb{R}^{|S| \times |S|}$

$R^\pi$ vector with entries $r_s^\pi$
$P^\pi$ matrix with entries $P_{ss'}^\pi$

Intractable to solve in real world scenarios

**Optimal Policy** A policy $\pi$ is better then $\pi'$ (weakly dominates) $(\pi \geq \pi')$ $\iff$ $V_\pi(s) \geq V_{\pi'}(s)$ $\forall s \in S$

$V_\pi(s) = \max_\pi V(s) = V^*(s)$ $\forall s \in S$

- $\forall MDP$ exists at least one optimal policy $\pi_* = \arg\max_\pi \mathbb{E}_\pi [G_1 | s_1 \sim P(s)] = \arg\max_\pi \mathbb{E}_{s \sim P(s)} [V_\pi(s_1)]$

initial state distribution

**Optimal Value Functions** $V_*(s) = \max_\pi V_\pi(s)$ $\to$ largest expected return achievable from each $s$

$q_*(s,a) = \max_\pi q_\pi(s,a)$

Th: $\forall MDP$ $\exists \pi^*(s)$ a deterministic optimal policy : $\pi^*(a|s) = \begin{cases} 1 & \text{if } a = \arg\max_{a \in A} Q^*(s,a) \\ 0 & \text{otherwise} \end{cases}$ greedy

**Bellman Optimality Equation** non linear

- $V^*(s) = \max_{a \in A} q^*(s,a) = \max_{a \in A} R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V^*(s')$

$\to$ implies determinist optimal policy with max
$\to$ Independent from $\pi$

- $q^*(a,s) = R_s^a + \gamma \sum_{s \in S} P_{ss'}^a V^*(s') = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a \max_{a \in A} q^*(a',s')$

$\to$ These equations are non-linear, so in general $\nexists$ a closed solution $V_*$
· we need ways to approximate $\to$ iteratives approaches

# DYNAMIC PROGRAMMING

$\to$ can be used to solve problems that :

1) can be decomposed in sub-problems $\to$ bellman equations
2) can store and reuse subproblems solutions $\to$ done by value function
   (computation of $\mathbb{E}$ return over a trajectory)

· Assume full MDP knowledge and finite MDP (finite A, S)

① **Policy Iteration** $\to$ alternates policy evaluation and improvement until convergence $\to$ Each iteration is called GPI (generalized policy iteration)

$\pi_0 \xrightarrow{E} V_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} V_{\pi_1} \to \cdots \xrightarrow{E} V_{\pi_*} \xrightarrow{I} \pi_*$

i) Policy Evaluation : Input $S, A, R, P, \pi, \varepsilon$

$\to$ I could compute $V_\pi$ exactly but it is computationally heavy $\to$ iterative method
$O(|S|^3)$ vs. $O(|S|^2 \cdot |A|)$ per iteration

$V_0(s) = 0$ $\forall s \in S$
$K = 0$
while $K == 0 \lor \|V_K - V_{K-1}\| > \varepsilon$ :
$V_3$: $V_{K+1}(s) = \sum_{a \in A} \pi(a|s) [R_s^a + \gamma \sum_{s \in S} P_{ss'}^a V_K(s')]$ $= R_s^{\pi(s)} + \gamma \sum_{s \in S} P_{ss'}^{\pi(s)} V_K(s')$ assumed $\pi(a|s) = \pi(s)$
$K++$
return $V_K$

2) Policy Improvement $\pi'(s) \leftarrow \arg\max_{a \in A} R_s^a + \gamma \sum_{s \in S} P_{ss'}^a V_K(s')$ $\forall s \in S$ computed using old $\pi$

deterministic $\underbrace{}_{q_\pi(a,s)}$

- improve just by looking a step-ahead
$\to$ typically $V_K$ takes more to converge wrt $\pi$
$\to$ do not need $q_\pi$ actually

Th: Policy Improvement lets $\pi' \leftarrow$ greedy $(V_\pi)$ $\land$ $\pi \neq \pi'$ $\Rightarrow$ $V_{\pi'} > V_\pi$ as above

1) $q_\pi(s, \pi'(s)) = \max_{a \in A} q_\pi(s,a) \geq q_\pi(s, \pi(s)) = V_\pi(s)$

in case of equality $V_\pi(s) = \max_{a \in A} Q_\pi(s,a)$ $\to$ $V_\pi = V_{\pi'} = V_{\pi^*}$ $\Rightarrow$ $\pi = \pi' = \pi^*$

2) $V_\pi(s_t) \leq q_\pi(s_t, \pi'(s_t)) = \mathbb{E}[r_{t+1} + \gamma V_\pi(s_{t+1}) | s_t, a_t = \pi'(s_t)] = \mathbb{E}_{\pi'}[r_{t+1} + \gamma V_\pi(s_{t+1}) | s_t] \leq \mathbb{E}_{\pi'}[r_{t+1} + \gamma q_\pi(s_{t+1}, a_{t+1}) | s_t] = \mathbb{E}_{\pi'}[r_{t+1} + \gamma r_{t+2} + \gamma^2 V_\pi(s_{t+2}) | s_t] \leq \cdots \leq V_{\pi'}(s_t)$

Th: Policy Iteration Convergence
- $|S|$ and $|A|$ must be finite
- $R$ must be bounded
- $0 \leq \gamma < 1$

$\Rightarrow$ $\{V_K\}$ converges to $V_*$ $\land$ $\{\pi\} \to \pi_*$

② **Value Iteration** $\to$ do not model $\pi(s)$, computes it only at the end
· uses Bellman Optimality Equation and not Bellman Equation

$V_1 \to V_2 \to \cdots \to V_*$

$V_{K+1}(s) = \max_{a \in A} R_s^a + \gamma \sum_{s \in S} P_{ss'}^a V_K(s')$ and then compute $\pi_* \leftarrow$ greedy $(V_{\pi^*})$

ASYNCRONOUS DP
- some state may not need to have their values updated as often as others
- in place updates of $V$ and $\pi$
- Early stopping $\to$ do not wait for full convergence
} improve performances

PROBLEMS of DP
- high computational effort, scanning all $S$ and $A$ each iteration $O(|S|^2 \cdot |A|)$ $\to$ do not scale with $|S|$
- new approaches must be considered sample-based $\to$ model-free
$\to$ model-based

$\to$ DP method performs bootstrapping and not sampling

# MODEL FREE EVALUATION → algorithms that do not need to know transitions probabilities $p(s'|s,a)$

· episode → list of experiences made between initial and final state (or max # of experiences T)

· experience$_i$ = < State, Action, Reward, Next_State > → sample

## 1) MONTE-CARLO RL

$$V_\pi(s) = \mathbb{E}_\pi \left[ G_t^T \mid S_t = s \right] \qquad G_t^T = \sum_{k=t}^{T} \gamma^{k-t} r_{k+1}$$

In order to estimate $V_\pi$ we have two possibilities:

1) **First-Encounter:** $V_\pi(s)$ is the avg of returns following first visit to s across episodes

$$\hat{V}(s) = \frac{1}{N} \sum_{i=1}^{N} G_i \rightarrow V_\pi(s) \quad \text{as } N \to +\infty \quad \text{for large number law}$$

↳ across N episodes
↳ first encounter return in episode i

2) **Every-Encounter:** consider all the visit to S across a single episode

---

Input: $\pi$, num_episodes

$N(s) \leftarrow 0 \quad \forall s \in S$
$Returns(s) \leftarrow 0 \quad \forall s \in S$

for $e \leftarrow 1$ to num_episodes do
  generate using $\pi$ an episode $S_0, A_0, R_1, \ldots, S_{T-1}, A_{T-1}, R_T, S_T$
  $G \leftarrow 0$
  for $t \leftarrow T-1$ to 0 do:
    $G \leftarrow [\gamma]^* G + R_{t+1}$
    $\left[ \text{if } S_t \notin < S_0, \ldots, S_{t-1} > \text{ then} \right]^*$ → only present in first-encounter
      $Returns(S_t) \leftarrow Returns(S_t) + G$ → sum across episodes
      $N(S_t) \leftarrow N(S_t) + 1$
  $V(s) \leftarrow \frac{Returns(S_t)}{N(S_t)} \quad \forall s \in S$

$\left[ G(s,a) \leftarrow \frac{Returns(s,a)}{N(s,a)} \right]$ ↝ needs two different counters the if is updated with $<S_t, a_t> \notin < S_0, a_0, \ldots, S_{t-1}, a_{t-1} >$

return V

---

· I could use the **incremental mean** to avoid storing $Returns(S_t)$

$$V(S_t) \leftarrow V(S_t) + \frac{1}{N(S_t)} \left[ G - V(S_t) \right]$$

· In **non-stationary MDP** (transitions and rewards changes), V cannot converge and we will use **exponential decay** to forgot about the past $V(S_t) \leftarrow V(S_t) + \gamma (G_t - V(S_t)) \quad \boxed{\gamma \in (0,1)}$
↳ fixed regardless $N(S_t)$ → prioritize recent observations

**Offline learning** because went the end of episode to update V, since I need to compute $G_t$
**Unbiased** since we are computing correctly the expected value
· **High Variance** wrt TD learning

---

## 2) TEMPORAL-DIFFERENCE (TD) LEARN → MC methods needs to went the end of an episode to make an update, meanwhile TD methods only went the next timestep

$$TD(0): \quad V(S_t) \leftarrow V(S_t) + \alpha \left[ R_{t+1} + \gamma V(S_{t+1}) - V(S_t) \right]$$
↳ new estimation for $G_t$

· Converges to optimality with high probability $\iff 0 < \alpha \ll 1$
· samples order affects the result a lot
· TD **converges faster** than MC in practice → **online learning**

· given $\pi$, TD(0) **converges** to $V_\pi \iff$ 1) $\pi$ induces an irreducible Markov chain (reach any state from any other)
2) finite MDP
3) $\sum_{k=1}^{\infty} \alpha_k = +\infty \wedge \sum_{k=1}^{\infty} \alpha_k^2 < +\infty \quad \boxed{\alpha_k = \frac{1}{k}}$

---

Input ($\pi$, $\alpha \in (0,1]$, num_ep, $\gamma$)
Initialize $V(s) \in \mathbb{R}$ randomly except $V(terminal) = 0$
for $e \leftarrow 1$ to num_ep
  $S_0 \leftarrow$ initial state $\quad t \leftarrow 0$
  while $S_t \neq$ terminal
    $a_t \leftarrow \pi(S_t)$
    Take action $a_t$ and observe $(S_{t+1}, r_{t+1})$
    $V(S_t) \leftarrow V(S_t) + \alpha \left[ r_{t+1} + \gamma V(S_{t+1}) - V(S_t) \right]$
    $t \leftarrow t+1$

