



# Evolutionary Computation for Non-Convex Machine Learning

Ke Tang

Southern University of Science and Technology

Yang Yu

Nanjing University

---

# Outlines

---

- **The need of machine learning**
- **Evolutionary computation for supervised Learning –  
New case studies.**
- **Evolutionary Reinforcement Learning**

# The Need of Machine Learning

---

- Three key components of machine learning:
  - Data/model representation
  - Evaluation
  - Training algorithm
- Most modern machine learning problems are essentially searching for the model that is optimal with respect to some objective function (e.g., generalization).
- Optimization algorithms thus play a crucial role in machine learning.

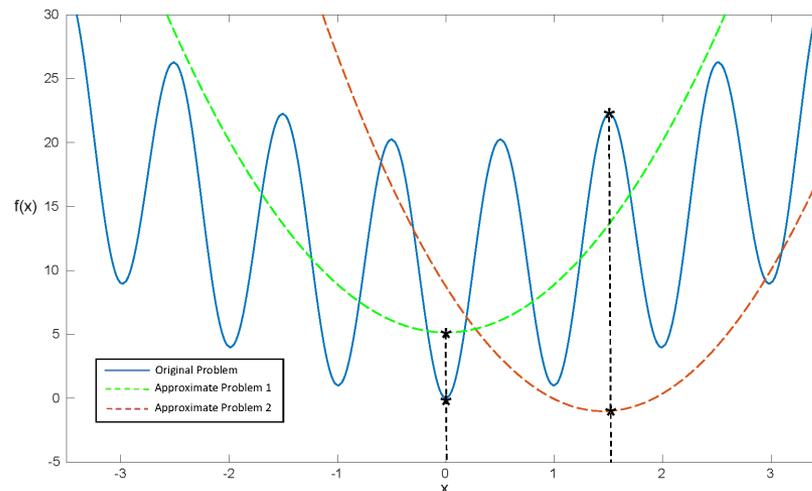
# The Need of Machine Learning

---

- Many machine learning tasks, when formulated as an optimization problem, cannot be well solved by traditional (e.g., convex) optimization techniques.

Plan A: Approximate the original (non-convex) optimization problem with a convex one.

Plan B: Seek an approximate solution to the original problem with heuristic methods, e.g., Evolutionary algorithms.



# The Need of Machine Learning

---

EAs have been applied to a large variety of learning problems in the past decades.

- Data representation
  - Feature selection
  - Feature extraction
  - Dimensionality Reduction
- Model training
  - Decision tree
  - Neural networks
  - Rule-based systems
  - Clustering
- Hyper-parameter tuning

# The Need of Machine Learning

---

Machine learning has its own characteristics that calls for specialized EAs.

- Huge problem size (i.e., the search space).
- Noisy fitness evaluation (since the generalization cannot be precisely measured)
- Expensive fitness evaluation
- Theoretical guarantee is more preferred than in other areas.

# Case Study (1)

---

- Subset selection: select a subset of size  $k$  from a total set of  $n$  variables for optimizing some criterion.

Formally stated: given all variables  $V = \{X_1, \dots, X_n\}$ , a criterion  $f$  and a positive integer  $k$ ,

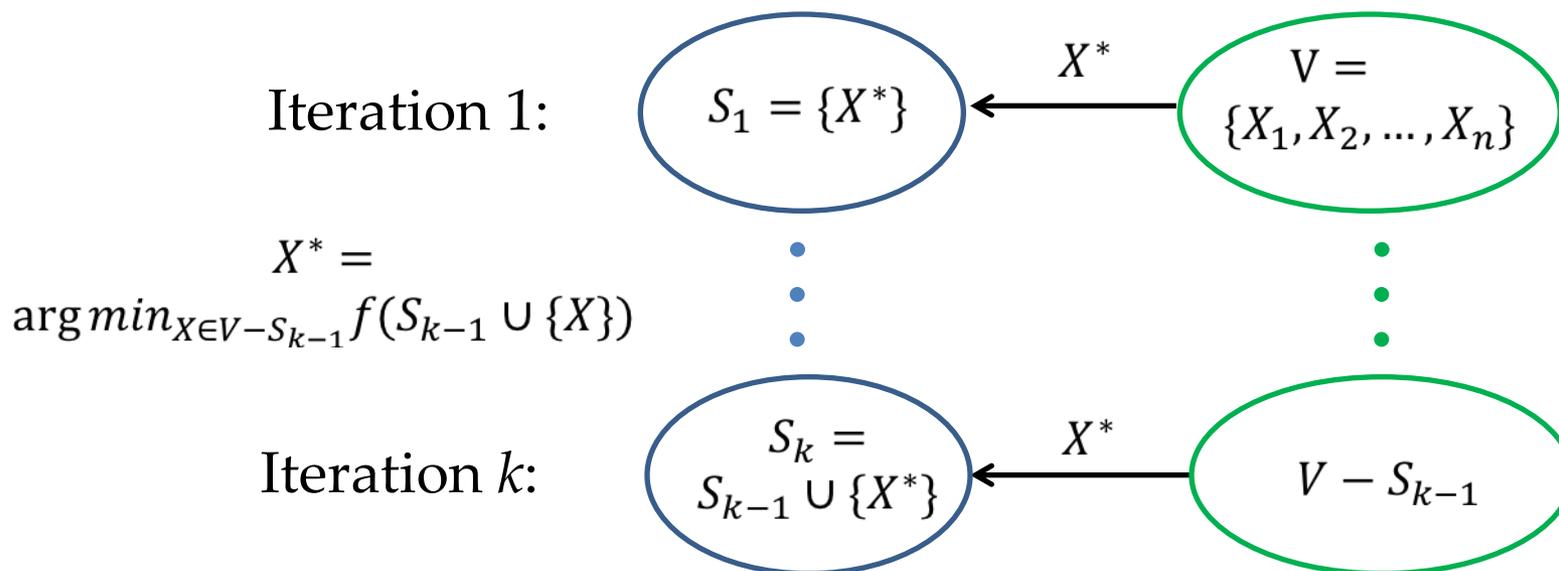
$$\arg \min_{S \subseteq V} f(S) \quad s.t. \quad |S| \leq k.$$

- NP-hard in general [Natarajan,1995; Davis et al., 1997] and arises in many learning problems:
  - Feature Selection
  - Sparse Learning
  - Compressed Sensing

# Case Study (1)

---

- Greedy algorithms [Gilbert et al., SODA'03; Tropp, TIT'04]
  - Process: iteratively select or abandon one variable that makes the criterion currently optimized
  - Weakness: get stuck in local optima due to the greedy behavior



# Case Study (1)

---

- Convex relaxation methods [Tibshirani, JRSSB'96; Zou & Hastie, JRSSB'05]
  - Process: replace the set size constraint with convex constraints, then find the optimal solutions to the relaxed problem.
  - Weakness: the optimal solution of the relaxed problem may be distant to the true optimum.

$$\begin{array}{l} \arg \min_{S \subseteq V} f(S) \quad s.t. \quad |S| \leq k \\ \updownarrow \\ \arg \min_{w \in R^n} g(w) \quad s.t. \quad |w|_0 \leq k \quad \text{non-convex} \\ \curvearrowright \\ \arg \min_{w \in R^n} g(w) \quad s.t. \quad |w|_1 \leq k \quad \text{convex} \end{array}$$

# Case Study (1)

---

- There have been numerous EAs for subset selection, while rigorous theoretical guarantee is few.
- Subset Selection as a bi-objective optimization problem

## POSS (Pareto Optimization for Subset Selection)

The basic idea:

$$\min_{S \subseteq V} f(S) \quad s.t. \quad |S| \leq k \quad \text{constrained}$$

↓

$$\min_{S \subseteq V} (f(S), |S|) \quad \text{bi-objective}$$

# Case Study (1)

---

## Algorithm 1 POSS

**Input:** all variables  $V = \{X_1, \dots, X_n\}$ , a given objective  $f$  and an integer parameter  $k \in [1, n]$

**Parameter:** the number of iterations  $T$

**Output:** a subset of  $V$  with at most  $k$  variables

**Process:**

1: Let  $s = \{0\}^n$  and  $P = \{s\}$ .

2: Let  $t = 0$ .

3: **while**  $t < T$  **do**

4:   Select  $s$  from  $P$  uniformly at random.

5:   Generate  $s'$  by flipping each bit of  $s$  with prob.  $\frac{1}{n}$ .

6:   Evaluate  $f_1(s')$  and  $f_2(s')$ .

7:   **if**  $\nexists z \in P$  such that  $z \prec s'$  **then**

8:      $Q = \{z \in P \mid s' \preceq z\}$ .

9:      $P = (P \setminus Q) \cup \{s'\}$ .

10:   **end if**

11:    $t = t + 1$ .

12: **end while**

13: **return**  $\arg \min_{s \in P, |s| \leq k} f_1(s)$

**Initialization:** randomly generate a solution, put it into the archive  $P$

**Reproduction:** pick a solution randomly from  $P$ , and randomly change it to make a new one

**Evaluation & Selection:** if the new solution is not dominated, put it into  $P$  and weed out bad solutions

**Output:** select the best feasible solution

- Chaoqian paper

# Case Study (1)

---

- Sparse regression is to find a sparse approximation solution to the regression problem.

Formally stated: given all observation variables  $V = \{X_1, \dots, X_n\}$ , a predictor variable  $Z$  and a positive integer  $k$ , define the mean squared error of a subset  $S \subseteq V$  as

$$MSE_{Z,S} = \min_{\alpha \in \mathbb{R}^{|S|}} E[(Z - \sum_{i \in S} \alpha_i X_i)^2]$$

Sparse regression is

$$\arg \min_{S \subseteq V} MSE_{Z,S} \quad s. t. \quad |S| \leq k.$$

Chao Qian, Yang Yu, and Zhi-Hua Zhou. Subset Selection by Pareto Optimization. In: Advances in Neural Information Processing Systems 28 (NIPS'15), Montreal, Canada, 2015, pp.1765-1773.

# Case Study (1)

---

## Previous theoretical bounds:

$u$ : the coherence,  $\gamma$ : the submodular ratio

- [Gilbert et al., 2003]:  $(1 + \Theta(uk^2)) \cdot OPT$  on  $MSE_{Z,S}$  for  $u \in O(1/k)$  by a two-phased approach
- [Tropp et al., 2003; Tropp, 2004]: improve the above bound
- [Das & Kempe, 2008]:  $(1 - \Theta(uk)) \cdot OPT$  on  $R_{Z,S}^2$  for  $u \in O(1/k)$  by the forward regression algorithm
- [Das & Kempe, 2011]:  $(1 - e^{-\gamma}) \cdot OPT$  on  $R_{Z,S}^2$  by the forward regression algorithm
- [Shalev-Shwartz et al., 2010; Yuan & Yan, 2013]: lower bounds on  $|S|$  for achieving  $OPT + \varepsilon$

strongest

# Case Study (1)

---

**Theorem 1.** For sparse regression, POSS with  $E[T] \leq 2ek^2n$  and  $I(\cdot) = 0$  (i.e., a constant function) finds a set  $S$  of variables with  $|S| \leq k$  and  $R_{Z,S}^2 \geq (1 - e^{-\gamma\phi,k}) \cdot OPT$ .

the best  
previous  
theoretical  
guarantee

POSS can do at least as well as previous methods.

**Theorem 2.** For the Exponential Decay subclass of sparse regression, POSS with  $E[T] \in O(k^2n^2 \log n)$  and  $I(\mathbf{s} \in \{0, 1\}^n) = \min\{i \mid s_i = 1\}$  can find the optimal solution.

**Proposition 1.** For Example 1 with  $n = 3$ ,  $r_2 = 0.03$ ,  $r_3 = 0.5$ ,  $Cov(Y_1, Z) = Cov(Y_2, Z) = \delta$  and  $Cov(Y_3, Z) = 0.505\delta$ , FR cannot find the optimal solution for  $k = 2$ .

POSS can do strictly better than previous methods.

# Case Study (1)

- POSS for Sparse Regression: Summary

Theory

**Theorem 1.** For sparse regression, POSS with  $\mathbb{E}[T] \leq 2ek^2n$  finds a set  $S$  of variables with  $|S| \leq k$  and  $R_{Z,S}^2 \geq (1 - e^{-\gamma_{\theta,k}}) \cdot OPT.$

the number of iterations

the best known polynomial-time approximation bound  
[Das & Kempe, ICML'11]

Experiment

| Data set           | OPT         | POSS        | FR           | FoBa         | OMP          | RFE          | MCP          |
|--------------------|-------------|-------------|--------------|--------------|--------------|--------------|--------------|
| housing            | .7437±.0297 | .7437±.0297 | .7429±.0300● | .7423±.0301● | .7415±.0300● | .7388±.0304● | .7354±.0297● |
| eunite2001         | .8484±.0132 | .8482±.0132 | .8348±.0143● | .8442±.0144● | .8349±.0150● | .8424±.0153● | .8320±.0150● |
| svmguide3          | .2705±.0255 | .2701±.0257 | .2615±.0260● | .2601±.0279● | .2557±.0270● | .2136±.0325● | .2397±.0237● |
| ionosphere         | .5995±.0326 | .5990±.0329 | .5920±.0352● | .5929±.0346● | .5921±.0353● | .5832±.0415● | .5740±.0348● |
| sonar              | –           | .5365±.0410 | .5171±.0440● | .5138±.0432● | .5112±.0425● | .4321±.0636● | .4496±.0482● |
| triazines          | –           | .4301±.0603 | .4150±.0592● | .4107±.0600● | .4073±.0591● | .3615±.0712● | .3793±.0584● |
| coil2000           | –           | .0627±.0076 | .0624±.0076● | .0619±.0075● | .0619±.0075● | .0363±.0141● | .0570±.0075● |
| mushrooms          | –           | .9912±.0020 | .9909±.0021● | .9909±.0022● | .9909±.0022● | .6813±.1294● | .8652±.0474● |
| clean1             | –           | .4368±.0300 | .4169±.0299● | .4145±.0309● | .4132±.0315● | .1596±.0562● | .3563±.0364● |
| w5a                | –           | .3376±.0267 | .3319±.0247● | .3341±.0258● | .3313±.0246● | .3342±.0276● | .2694±.0385● |
| gisette            | –           | .7265±.0098 | .7001±.0116● | .6747±.0145● | .6731±.0134● | .5360±.0318● | .5709±.0123● |
| farm-ads           | –           | .4217±.0100 | .4196±.0101● | .4170±.0113● | .4170±.0113● | –            | .3771±.0110● |
| POSS: win/tie/loss |             | –           | 12/0/0       | 12/0/0       | 12/0/0       | 11/0/0       | 12/0/0       |

