

Background: What's wrong

Extreme Value Theory: What's right

1. Balancing Exploitation vs Exploration

$h$  : an estimate of distance/cost to the goal

misguiding heuristics

How should we balance trusting and doubting the estimate?

2. Monte-Carlo Tree Search (Kocsis & Szepesvári, 2006)

Selection Expansion Evaluation Backprop repeat

Using "tree policy": Multi-Armed Bandit at each node

Policy UCB1: Select the max of

$\hat{h}_i + c\sqrt{(2\log T)/t_i}$

average of samples  $\{h_1, h_2, h_3, \dots\}$  num. samples for parent/child

statistics of leaves

3. How to evaluate? (Schulte et al, 2014)

For PDDL, domain-independent heuristics makes sense

$h(s) = FF(s)$   $h(s) = \frac{1}{N} \sum r_i$

Heuristics Monte-Carlo Simulation

Games

(define (...) ...) Planning

guaranteed to terminate in a fixed number of steps

difficult to design domain-dependent

Infeasible in planning (terminal nodes, a.k.a. goals, are rare!)

4. GBFS is just "min backprop" MCTS! (Schulte et al, 2014)

So, let's call it THTS! (Trial-based Heuristic Tree Search)

std. MCTS  $avg(h_1, h_2, h_3, h_4)$  GBFS  $min(h_1, h_2, h_3, h_4)$

$avg(h_1, h_2)$   $avg(h_3, h_4)$   $min(h_1, h_2)$   $min(h_3, h_4)$

$h_1$   $h_2$   $h_3$   $h_4$  = priority queue

5. UCB1 is WRONG (Wissow & Asai, 2023)

UCB1 : designed for 0/1 rewards (games are like that! : win = 1, loss = 0) BUT Heuristics have no upper bound !!! ([0, 1] is an overspecification)

Algorithm Assumption

UCB1 known finite support distributions, like [0, 1]

UCB1-Normal Gaussian + assumptions (may not hold)

UCB1-Normal2 Gaussian + different assumptions (more likely to hold in classical planning)

Solution:

- Use an unbounded distribution! Gaussian:  $[-\infty, \infty]$
- Backpropagates both (mean, variance):  $N(\mu, \sigma)$
- As seen here, quite powerful vs UCB1-based MCTS in classical planning

6. Gaussian is better but is STILL WRONG (this paper)

Using the average is SO wrong

- Gaussian: no bounds at all  $h \notin [-\infty, \infty] = R$
- Heuristics have unknown bounds !!!  $h \in [0, \infty]$  (underspecification)  $h_{add} \in [h_{max}, \infty], l_{mcut} \in [0, h^+]$  ...

- We are interested in good nodes. Why do we use the average?
  - Average takes ALL bad nodes into consideration
- Why change it when it's not broken? GBFS uses the minimum.
  - Schulte et.al. proposed min-backup, but it is not good
  - Why not good?  $\rightarrow$  min-backup lacks statistical justifications. (bandit theory)
  - What's the statistical theory of minimum/maximum (extremum)?
- Dead-ends ( $h=\infty$ ) break the average. Average of  $[2, 3, \infty, 7, 5]$  is  $\infty$ .
  - Existing work removes the dead-ends because otherwise it doesn't work
  - GBFS implicitly does it (minimum can discard  $\infty$  naturally)
  - "Removing them just to make the algorithm work" is ad-hoc and wrong

- Statistical theory of the maximum ( $\neq$  average)
- Used in safety-critical applications: e.g. Maximum water level
- There are two types:
  - Method of block maxima: (block) e.g. Predict next monthly maximum from several monthly maxima
  - Peaks-over-Threshold: Predict the exceedance over the threshold