· **Biased** since use **bootstrapping** (use of estimated values in update) meanwhile MC samples
↳ biased towards the current estimation

$$V_\pi(s) = \mathbb{E}_\pi \left[ G_t \mid S_t = s \right] \quad \leftarrow \text{MC based}$$
↳ true
$$= \mathbb{E}_\pi \left( R_{t+1} + \gamma V_\pi(S_{t+1}) \mid S_t = s \right) \quad \leftarrow \text{TD based}$$
↳ approximation

· **Low variance** since a lot of updates (shorter stochastic chains)
eg. MC: $Var \left[ G_t \mid S_t \right] = Var \left[ R_{t+1} + \gamma R_{t+2} + \ldots \right]$

# DIFFERENCES TD vs. MC

- MC converges to MSE solution $\quad V_\pi = \arg\min_V \sum_{c=1}^{E} \sum_{t=0}^{T_c-1} (G_t - V(s_t))^2$, meanwhile TD(0) to maximum likelihood MDP

- MC doesn't respect Markovian Property → update $V_\pi$ on all the episode

  meanwhile TD(0) does → update $V_\pi$ only on $s_t$ and $s_{t+1}$



MC: $V_\pi(A) = 0$

TD: $V_\pi(A) = \frac{3}{4}$ → as if I first estimate $\hat{P}^a_{ss'}$ and $\hat{R}^a_s$ and then compute $\mathbb{E}[G_t \mid s_t = A]$

| A,0, B,0 | B,1 |
| B,1 | B,1 |
| B,1 | B,1 |
| B,1 | B,0 |

→ because TD(0) assumes and uses the MDP structure in one-step update

$$\begin{cases} \hat{P}^a_{ss'} = \frac{1}{N(s,a)} \sum_{c=1}^{E} \sum_{t=0}^{T_c-1} \mathbb{I}(s_t^c, a_t^c, s_{t+1}^c) \\ \hat{R}^a_s = \frac{1}{N(s,a)} \sum_{c=1}^{E} \sum_{t=0}^{T_c-1} \mathbb{I}(s_t, a_t^c) \, r_t^c \end{cases}$$

# MODEL FREE OPTIMIZATION

## ON-POLICY OPTIMIZATION → There are two types of policies:

1) **Target policy**: used to optimize decision making (typically greedy)
2) **Behavioural policy**: used to navigate the env (could be random or $\varepsilon$-greedy)

In on-policy algorithms the two are equal, meanwhile different in off-policy ones

→ Problems of Greedy approaches for updating the policy in on-policy algorithms: $\pi_{k+1}(a|s) \leftarrow \arg\max_{a \in A} R^a_s + \gamma \sum_{s' \in S} P^a_{ss'} V_{\pi_k}(s')$

1) $R^a_s$ and $P^a_{ss'}$ are unknown → use Q instead of V: $\pi_{k+1}(a|s) = \arg\max_{a \in A} Q_{\pi_k}(s,a)$

2) There is no exploration of the possible actions since only greedy choices are made
   - $\varepsilon$-greedy select with $p = 1-\varepsilon$ the greedy choice and with $p = \varepsilon$ a choice at random
   - Exploration (get infos about the word) - Exploitation (getting rewards) tradeoff

not deterministic

greedy action

$$\pi_{k+1}(a|s) = \begin{cases} \frac{\varepsilon}{|A(s)|} + (1-\varepsilon) & \text{if } a = \arg\max_{a \in A} Q(a,s) \\ \frac{\varepsilon}{|A(s)|} & \text{otherwise} \end{cases}$$

$m = |A(s)|$

**Th: $\varepsilon$-greedy policy improvement**: $\varepsilon$-greedy generates a policy $\pi_{k+1}$ that $V_{\pi_{k+1}} \geq V_{\pi_k}$

$\max_a Q_{\pi_k}(s,a) \geq \sum_a w_a \cdot Q_{\pi_k}(s,a)$ with $\sum w_a = 1 \wedge w_a \geq 0$

Proof: $\mathbb{E}_{\pi_{k+1}}[Q_{\pi_k}(s,a)|s] = \sum_{a \in A} \pi_{k+1}(a|s) Q_{\pi_k}(s,a) = \frac{\varepsilon}{m} \sum_{a \in A} Q_{\pi_k}(s,a) + (1-\varepsilon) \max_{a \in A} Q_{\pi_k}(s,a) \geq \frac{\varepsilon}{m} \sum_{a \in A} Q_{\pi_k}(s,a) + (1-\varepsilon) \sum_{a \in A} \frac{\pi_k(a|s) - \frac{\varepsilon}{m}}{1-\varepsilon} Q_{\pi_k}(s,a) = \sum_{a \in A} \pi_k(a|s) Q_{\pi_k}(s,a) = V_{\pi_k}(s)$

- recursively apply this inequality for policy improvement theorem $V_{\pi_{k+1}} \geq V_{\pi_k}$

## GLIE PROPERTY (Greedy in the Limith with Infinite Exploration):

$\pi$ has GLIE property $\overset{\text{sufficient}}{\Longrightarrow}$ Q-style update rule converges $(\pi^*, Q^*)$

**DEF** a stochastic policy has GLIE property if:

1) $\lim_{K \to +\infty} N_K(s,a) = +\infty$
2) $\lim_{K \to +\infty} \pi_K(a|s) = \mathbb{I}(a = \arg\max_{a' \in A} Q_K(s,a'))$

- each state-action pair is visited infinitely many times
- policy converges to the greedy one → guarantees that GLIE converges to optimum $\pi^*$

→ If $\varepsilon$-greedy has $\varepsilon_K \in \Theta(\frac{1}{K})$ ⇒ it has GLIE property since:

1) $\sum_{K=1}^{\infty} \varepsilon_K = \sum_{K=1}^{\infty} \frac{1}{K} = +\infty$ → infinite exploration
2) since $\varepsilon_K \to 0$, tends to greedy selection

## 1) ON-POLICY MC

```
Initialize Q(s,a) ∈ ℝ arbitrary ∀s,a ∈ S×A,  π_ε(s) ∈ A(s) arbitrarily   (Q ≈ q*)
for e ← 1 to max-ep
    generate an episode following π_ε:  s₀, a₀, r₁, ...  → π both behavioural
                                                            and target policy
    G ← 0
    for t ← T-1 to 0:
        G ← γG + R_{t+1}
        N(s_t, a_t) ← N(s_t, a_t) + 1
        Q(s_t, a_t) ← Q(s_t, a_t) + 1/N(s_t,a_t) [ G - Q(s_t, a_t) ]
    ε ← 1/e
    π_ε(a|s) ← ε-greedy  with Q  ∀a,s
```

## 2) SARSA → similar to TD evaluation on Q instead of V

- sample$_i$ - $\langle S, A, R, S', A' \rangle$ → transitions of $\langle S, A \rangle$ pairs

  Update Rule: $Q(s,a) \leftarrow Q(s,a) + \gamma [ r + \gamma Q(s',a') - Q(s,a) ]$

→ converges (as all the weighted averages methods) to $q^*(s,a)$

just sufficient exploration is needed

$\Longleftrightarrow \sum_{t=1}^{\infty} \gamma_t = +\infty \wedge \sum_{t=1}^{\infty} \gamma_t^2 < +\infty \wedge$ GLIE policy $\pi_\varepsilon(a|s)$

- **EXPECTED SARSA** $\quad Q(s,a) \leftarrow Q(s,a) + \gamma [ r + \gamma \sum_{a' \in A} \pi(a'|s') Q(s',a') - Q(s,a) ]$

$\mathbb{E}_\pi[Q(s',a_{t+1})|s']$

→ fails when $|A|$ is too large to compute

→ prevents the randomness of $a' \sim \pi$ to increase the variance by taking on $\mathbb{E}_{a \sim \pi}$
- has the same bias but lower variance wrt SARSA

↳ more robust updates ⇒ faster convergence

---

Input $(\alpha \in (0,1], \text{ num-ep}, \gamma)$   (estimate $q^*$)

Initialize $Q(s,a) \in \mathbb{R}$ randomly except $Q(\text{terminal}) = 0$

```
for e ← 1 to num-ep
    S₀ ← initial state   t ← 0   a₀ ← π(s₀, Q)
    while s_t ≠ terminal
        Take action a_t and observe (s_{t+1}, r_{t+1})
        a_{t+1} ← π(s_{t+1}, Q) → target policy in this iteration and    [ε-g]
                                   exploratory on the next one
        Q(s_t, a_t) ← Q(s,a) + α( R_{t+1} + γ Q(s_{t+1}, a_{t+1}) - Q(s,a) )
        t ← t+1
```

# OFF-POLICY OPTIMIZATION

## 1) MC OPTIMIZATION

Importance Sampling $\quad \mathbb{E}_p[f(x)] \approx \frac{1}{N} \sum_{i=1}^{N} f(x_i) \quad$ with $x_i \sim f(x)$

$\frac{1}{N} \sum_{i=1}^{N} f(x_i) = S \xrightarrow{\text{large } N} N(\mu, \sigma^2) = \begin{cases} \mu = \mathbb{E}_p[f(x)] \\ \sigma^2 = \frac{1}{N} V_p[f(x)] \end{cases} \rightarrow$ unbiased

$\uparrow$ law of big numbers

$\uparrow$ random variable since changing sample change s

$\mathbb{E}_p[f(x)] = \sum_x p(x) f(x) = \sum_x Q(x) \frac{P(x)}{Q(x)} f(x) = \mathbb{E}_q\left[\frac{P(x)}{Q(x)} f(x)\right] \qquad$ — requirements: $\underline{P(x) > 0, \ q(x) > 0}$  otherwise cannot sample from same areas

$\uparrow$ target policy $\qquad\qquad\qquad \uparrow$ behavioural policy

I can estimate $\mathbb{E}_p$ with samples from $x \sim q \qquad \mathbb{E}_q\left[\underset{w(x)}{\frac{P(x)}{Q(x)}} f(x)\right] \approx \frac{1}{N} \sum_{i=1}^{N} \frac{P(x_i)}{Q(x_i)} f(x_i) = r \rightarrow$ unbiased but different variance $\quad V_q[r] = \frac{1}{N} V_q\left[\frac{P(x)}{Q(x)} f(x)\right]$

There are two samples: 
i) $G_t^\mu$ = return from $\mu$ (behavioural)
ii) $G_t^\pi$ = return of target policy computed by $\mu$ experiences $\quad Q(s_t, a_t) \leftarrow Q(s_t, a_t) + r[G_t^\pi - Q(s_t, a_t)]$

$P(A_t, S_{t+1}, A_{t+1}, \dots, S_T \mid S_t, A_{t:T-1} \sim \pi) = \prod_{K=t}^{T-1} \pi(A_K \mid S_K, \text{past}) \ p(S_{K+1} \mid S_K, A_K, \text{past})$

$\uparrow \ P(x_1, \dots, x_n \mid Y) = P(x_1 \mid Y) \cdot P(x_2 \mid x_1, Y) \cdots$ chain rule