significantly better than all the compared methods on all data sets

# Case Study (2)

---

## Algorithm 1 POSS

---

**Input:** all variables  $V = \{X_1, \dots, X_n\}$ , a given objective  $f$  and an integer parameter  $k \in [1, n]$

**Parameter:** the number of iterations  $T$

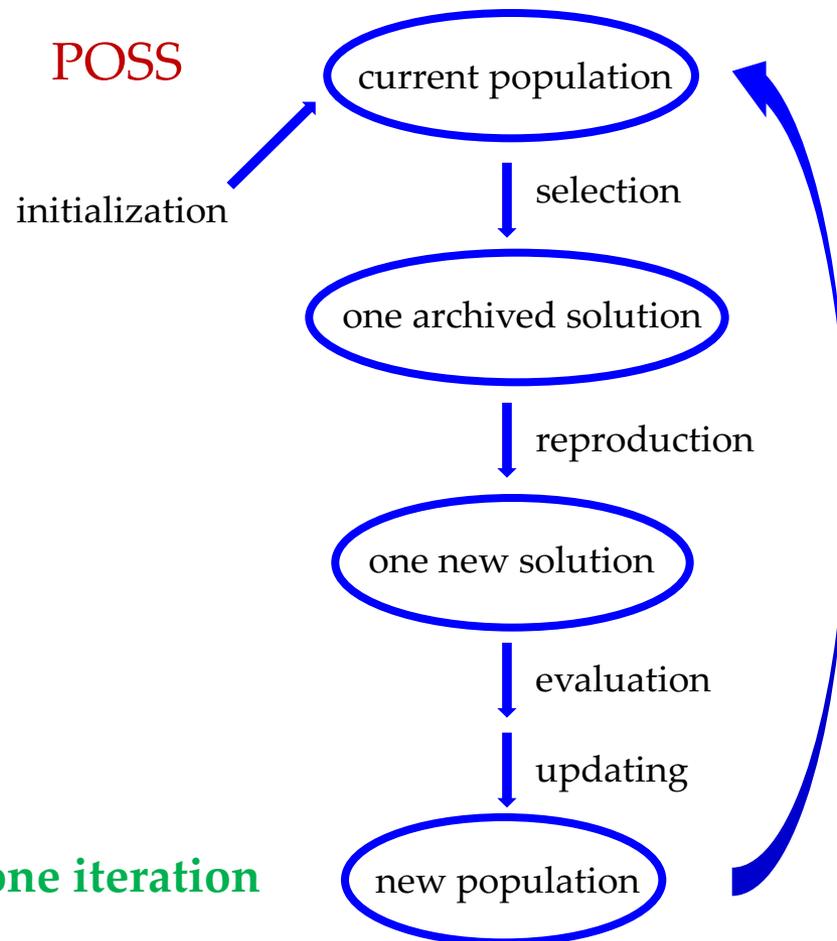
**Output:** a subset of  $V$  with at most  $k$  variables

**Process:**

- 1: Let  $s = \{0\}^n$  and  $P = \{s\}$ .
  - 2: Let  $t = 0$ .
  - 3: **while**  $t < T$  **do**
  - 4:   Select  $s$  from  $P$  uniformly at random.
  - 5:   Generate  $s'$  by flipping each bit of  $s$  with prob.  $\frac{1}{n}$ .
  - 6:   Evaluate  $f_1(s')$  and  $f_2(s')$ .
  - 7:   **if**  $\nexists z \in P$  such that  $z \prec s'$  **then**
  - 8:      $Q = \{z \in P \mid s' \preceq z\}$ .
  - 9:      $P = (P \setminus Q) \cup \{s'\}$ .
  - 10:   **end if**
  - 11:    $t = t + 1$ .
  - 12: **end while**
  - 13: **return**  $\arg \min_{s \in P, |s| \leq k} f_1(s)$
- 

generate only one new solution in one iteration

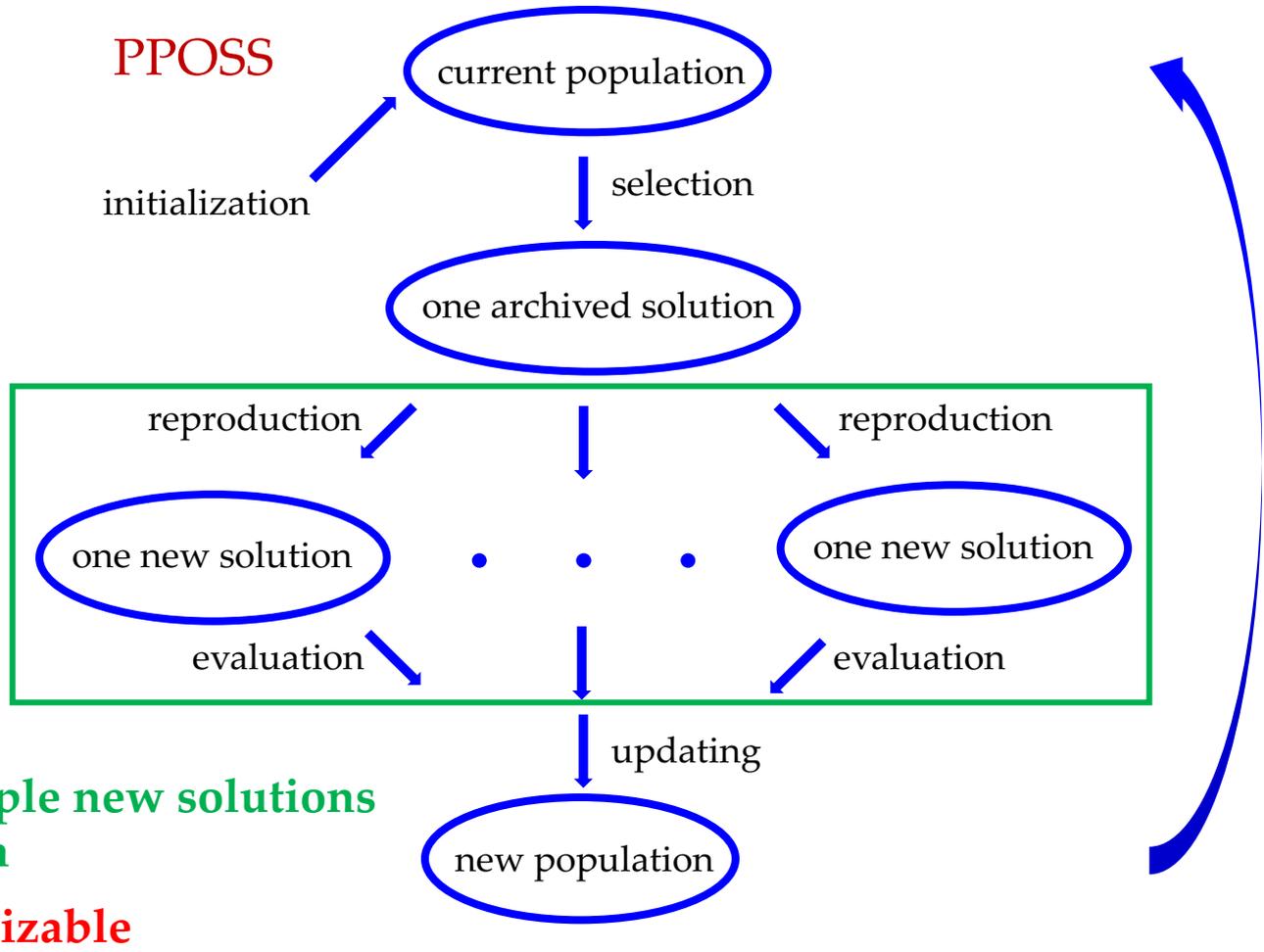
sequential



A sequential algorithm that cannot be readily parallelized  
restrict the application to large-scale real-world problems

# Case Study (2)

Theoretical guarantee still holds?



# Case Study (2)

- Theoretical results

$$f(S_2) \geq f(S_1) \text{ for any } S_1 \subseteq S_2$$

**Theorem 1.** For maximizing a **monotone** function under the set size constraint, the expected number of iterations until PPOSS finds a solution  $s$  with  $|s| \leq k$  and  $f(s) \geq (1 - e^{-\gamma_{\min}}) \cdot OPT$ , where  $\gamma_{\min} = \min_{s:|s|=k-1} \gamma_{s,k}$ , is

(1) if  $N = o(n)$ , then  $\mathbb{E}[T] \leq 2ek^2n/N$ ;

(2) if  $N = \Omega(n^i)$  for  $1 \leq i \leq k$ , then  $\mathbb{E}[T] = O(k^2/i)$ ;

(3) if  $N = \Omega(n^{\min\{3k-1, n\}})$ , then  $\mathbb{E}[T] = O(1)$ .

the same approximation bound

- When the number of processors is less than the number of variables, the number of iterations can be reduced linearly w.r.t. the number of processors
- With increasing number of processors, the number of iterations can be continuously reduced, eventually to a **constant**

The best previous known bound

- submodular [Nemhauser & Wolsey, MOR'78]
- sparse regression (non-submodular) [Das & Kempe, ICML'11]

# Case Study (2)

Sparse  
regression

$$\arg \min_{S \subseteq V} MSE_{Z,S} \quad s.t. \quad |S| \leq k$$

$$R_{Z,S}^2 = (Var(Z) - MSE_{Z,S}) / Var(Z)$$

the larger  
the better

data set

| data set          | #inst | #feat | data set         | #inst | #feat | data set        | #inst | #feat |
|-------------------|-------|-------|------------------|-------|-------|-----------------|-------|-------|
| <i>housing</i>    | 506   | 13    | <i>sonar</i>     | 208   | 60    | <i>clean1</i>   | 476   | 166   |
| <i>eunite2001</i> | 367   | 16    | <i>triazines</i> | 186   | 60    | <i>w5a</i>      | 9888  | 300   |
| <i>svmguide3</i>  | 1284  | 21    | <i>coil2000</i>  | 9000  | 86    | <i>gisette</i>  | 7000  | 5000  |
| <i>ionosphere</i> | 351   | 34    | <i>mushrooms</i> | 8124  | 112   | <i>farm-ads</i> | 4143  | 54877 |

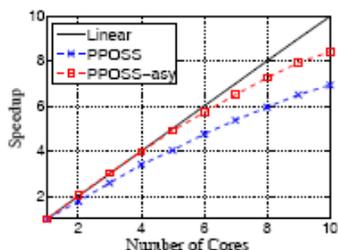
the sparsity  $k = 8$

the number of cores  $N = 1 \rightarrow 10$

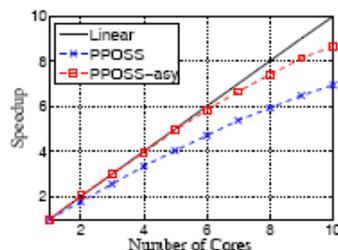
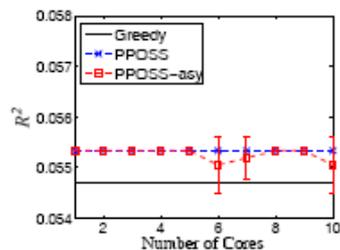
For PPOSS with each  $N$  value on each data set, the run is repeated for 10 runs independently, and the average results are reported

# Case Study (2)

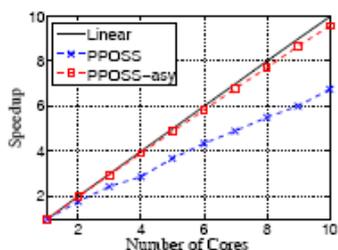
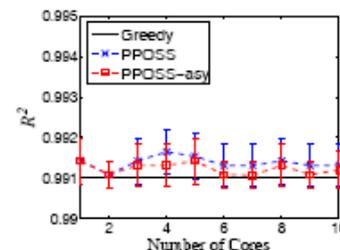
Compare the **speedup** as well as the solution quality measured by  $R^2$  values with different number of cores



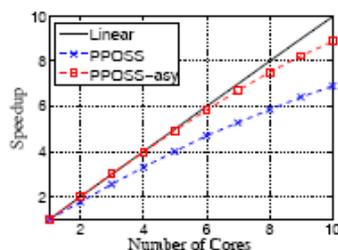
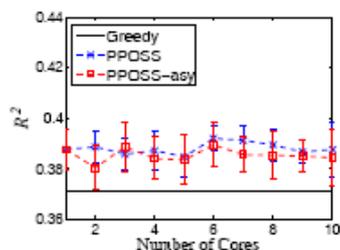
(a) on *coil2000* (9000 instances, 86 features)



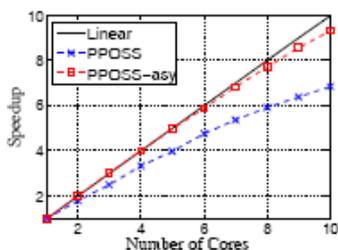
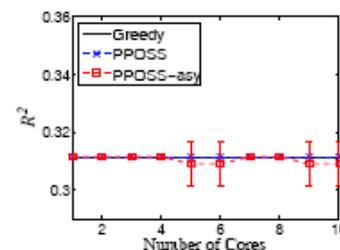
(b) on *mushrooms* (8124 instances, 112 features)



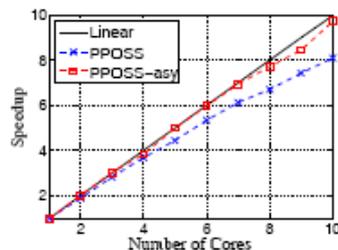
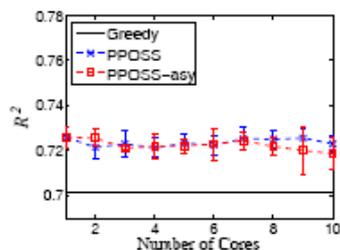
(c) on *clean1* (476 instances, 166 features)



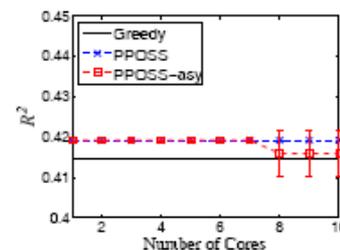
(d) on *w5a* (9888 instances, 300 features)



(e) on *gisette* (7000 instances, 5000 features)

