

ENFrame at a glance

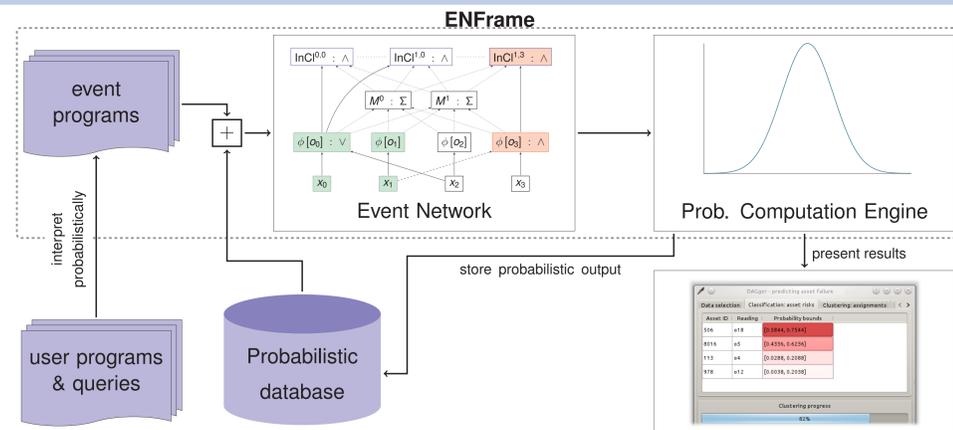
- Users write code in a subset of Python + relational queries.
- User code is oblivious to the probabilistic nature of the data: ENFrame interprets the code and runs it on probabilistic data.
- We already tested ENFrame on several algorithms, e.g., k -medoids clustering and k -nearest neighbour classification applied to probabilistic data representing query results.

Key Technical Features of ENFrame

- Language to express probabilistic events that capture arbitrary correlations in the input data and in the output as induced by program traces.
- Sequential and parallel algorithms for exact and approximate probabilistic computation of user programs.
- Quality metric for probabilistic computation using, e.g., ENFrame, naïve and sampling-based methods.

Event Language, Event Programs, and Event Network

- **Event program:**
 - ENFrame's probabilistic interpretation of the user program is captured by *events*.
- **Event language:**
 - Event = random variables, conditioned values, Boolean formulas over events, arithmetic operations over events.
 - Variables of type T from user program become discrete random variables with pdfs over values in T .
 - **c-values:** values (numbers/vectors) conditioned on events.
 - For a number v and variable Φ :
c-value $\Phi \otimes v$ evaluates to v if Φ is true, or 0 otherwise.
 - c-values can be summed and compared:
 $\Phi_1 \otimes v_1 + \dots + \Phi_n \otimes v_n \leq \Psi_1 \otimes w_1 + \dots + \Psi_m \otimes w_m$
 - Application example: sum of distances between objects.
- **Event Network:** Joint representation of interconnected events.

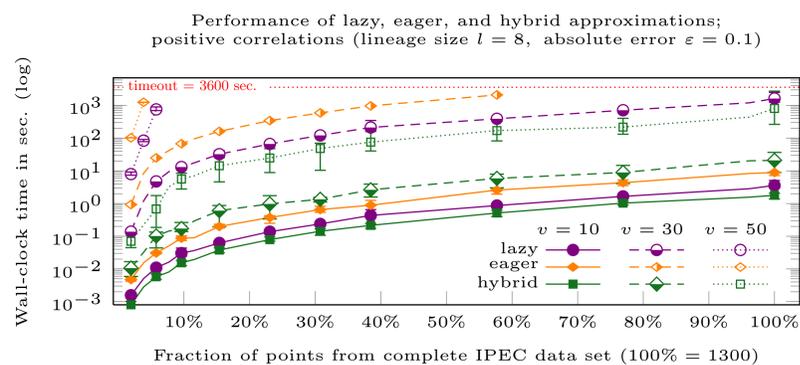


Exact and Approximate Probability Computation

- ENFrame can compute, e.g., the probability that two objects are in the same cluster (k -medoids), or that an object is assigned to a class (k -NN).
- Main idea: exhaustive/partial structural decomposition of an entire event network.
 - Outcome equivalent to clustering in every possible world.
 - Multiple pruning strategies for approximate probability computation with error guarantees.
 - Efficient parallel computation on multi-cores.