$= \prod_{K=t}^{T-1} \pi(A_K \mid S_K) \ p(S_{K+1} \mid S_K, A_K) \rightarrow$ Markovian property

$\Rightarrow w_{t:T-1} = \dfrac{\prod_{K=t}^{T-1} \pi(A_K \mid S_K) \ p(S_{K+1} \mid S_K, A_K)}{\prod_{K=t}^{T-1} \mu(A_K \mid S_K) \ p(S_{K+1} \mid S_K, A_K)} = \prod_{K=t}^{T-1} \dfrac{\pi(A_K \mid S_K)}{\mu(A_K \mid S_K)}$

$\uparrow$ importance ratio $\qquad\qquad \uparrow$ unknown

$\boxed{\mathbb{E}_\mu(w_{t:T-1} \ G_t^\mu \mid S_t = s) = \mathbb{E}_\pi(G_t^\pi \mid S_t = s) = V_\pi(s)}$

$\rightarrow$ unfortunately $w_t$ and the episodes increases the variance
- $w_t \neq 1 \rightarrow$ when $\pi$ and $\mu$ behave differently
- $\mu$ and $\pi$ should be as close as possible (like $\varepsilon$-greedy and greedy)

using $\mu$ sample on episode $s_0, a_0, r_1, \dots$

$w_T \leftarrow 1, \ G^\mu = G^\pi \leftarrow 0$

for $t \leftarrow T-1$ to $0$

$\quad w_t \leftarrow w_{t+1} \cdot \dfrac{\pi(A_t \mid S_t)}{\mu(A_t \mid S_t)}$

$\quad\quad\quad\quad \uparrow$ typically deterministic

$\quad\quad\quad\quad \downarrow$ typically stochastic

$\quad G^\mu \leftarrow R_{t+1} + \gamma G^\mu$

$\quad G^\pi \leftarrow w_t G^\mu$

$\quad Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [G^\pi - Q(s_t, a_t)]$

$\quad \pi(s_t) \leftarrow \underset{a \in A}{\arg\max} \ Q(s_t, a)$

## 2) TD optimization

$V(s_t) \leftarrow V(s_t) + \nu\left[\underset{G_t^\mu}{\boxed{\frac{\pi(a \mid s)}{\mu(a \mid s)} \left(r_{t+1} + \gamma V(s_{t+1})\right)}} - V(s_t)\right]$

$\uparrow$ depends on $\mu$ for $a$ (influence only $r_{t+1}$) $\qquad \uparrow$ reward $r_{t+1}$ just depends on current $\langle a, s \rangle$ pairs not on all trajectory

$V_\pi(s) = \mathbb{E}_\mu\left(w_t(r_{t+1} + \gamma G_{t+1}^\mu) \mid S_t = s\right)$

$\rightarrow$ scaling by just one factor lower significantly the ==variance==

## 3) Q-LEARNING $\quad \rightarrow$ do not use importance sampling

- $\pi(s) = \underset{a \in A}{\arg\max} \ Q(a, s)$

- $\mu(a \mid s) = \begin{cases} \frac{\varepsilon}{|A(s)|} + (1 - \varepsilon) & \text{if } a = \underset{a \in A}{\arg\max} \ Q_\pi(a, s) \\ \frac{\varepsilon}{|A(s)|} & \text{otherwise} \end{cases}$

$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + r\left[\boxed{r_{t+1}^\mu + \gamma Q(s_{t+1}, a_{t+1}^\pi)} - Q(s_t, a_t)\right]$

$\uparrow$ doesn't depend on $\mu \qquad\qquad \uparrow$ importance sample is not needed because $a_t$ is given

$= r_{t+1}^\mu + \gamma Q(s_{t+1}, \underset{a \in A}{\arg\max} \ Q(s_{t+1}, a'))$

$= r_{t+1}^\mu + \gamma \underset{a \in A}{\max} \ Q(s_{t+1}, a')$

- optimal convergence with GLIE and $\sum \gamma_t = +\infty \wedge \sum \nu_t^2 < +\infty$
- Q-learning is a ==tabular method== since I can express $Q$ by a $\mathbb{R}^{|S| \times |A|}$ matrix

## SARSA vs Q-Learning (on-policy vs off-policy)



SARSA

SARSA waves discourage $Q_t$

$Q(s_t, a_t) \uparrow$

$s_t \ a_t \ s_{t+1} \ a_{t+1}^\pi \quad$ Q-L

$Q(s_{t+1}, a_{t+1}^\mu) \downarrow$

CLIFF R = -100

- Q-learning is learning the actual ==optimal== path (greedy) $\rightarrow$ ==optimistic==

- SARSA is learning the safest path from $\varepsilon$-greedy $\rightarrow$ more ==conservatory== in dangerous situations

  - it learns a ==sub-optimal== policy

# N-STEPS BOOTSTRAPPING → in between MC methods with $n \to +\infty$ and TD(0)

$$G_t = \overbrace{R_{t+1} + \gamma R_{t+2} + \cdots + \gamma^n R_{t+n+1}}^{n+1 \text{ rewards}} + \underbrace{\gamma^{n+1} R_{t+n+2} + \cdots + \gamma^{T-t-1} R_T}_{\gamma^{n+1} V(S_{t+n+1}) \quad \text{bootstrapping}}$$

## 1) TD(n)
$$G_t^{(n)} = \overbrace{\sum_{k=0}^{n-1} \gamma^k R_{t+k+1}}^{n \text{ rewards}} + \gamma^n V(S_{t+n})$$

→ typically the best result is in between $n=1$ and $n=+\infty$

$$G_t^{(1)} = R_{t+1} + \gamma V(S_{t+1})$$
↳ will be the one of $t+n$ steps ahead

· The first time $R_{t+n+1}$ is available is $n$ timesteps forward (online but delay)

$$V(S_t) \leftarrow V(S_t) + v\left[G_t^{(n)} - V(S_t)\right]$$

## 2) TD($\lambda$) → weighted mean of $G_t^{(n)}$ for different $n$
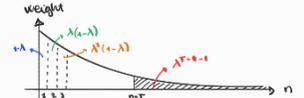
offline

$$G_t^{(\lambda)} = (1-\lambda) \sum_{n=1}^{+\infty} \lambda^{n-1} G_t^{(n)}$$
→ or long N (slide)

→ $\lambda$ controls bias-variance trade-off

$\lambda \in (0,1) \Rightarrow \sum_{n=1}^{+\infty} (1-\lambda)\lambda^{n-1} = 1$

short term rewards    long term rewards

weight

$\lambda=0 \Rightarrow G_t^{(\lambda)} = G_t^{(1)}$    TD(0)

$\lambda=1 \Rightarrow G_t^{(\lambda)} = (1-\lambda)\sum_{n=1}^{T-t-1} \lambda^{n-1} G_t^{(n)} + (1-\lambda)\sum_{n=T-t}^{\infty} \lambda^{n-1} G_t^{(n)} = G_t(1-\lambda)\sum_{m=0}^{+\infty} \lambda^{m+T-t-1} = G_t(1-\lambda)\lambda^{T-t-1} \sum_{m=0}^{+\infty} \lambda^m$

episodic (end at T)    post termination contributions    $= G_t$ because it seen all the timesteps $t:T$    geometric serie

$m = n-T+t$
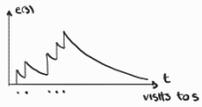
$= G_t(1-\lambda)\lambda^{T-t-1} \cdot \frac{1}{1-\lambda} = G_t \cdot \lambda^{T-t-1} = G_t$    MC

## 3) Backward TD($\lambda$) → implement TD($\lambda$) in an online way

· Associated with each state there is what is called an eligibility trace $e_t(s) \in \mathbb{R}^+$

→ $e(s)$ records which states have been recently visited

At each step update $e_t(s)$ $\forall s \in S$ in this way: $e_t(s) = \begin{cases} \gamma\lambda e_{t-1}(s) & \text{if } s \neq s_t,\ t\neq 0 \\ \gamma\lambda e_{t-1}(s)+1 & \text{if } s=s_t,\ t\neq 0 \\ 0 & \text{if } t=0 \end{cases}$

decay factor

We define $\delta_t = r_{t+1} + \gamma V_t(s_{t+1}) - V_t(s_t)$    The update policy is $V(s) \leftarrow V(s) + v\,\delta_t\,e_t(s)$    $\forall s \in S$

error    ↳ constant $\forall s \in S$

· BTD(0) has $e_t(s) = \mathbb{I}(s=s_t)$, thus equivalent to TD

· meanwhile with BTD(1) you get MC, but better → online + extendable to continuing tasks

→ can be shown that TD($\lambda$) and BTD($\lambda$) are equivalent (with $\gamma\lambda < 1$)

$v\left[\sum_{k=t}^{\infty} (\gamma\lambda)^{k-t} d_k\right] = v\left[G_t^\lambda - V(s_t)\right]$

total update of $V(s_t)$ from BTD($\lambda$)    TD($\lambda$) update of $V(s_t)$

→ At the end of the episode they produced the same updates

```
∀ episode
  e(s) ← 0        ∀s ∈ S
  So ← initialize   t ← 0
  while St ≠ terminal
      a ← π(St)
      take action a and observe (r, St+1)
      δ ← r + γV(St+1) − V(St)
      e(St) ← e(St) + 1
      ∀s ∈ S,
          V(s) ← V(s) + αδe(s)
          e(s) ← γλe(s)
      t ← t+1      ↳ exponential decay
```

## 4) SARSA (n)
Defining $Q_t^{(n)} \leftarrow \sum_{k=0}^{n-1} \gamma^k r_{t+k+1} + \gamma^n Q(s_{t+n+1}, a_{t+n+1})$

We got the update $Q(s,a) \leftarrow Q(s,a) + v\left[Q_t^{(n)} - Q(s,a)\right]$

## 5) SARSA ($\lambda$)
Defining $Q_t^{(\lambda)} \leftarrow (1-\lambda)\sum_{n=1}^{+\infty} \lambda^{n-1} G_t^{(n)}$

we got the update $Q(s,a) \leftarrow Q(s,a) + v\left[Q_t^{(\lambda)} - Q(s,a)\right]$

## 6) Backward SARSA ($\lambda$)
$Q(S,A) \leftarrow Q(S,A) + v\,\delta_t E_t(S,A)$    $\forall s,a \in S \times A(s)$

$\delta_t = r_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)$

$e_t(s,a) = \gamma\lambda e_{t-1}(s,a) + \mathbb{I}(s_t=s, a_t=a)$

# Approximate Solutions Methods

· Tabular methods have difficulties for large spaces (images) or continuous ones
  - they lack the ability to recognize similar states and generalize

No more tabular version of $V$ and $Q$ but we define two **parametric functions**

$$\hat{v}: S \times \Theta \to \mathbb{R} \qquad \hat{Q}: S \times \Theta \to \mathbb{R}^{|A|}$$