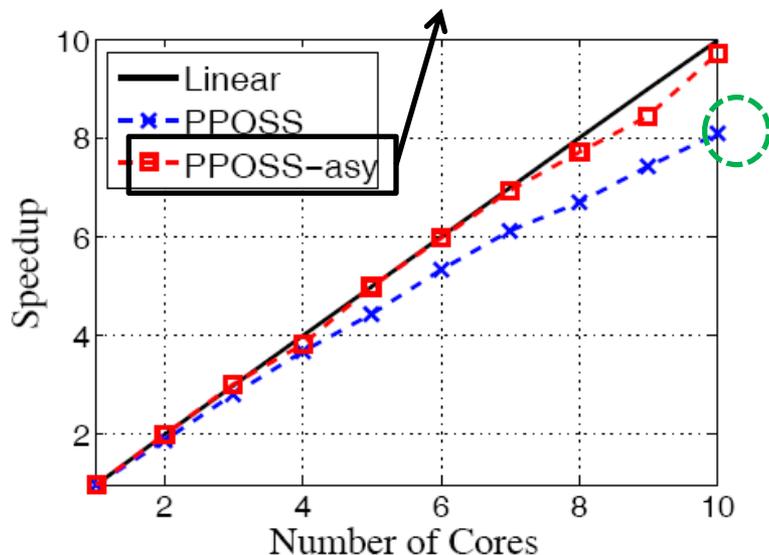


(f) on *farm-ads* (4143 instances, 54877 features)

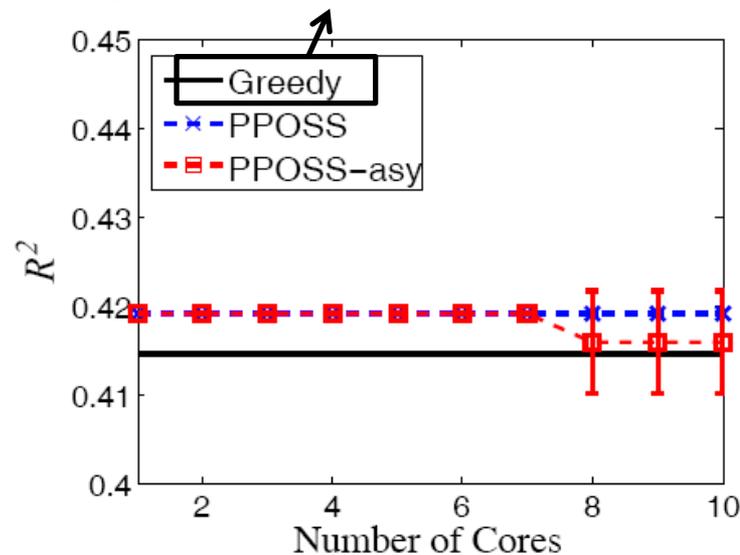


# Case Study (2)

the asynchronous version of PPOSS



the best previous algorithm [Das & Kempe, ICML'11]



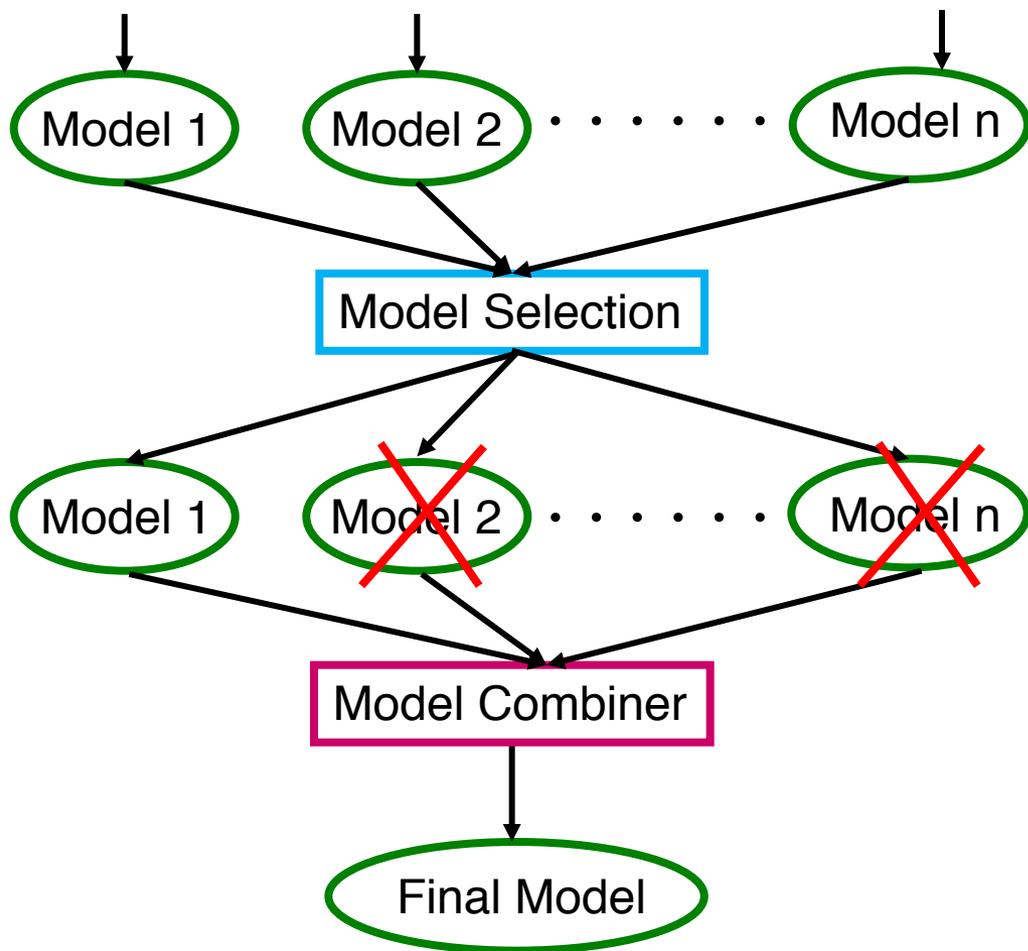
(f) on *farm-ads* (4143 instances, 54877 features)

**PPOSS (blue line):** achieve speedup around 8 when the number of cores is 10; the  $R^2$  values are stable

**PPOSS-asy (red line):** achieve better speedup (avoid the synchronous cost); the  $R^2$  values are slightly worse (the noise from asynchronization)

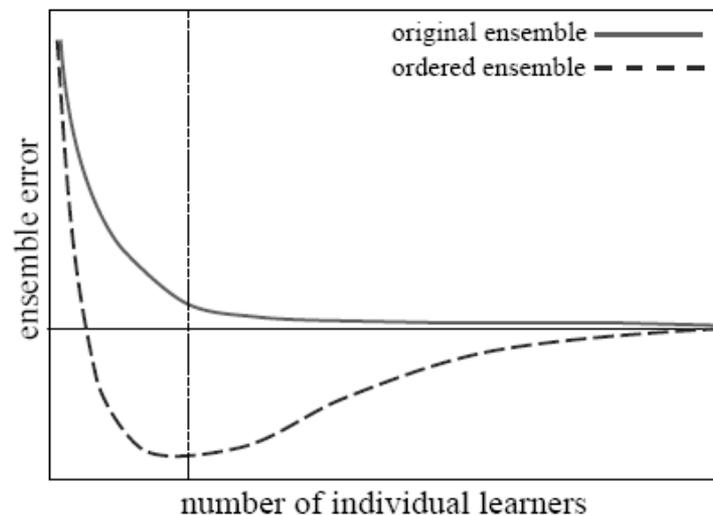
# Case Study (3)

- Ensemble Pruning is also a subset selection problem.



**Ensemble pruning:** select a subset of learners to combine

- reduce storage & improve efficiency
- better performance than the complete ensemble



[Martinez-Munoz and Suarez, ICML'06]

# Case Study (3)

- PEP: Pareto Ensemble Pruning

| Data set                 | Test Error       |                  |            |                  |                  |                  |            |                  |                  |
|--------------------------|------------------|------------------|------------|------------------|------------------|------------------|------------|------------------|------------------|
|                          | PEP              | Bagging          | BI         | RE               | Kappa            | CP               | MD         | DREP             | EA               |
| australian               | .144±.020        | <b>.143±.017</b> | .152±.023● | .144±.020        | <b>.143±.021</b> | .145±.022        | .148±.022  | .144±.019        | <b>.143±.020</b> |
| breast-cancer disorders  | <b>.275±.041</b> | .279±.037        | .298±.044● | .277±.031        | .287±.037        | .282±.043        | .295±.044● | <b>.275±.036</b> | <b>.275±.032</b> |
| heart-statlog            | <b>.304±.039</b> | .327±.047●       | .365±.047● | .320±.044●       | .326±.042●       | .306±.039        | .337±.035● | .316±.045        | .317±.046●       |
| house-votes              | .197±.037        | .195±.038        | .235±.049● | <b>.187±.044</b> | .201±.038        | .199±.044        | .226±.048● | .194±.044        | .196±.032        |
| ionosphere               | .045±.019        | <b>.041±.013</b> | .047±.016  | .043±.018        | .044±.017        | .045±.017        | .048±.018● | .045±.017        | <b>.041±.012</b> |
| kr-vs-kp                 | .088±.021        | .092±.025        | .117±.022● | .086±.021        | <b>.084±.020</b> | .089±.021        | .100±.026● | .085±.021        | .093±.026        |
| letter-ah                | <b>.010±.003</b> | .015±.007●       | .011±.004  | <b>.010±.004</b> | <b>.010±.003</b> | .011±.003        | .011±.005  | .011±.003        | .012±.004        |
| letter-br                | .013±.005        | .021±.006●       | .023±.008● | .015±.006●       | <b>.012±.006</b> | .015±.006        | .017±.007● | .014±.005        | .017±.006●       |
| letter-oq                | <b>.046±.008</b> | .059±.013●       | .078±.012● | .048±.012        | .048±.014        | .048±.012        | .057±.014● | .048±.009        | .053±.011●       |
| optdigits                | .043±.009        | .049±.012●       | .078±.017● | .046±.011        | .042±.011        | .042±.010        | .046±.011  | <b>.041±.010</b> | .044±.011        |
| satimage-12v57           | <b>.035±.006</b> | .038±.007●       | .095±.008● | .036±.006        | <b>.035±.005</b> | .036±.005        | .037±.006● | <b>.035±.006</b> | <b>.035±.006</b> |
| satimage-2v5             | <b>.028±.004</b> | .029±.004        | .052±.006● | .029±.004        | <b>.028±.004</b> | .029±.004        | .029±.004  | .029±.004        | .029±.004        |
| sick                     | <b>.021±.007</b> | .023±.009        | .033±.010● | .023±.007        | .022±.007        | <b>.021±.008</b> | .026±.010● | .022±.008        | <b>.021±.008</b> |
| sonar                    | <b>.015±.003</b> | .018±.004●       | .018±.004● | .016±.003        | .017±.003●       | .016±.003●       | .017±.003● | .016±.003        | .017±.004●       |
| spambase                 | <b>.248±.056</b> | .266±.052        | .310±.051● | .267±.053●       | .249±.059        | .250±.048        | .268±.055● | .257±.056        | .251±.041        |
| tic-tac-toe              | <b>.065±.006</b> | .068±.007●       | .093±.008● | .066±.006        | .066±.006        | .066±.006        | .068±.007● | <b>.065±.006</b> | .066±.006        |
| vehicle-bo-vs            | .131±.027        | .164±.028●       | .212±.028● | .135±.026        | .132±.023        | .132±.026        | .145±.022● | <b>.129±.026</b> | .138±.020        |
| vehicle-b-v              | <b>.224±.023</b> | .228±.026        | .257±.025● | .226±.022        | .233±.024●       | .234±.024●       | .244±.024● | .234±.026●       | .230±.024        |
| vote                     | <b>.018±.011</b> | .027±.014●       | .024±.013● | .020±.011        | .019±.012        | .020±.011        | .021±.011● | .019±.013        | .026±.013●       |
| count of the best        | <b>12</b>        | 2                | 0          | 2                | 7                | 1                | 0          | 5                | 5                |
| PEP: count of direct win |                  | <b>17</b>        | <b>20</b>  | <b>15.5</b>      | 12.5             | <b>17</b>        | <b>20</b>  | 12.5             | <b>15.5</b>      |

**Better than any other method on more than 60% (12.5/20) data sets, and never significantly worse**

# Case Study (3)

| Data set                 | Ensemble Size   |                 |             |                  |                 |                 |           |
|--------------------------|-----------------|-----------------|-------------|------------------|-----------------|-----------------|-----------|
|                          | PEP             | RE              | Kappa       | CP               | MD              | DREP            | EA        |
| australian               | 10.6±4.2        | 12.5±6.0        | 14.7±12.6   | 11.0±9.7         | <b>8.5±14.8</b> | 11.7±4.7        | 41.9±6.7● |
| breast-cancer            | 8.4±3.5         | 8.7±3.6         | 26.1±21.7●  | 8.8±12.3         | <b>7.8±15.2</b> | 9.2±3.7         | 44.6±6.6● |
| disorders                | 14.7±4.2        | <b>13.9±4.2</b> | 24.7±16.3●  | 15.3±10.6        | 17.7±20.0       | <b>13.9±5.9</b> | 42.0±6.2● |
| heart-statlog            | <b>9.3±2.3</b>  | 11.4±5.0●       | 17.9±11.1●  | 13.2±8.2●        | 13.6±21.1       | 11.3±2.7●       | 44.2±5.1● |
| house-votes              | <b>2.9±1.7</b>  | 3.9±4.0         | 5.5±3.3●    | 4.7±4.4●         | 5.9±14.1        | 4.1±2.7●        | 46.5±6.1● |
| ionosphere               | <b>5.2±2.2</b>  | 7.9±5.7●        | 10.5±6.9●   | 8.5±6.3●         | 10.7±14.6●      | 8.4±4.3●        | 48.8±5.1● |
| kr-vs-kp                 | <b>4.2±1.8</b>  | 5.8±4.5         | 10.6±9.1●   | 9.6±8.6●         | 7.2±15.2        | 7.1±3.9●        | 45.9±5.8● |
| letter-ah                | <b>5.0±1.9</b>  | 7.3±4.4●        | 7.1±3.8●    | 8.7±4.7●         | 11.0±10.9●      | 7.8±3.6●        | 42.5±6.5● |
| letter-br                | <b>10.9±2.6</b> | 15.1±7.3●       | 13.8±6.7●   | 12.9±6.8         | 23.2±17.6●      | 11.3±3.5        | 38.3±7.8● |
| letter-oq                | <b>12.0±3.7</b> | 13.6±5.8        | 13.9±6.0    | 12.3±4.9         | 23.0±15.6●      | 13.7±4.9        | 39.3±8.2● |
| optdigits                | 22.7±3.1        | 25.0±9.3        | 25.2±8.1    | <b>21.4±7.5</b>  | 46.8±23.9●      | 25.0±8.0        | 41.4±7.6● |
| satimage-12v57           | <b>17.1±5.0</b> | 20.8±9.2●       | 22.1±10.3●  | 21.2±10.0●       | 37.6±24.3●      | 18.1±4.9        | 42.7±5.2● |
| satimage-2v5             | <b>5.7±1.7</b>  | 6.8±3.2         | 7.6±4.2●    | 10.9±7.0●        | 26.2±28.1●      | 7.7±3.5●        | 44.1±4.8● |
| sick                     | <b>6.9±2.8</b>  | 7.5±3.9         | 10.9±6.0●   | 11.5±10.0●       | 8.3±13.6        | 11.6±6.7●       | 44.7±8.2● |
| sonar                    | 11.4±4.2        | <b>11.0±4.1</b> | 20.6±9.3●   | 13.9±7.1         | 20.6±20.7●      | 14.4±5.9●       | 43.1±6.4● |
| spambase                 | 17.5±4.5        | 18.5±5.0        | 20.0±8.1    | 19.0±9.9         | 28.8±17.0●      | <b>16.7±4.6</b> | 39.7±6.4● |
| tic-tac-toe              | 14.5±3.8        | 16.1±5.4        | 17.4±6.5    | 15.4±6.3         | 28.0±22.6●      | <b>13.6±3.4</b> | 39.8±8.2● |
| vehicle-bo-vs            | 16.5±4.5        | 15.7±5.7        | 16.5±8.2    | <b>11.2±5.7○</b> | 21.6±20.4       | <b>3.2±5.0○</b> | 41.9±5.6● |
| vehicle-b-v              | <b>2.8±1.1</b>  | 3.4±2.1         | 4.5±1.6●    | 5.3±7.4          | <b>2.8±3.8</b>  | 4.0±3.9         | 48.0±5.6● |
| vote                     | <b>2.7±1.1</b>  | 3.2±2.7         | 5.1±2.6●    | 5.4±5.2●         | 6.0±9.8         | 3.9±2.5●        | 47.8±6.1● |
| count of the best        | <b>12</b>       | 2               | 0           | 2                | 3               | 3               | 0         |
| PEP: count of direct win |                 | <b>17</b>       | <b>19.5</b> | <b>18</b>        | <b>17.5</b>     | <b>16</b>       | <b>20</b> |