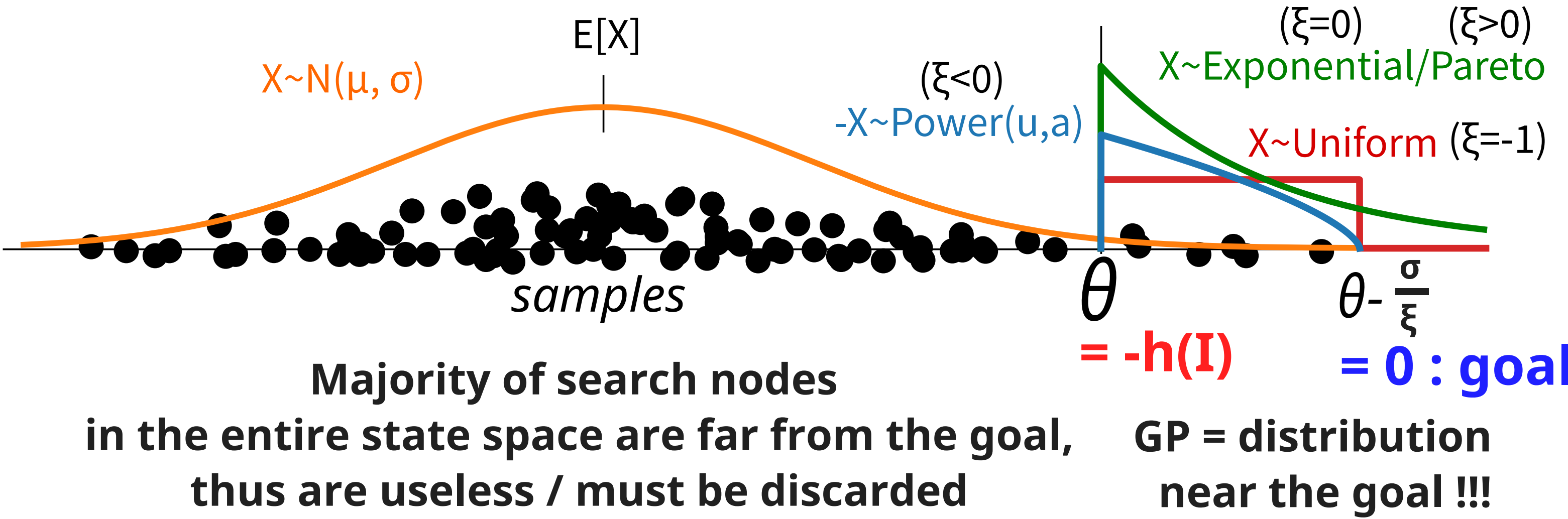
What to model?	Limit Theorem ( $N \rightarrow \infty$ )	converges to
Average	Central Limit Theorem	Gaussian Distribution
Block maxima	Fisher–Tippett–Gnedenko	Extreme Value Distribution (EVD)
Exceedance	Pickands–Balkema–de Haan	Generalized Pareto (GP) Distribution

$$GP(x \mid \theta, \sigma, \xi) = \begin{cases} \frac{1}{\sigma} \left(1 + \xi \frac{x-\theta}{\sigma}\right)^{-\frac{\xi+1}{\xi}} & (\xi \neq 0) \\ \frac{1}{\sigma} \exp\left(-\frac{x-\theta}{\sigma}\right) & (\xi = 0) \end{cases} \quad (x > \theta)$$

Note: It models the maximum, so we should invert the signs for minimization

Theoretical reason for removing the dead-ends ( $x = -\infty$ )

- We predict the exceedance above  $\theta = -h(I)$  Initial heuristic value
- fit  $N(\mu, \sigma) =$  based on all samples (incl. bad  $h$ ), compute the average & the variance
- fit  $GP(\theta, \sigma, \xi) =$  discard samples below the threshold  $\theta$ , compute the maximum and shape  $\xi$



Backprop GP, not Gaussian

We focus on GP's special cases: Uniform and Power. (see paper for why)