→ abilities to generalize from unseen states
→ needs to be **differentiable** for SGD

Two problems
1) **Non-stationary data** → environment may change (thus the policy should change too)
2) **Non-IID data** → correlated events that could cause some problem with the convergence of SGD

MSE: $\mathcal{L}(\Theta) = \mathbb{E}_{\mu}\left[ \left( v_\pi(s) - \hat{v}(s,\Theta) \right)^2 \right]$   $\mu$ : state probability distribution (how much we care about the error in state s)

( unknown → over $v_\pi(s)$ )

**LINEAR APPROXIMATORS**   $\hat{v}(s,\Theta) = \phi(s)^T \Theta$

( state space encoding function → $\phi(s)$ )

- $\mathcal{L}(\Theta)$ is quadratic   · $\nabla_\Theta \mathcal{L} = -2(v_\pi(s) - \hat{v}(s,\Theta)) \cdot \frac{\partial \hat{v}(s,\Theta)}{\partial \Theta} = -2(v_\pi(s) - \hat{v}(s,\Theta))\phi(s)$   ⟹   $\Theta_{t+1} \leftarrow \Theta_t + \gamma\left[ v_\pi(s) - \hat{v}(s,\Theta) \right]\phi(s)$

( gradient descend )

- linear approximator is enough thanks to **Representer Theorem** (difficult part is to find the kernel $\phi$)

Since $v_\pi(s)$ is unknown we need to replace it with the empirical target $G_t$

**MC** : $G_t$ is unbiased but noisy
$\mathbb{E}[G_t | S = S_t] = v_\pi(S_t)$
- strong convergence property

```
loop :
   generate with π on episode    S₀, A₀, R₁, ....
   for each step t :
      Θ ← Θ + α[ G_t - v̂(S_t,Θ)] φ(s)
```

**TD(0)** : $G_t$ is reliable but biased
- weak convergence property
  (for biased GD but faster)

· TD(0) under linear approximations converges

```
loop (episodes) :
   S₀ ← initialize
   loop (steps) :
      A ~ π(a|s)
      take action A and observe R, S'
      Θ ← Θ + α[ R + γv̂(S',Θ) - v̂(s,Θ)] φ(s)
      S ← S'
```

→ Since the target depends on current $\Theta$ the GD will be **biased**
⟹ **semi-gradient methods** because the gradient is not the true gradient but an approximation
· All the bootstrapping are affected

The eligibility for TD($\lambda$) can be defined as: $E_t = \gamma\lambda E_{t-1} + \phi(s_t)$

( instead of iteratives methods )

- We can use **batched episodes** to generate more **stable** gradients → solves iid problems with dataset sampling from **experience replay buffer**   ( $(s,a,r,s') \sim D$ )
  · I could add a **temporal weight** to overcome the non-stationary data

**DNN APPROXIMATORS**   Deep Q-Network (DQN) is a DNN to approximate Q   ( → cannot be done with Deep-SARSA since it depends on $a' \sim \pi_{old}$ )
- it uses an **experience replay buffer** → breaks correlation between consecutives experience
- uses $\varepsilon$ greedy policy

· It uses a separate NN to estimate Q-values target ($\Theta^-$ parameters) → prevents rapid oscillations caused by moving targets   ( → diverges otherwise )
  · it is updated **periodically** (not on every epoch $\Theta^- = \Theta$) → it reduces the bias
  $\Theta^- = \tau\Theta + (1-\tau)\Theta^-$   if $\tau=1$ ⟹ hard update, otherwise soft

( over replay buffer )   ( → semi-gradient update )

Simple DQN: $\mathcal{L}(\Theta) = \mathbb{E}_D\left[ \left( r + \gamma\max_{a' \in A} \hat{Q}(s',a',\Theta^-) - \hat{Q}(s,a,\Theta) \right)^2 \right]$

· Problem: **upward positive bias** from taking the maximum of a noisy esteem $\hat{Q}(s,a,\Theta^-)$   → problem is solved if the noise on $\hat{Q}_\Theta$ and $\hat{Q}_{\Theta^-}$ are decorrelated
  ⟹ splits action selection ($\hat{Q}_\Theta$) and action evaluation ($\hat{Q}_{\Theta^-}$) and then vice-versa

Double DQN: $\mathcal{L}(\Theta) = \mathbb{E}_D\left[ \left( r + \gamma\hat{Q}(s', \underset{a' \in A}{argmax}\ \hat{Q}(s',a',\Theta), \Theta^-) - \hat{Q}(s,a,\Theta) \right)^2 \right]$   · randomly select if updating $\Theta$ or $\Theta^-$ at each iteration

$\mathcal{L}(\Theta) = \mathbb{E}_D\left[ \left( r + \gamma\hat{Q}(s', \underset{a' \in A}{argmax}\ \hat{Q}(s',a',\Theta^-), \Theta) - \hat{Q}(s,a,\Theta^-) \right)^2 \right]$

**Duelling Networks**  $\hat{Q}(s,a,\theta) = \hat{V}(s,\theta_V) + \left[ \hat{A}(s,a,\theta_A) - \max_{a' \in A} \hat{A}(s,a',\theta_A) \right]$   $\theta \cdot (\theta_V, \theta_A)$

Uses two fuction :   i) Value function $\hat{V}(s,\theta)$ → expected reward for being in a state $s$

2) Advantage function $\hat{A}(s,a,\theta_A)$ → quantifies the advantage to choose a particular action in state $s$
   $A(s,a) = Q(s,a) - V(s)$

· In this way the model can learn independently what states are the best and what actions are advantageus
and then combine them in the Q value   $\hat{Q}(s,a,\theta_V,\theta_A) = \hat{V}(s,\theta_V) + \hat{A}(s,a,\theta_A)$

→ more stable and robust   1) can estimate Q-values from unseen action better, because it has a stronger starting point then $\hat{q}(s,a)$

2) faster convergence when a lot of similar actions since $V(s)$ will be shared and fixed and only $\hat{A}$ will change



→ once experienced a positive reward going left ⇒ bias towards B

- no guarantees of convergence, higher variance, $\pi$ diff

# POLICY GRADIENT METHODS
→ learn a parametrized policy $\pi(a|s,\theta)$ not a value function to get the policy → more direct

● typically $\pi^*$ approximation converges faster then $V^*$   ● can be used in continuous A space

● wrt $\varepsilon$-greedy the value of $\pi_\theta$ change continously with $d\theta$ changes → stronger convergence

I would live to $\theta_{t+1} \leftarrow \theta_t + \alpha \nabla_\theta \pi_\theta(a^*|s)$   · Unfortunately $a^*$ is unknown → weight the updates   $\theta_{t+1} \leftarrow \theta_t + \alpha \hat{Q}(s,a) \dfrac{\nabla_\theta \pi_{\theta_t}(a|s)}{\pi_{\theta_t}(a|s)}$   correction in order to not update too much more probable $a$

optimal action in state $s$

· $\theta_t + \alpha \hat{Q}(s,a) \nabla_\theta \ln \pi_{\theta_t}(a|s)$

# POLICY GRADIENT THEOREM   $\hat{\pi} \in C'(\mathbb{R}^n)$   In any policy gradient method where   $L(\theta) \propto \mathbb{E}_{\hat{\pi}}\left( V_{\hat{\pi}}(s) \right) = \sum_s d_{\hat{\pi}}(s) V_{\hat{\pi}}(s) = \sum_s d_{\hat{\pi}}(s) \sum_a \pi(a|s) q_{\hat{\pi}}(s,a)$   ⟹

continously differentiable   perform gradient ascend   true value function wrt $\hat{\pi}$

$\nabla_\theta J(\theta) \propto \sum_s d_{\hat{\pi}}(s) \sum_a q_{\hat{\pi}}(s,a) \nabla_\theta \hat{\pi}(a|s,\theta) = \mathbb{E}_{\hat{\pi}}\left[ \sum_a \nabla_\theta \hat{\pi}(a|s,\theta) q_{\hat{\pi}}(s,a) \right] = \mathbb{E}_{\hat{\pi}}\left[ \sum_a \hat{\pi}(a|s,\theta) \nabla_\theta \ln \hat{\pi}(a|s,\theta) q_{\hat{\pi}}(s,a) \right]$

on policy state distribution   $\nabla_\theta \ln \hat{\pi} \cdot \frac{1}{\hat{\pi}} \cdot \nabla_\theta \hat{\pi}$   $s \sim d_{\hat{\pi}}(s)$

$d_{\hat{\pi}}(s) = \lim_{t \to \infty} P(s_t = s | s_0, \hat{\pi})$

· $\mathbb{E}_{\hat{\pi}}\left[ \nabla_\theta \ln \hat{\pi}(a|s,\theta) q_{\hat{\pi}}(s,a) \right]$
$s \sim d_{\hat{\pi}}(s), a \sim \hat{\pi}$
$(s,a) \sim d_{\hat{\pi}}(s,a)$

$\nabla_\theta J(\theta)$   depends both on :   1) action selection $(\pi_\theta)$
2) state distribution $(d_{\hat{\pi}})$ → it's change is a function of the environment, and thus unknown   ● we do not need $\nabla d_{\hat{\pi}}$ for previous theorem
and neither $\nabla q_{\hat{\pi}_\theta}(s,a)$

MC policy gradient   $G_t$ is an unbiased sample of $Q_{\hat{\pi}}(s,a)$ → $G_t = \mathbb{E}_{\hat{\pi}}\left[ Q_{\hat{\pi}}(s,a) | s,a \right]$

By sampling the gradient policy theorem   $\nabla_\theta L(\theta) \propto \dfrac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \nabla_\theta \ln \hat{\pi}_\theta(a_i|s_i,\theta) G_t^i$   ⟹   $\theta_{t+1} \leftarrow \theta_t + \alpha \nabla_\theta L(\theta)$

Maximum Likelihood gradient   · weights more successful trajectories

· Weighted Maximum Likelihood   ↳ collect $N$ episodes   → on-policy approach (otherwise importance sampling, but higher variance)

# ACTOR CRITIC METHODS   We want to reduce MC policy gradient variance

Two entities :   1) Actor $\hat{\pi}$ that generates actions
2) Critic $\hat{Q}$ that evaluates the quality of the action taken   ⟨ estimated by TD, MC,... methods
or using DNN   if $w \neq \theta$ two distinct NN, no gradient conflicts

↳ Using a function to estimate G instead of the environment (stochastic) has lower variance but higher bias

# COMPATIBLE FUNCTION APPROXIMATION THEOREM

If   1) $\nabla_w \hat{Q}(s,a,w) = \nabla_\theta \ln \hat{\pi}(a|s,\theta)$   [compatibility hyp.]   holds when :   $Q(s,a,w) = \nabla_\theta \ln \hat{\pi}(a|s,\theta)^T w$