**Better than any other method on more than 80% (16/20) data sets; never significantly worse, except two losses**

# Case Study (4)

---

- Many real-world classification problems are cost-sensitive.
- The optimal classifier for a binary classification problem

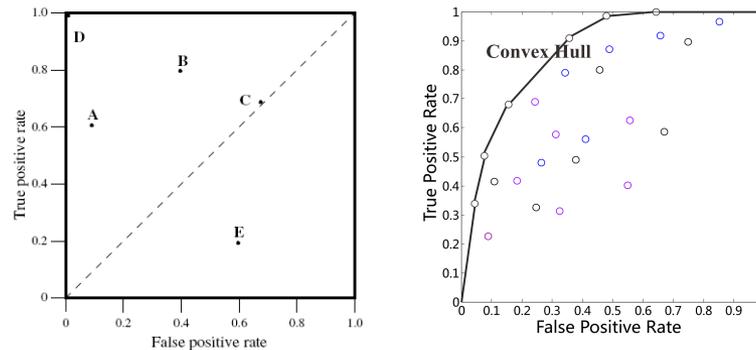
$$h_{opt} = \operatorname{argmin}_{h \in \Omega} C_{10} \cdot FNR(h) + C_{01} \cdot FPR(h) \quad (1)$$

| Cost matrix |             |             | Confusion matrix |             |             |
|-------------|-------------|-------------|------------------|-------------|-------------|
|             | Predicted + | Predicted - |                  | Predicted + | Predicted - |
| +           | 0           | $C_{10}$    | +                | TPR         | FNR         |
| -           | $C_{01}$    | 0           | -                | FPR         | TNR         |

- A good classifier could be used via *cost-sensitive learning*.

# Case Study (4)

- However, costs are often subject to great **uncertainty**.
  - Very difficult to specify the exact cost values *before training*.
  - The costs may change over time.
- Alternative: seeking **a group of classifiers** that maximize the the Receiver Operating Characteristic (ROC) convex hull.



*“The optimal classifier for any cost values must be a vertex or on the edge of the **convex hull** of all (achievable) classifiers in the ROC space.” [Provost and Fawcett, 2001]*

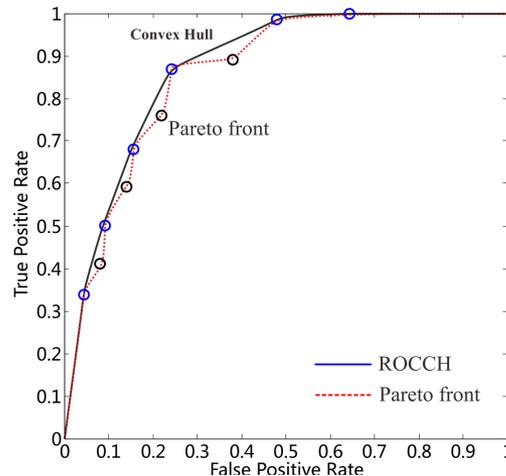
# Case Study (4)

---

- **New Learning Target:** To obtain a set of classifiers such that their ROCCH is maximized.
- This is a set-oriented optimization problem can hardly be solved with existing approaches.
- Evolutionary Algorithms provides a natural way to search for a set (population) of classifiers.

# Case Study (4)

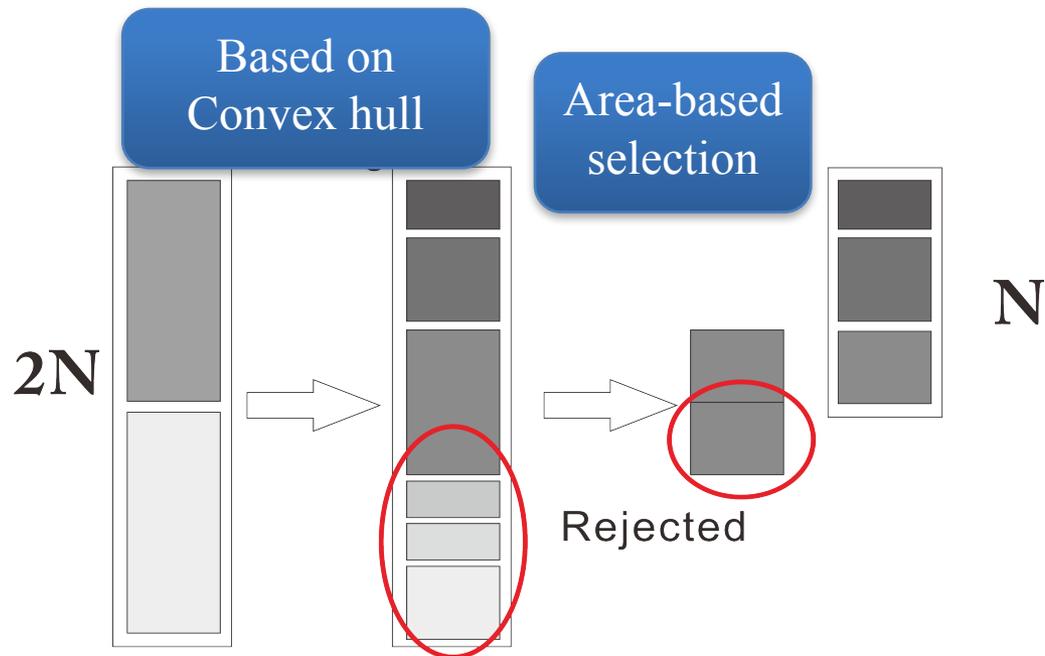
- Multi-objective evolutionary algorithms (MOEAs) are off-the-shelf tools for this problem
  - Maximize TPR
  - Minimize FPR
- Direct application of an MOEA is OK, while not ideal.
  - A Pareto optimal solution is not necessarily a vertex on the convex hull.
  - Many-to-one mapping between the hypothesis and ROC spaces.



# Case Study (4)

---

- Approach: Convex Hull-based MOEA (CH-MOEA)
- Features of CH-MOEA:
  - Convex hull-based sorting
  - Redundancy elimination



# Case Study (4)

- Convex Hull-based Sorting

---

## Algorithm 2 *DeltaArea* ( $T$ )

---

**Require:**  $T \neq \emptyset$

1:  $T$  is a solution set

**Ensure:** *DeltaArea*

2:  $Q = T$

3:  $m = \text{sizeof}(Q)$

4:  $\mathbf{E}$  is performance of  $Q$

5:  $\mathbf{DeltaH}_1, \dots, \mathbf{DeltaH}_m \leftarrow 0$

6: **if**  $m < 3$  **then**

7:   Set  $\mathbf{DeltaH}_1, \dots, \mathbf{DeltaH}_m \leftarrow \infty$

8: **else**

9:   Set  $\mathbf{DeltaH}_1, \mathbf{DeltaH}_m \leftarrow \infty$

10:   **for**  $2 \leq i \leq \text{sizeof}(Q) - 1$  **do**

11:      $\mathbf{DeltaH}_i = 0.5 \times \det((\mathbf{E}_i - \mathbf{E}_{i-1}) \circ (\mathbf{E}_{i+1} - \mathbf{E}_{i-1}))$

12:   **end for**

13:   **while**  $\text{sizeof}(Q) > 2$  **do**

14:      $r \leftarrow \text{argmin}\{\mathbf{DeltaH}\}$

15:      $Q \leftarrow Q \setminus \{Q_r\}$

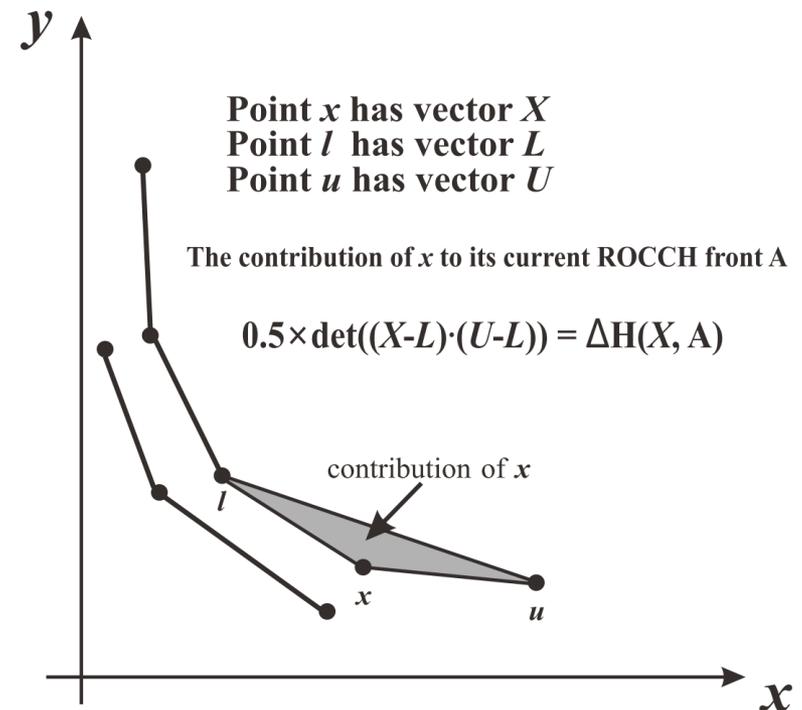
16:     Update( $\mathbf{DeltaH}_{r-1}, \mathbf{DeltaH}_{r+1}$ )

17:   **end while**

18: **end if**

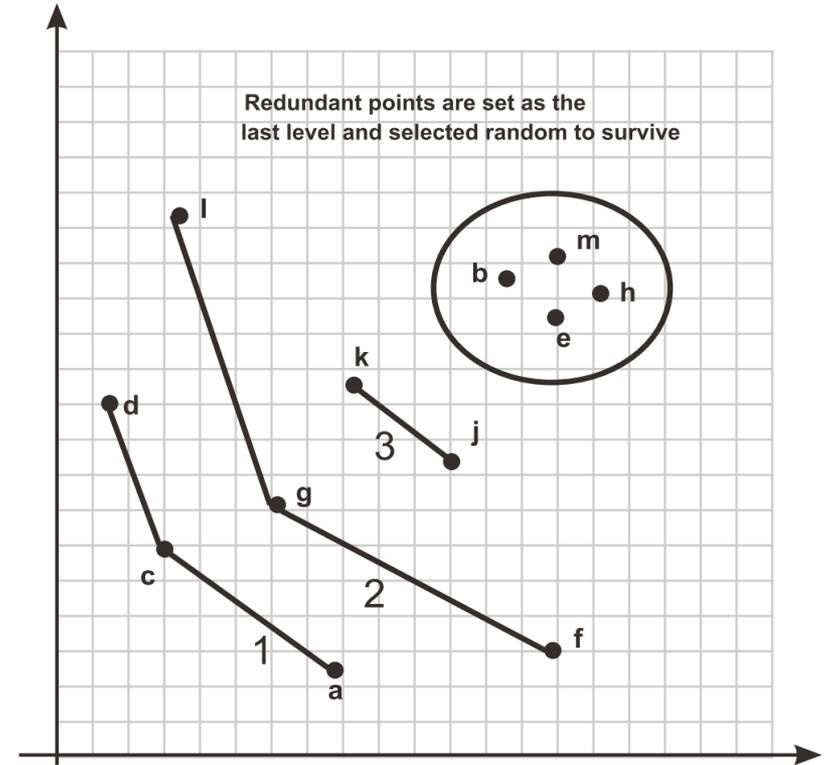
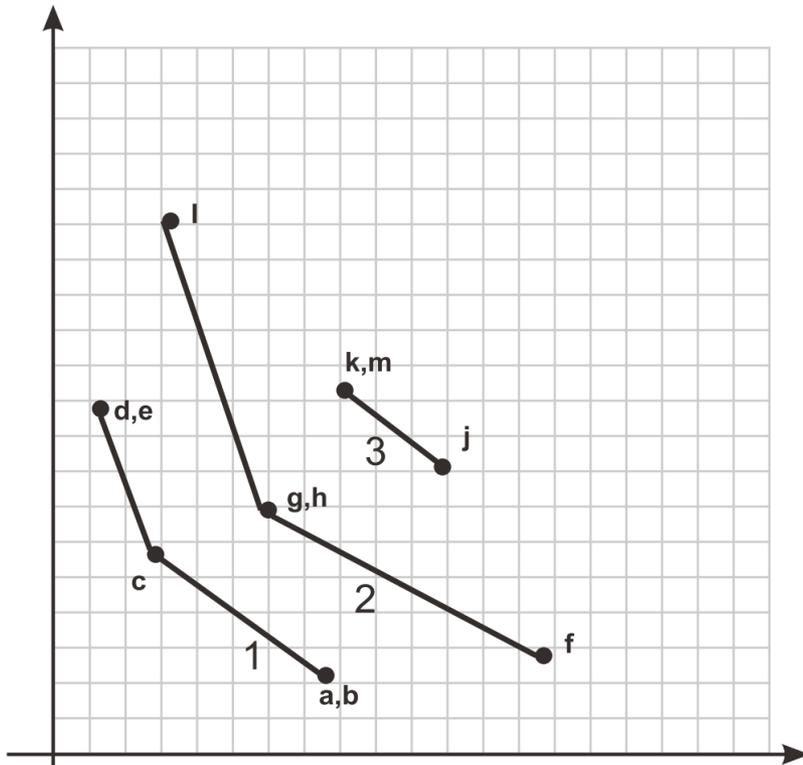
19: **Return** ( $\mathbf{DeltaH}$ )

---



# Case Study (4)

- Redundancy Elimination



# Case Study (4)

---

- CH-MOEA can be combined with any base learners to build either homogeneous and heterogeneous ensembles
  - Neural Network
  - Decision Tree
  - SVM
  - ...
- Different types of base learners need different search operators.
- We implemented CH-MOEA with Genetic Programming (CH-MOGP).