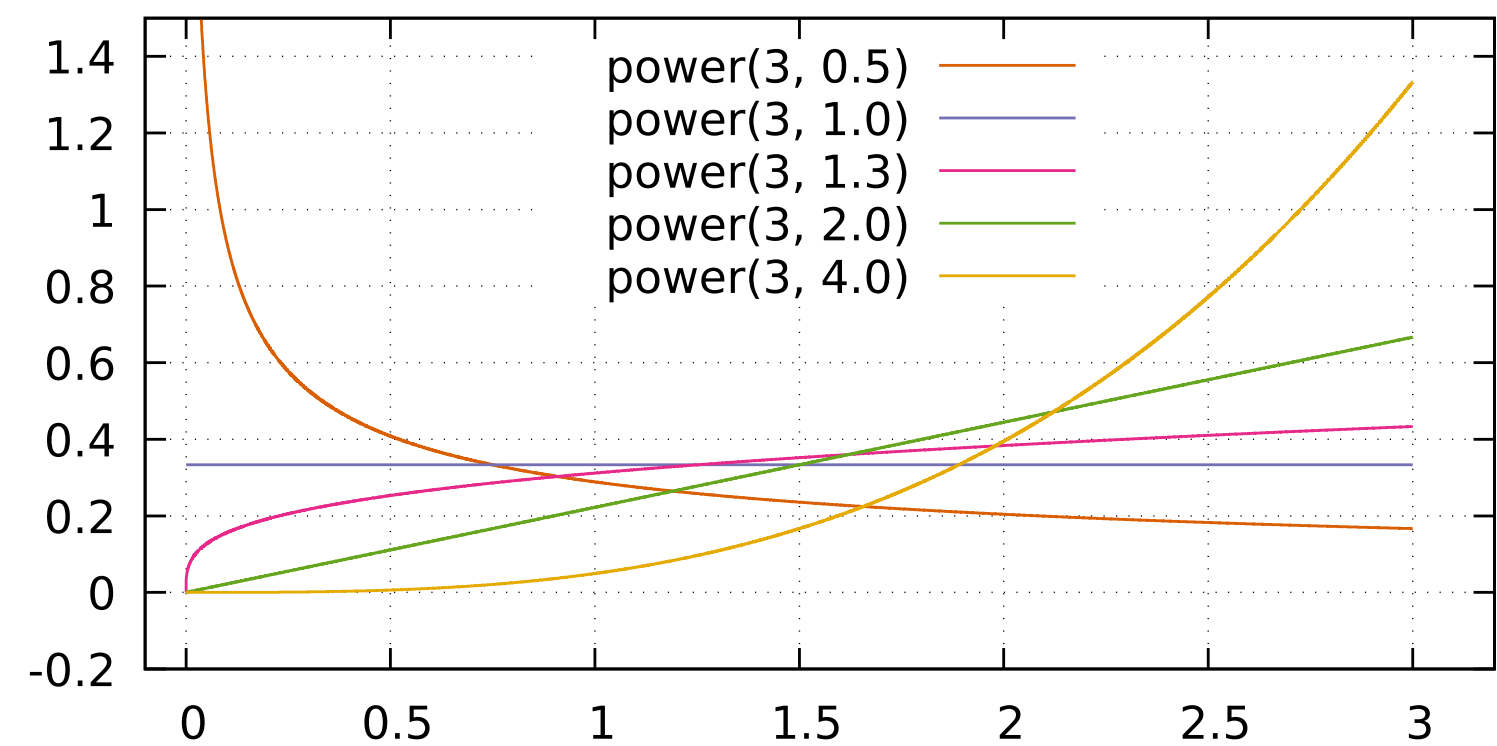
(GP with  $\xi < 0$ )  $Pow(x|u, a) = \frac{ax^{a-1}}{u^a} \cdot (0 < x < u, 0 < a)$

(GP with  $\xi = -1$ )  $U(x|l, u) = \frac{1}{u-l} \cdot (l < x < u)$

Backpropagate min  $h_i$ , max  $h_i$ , mean log  $h_i$ .

Power:  $\hat{u} = \max_i x_i$  and  $\hat{a} = (\log \hat{u} - \frac{1}{N} \sum_i \log x_i)^{-1}$   $x_i$ : heuristic of leaf  $i$

Uniform:  $\hat{u} = \max_i x_i$  and  $\hat{l} = \min_i x_i$



Power's shape parameter  $a$  : Rarity of  $h$  near 0

$a$  is estimated from backpropagated heuristics

Small  $a \Leftrightarrow$  nodes with small  $h$  are common

Large  $a \Leftrightarrow$  nodes with small  $h$  are rare

Multi-Armed Bandit for GP (with regret bounds!)

LCB1-Uniform $_i = \frac{\hat{u}_i + \hat{l}_i}{2} - (\hat{u}_i - \hat{l}_i) \sqrt{6t_i \log T}$

LCB1-Power $_i = \frac{\hat{u}_i \hat{a}_i}{\hat{a}_i + 1} - \hat{u}_i \sqrt{6t_i \log T}$

Results

Num. solved on 24 IPC domains w/  $10^4$  evaluations

	$h =$	$h^{FF}$	$h^{add}$	$h^{max}$	$h^{GC}$	$h^{FF+PO}$	$h^{FF+DE}$	$h^{FF+DE+PO}$
GBFS		538	518	224	354	-	489	-
Softmin-Type(h)		576	542.6	297.2	357.6	-	578	-
GUCT	Uses UCB1	412	397.8	228.4	285.2	454	389.2	439.4
-Normal	Uses UCB1-Normal	283.4	265	212	233.4	372.4	289	381.6
*-Normal	* backprop min	318.8	300	215.2	246.2	378.05	304.4	386.7
-Normal2		581.8	535.8	316.6	379	621	518	578
*-Normal2		567.2	533.8	263	341	618	511.4	567.8
-Power		596	541.8	450.6	463.2	623.4	413.6	583
-Uniform		594.8	543.8	450.6	463.8	626.4	416.4	583

FD/C++ implementation is on the way and showing promising results