2) $\mathbb{E}_{\hat{\pi}}\left[ (Q_{\hat{\pi}} - \hat{Q})^2 \right] \xrightarrow{t \to \infty} 0$   [unbiased property]
↳ true function

⇒ The policy gradient theorem still holds true   $\nabla_\theta L(\theta) \propto \mathbb{E}_{\hat{\pi}}\left[ \nabla_\theta \ln \hat{\pi}(a|s,\theta) \hat{Q}(s,a,w) \right]$
$(s,a) \sim d_{\hat{\pi}}(s,a)$

The variance can be further reduced using a baseline $B(\cdot)$ → represents empirical knowledge about the model

$\nabla_\theta L(\theta) \propto \mathbb{E}_{\hat{\pi}}\left[ \sum_a \nabla_\theta \hat{\pi}(a|s,\theta) (Q_{\hat{\pi}} - B(\cdot)) \right]$   → keep only the variance about the action, not future states

· If $B(\cdot)$ depends only on $s$   ⇒   doesn't change the gradient since   $\mathbb{E}_{\hat{\pi}}\left[ \sum_a \nabla_\theta \hat{\pi}(a|s,\theta) B(s) \right] = \mathbb{E}_{\hat{\pi}}\left[ B(s) \nabla_\theta \sum_a \hat{\pi}(a|s,\theta) \right] = \mathbb{E}_s\left[ B(s) \nabla_\theta 1 \right] = 0$

· common choice   $B(s) = V_{\hat{\pi}}(s)$   → reducing the variance of gradient updates (no bias introduced) → faster convergence
- becomes the Advantage function

↳ Without baseline if $Q_{\hat{\pi}}(s,a) > 0 \ \forall(s,a)$   ⇒   the network couldn't learn to do less of an action but just more the others

**Advantage function**    $A_{\bar\pi}(s,a) = Q_{\bar\pi}(s,a) - V_{\bar\pi}(s)$    <mark>measures the value added by action $a$ to state $s$</mark>

· advantage policy gradient reformulation    $\nabla_\theta \mathcal{L}(\theta) \propto \mathbb{E}_{\bar\pi}\left[\nabla_\theta \ln \bar\pi(a|s,\theta) A_{\bar\pi}(s,A)\right]$

→ In actor-critic scenario: $\hat A(s,a,\omega,\rho) = \hat Q(s,a,\omega) - \hat V(s,\rho)$

· Thus we should learn $(\theta,\omega,\rho)$ but in reality:    $\qquad A_{\bar\pi}(s,a) = \mathbb{E}_{\gamma}\left[r + \delta V_{\bar\pi}(s') - V_{\bar\pi}(s) \mid s,a\right]$

$Q_{\bar\pi}(s,a) = \mathbb{E}_{\gamma}[r + \delta V_{\bar\pi}(s')|s,a]$

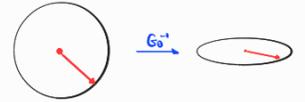$\hookrightarrow A_{\bar\pi}$ could be biased if $V_{\bar\pi}$ estimate is not perfect

## NATURAL POLICY GRADIENT

Policy gradients could perform poorly because of plateaus or discontinuities in policy space → <mark>small $\theta$ change causes big $\bar\pi$ changes</mark>

→ critic $\hat Q$ not valid anymore

Solution: make changes to $\theta$ proportional to policy variation

$\begin{cases} d^* = \arg\min_d \mathcal{L}(\theta_t + d) \\ \text{s.t. } D_{KL}(\bar\pi(\theta_t) \| \bar\pi(\theta_t + d)) \leq \varepsilon \end{cases}$    $D_{KL}(\bar\pi(\theta_t) \| \bar\pi(\theta_{t+1})) = \int_{s\times A} \bar\pi(a|s,\theta_t) \ln \frac{\bar\pi(a|s,\theta_t)}{\bar\pi(a|s,\theta_{t+1})} ds\,da$

trust region radius

$\hookrightarrow$ limits the loss of information (penalize $d$ that alters a lot $\bar\pi_\theta$)

Using the tylor approximation of second order of $D_{KL}$ we can find that    <span style="background:lightblue">$\theta_{t+1} - \theta_t \propto \vec{G_\theta} \nabla_\theta \mathcal{L}(\theta)$</span>  $\overset{\text{policy th}}{=}$  $\vec{G_\theta} \nabla_\theta \ln \bar\pi(a|s,\theta) Q_{\bar\pi}(s,a)$

With $G_\theta$ the <mark>Fisher information matrix</mark>    $G_\theta = \mathbb{E}_{\bar\pi}\left[\nabla_\theta \ln \bar\pi(a|s,\theta) \nabla_\theta \ln \bar\pi(a|s,\theta)^T\right]$

→ It's like a covariance matrix, where big entry in $G_\theta$ means small change in $\theta$ ⇒ big change in $\bar\pi$

→ Its inverse helps to normalize the <span style="background:lightgreen">parameters search space</span> by making all equally important

If <mark>compatibility</mark> holds we know that    $Q(s,a,\omega) = \nabla \ln \bar\pi(a|s,\theta)^T \omega$

$\Rightarrow \quad \Delta\theta_t \propto \vec{G_\theta} \nabla_\theta \ln \bar\pi(a|s,\theta) \nabla_\theta \ln \bar\pi(a|s,\theta)^T \omega = \vec{G_\theta} G_\theta \omega = \omega$

## MODEL-BASED RL

→ learn faster with fewer data

→ learns the model and from it get the value function from simulated experience (<span style="background:orange">learning</span>)

· needs to learn $P_{ss'}^a$ and $R_s^a$

· The simplest one is table lookup:    $P_{ss'}^a = \frac{1}{N(s,a)} \sum_{t=1}^{T-1} \mathbb{I}(s_t, a_t, s_{t+1})$    $R_s^a = \frac{1}{N(s,a)} \sum_{t=1}^{T} \mathbb{I}(s_t, a_t) r_t$

model variables

$M_\eta = \langle P_\eta, R_\eta \rangle$    <mark>- Inaccurate models lead to inaccurate policies</mark> → <span style="background:salmon">two source of errors</span> ($M_\eta$ and $\bar\pi$)

→ If the uncertainty is too high, go back to model-free or use an integrated approach

<span style="background:lightgreen">Integrated Approach</span> :    1) $M_\eta$ learning from real experience

2) $\hat V, \hat Q, \bar\pi$ learning from both $M_\eta$ and real experience

## BACKWARD SEARCH

<span style="color:blue">Dyna-Q</span>    → <span style="color:red">integrated approach + deterministic environment</span>

It has an extension <span style="background:orange">Dyna-Q+</span> to add an exploration bonus

→ Dyna-Q fails if the env is not static ($r$ changes)

$(r + K\sqrt{\tau_{sa}}, s') \leftarrow \text{Model}(s,a)$

$\hookrightarrow$ #timesteps last tried $(s,a)$

· Increasing the reward induces the policy to follow those a tried a long time ago

· In the slides is said that $(s,a) \sim$ weighed distribution
   $W_s(\tau) = R_s^a + K\sqrt{\tau}$

<span style="color:blue">Dyna-2</span>    → Uses two $Q$ functions:

1) <span style="background:lightgreen">$Q^T$</span> (transient) → adjusted from simulation

2) <span style="background:lightgreen">$Q^\rho$</span> (persistent) → learned only from real data

· Takes action combining both $Q^T$ and $Q^\rho$    $a_{t+1} \leftarrow \varepsilon\text{-greedy}(Q^T(s_t,\cdot) + Q^\rho(s_t,\cdot))$

---

> **Input** $S, A, \gamma$
> Initialize $Q(s,a)$ and $M(s,a)$    $\forall (s,a) \in S \times A$
> loop:
>   $S \leftarrow$ current state
>   $A \leftarrow \varepsilon\text{-greedy}(s,Q)$
>   execute $A$ and observe $(R,S')$
>   $Q(s,a) \leftarrow Q(s,a) + \alpha\left[R + \gamma\max_a Q(s',a) - Q(s,a)\right]$    # TD(0)
>   model → $M(s,a) \leftarrow (r,s)$    # Assuming deterministic env.
>   <span style="background:lightgreen">for n times:</span>
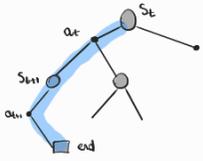>     choose $s \in S$ from $H$    → #history
>     choose $a \in A$ from $M(a)$
>     $(r,s') \leftarrow M(s,a)$
>     $Q(s,a) \leftarrow Q(s,a) + \alpha\left[R + \gamma\max_a Q(s',a) - Q(s,a)\right]$

# FORWARD SEARCH

builds a search tree with current state $S_t$ as root → uses $M_\eta$ as look ahead

→ simulates episodes from now with the model

→ apply model-free methods to simulated episodes



## NAIVE MC Search

· you are not learning $\hat{\pi}$ from the simulations
→ inefficient

```
Input   M, π, S_t

for all a ∈ A do:
        simulate K episodes  {s_t, a, R^κ_{t+1}, S^κ_{t+1}, a^κ_{t+1}, ----, S^κ_T} ~ M_η, π
        Q(s_t, a) = 1/K Σ^K_{a=1} G^κ_t    → mean return on simulation

a_t ← argmax_a Q(s_t, a)
```
completely from simulations

## MC TREE-SEARCH   → Builds a search tree with visited states and actions

starting point where       initially empty        · strategy to select which node to visit often : ε-greedy
I want to take an action

1) Selection: start from root $S^0_0$ and traverse the tree using a tree policy → like UCB to balance exploration and exploitation

stops when:   a) reach a node with an action that is not expanded

b) using tree policy the next node to visit is not on the tree

c) terminal state is reached   → in this case skip to back propagation

2) Expansion:   case a) → add action node and a sampled state from $M_\eta$
case b) → action was already present, add only s ↝ stochastic result state after a

greedy

3) Simulation: from the state $S_t$ simulate using a default policy until terminal state reached

$\{S_t, A^κ_t, R^κ_{t+1}, ----, S^κ_T\} ~ M_\eta, \pi$

4) BackPropagation: update every (s,a) along the path from selected node to root

a) Increment visit count $N(s,a)$ by one

b) $G^κ_t ← r^κ_t + \gamma G^κ_{t+1}$   (initial G is the return of selected node along the generated path)

c) Update Q estimate   $Q(s,a) = \frac{1}{N(s,a)} \sum_K \sum^{T_n}_{u=1} \mathbb{I}(S^κ_u = s, A^κ_u = a) G^κ_u$   → could be done more efficiently with incremental mean

# of simulations done

Repeat the four phases until there is computation time, then take $a_t$ and discard the tree   (some implementation could reuse it)

→ Exhaustive search is impractible in complex spaces
→ Highly parallelizable
→ The two policies improves

## TD-SEARCH   → uses bootstrapping instead of sampling simulation. apply SARSA from $(S_t, a_t)$

$\Delta Q(S_t, a_t) = \alpha [R_t + \gamma Q(S_{t+1}, a_{t+1}) - Q(S_t, a_t)]$   · could be use to simulate n steps ahead like TD(n)

From $S_t$ update Q at each iteration → lower variance
· Dyna2 learns $Q^T$ in this way

# EXPLORATION-EXPLOITATION TRADE-OFF

In $\varepsilon$-greedy the exploration is random (noise + greedy).