# Case Study (4)

---

- Empirical studies
  - Which MOEA framework performs the best for our problem?
  - Is CH-MOGP competitive in comparison to non-evolutionary methods?
- Compared methods
  - NSGA-II
  - MOEA/D
  - SMS-EMOA (an indicator based MOEA)
  - C4.5
  - PRIE ([Fawcett, 2008], a state-of-the-art heuristic approach for ROCCH maximization)
  - Naïve Bayes

All evolutionary approaches adopt the same base learner and reproduction operator

# Case Study (4)

- CH-MOGP outperformed state-of-the-art MOEAs

Performance on Nineteen UCI Datasets

| Dataset             | CH-MOGP      | SMS-EMOA     | NSGA-II      | MOEA/D       | Dataset            | CH-MOGP      | SMS-EMOA     | NSGA-II      | MOEA/D       |
|---------------------|--------------|--------------|--------------|--------------|--------------------|--------------|--------------|--------------|--------------|
| <i>australian</i>   | 91.49 ± 2.72 | 91.67 ± 2.48 | 91.16 ± 2.41 | 90.29 ± 2.75 | <i>bands</i>       | 77.00 ± 4.05 | 76.38 ± 4.09 | 75.54 ± 3.56 | 71.85 ± 3.82 |
| <i>bcw</i>          | 97.94 ± 1.20 | 97.73 ± 1.56 | 97.84 ± 1.41 | 97.48 ± 1.48 | <i>crx</i>         | 91.30 ± 2.45 | 91.16 ± 2.33 | 91.14 ± 2.36 | 89.88 ± 2.51 |
| <i>german</i>       | 73.10 ± 3.24 | 73.32 ± 3.33 | 72.39 ± 3.07 | 71.45 ± 2.85 | <i>house-votes</i> | 97.94 ± 1.56 | 97.69 ± 1.59 | 97.74 ± 1.71 | 97.15 ± 1.75 |
| <i>ionosphere</i>   | 91.07 ± 4.95 | 90.51 ± 4.52 | 90.45 ± 4.53 | 89.89 ± 4.83 | <i>kr-vs-kp</i>    | 98.40 ± 0.89 | 98.63 ± 0.75 | 98.39 ± 0.79 | 96.67 ± 1.43 |
| <i>mammographic</i> | 89.75 ± 2.01 | 89.48 ± 1.94 | 89.41 ± 1.87 | 87.50 ± 2.23 | <i>monks-1</i>     | 99.70 ± 1.68 | 97.62 ± 3.71 | 99.62 ± 1.35 | 96.51 ± 5.69 |
| <i>monks-2</i>      | 91.05 ± 8.00 | 89.28 ± 5.58 | 90.53 ± 5.19 | 73.26 ± 9.14 | <i>monks-3</i>     | 99.81 ± 0.43 | 99.74 ± 0.45 | 99.45 ± 2.87 | 99.07 ± 0.88 |
| <i>parkinsons</i>   | 86.79 ± 6.86 | 85.11 ± 6.68 | 84.90 ± 7.54 | 83.94 ± 6.72 | <i>pima</i>        | 80.08 ± 3.38 | 79.85 ± 3.38 | 79.29 ± 3.70 | 76.93 ± 3.10 |
| <i>sonar</i>        | 79.42 ± 5.87 | 78.04 ± 5.91 | 77.79 ± 7.34 | 75.75 ± 5.66 | <i>spect</i>       | 77.38 ± 7.36 | 76.27 ± 7.14 | 76.91 ± 8.46 | 74.88 ± 6.43 |
| <i>tic-tac-toe</i>  | 83.40 ± 10.4 | 79.56 ± 11.1 | 79.07 ± 13.4 | 70.85 ± 10.4 | <i>transfusion</i> | 71.62 ± 4.62 | 71.48 ± 4.47 | 71.49 ± 4.84 | 68.77 ± 4.63 |
| <i>wdbc</i>         | 96.78 ± 1.92 | 96.49 ± 2.25 | 96.70 ± 2.11 | 95.90 ± 2.19 |                    |              |              |              |              |

Performance on Three Large-scaled UCI Datasets

| Dataset      | CH-MOGP      | SMS-EMOA     | NSGA-II      | MOEA/D       | Dataset        | CH-MOGP      | SMS-EMOA     | NSGA-II      | MOEA/D       |
|--------------|--------------|--------------|--------------|--------------|----------------|--------------|--------------|--------------|--------------|
| <i>adult</i> | 84.58 ± 1.40 | 82.53 ± 2.15 | 84.01 ± 1.38 | 77.04 ± 2.54 | <i>magic04</i> | 83.02 ± 1.04 | 81.76 ± 1.57 | 82.01 ± 1.19 | 76.39 ± 3.07 |
| <i>skin</i>  | 97.10 ± 1.11 | 95.46 ± 1.85 | 96.57 ± 1.25 | 93.20 ± 2.37 |                |              |              |              |              |

Results for 19 UCI Dataset

| Ratio           | $\frac{1}{15}$ | $\frac{1}{10}$ | $\frac{1}{4}$ | $\frac{1}{3}$ | $\frac{1}{2}$ | $\frac{2}{3}$ | 1      |
|-----------------|----------------|----------------|---------------|---------------|---------------|---------------|--------|
| <i>NSGA-II</i>  | 4-15-0         | 4-15-0         | 2-17-0        | 4-15-0        | 5-14-0        | 5-14-0        | 4-15-0 |
| <i>SMS-EMOA</i> | 11-8-0         | 11-8-0         | 6-13-0        | 5-14-0        | 4-15-0        | 4-15-0        | 5-14-0 |
| <i>MOEA/D</i>   | 19-0-0         | 19-0-0         | 19-0-0        | 19-0-0        | 19-0-0        | 19-0-0        | 19-0-0 |

Results for 3 Large-Scaled UCI Dataset

| Ratio           | $\frac{1}{15}$ | $\frac{1}{10}$ | $\frac{1}{4}$ | $\frac{1}{3}$ | $\frac{1}{2}$ | $\frac{2}{3}$ | 1     |
|-----------------|----------------|----------------|---------------|---------------|---------------|---------------|-------|
| <i>NSGA-II</i>  | 0-3-0          | 1-2-0          | 1-2-0         | 1-2-0         | 2-1-0         | 2-1-0         | 2-1-0 |
| <i>SMS-EMOA</i> | 3-0-0          | 3-0-0          | 3-0-0         | 3-0-0         | 3-0-0         | 3-0-0         | 3-0-0 |
| <i>MOEA/D</i>   | 3-0-0          | 3-0-0          | 3-0-0         | 3-0-0         | 3-0-0         | 3-0-0         | 3-0-0 |

# Case Study (4)

| Dataset             | CH-MOGP      | C4.5         | NB           | PRIE         |
|---------------------|--------------|--------------|--------------|--------------|
| <i>australian</i>   | 91.97 ± 2.53 | 85.52 ± 4.05 | 89.47 ± 2.78 | 91.75 ± 2.36 |
| <i>bands</i>        | 78.50 ± 3.56 | 74.56 ± 4.59 | 73.91 ± 4.68 | 76.07 ± 4.81 |
| <i>bcw</i>          | 98.17 ± 1.06 | 95.05 ± 2.55 | 98.92 ± 0.62 | 98.16 ± 1.09 |
| <i>crx</i>          | 91.82 ± 2.27 | 85.51 ± 3.94 | 87.88 ± 3.16 | 90.65 ± 2.77 |
| <i>german</i>       | 74.27 ± 2.79 | 65.36 ± 4.74 | 78.42 ± 2.94 | 75.95 ± 3.25 |
| <i>house-votes</i>  | 98.23 ± 1.26 | 96.35 ± 2.04 | 98.05 ± 1.04 | 97.80 ± 1.49 |
| <i>ionosphere</i>   | 92.42 ± 3.66 | 88.20 ± 5.65 | 93.57 ± 3.18 | 93.68 ± 4.23 |
| <i>kr-vs-kp</i>     | 99.40 ± 0.26 | 99.71 ± 0.23 | 93.21 ± 1.00 | 98.26 ± 0.44 |
| <i>mammographic</i> | 90.20 ± 1.76 | 87.66 ± 2.21 | 89.77 ± 1.96 | 89.70 ± 2.02 |
| <i>monks-1</i>      | 100.0 ± 0.00 | 77.13 ± 6.90 | 73.18 ± 4.58 | 70.93 ± 5.59 |
| <i>monks-2</i>      | 95.68 ± 4.61 | 94.17 ± 5.93 | 52.38 ± 7.04 | 51.25 ± 6.16 |

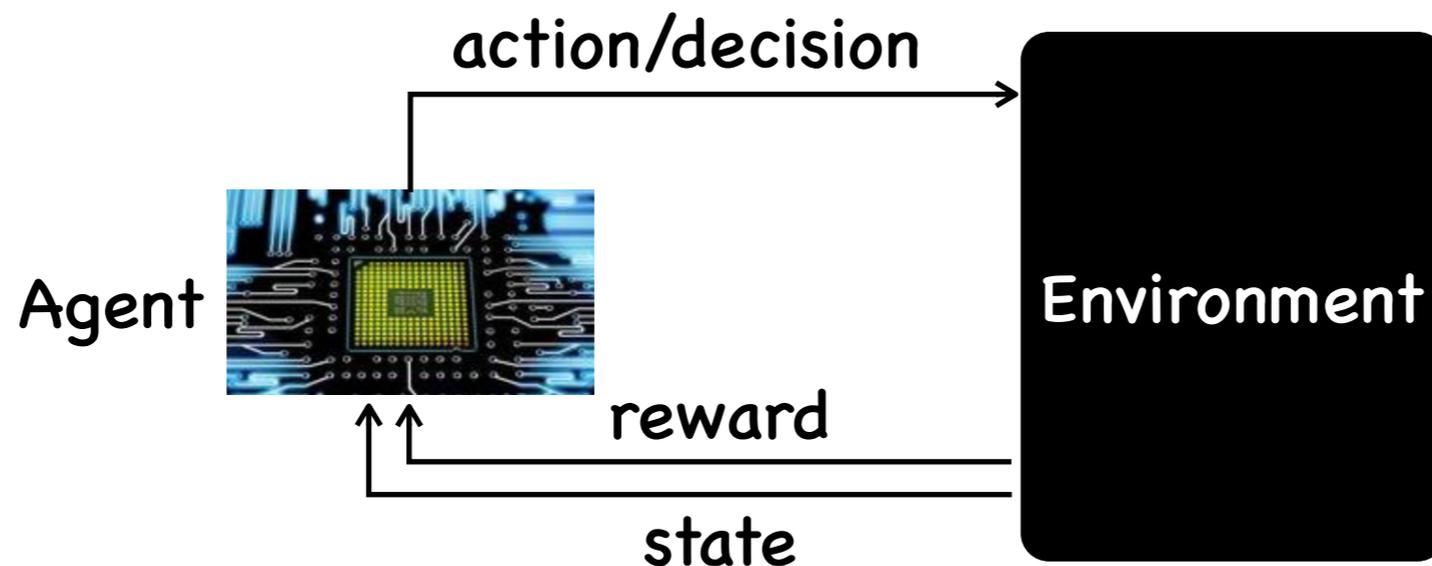
| Dataset            | CH-MOGP      | C4.5         | NB           | PRIE         |
|--------------------|--------------|--------------|--------------|--------------|
| <i>monks-3</i>     | 100.0 ± 0.00 | 100.0 ± 0.00 | 95.94 ± 2.17 | 99.60 ± 0.27 |
| <i>parkinsons</i>  | 86.10 ± 6.66 | 78.91 ± 9.76 | 85.91 ± 6.11 | 88.24 ± 5.83 |
| <i>pima</i>        | 80.74 ± 3.12 | 75.23 ± 4.93 | 81.40 ± 3.01 | 79.58 ± 2.92 |
| <i>sonar</i>       | 81.44 ± 5.15 | 73.85 ± 7.84 | 80.12 ± 7.03 | 69.92 ± 8.64 |
| <i>spect</i>       | 78.56 ± 7.44 | 76.88 ± 8.91 | 84.09 ± 6.03 | 83.51 ± 7.01 |
| <i>tic-tac-toe</i> | 90.07 ± 8.88 | 84.91 ± 13.9 | 61.50 ± 14.7 | 70.41 ± 12.5 |
| <i>transfusion</i> | 72.19 ± 4.89 | 71.08 ± 5.08 | 70.93 ± 4.94 | 70.87 ± 5.39 |
| <i>wdbc</i>        | 97.32 ± 1.40 | 92.74 ± 3.16 | 98.14 ± 1.33 | 96.58 ± 1.94 |
| <i>adult</i>       | 88.97 ± 0.37 | 88.89 ± 0.53 | 85.27 ± 0.37 | 90.37 ± 0.25 |
| <i>magic04</i>     | 87.16 ± 0.74 | 86.76 ± 0.83 | 75.70 ± 0.74 | 85.37 ± 0.76 |
| <i>skin</i>        | 99.49 ± 0.11 | 99.93 ± 0.02 | 94.17 ± 0.07 | 98.15 ± 0.08 |