Experimental Evaluation

Comparison of approximation algorithms for k -medoids ($k = 3$), input data with positive correlations and varying number of variables (v). Further experiments in [vSOF14].



Example: User Program and Event Program for k -medoids Clustering

```
#Initialisation phase: Select k cluster medoids (centres)
1: (O, n) = loadData()      # list and number of objects
2: (k, iter) = loadParams() # number of clusters and iterations
3: M = init()              # initialise medoids

4: for it in range(0,iter): # clustering iterations
    #Assignment phase: assign objects to closest medoid
5:   InCl = [None] * k
6:   for i in range(0,k):
7:     InCl[i] = [None] * n
8:     for l in range(0,n):
9:       InCl[i][l] = reduce_and(
10:        [(dist(O[l],M[i]) <= dist(O[l],M[j])) for j in range(0,k)])
11:   InCl = breakTies2(InCl) # each object is in exactly one cluster

    #Update phase: Select new cluster medoids
12:   DistSum = [None] * k
13:   for i in range(0,k):
14:     DistSum[i] = [None] * n
15:     for l in range(0,n):
16:       DistSum[i][l] = reduce_sum(
17:        [dist(O[l],O[p]) for p in range(0,n) if InCl[i][p]])

18:   Centre = [None] * k
19:   for i in range(0,k):
20:     Centre[i] = [None] * n
21:     for l in range(0,n):
22:       Centre[i][l] = reduce_and(
23:        [DistSum[i][l] <= DistSum[i][p] for p in range(0,n)])
24:   Centre = breakTies1(Centre) # enforce one Centre per cluster

25: M = [None] * k
26: for i in range(0,k):
27:   M[i] = reduce_sum([O[l] for l in range(0,n) if Centre[i][l]])
```

$$\forall i \text{ in } 0..n-1 : O^i \equiv \Phi(o_i) \otimes \vec{o}_i$$

$$M_{-1}^0 \equiv \Phi(o_{\pi(0)}) \otimes \vec{o}_{\pi(0)}; \dots; M_{-1}^{k-1} \equiv \Phi(o_{\pi(k-1)}) \otimes \vec{o}_{\pi(k-1)}$$

$\forall \text{it in } 0..\text{iter} - 1 :$

$$\forall i \text{ in } 0..k-1 :$$

$$\forall l \text{ in } 0..n-1 :$$

$$\text{InCl}_{\text{it}}^{i,l} \equiv \bigwedge_{j=0}^{k-1} [\text{dist}(O^l, M_{\text{it}-1}^j) \leq \text{dist}(O^l, M_{\text{it}-1}^i)]$$

Encoding of breakTies2 omitted

$$\forall i \text{ in } 0..k-1 :$$

$$\forall l \text{ in } 0..n-1 :$$

$$\text{DistSum}_{\text{it}}^{i,l} \equiv \sum_{p=0}^{n-1} \text{InCl}_{\text{it}}^{i,p} \otimes \text{dist}(O^l, O^p)$$

$$\forall i \text{ in } 0..k-1 :$$

$$\forall l \text{ in } 0..n-1 :$$

$$\text{Centre}_{\text{it}}^{i,l} \equiv \bigwedge_{p=0}^{n-1} [\text{DistSum}_{\text{it}}^{i,l} \leq \text{DistSum}_{\text{it}}^{i,p}]$$

Encoding of breakTies2 omitted

$$\forall i \text{ in } 0..k-1 :$$

$$M_{\text{it}}^i = \sum_{l=0}^{n-1} \text{Centre}_{\text{it}}^{i,l} \wedge O^l$$

Challenges Currently under Microscope

- Which additional event language constructs are needed to capture further data analysis tasks?
- Trade-off: functionality (event-based result explanation, sensitivity analysis) vs. performance (coarse events compiled to C++ code)?