Two strategies available:
1) Optimism → estimates action-state uncertainty and explore the more uncertain
2) Information State → uses hystory knoledge to explore

MULTI-ARMED BANDIT SCENARIO
- deterministic starting state $S_0$
- multiple actions that brings to a different terminal state (no state space) with rewards that depends on each probability distribution
  → we want to find $R_{S_0}^a$ (expected reward) in order to maximize $\sum_{t=1}^{k} r_t$
  ↳ cumulative return over actions

- using greedy we are expoiting our current knowledge to maximize one-step-ahead return, but we are not exploring
- using $\varepsilon$-greedy may produce greater total reward in the long run

$$V^* = Q(a^*) = \max_a Q(a)$$
true action value (unknown)

regret: expected error doing suboptimal $a$
$$I_t = \mathbb{E}_{a \sim \pi}[V^* - Q(a_t)]$$
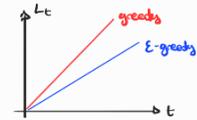
total regret
$$L_T = \sum_{t=1}^{T} I_t$$

minimize $\sum_{t=1}^{T} I_t \iff$ maximize $\sum_{t=1}^{T} r_t$

number of selection of a across $T$ timesteps (depends on $\pi$)
$$L_T = \sum_{t=1}^{T} \mathbb{E}(V^* - Q(a_t)) = \sum_a \mathbb{E}_\pi [N_T(a)] (V^* - Q(a))$$


greedy
$\varepsilon$-greedy

greedy (no exploration) or $\varepsilon$-greedy (random exploration) has linear regret

- in greedy $\mathbb{E}[N_T(a)]$ can be $T$ for a single $a^*$ ⟹ $L_T = T(V^* - Q(a^*))$
- in $\varepsilon$-greedy $\mathbb{E}[N_T(a)] \geq \frac{\varepsilon}{|A|} T$ ⟹ $L_T \geq \frac{\varepsilon}{|A|} T \sum_a (V^* - Q(a))$

Hoeffding Inequality: Let $x_1, \dots, x_n$ random variables in $[0,1]$, let $\bar{x}_t = \frac{1}{t} \sum_{k=1}^{t} x_k$ ⟹
$$P(\mathbb{E}(X) - \bar{x}_t \geq u) \leq e^{-2tu^2} \quad \forall u \geq 0$$

estimate $\hat{Q}_t(a) = \frac{1}{N_t(a)} \sum_{s=1}^{t} r_s \cdot \mathbb{I}(a_s = a)$

(UCB) Upper Confidence Bound
rewards in $[0,1]$
Call the upper bound on $Q(a) \to \hat{U}_t(a)$ ⟹ $P[Q(a) > \hat{Q}_t(a) + \hat{U}_t(a)] \leq e^{-2N_t(a)\hat{U}_t(a)^2} = p$

→ $\hat{U}_t(a) = \sqrt{\frac{-\log p}{2N_t(a)}}$ ⟿ less time chosen greater the uncertainty

OPTIMISM Between all the high probable $(1-p)$ models choose the most optimistic to improve the policy $\hat{Q}_t(a) + \hat{U}_t(a)$

- If we reduce $p_t = t^{-\alpha}$ (fixed $\alpha > 0$) we have a sublinear policy
$$a = \arg\max_a \left[ \hat{Q}_t(a) + \sqrt{\frac{\alpha \log t}{2N_t(a)}} \right] = \pi(s_0^i)$$
increases with uncertainty and time

→ Guarantees asymptotic convergence to optimal action
- It has sublinear regret → because it doesn't use random exploration but it's based on uncertain actions that could be optimal

# STATE of THE ART MODELS

## PROXIMAL POLICY GRADIENT

We want to limit the variation of $\pi$ across weight updates

· Is heavy to put a constraint on $D_{KL}$ or to compute $G_\theta^{-1}$

To keep $D_{KL}$ small at every step we can:

1) **PPO-Penalty** (proximal policy optimization)

→ how much $P$ of taking $a_t | s_t$ would increase with $\theta_t$

$$\mathcal{L}(\theta) = \mathbb{E}_t \left[ \frac{\pi_\theta(a_t | s_t)}{\pi_{\theta_{k-1}}(a_t | s_t)} \hat{A}_w(a_t, s_t) - \beta_n D_{KL}(\pi_\theta \| \pi_{\theta_{k-1}}) \right]$$

↳ aver. all episode

↳ penalize big $D_{KL}$

→ $(a_t, s_t)$ were collected via $\pi_{\theta_{k-1}}$ → importance sampling

$$\beta_{n+1} = \begin{cases} \frac{\beta_n}{2} & \text{if } d < \frac{2}{3}\bar{d} \quad \text{← user-defined th.} \\ 2\beta_n & \text{if } d > \frac{3}{2}\bar{d} \\ \beta_n & \text{otherwise} \end{cases}$$
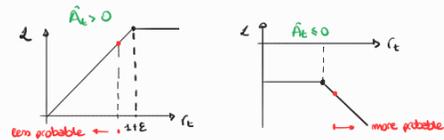
2) **PPO-Clipping**
$$\mathcal{L}(\theta) = \mathbb{E}_t \left[ \min\left\{ \frac{\overbrace{\pi_\theta(a_t | s_t)}^{r_t}}{\pi_{\theta_{k-1}}(a_t | s_t)}, g(\varepsilon, A_t) \right\} \hat{A}_w(a_t, s_t) \right]$$

$$g(\varepsilon, A_t) = \begin{cases} 1+\varepsilon & \text{if } \hat{A}_t > 0 \\ 1-\varepsilon & \text{otherwise} \end{cases}$$

We are clipping $r_t$ :    i) if $\hat{A}_t$ is good ($\hat{A}_t > 0$)  by  $1+\varepsilon$       → prevents too big changes ($r_t$ really high)

ii) if $\hat{A}_t$ is bad  by  $1-\varepsilon$

· There is no reward of having $r_t \notin [1-\varepsilon, 1+\varepsilon]$



## DEEP DETERMINISTIC POLICY GRADIENT

**DPG Theorem** : Let $\mathcal{L}(\theta) = \mathbb{E}_{s \sim d(s)}\left[ R(s, \overset{\text{deterministic}}{\pi_\theta(s)}) \right]$  with  $\pi_\theta, R, P \in C'$    $\Rightarrow$    $\nabla_\theta \mathcal{L}(\theta) = \mathbb{E}_{s \sim d(s)}\left[ \nabla_\theta \pi_\theta(s) \nabla_a Q(s,a)\big|_{a=\pi_\theta(s)} \right]$

reward ↗           more efficient specifically in continuous-action space

**Deterministic compatible Value Approximation Theorem:**  If we use $\hat{Q}_w(s,a) \in C'$ a parametrized version of $Q(s,a)$ → DPG theorem holds  $\Longleftrightarrow$

1) $\nabla_a \hat{Q}_w(s,a)\big|_{a=\pi_\theta(s)} = \nabla_\theta \pi_\theta(s)^T w$

2) $\mathbb{E}_{s \sim d(s)}\left[ (\nabla_a \hat{Q}_w(s,a) - \nabla_a Q(s,a))^2 \big|_{a=\pi_\theta(s)} \right] \xrightarrow{t \to +\infty} 0$

}  guarantees to have no bias substituting $Q$ with $\hat{Q}_w$

From 1) derives that   $Q_w(s,a) = a^T \nabla_\theta \pi_\theta(s)^T w$

→ $\delta_t = r_t + \gamma Q_{w_t}(s_{t+1}, \pi_\theta(s_{t+1})) - Q_{w_t}(s_t, a_t)$    → TD(0) error

① $\theta_{t+1} = \theta_t + \lambda_\theta \nabla_\theta \mathcal{L} = \theta_t + \lambda_\theta \nabla_\theta \pi_\theta(s) \nabla_a Q(s,a)\big|_{a=\pi(s)}$   $\Rightarrow$   $\theta_{t+1} = \theta_t + \lambda_\theta \nabla_\theta \pi_\theta(s) \nabla_\theta \pi_\theta(s)^T w_t$

↑ learning rate $\theta$

② $w_{t+1} = w_t + \lambda_w \delta_t \nabla_w Q(s_t, a_t)$       $\Rightarrow$     $w_{t+1} = w_t + \lambda_w \delta_t \nabla_\theta \pi_\theta(s) a$

**Natural Policy Modification** :   $\theta_{t+1} = \theta_t + \lambda_\theta w_t$     when we use the Fisher matrix to weight $\nabla_\theta \mathcal{L}(\theta)$

## TWIN DELAYED DDPG (TD3)    → DDPG is sensitive to hyperparameters and tends to overestimate $Q$

① TARGET POLICY SMOOTHING  → add noise to $\pi$, making it difficult to exploit $Q$ errors (smoothing)

$a = \text{clip}\left[ \pi_\theta(s) + \text{clip}(\varepsilon, -c, +c), a_\ell, a_h \right]$      with  $a_\ell, a_h \in \mathbb{R}$   $a_\ell < a_h$,   $c \in \mathbb{R}^+$,   $\varepsilon \sim N(0,1)$

② CLIPPED DOUBLE-Q LEARNING    uses two $Q_{w_1}, Q_{w_2}$ trained independently  → Difficult to have the same overestimation

$y(r, s') = r + \gamma \underset{i=1,2}{\min} Q_{w_i}(s', \overset{\text{clipped}}{a(s')})$

$L(w_1) = \mathbb{E}_D\left[ (Q_{w_1}(s,a) - y(r,s'))^2 \right]$     → same for $w_2$ (start from different $\bar{w}_1 \neq \bar{w}_2$)

$\pi_{\theta'} = \underset{\pi_\theta}{\arg\max} \underbrace{\mathbb{E}_D\left[ Q_{w_1}(s, \pi_\theta(s)) \right]}_{-\mathcal{L}(\theta)}$  ↳ uses DPG to update $\theta$, but not compatibility $\Rightarrow$ BIAS

③ DELAYED POLICY UPDATES  → updates policy $\pi_\theta$ once every $d$ critic update $(w_1, w_2)$

↑ typically $d=2$

· this allows to train $\pi_\theta$ with a lower variance value estimate → more stable

# SOFT ACTOR CRITIC (SAC) → We can use entropy regularization to enforce exploration

· Given a density function $p$ and a random variable $x \sim p$ ⟹ $H(p) \doteq -\mathbb{E}_x[\ln(p(x))]$