**CH-MOEA outperformed other state-of-the-art methods in terms of solution quality**

# Evolutionary Reinforcement Learning

# What is reinforcement learning

learning a strategy to interact with the environment for maximizing the long-term reward



**Agent:** Policy:  $\pi : S \times A \rightarrow \mathbb{R}$ ,  $\sum_{a \in A} \pi(a|s) = 1$

Policy (deterministic):  $\pi : S \rightarrow A$

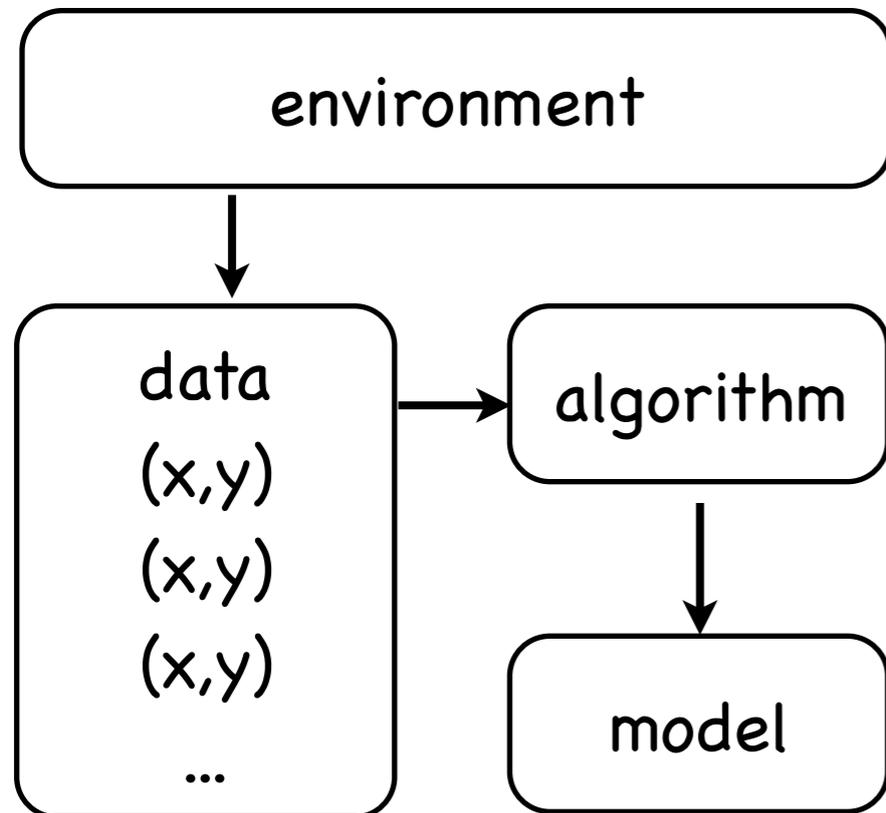
**Agent's goal:** learn a policy to maximize long-term total reward

T-step:  $\sum_{t=1}^T r_t$

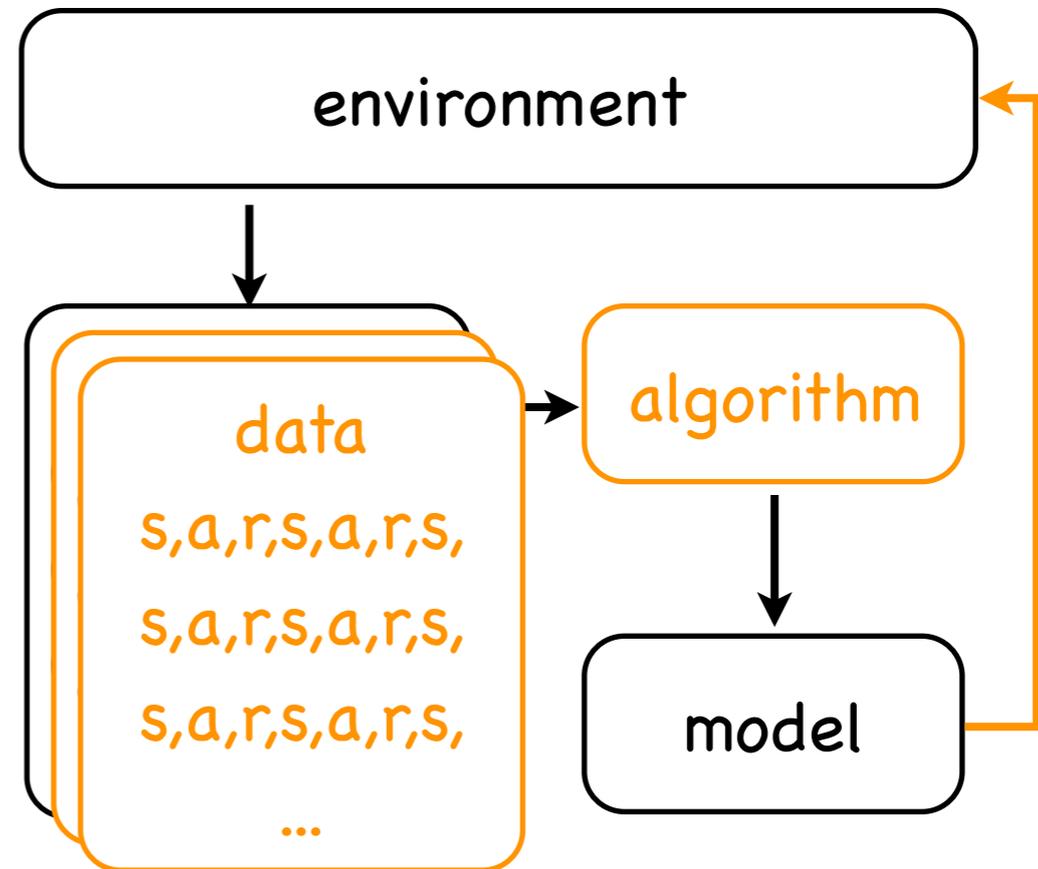
discounted:  $\sum_{t=1}^{\infty} \gamma^t r_t$

# Compare RL with SL

## supervised learning



## reinforcement learning

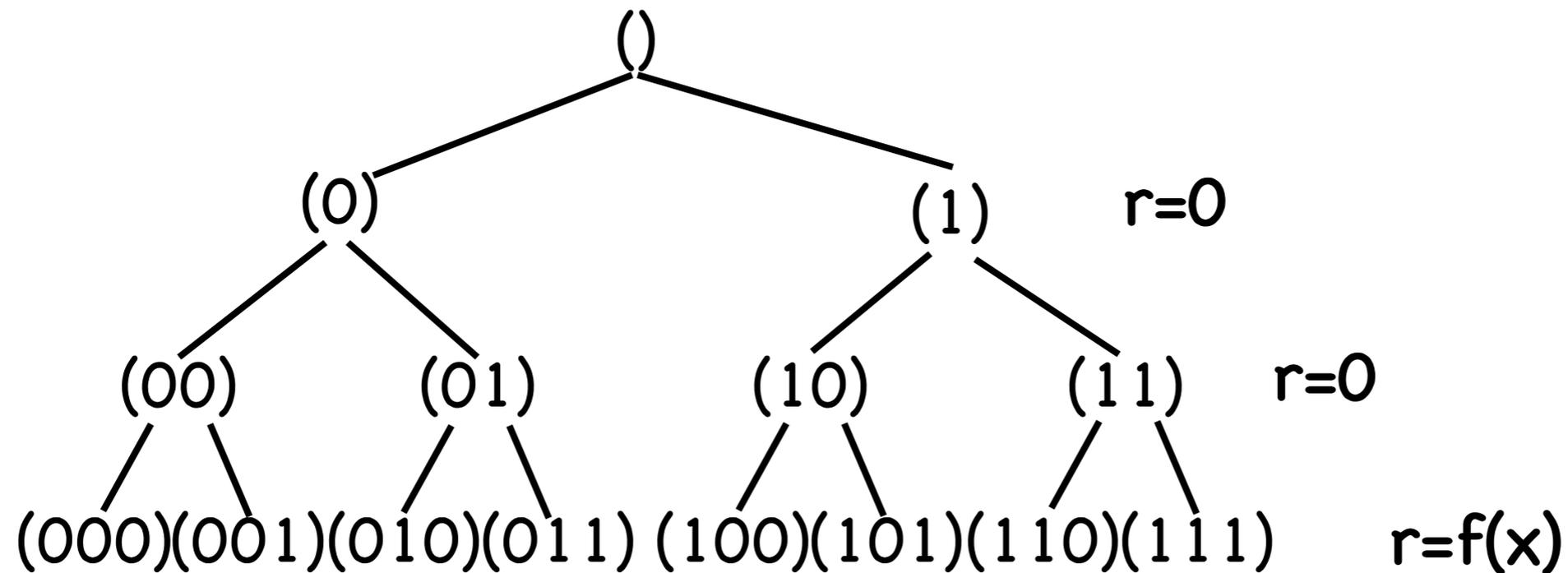


SL searches for a model

RL searches for the right output and a model

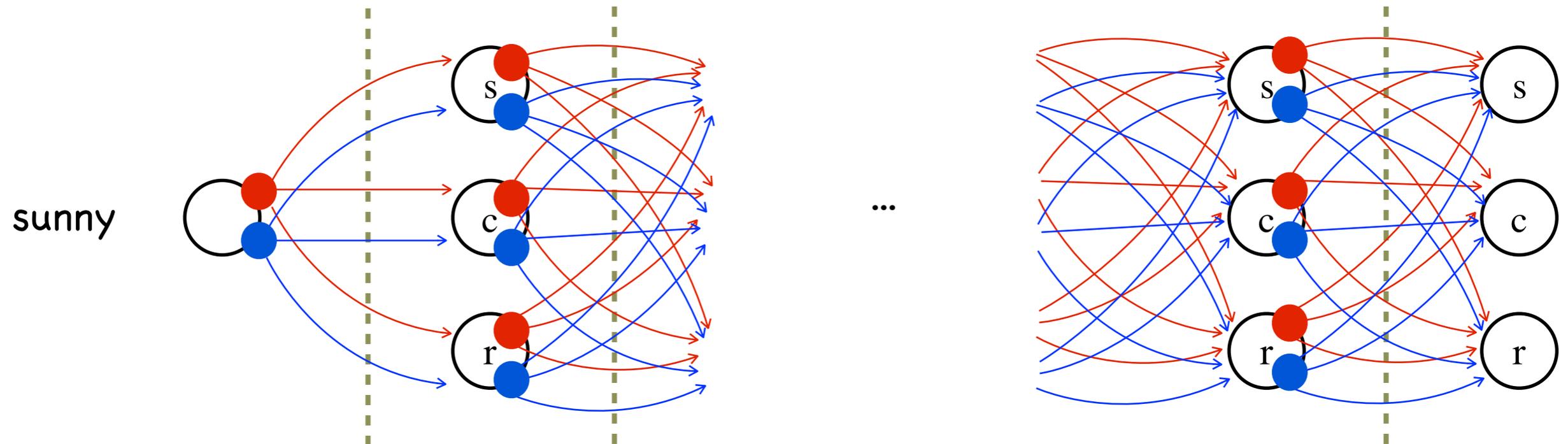
# Hardness of RL

general binary space problem  $\max_{x \in \{0,1\}^n} f(x)$



**solving the optimal policy is NP-hard!**

# Value-based methods



## dynamic programming

$$V^\pi(s) = \sum_a \pi(a|s) \sum_{s'} P(s'|s, a) (R(s, a, s') + V^\pi(s'))$$

$$Q^\pi(s, a) = \sum_{s'} P(s'|s, a) (R(s, a, s') + V^\pi(s'))$$

# Value-based methods

overall idea:

how is the current policy **policy evaluation**

improve the current policy **policy improvement**

policy iteration:

**policy evaluation:** **backward calculation**

$$V^\pi(s) = \sum_a \pi(a|s) \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V^\pi(s'))$$

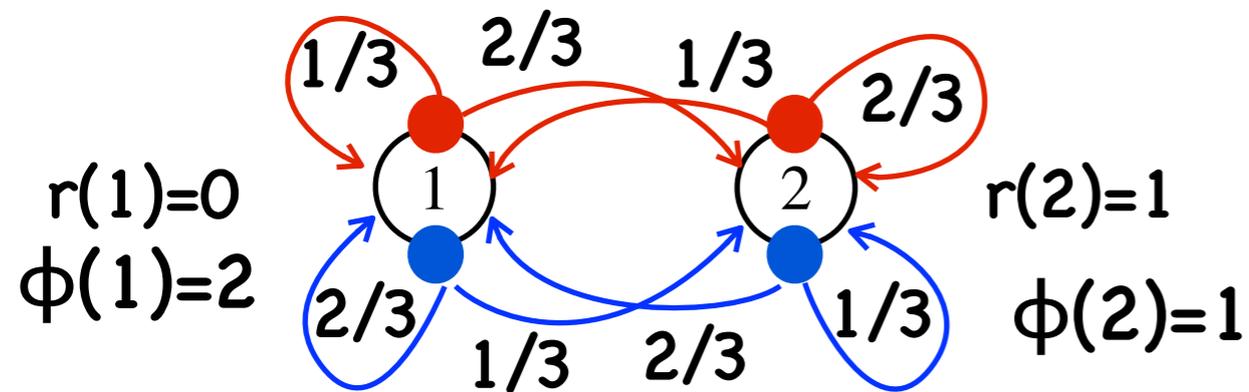
**policy improvement:** **from the Bellman optimality equation**

$$V(s) \leftarrow \max_a Q^\pi(s, a)$$

# Value-based methods

## policy degradation in value-based methods

[Bartlett. An Introduction to Reinforcement Learning Theory: Value Function Methods. Advanced Lectures on Machine Learning, LNAI 2600]



optimal policy: red  
 $V^*(2) > V^*(1) > 0$

let  $\hat{V}(s) = w\phi(s)$ , to ensure  $\hat{V}(2) > \hat{V}(1)$ ,  $w < 0$

as value function based method minimizes  $\|\hat{V} - V^*\|$   
results in  $w > 0$

sub-optimal policy, better value  $\neq$  better policy

# Policy search

## parameterized policy

$$\pi(a|s) = P(a|s, \theta)$$

Gibbs policy (logistic regression)

$$\pi_{\theta}(i|s) = \frac{\exp(\theta_i^{\top} \phi(s))}{\sum_j \exp(\theta_j^{\top} \phi(s))}$$

Gaussian policy (continuous !)

$$\pi_{\theta}(a|s) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(\theta^{\top} s - a)^2}{\sigma^2}\right)$$

# Policy search

direct objective functions

episodic environments: over all trajectories

$$J(\theta) = \int_{Tra} p_{\theta}(\tau) R(\tau) d\tau$$

where  $p_{\theta}(\tau) = p(s_0) \prod_{i=1}^T p(s_i | a_i, s_{i-1}) \pi_{\theta}(a_i | s_{i-1})$

is the probability of generating the trajectory

continuing environments: over stationary distribution

$$J(\theta) = \int_S d^{\pi_{\theta}}(s) \int_A \pi_{\theta}(a | s) R(s, a) ds da$$

$d^{\pi_{\theta}}$  is the stationary distribution of the process

# Policy search by gradient: policy gradient

$$J(\theta) = \int_{Tra} p_{\theta}(\tau) R(\tau) d\tau$$

**logarithm trick**  $\nabla_{\theta} p_{\theta} = p_{\theta} \nabla_{\theta} \log p_{\theta}$

**as**  $p_{\theta}(\tau) = p(s_0) \prod_{i=1}^T p(s_i | a_i, s_{i-1}) \pi_{\theta}(a_i | s_{i-1})$

$$\nabla_{\theta} \log p_{\theta}(\tau) = \sum_{i=1}^T \nabla_{\theta} \log \pi_{\theta}(a_i | s_{i-1}) + \text{const}$$

**gradient:**  $\nabla_{\theta} J(\theta) = \int_{Tra} p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau) R(\tau) d\tau$

$$= E\left[\sum_{i=1}^T \nabla_{\theta} \log \pi_{\theta}(a_i | s_i) R(s_i, a_i)\right]$$

use samples to estimate the gradient (unbiased estimation)

# Policy search v.s. value-based methods

## Policy search advantages:

- effective in high-dimensional and continuous action space
- learn stochastic policies directly
- avoid policy degradation

## disadvantages:

- converge only to a local optimum
- high variance

# Policy gradient: variance control

## actor-critic

$$\nabla_{\theta} J(\theta) \approx E[\nabla_{\theta} \log \pi_{\theta}(a|s) Q_w(s, a)]$$

if  $w$  is a minimizer of  $E[(Q^{\pi_{\theta}}(s, a) - Q_w(s, a))^2]$

Learn policy (actor) and Q-value (critic) simultaneously

## baseline

$$\nabla_{\theta} J(\theta) = E[\nabla_{\theta} \log \pi_{\theta}(a|s) (Q^{\pi}(s, a) - b(s))]$$

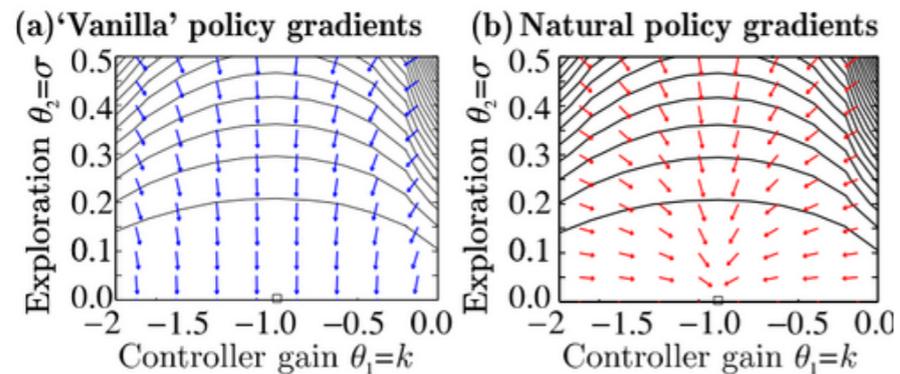
**advantage function:**  $A^{\pi}(s, a) = Q^{\pi}(s, a) - V^{\pi}(s)$

$$\nabla_{\theta} J(\theta) = E[\nabla_{\theta} \log \pi_{\theta}(a|s) A^{\pi}(s, a)]$$

learn policy, Q and V simultaneously

# Policy gradient: other gradients

## nature policy gradient



[Kakade. A Natural Policy Gradient. NIPS'01]

## functional policy gradient

$$\pi_{\Psi}(a|\mathbf{s}) = \frac{\exp(\Psi(\mathbf{s}, a))}{\sum_{a'} \exp(\Psi(\mathbf{s}, a'))}$$

$$\Psi_t = \sum_{i=1}^t h_t$$

[Yu et al. Boosting nonparametric policies. AAMAS'16]

## parameter-level exploration

$$\theta \sim \mathcal{N}$$

[Sehnke et al. Parameter-exploring policy gradients. Neural Networks'10]

## asynchronous gradient update

[Mnih et al. Asynchronous Methods for Deep Reinforcement Learning . ICML'16]

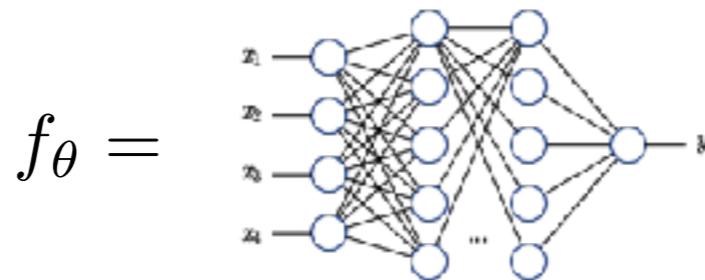
# Optimization difficulty

the non-convexity

$$J(\theta) = \int_{Tra} p_{\theta}(\tau) R(\tau) d\tau$$

where  $p_{\theta}(\tau) = p(s_0) \prod_{i=1}^T p(s_i | a_i, s_{i-1}) \pi_{\theta}(a_i | s_{i-1})$

$$\pi_{\theta}(i | s) = \frac{\exp(f_{\theta}(i; \phi(s)))}{\sum_j \exp(f_{\theta}(j; \phi(s)))}$$

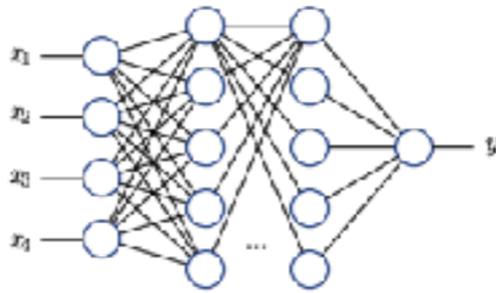


too many local minima

# EARL - EA for RL

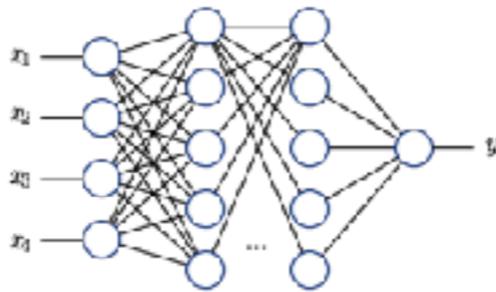
value-function representation

$$Q(s, a) =$$

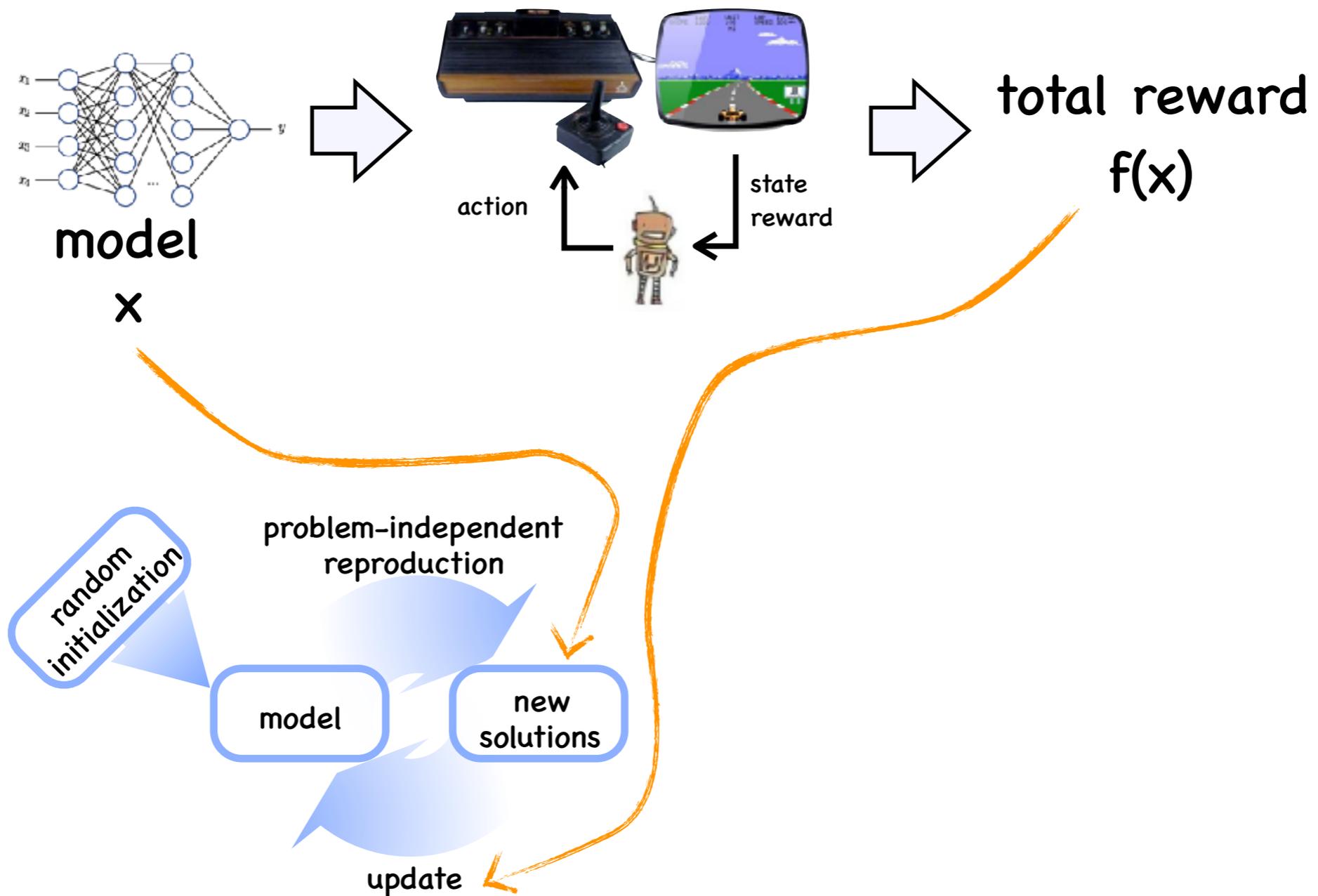


policy representation

$$\pi(s) =$$



# Parameter search by EAs



# Issue: low-efficiency

Does not utilize problem structure

1. NN parameters  $\rightarrow$  NN structure

2. Parameters  $\rightarrow$  dynamic programming

## NEAT+Q

generations

evaluation  
with TD

reproduction

---

```
1: // S: set of all states, A: set of all actions, c: output scale, p: population size
2: //  $m_n$ : node mutation rate,  $m_l$ : link mutation rate, g: number of generations
3: // e: number of episodes per generation,  $\alpha$ : learning rate,  $\gamma$ : discount factor
4: //  $\lambda$ : eligibility decay rate,  $\epsilon_{td}$ : exploration rate
5:
6:  $P[] \leftarrow$  INIT-POPULATION( $S, A, p$ ) // create new population P with random networks
7: for  $i \leftarrow 1$  to  $g$  do
8:   for  $j \leftarrow 1$  to  $e$  do
9:      $N, s, s' \leftarrow$  RANDOM( $P[]$ ), null, INIT-STATE( $S$ ) // select a network randomly
10:    repeat
11:       $Q[] \leftarrow c \times$  EVAL-NET( $N, s'$ ) // compute value estimates for current state
12:
13:      with-prob( $\epsilon_{td}$ )  $a' \leftarrow$  RANDOM( $A$ ) // select random exploratory action
14:      else  $a' \leftarrow$  argmax $_k Q[k]$  // or select greedy action
15:      if  $s \neq$  null then
16:        BACKPROP( $N, s, a, (r + \gamma \max_k Q[k])/c, \alpha, \gamma, \lambda$ ) // adjust weights toward target
17:
18:       $s, a \leftarrow s', a'$ 
19:       $r, s' \leftarrow$  TAKE-ACTION( $a'$ ) // take action and transition to new state
20:       $N.fitness \leftarrow N.fitness + r$  // update total reward accrued by N
21:    until TERMINAL-STATE?( $s$ )
22:     $N.episodes \leftarrow N.episodes + 1$  // update total number of episodes for N
23:     $P'[] \leftarrow$  new array of size  $p$  // new array will store next generation
24:    for  $j \leftarrow 1$  to  $p$  do
25:       $P'[j] \leftarrow$  BREED-NET( $P[]$ ) // make a new network based on fit parents in P
26:      with-probability  $m_n$ : ADD-NODE-MUTATION( $P'[j]$ ) // add a node to new network
27:      with-probability  $m_l$ : ADD-LINK-MUTATION( $P'[j]$ ) // add a link to new network
28:     $P[] \leftarrow P'[]$ 
```