Define a new loss with exploratory bonus: $\pi^* = \arg\max_{\pi} \mathbb{E}_b \left[ \underbrace{\sum_{t=0}^{\infty} \gamma^t (R(a_t, a_t)}_{G_t} \cdot \boxed{\nu H(\pi(\cdot|s_t))} \right]$   $\nu > 0$

↳ entropy bonus → more exploration

· Approximate   $G_t \sim \min_{i=1,2} \hat{Q}_{w_i}(a_t, s_t)$  ⟹  $\arg\max_{\pi_\theta} \mathbb{E}_{\substack{s \sim D \\ a \sim \pi_\theta}} \left[ \min_{j=1,2} Q_{w_j}(s,a) - \boxed{\nu \ln \pi_\theta(a|s)} \right]$

↳ not used in TD3

↳ typically decreased over time $\nu_t$

<span>reparametrization trick</span> since is difficult to compute directly $\pi_\theta^*$  I used  divide it in two:

$a_\theta(s, \varepsilon) = \tanh(\mu_\theta(s) + \sigma_\theta(s) \odot \varepsilon) \in (-1, +1)$   with   $\varepsilon \sim N(0,1)$   now $\pi_\theta \overset{\approx}{\sim} N(\mu_\theta, \sigma_\theta)$

· Decoupled randomness of $\pi_\theta(a|s)$ from model parameters

⟹   $\pi^* = \arg\max_{\theta} \mathbb{E}_{\substack{s \sim D \\ \varepsilon \sim N}} \left[ \min_{j=1,2} Q_{w_j}(s, a_\theta(s,\varepsilon)) - \nu \ln \pi_\theta(a_\theta(s,\varepsilon)|s) \right]$   ⟹ $\pi_\theta \sim N(\mu_\theta, \sigma_\theta)$

· Can be written the <span>entropy bellman equation</span> in order to learn $Q$

$Q_\pi(s,a) = \mathbb{E}_{\substack{s' \sim p \\ a' \sim \pi(\cdot|s')}} \left[ R(s,a) + \gamma \left( Q_\pi(s',a') - \nu \ln(\pi(a'|s')) \right) \Big| s, a \right]$

⟹  $\mathcal{L}(w_i) = \mathbb{E}_\pi \left[ \delta_t^2 \right]$   → similar to TD3 but with entropy correction

Differences   1) TD3 uses a deterministic policy
2) SAC uses entropy regularization
3) SAC doesn't use policy smoothing

# ALPHA THEORY

Mathematical field in which the Axiom of Archimede lacks

· let $F$ be a totally ordered field $\rightarrow$ $\forall x, y \in F$, $0 < x < y$, $\exists n \in \mathbb{N}$ : $y < nx$

$\not\exists n \in \mathbb{N}$ : $\alpha < nx$ $\quad \forall x \in \mathbb{R} \subseteq \mathbb{E}$

**Axiom 1** (existence) : $\exists$ an ordered field $\mathbb{E} \supset \mathbb{R}$ whose number are called **euclidian number**

$\left[ \text{and a } \textbf{function} \quad \lim_{n \to \alpha} : \mathbb{R}^{\mathbb{N}} \to \mathbb{E} : \quad \text{axiom 3} \right]$

sequence

**DEF**: given $\mathcal{E} \in \mathbb{E}$ then:

1) $\mathcal{E}$ is infinite $\iff$ $\forall n \in \mathbb{N}$, $|\mathcal{E}| > n$

2) $\mathcal{E}$ is finite $\iff$ $\exists n \in \mathbb{N}$ : $\frac{1}{n} < |\mathcal{E}| < n$

3) $\mathcal{E}$ is infinitesimal $\iff$ $\forall n \in \mathbb{N}$. $|\mathcal{E}| < \frac{1}{n}$

**Axiom 2** : $\exists$ a function $\quad$ num $: \underbrace{\mathbb{U}}_{\text{set of sets}} \to \mathbb{E}$

· num$(A) = |A|$ if $A$ is finite

· num$(A) <$ num$(B)$ if $A \subset B$

· num$(A \cup B) =$ num$(A) +$ num$(B) -$ num$(A \cap B)$ $\qquad \longrightarrow$ num$(\mathbb{N} \cup \{a\}) = \alpha + 1$ $\quad$ but $\quad |\mathbb{N} \cup \{a\}| = |\mathbb{N}|$

· num$(A \times B) =$ num$(A) \cdot$ num$(B)$ $\quad \longrightarrow$ contrast with cantour $\quad$ num$(\mathbb{N} \times \mathbb{N}) = \alpha^2$ $\quad$ but $\quad |\mathbb{N} \times \mathbb{N}| = |\mathbb{N}|$

**DEF**: $\alpha \doteq$ num$(\mathbb{N})$

$\eta \doteq \alpha^{-1}$ $\qquad \lim_{n \to \alpha} \frac{1}{n} = \eta$

**Axiom 3** : every sequence $\varphi, \varphi' : \mathbb{N} \to \mathbb{R}$ has a **unique** $\alpha$ limit which satisfy the following properties

· If $\varphi(n) = n$ then $\mathbb{N} \Rightarrow \lim_{n \to \alpha} \varphi(n) = \alpha$

· If $\varphi(n) = r$ $\forall n \in \mathbb{N}$ $\Rightarrow \lim_{n \to \alpha} \varphi(n) = r$

· $\lim_{n \to \alpha} \varphi(n) + \lim_{n \to \alpha} \varphi'(n) = \lim_{n \to \alpha} \left[ \varphi(n) + \varphi'(n) \right]$

· $\lim_{n \to \alpha} \varphi(n) \cdot \lim_{n \to \alpha} \varphi'(n) = \lim_{n \to \alpha} \left[ \varphi(n) \cdot \varphi'(n) \right]$

**IMPLICATIONS** : 1) if $\varphi(n) \in \psi(n) \Rightarrow \varphi[\alpha] \in \psi[\alpha]$ $\quad \in$ operator preservance

2) Given $A$ non empty : $A^* = \left\{ \varphi[\alpha] \mid \varphi : \mathbb{N} \to A \right\}$ $\quad$ eg. $f^* : A^* \to B^*$

3) Given $\mathbb{R}^* = \left\{ \varphi[\alpha] \mid \varphi : \mathbb{N} \to \mathbb{R} \right\}$

**Axiom** (Internal set) : If $\psi(n)$ is a sequence of non-empty set $\Rightarrow \lim_{n \to \alpha} \psi(n) = \left\{ \lim_{n \to \alpha} \varphi(n) \mid \varphi(n) \in \psi(n) \, \forall n \right\}$

internal set

sequence that picks on element $\forall n \in \mathbb{N}$ from $\psi$

eg. $\psi(n) = \left\{ 0, \frac{1}{n} \right\}$ then $\varphi(n) = \frac{1}{n} \in \psi(n)$

**Axiom** (Extension) : let $\varphi : \mathbb{N} \to \mathbb{R}$ and $f : \mathbb{R} \to \mathbb{R}$ be a function for which $f \circ \varphi$ is defined $\Rightarrow$ $f$ can be **uniquely** **extended** at $\mathcal{E} = \lim_{n \to \alpha} \varphi(n)$ as follows.

$f^*(x) \cdot f(x)$ $\quad \forall x \in \mathbb{R}$

$f^*(\mathcal{E}) = \lim_{n \to \alpha} f(\varphi(n))$ $\qquad f^* : \mathbb{E} \to \mathbb{E}$

$\hookrightarrow$ can be demonstrated that it doesn't depend on $\varphi$ definition

· define $\varphi(n) = x$ $\forall x \in \mathbb{R}$ $\Rightarrow f^*(x) = \lim_{n \to \alpha} f(x) = f(x)$

$\hookrightarrow$ constant sequence $\psi(n) = f(x)$

**DEF** An elementary formula is a finite string of symbols in first order logic $\quad$ (alphabet, $\underbrace{\neg, \wedge, \vee, \Rightarrow, \leftrightarrow}_{\text{logic connectives}}, \underbrace{\exists, \forall}_{\text{quantifiers}}, =, \in$)

**Th** (Transfer Principle) let $\sigma(x_1, \dots, x_n)$ an elementary formula and let $\varphi_1, \dots, \varphi_k$ be arbitrary sequences

$\sigma(\varphi_1(n), \dots, \varphi_k(n))$ holds almost every where $\iff$ $\sigma(\lim_{n \to \alpha} \varphi_1(n), \dots, \lim_{n \to \alpha} \varphi_k(n))$

$\hookrightarrow$ does not hold in a measure zero set

$\exists n^* \in \mathbb{N}$ : $\forall n \geqslant n^*$ $\sigma(\varphi_1(n), \dots \varphi_k(n))$ holds

$\rightarrow$ guarantees that some properties are preserved by the $\alpha$-limit $\quad$ eg. continuity, differentiability, ....

**DEF** (completeness) $\quad X$ is complete $\iff$ every non-empty subset of $X$ having on upper-bound must have a supremum in $X$

· Transfer principle doesn't transfer $\mathbb{R}$ completeness

eg. $\mu(0) = \{ y \in \mathbb{E} \mid y$ is infinitesimal $\}$ is bounded by $1$ but doesn't have a sup in $\mathbb{E}$

$\mathbb{R}$ is too large to fit in computers

**Algorithmic Field** is the set of all numbers that can be represented exactly in a machine (also number not preset in the original field it approximate like nan, $\pm \infty$) $(\hat{R}, \hat{E})$

1) **Fixed length representation** - faster code with possibility of hardware accellerators

- deterministic time consumption

2) **Symbolic representation** - variable size with "infinite" precision $\rightarrow$ programs will be slower

**Algorithmic Number** (AN): $\varepsilon \in \mathbb{E}$ is an AN $\iff$ can be represented by a finite sum $\quad \varepsilon = \sum_{k=0}^{\varepsilon} r_k \alpha^{s_k} \quad$ with $r_k \in \mathbb{R}, \; s_k \in \mathbb{Q}, \; s_k > s_{k+1}$

Th (AN normal form): Any AN $\varepsilon$ can be represented by $\quad \varepsilon = \alpha^p P(\eta^{\pm}) \quad$ with $p \in \mathbb{Z}, \; m \in \mathbb{N}, \; P(x)$ is a polynomial with real coefficients: $P(0) \neq 0$

· parallelism $\quad 1.37 \cdot 10^3 \qquad 10^3$ is $\alpha^p \quad$ and $\quad 1.37 = 1 \cdot 10^0 + 3 \cdot 10^{-1} + 7 \cdot 10^{-2}$ is $P(\eta^{\pm})$ polynomial with negative coefficients in $\alpha$

Two problems arise with AN:

    1) The inverse of a AN is not always an AN (not closed under inversion) e.g. $(1+\alpha)^{-1} \notin \hat{\mathbb{E}} \qquad \frac{1}{1+\alpha} = \frac{1}{\alpha(1+\eta)} \overset{\text{taylor}}{=} \eta \left[ 1 - \eta + \eta^2 - \eta^3 + \cdots \right] \qquad \frac{1}{1-x} = 1+x+x^2+\cdots$

        · $\mathcal{J}^{-1}$ has infinite # of digits

    2) They have variable length coding

To solve those we need **truncation** of $P(\cdot)$ $\qquad P(x) = p_0 x^0 + \cdots + p_m x^m \quad \Rightarrow \quad t_n[P(x)]: \begin{cases} P(x) & n \geq m \\ p_0 x^0 + \cdots + p_n x^n & n < m \end{cases}$

ensures closeness under basic operations

AN with $m=1$

**Bounded AN** (BAN): $\varepsilon \in \mathbb{E}$ is a BAN $\iff$ can be represented by $\varepsilon = \alpha^p P(\eta) \qquad$ [BAN $\subset$ AN] $\qquad 10^n P(10^{-1})$ notazione scientifica

    integer polynomium $\quad \alpha^2 [1 + 0.7\eta + 3.41\eta^2]$

    – sufficient for elementary operations (needed in ex) but not for non-linear ones e.g. $x^2 = \alpha \Rightarrow x = \sqrt{\alpha} \notin$ BAN

    – $m=1$ because operation between BANs are a lot easier then ANs

**LMOP** ( lexicographic M-O problems) $\quad \begin{cases} \text{lexmin}_x \; f_1, \cdots, f_n \\ x \in D \end{cases}$

$D \subseteq \mathbb{R}^n$

$f_i: \mathbb{R}^m \to \mathbb{R}$

$\boxed{\begin{array}{l} \Omega = D \\ \text{for } i=1,\cdots,n \text{ do:} \\ \quad \text{if } \Omega = \{\emptyset\} \lor |\Omega| = 1 \text{ then} \\ \qquad \text{break} \\ \quad \Omega = \{x \mid \text{argmin}_{x \in \Omega} f_i(x)\} \\ \text{return } \Omega \end{array}}$

It can be reformulated with $n$ optimization problems

    PRE-EMPTIVE APPROACH

$\begin{cases} \min_x \; f_i(x) \\ f_j(x) = \min_{x \in D} f_j(x) \quad \forall j = \{1 \cdots, i-1\} \\ x \in D \end{cases}$

$\rightarrow$ after each run a constrain is added

    - it gets more and more complex

    - bad use of hardware (cache) and software (unless the warm start)

    - nature of the problem may change. e.g. convex to non-convex with new constraints

    SCALARIZATION $\quad \{w_i\}_{i=1}^n$ : $\frac{w_{i+1}}{w_i} \ll 1$

$\begin{cases} \min_x \; \sum_{i=1}^n w_i f_i(x) \\ x \in D \end{cases}$

    - no guarantees of problem equivalence

    - choiche of weights are error-prone and arbitrary $\rightarrow$ may induce numeric instabilities

    NON-STANDARD SCALARIZATION $\quad w_i = \alpha^{1-i} \rightarrow$ the choice is not unique, it is sufficient to have infinitesimal $\frac{w_{i+1}}{w_i}$

$\begin{cases} \min_x \; \sum_{i=1}^n \alpha^{1-i} f_i(x) \\ x \in D \end{cases}$

    - the problems are equivalent and here we use a single objective

    - the transfer principle guarantees the equivalence between optimal solutions of the original LMOP and its scalarization

Th (Equivalence) Given a LMOP with $f_1, \cdots, f_n$ real functions, with priority induced by natural order

    $\Rightarrow \exists$ an equivalent scalar program over the same domain, whose objective is non-standard with the following form: $F(x) = \beta_1 f_1(x) + \cdots + \beta_n f_n(x)$

    with $\beta_1, \cdots, \beta_n \in \mathbb{E}$ : $\frac{\beta_{i+1}}{\beta_i} \approx 0 \quad \forall i = \{1, \cdots, n-1\}$

Proof: Let $\Omega$ be the domain of the two problems, and $w \in \Omega$ be a global **maximum** of LMOP
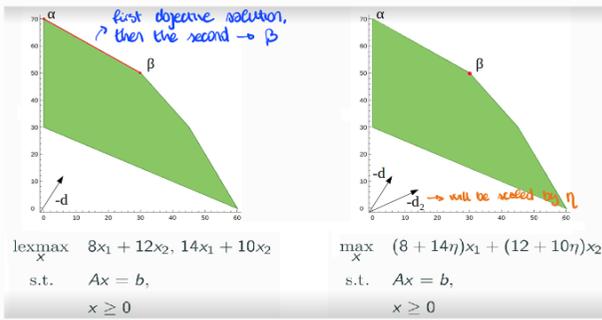
    $\Rightarrow \nexists w' \in \Omega: \; f_1(w') > f_1(w) \quad \lor \quad \left[ \exists i \in \{2, \cdots, n\}: \; f_j(w') = f_j(w) \quad \forall j = \{1, \cdots, i-1\} \; \land \; f_i(w') > f_i(w) \right]$

    This is true $\iff \nexists w' \in \Omega: F(w') > F(w)$ by the definition of $F$ itself

**LMOLP** (linear) $\quad \begin{cases} \text{lexmin}_x \; c_1^T x, \cdots, c_n^T x \\ Ax = b \\ x \geq 0 \end{cases} \iff \begin{cases} \min_x \; \sum_{i=1}^n c_i^T x \, \alpha^{1-i} \\ Ax = b \\ x \geq 0 \end{cases}$ could use the **NA-L-Simplex** algorithm

    $\rightarrow$ simplex with non-standard cost function (same property e.g. big-O)

    $\rightarrow$ transfer principle guarantees all the property of LP

$$\underset{X}{\text{lexmax}} \quad 8x_1 + 12x_2, \ 14x_1 + 10x_2$$
$$\text{s.t.} \quad Ax = b,$$
$$\quad x \geq 0$$

$$\underset{X}{\text{max}} \quad (8+14\eta)x_1 + (12+10\eta)x_2$$
$$\text{s.t.} \quad Ax = b,$$
$$\quad x \geq 0$$

## SIMPLEX PROBLEM — What if a starting feasible solution is unknown

AUXILIARY PROBLEM

$$Ax = b \quad \Rightarrow \quad \overset{mxn}{Ax} + \overset{mxm}{Is} = b \quad \Rightarrow \quad \overset{mx(n+m)}{[A \mid I]} \binom{x}{s} = b$$

$(0,b)$ is feasible of the new problem

such that with other variables $=0$ $\quad A_B x_B = b$

setting $\bar{x} = \bar{0}$ we have that $B = \{n+1, ..., n+m\}$ is a valid base (column indexes of $(A|I)$ to get a $m \times m$ invertible matrix)

In fact $I^{-1} = I \quad \Rightarrow \quad s = b$

Two methods to make the original and auxiliary problem the same:

### 1) Two-Phase Method → I need to get $\bar{s} = \bar{0}$ to be feasible in the first problem

First-Phase:
$$\begin{cases} \underset{x,s}{\text{min}} \ \sum_{i=1}^{m} s_i \\ Ax + Is = b \\ x, s \geq 0 \end{cases}$$

If $s^* = \bar{0}$ move to phase 2. otherwise P is infeasible
- The starting base will be $B'$

Second-Phase: Solve P starting from previous $x^*$. $B = B'[1:n] \rightsquigarrow$ original variables

### 2) Big-M Method → try to make the two thing at once

$$\begin{cases} \underset{x,s}{\text{min}} \ \overset{original}{d^T x} + M \overset{penalty}{\sum_{i=1}^{m} s_i} \\ Ax + Is = b \\ x, s \geq 0 \end{cases}$$

Th: $\exists M > 0$ sufficiently large: $\bar{x}' = [\bar{x}, 0]$ is the solution to the easier problem $\iff \bar{x}$ is solution to the original one

### 3) Infinite Big-M method (I-Big-M)
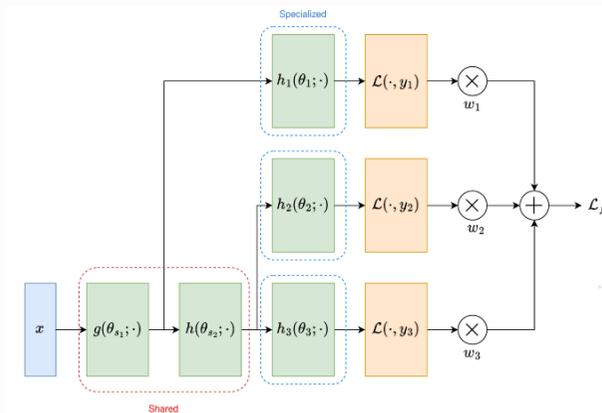
M prioritize the zeroing of s over the minimization of $d^T x$
→ similar to lexicographic problem

Set $M = \infty$

- M is always large enough for equivalence → do not need to retry
- No numerical instabilities by very large M

## HIERARCHICAL CLASSIFICATION → classification based on trees (hierarchy)

### 1) Branching DNN



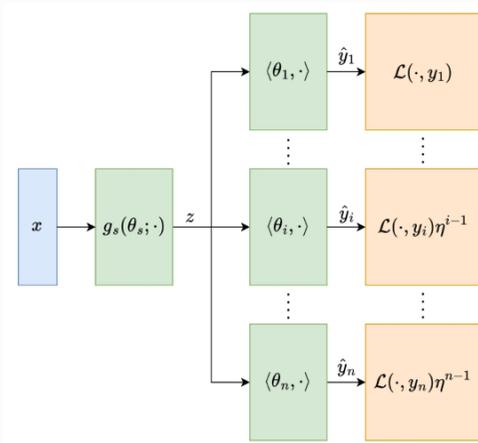$$\mathcal{L}_{B\text{-}DNN}(\Theta) = \sum_{i=1}^{n} w_i(t) \mathcal{L}(f_i(\Theta, x), y_i)$$

$$(w_1, w_2, w_3)|_{t=0} = (1,0,0)$$
$$(w_1, w_2, w_3)|_{t=finish} = (0,0,1)$$

## 2) Lexicographic DNN



$$\text{lexmin} \quad \mathcal{L}(\ell_1(\Theta, x), y_1), \ldots, \mathcal{L}(\ell_n(\Theta, x), y_n)$$

$$\iff \quad \min \sum_{i=1}^{n} L(\ell_i(\Theta, x), y_i) \, \eta^{i-1}$$

→ Cannot directly backpropagate otherwise updates of less important losses may worsen more important one
- perform a particular projection on gradient step to avoid it

Differences    i) Different depth splitting
               2) Different losses

→ allows a NN to obtain a classification if unseen data
→ prioritized learning