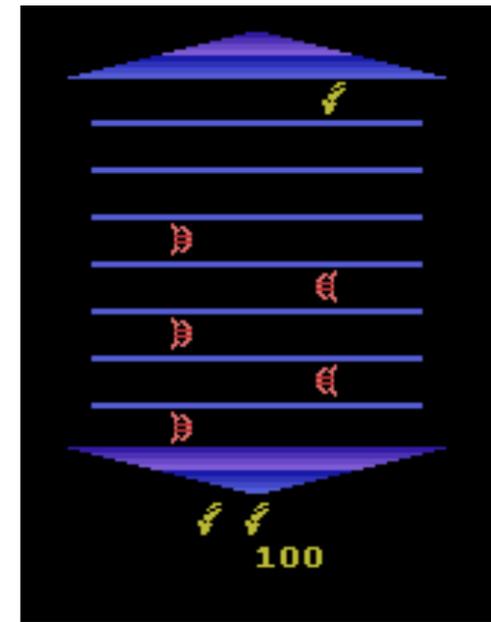
---

# Some comparison

on Atari games



freeway



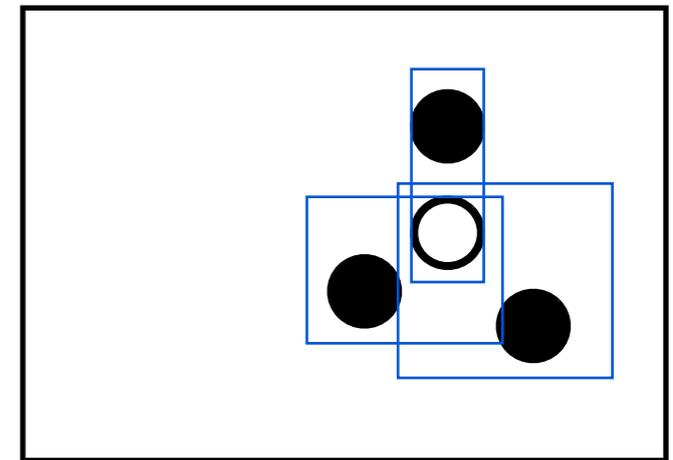
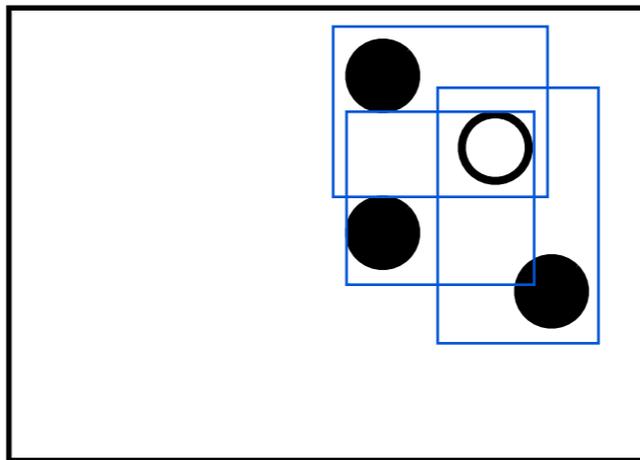
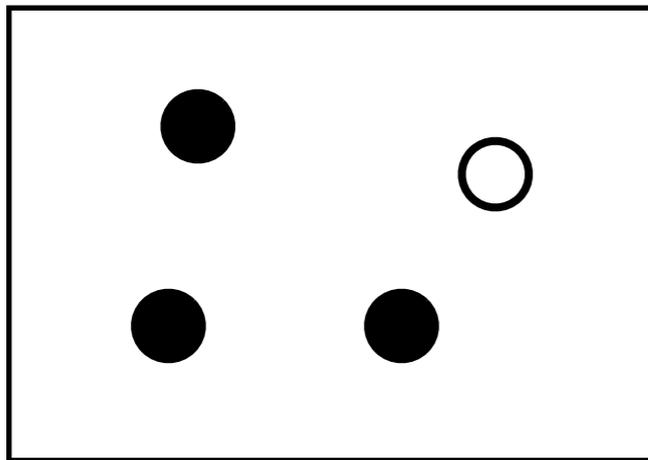
asterix

|                          | Freeway | Asterix |
|--------------------------|---------|---------|
| Sarsa( $\lambda$ )-BASS  | 0       | 402     |
| Sarsa( $\lambda$ )-DISCO | 0       | 301     |
| Sarsa( $\lambda$ )-RAM   | 0       | 545     |
| Random                   | 0       | 156     |
| HyperNEAT-GGP (Average)  | 27.4    | 870     |
| HyperNEAT-GGP (Best)     | 29      | 1000    |

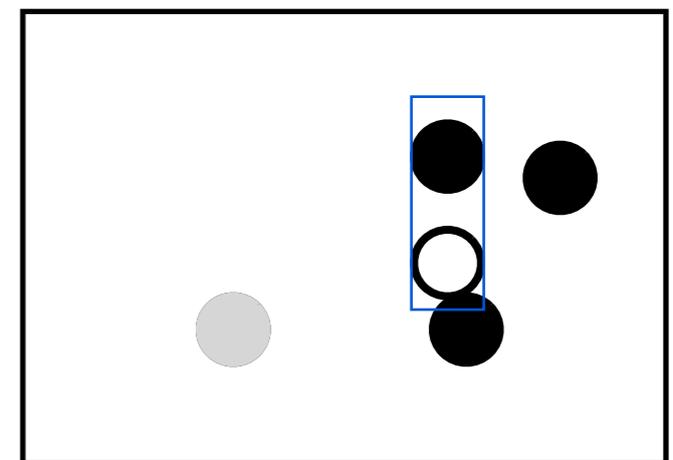
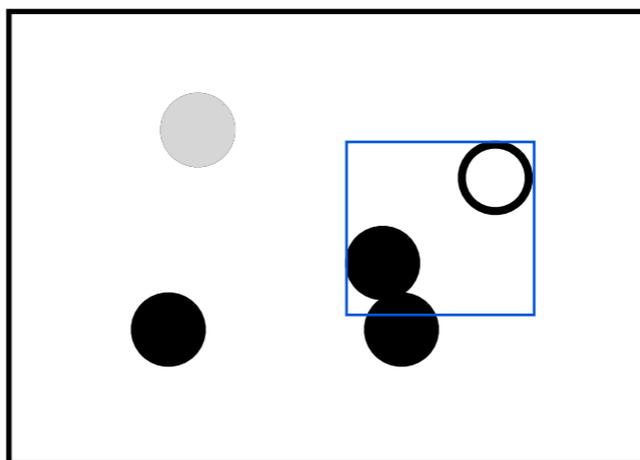
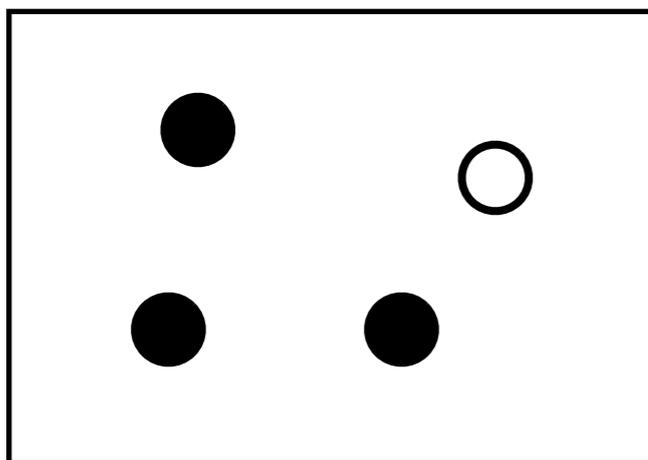
# Issue: low-efficiency

Batch sampling  
Online update

batch mode:



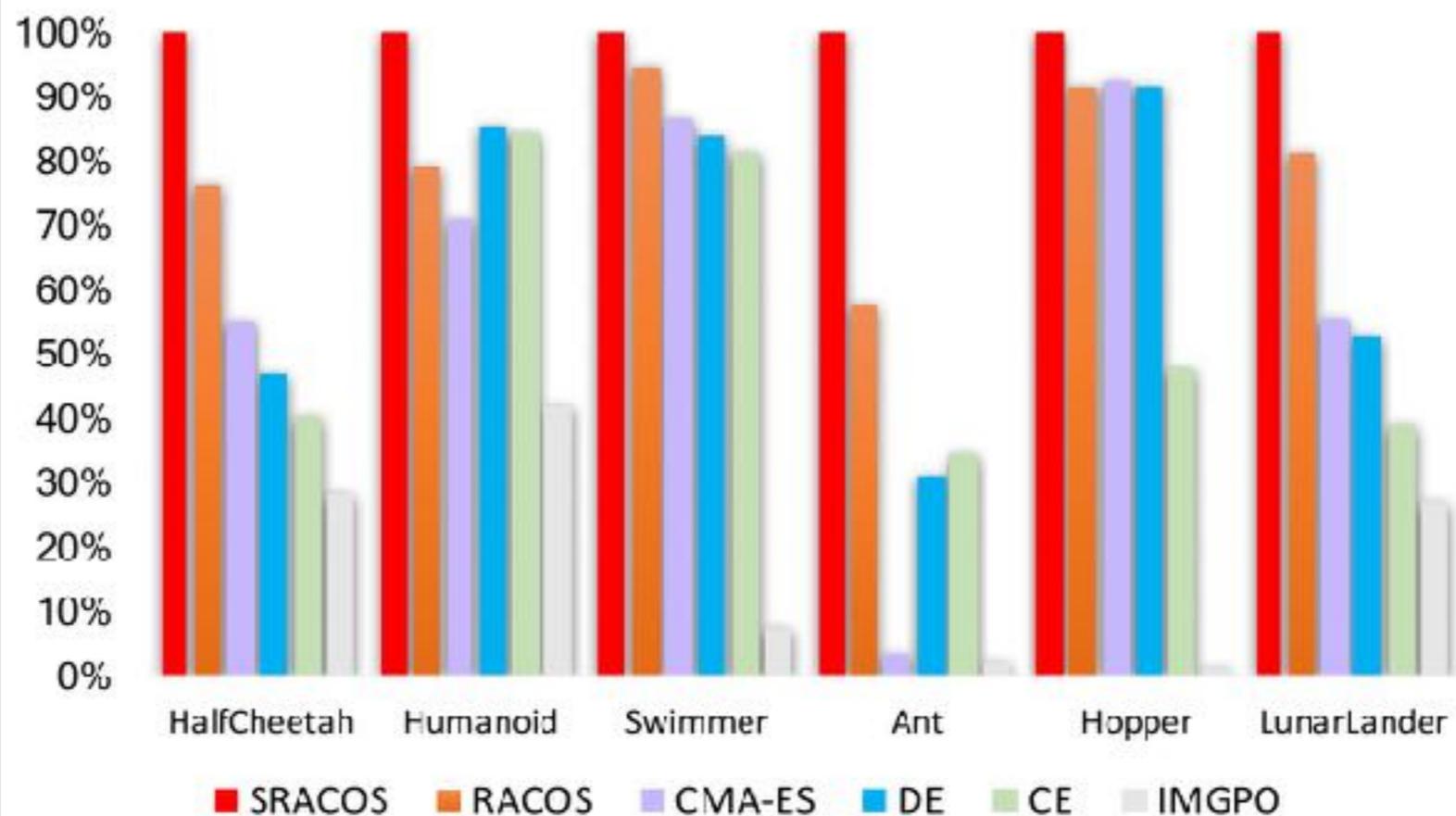
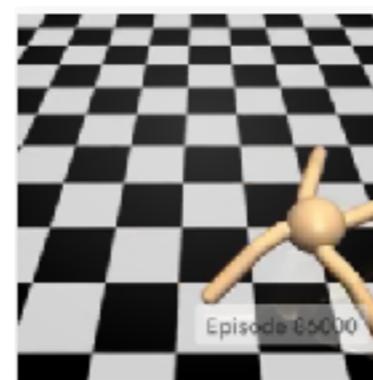
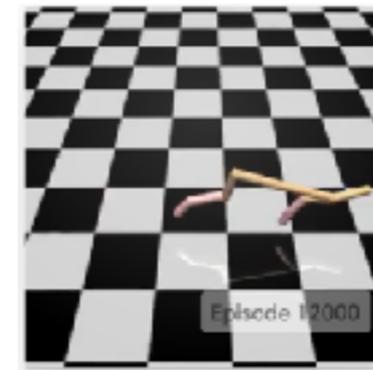
sequential mode:



# Issue: low-efficiency

Batch sampling      Online update

| Task name   | $d_{\text{State}}$ | #Actions | NN nodes | #Weights | Horizon |
|-------------|--------------------|----------|----------|----------|---------|
| Acrobot     | 6                  | 1        | 5, 3     | 48       | 2,000   |
| MountainCar | 2                  | 1        | 5        | 15       | 10,000  |
| HalfCheetah | 17                 | 6        | 10       | 230      | 10,000  |
| Humanoid    | 376                | 17       | 25       | 9825     | 50,000  |
| Swimmer     | 8                  | 2        | 5, 3     | 61       | 10,000  |
| Ant         | 111                | 8        | 15       | 1785     | 10,000  |
| Hopper      | 11                 | 3        | 9, 5     | 159      | 10,000  |
| LunarLander | 8                  | 1        | 5, 3     | 58       | 10,000  |



# Issue: low-efficiency

## Parallel: evolutionary strategies

centralized

- 
- 1: **Input:** Learning rate  $\alpha$ , noise standard deviation  $\sigma$ , initial policy parameters  $\theta_0$
  - 2: **for**  $t = 0, 1, 2, \dots$  **do**
  - 3:   Sample  $\epsilon_1, \dots, \epsilon_n \sim \mathcal{N}(0, I)$
  - 4:   Compute returns  $F_i = F(\theta_t + \sigma\epsilon_i)$  for  $i = 1, \dots, n$
  - 5:   Set  $\theta_{t+1} \leftarrow \theta_t + \alpha \frac{1}{n\sigma} \sum_{i=1}^n F_i \epsilon_i$
  - 6: **end for**
- 



parallel

- 
- 1: **Input:** Learning rate  $\alpha$ , noise standard deviation  $\sigma$ , initial policy parameters  $\theta_0$
  - 2: **Initialize:**  $n$  workers with known random seeds, and initial parameters  $\theta_0$
  - 3: **for**  $t = 0, 1, 2, \dots$  **do**
  - 4:   **for** each worker  $i = 1, \dots, n$  **do**
  - 5:     Sample  $\epsilon_i \sim \mathcal{N}(0, I)$
  - 6:     Compute returns  $F_i = F(\theta_t + \sigma\epsilon_i)$
  - 7:   **end for**
  - 8:   Send all scalar returns  $F_i$  from each worker to every other worker
  - 9:   **for** each worker  $i = 1, \dots, n$  **do**
  - 10:     Reconstruct all perturbations  $\epsilon_j$  for  $j = 1, \dots, n$
  - 11:     Set  $\theta_{t+1} \leftarrow \theta_t + \alpha \frac{1}{n\sigma} \sum_{j=1}^n F_j \epsilon_j$
  - 12:   **end for**
  - 13: **end for**
- 

time steps of ES to the top  
performance of TRPO

---

| ENVIRONMENT            | 25%  | 50%  | 75%  | 100% |
|------------------------|------|------|------|------|
| HALFCHEETAH            | 0.15 | 0.49 | 0.42 | 0.58 |
| HOPPER                 | 0.53 | 3.64 | 6.05 | 6.94 |
| INVERTEDDOUBLEPENDULUM | 0.46 | 0.48 | 0.49 | 1.23 |
| INVERTEDPENDULUM       | 0.28 | 0.52 | 0.78 | 0.88 |
| SWIMMER                | 0.56 | 0.47 | 0.53 | 0.30 |
| WALKER2D               | 0.41 | 5.69 | 8.02 | 7.88 |

---

# More issues to be solved

Noisy evaluation: too many repetitions

Large model: too large search space

Long-term reward: too complex objective function

---

Thanks for your time!

Questions/comments?