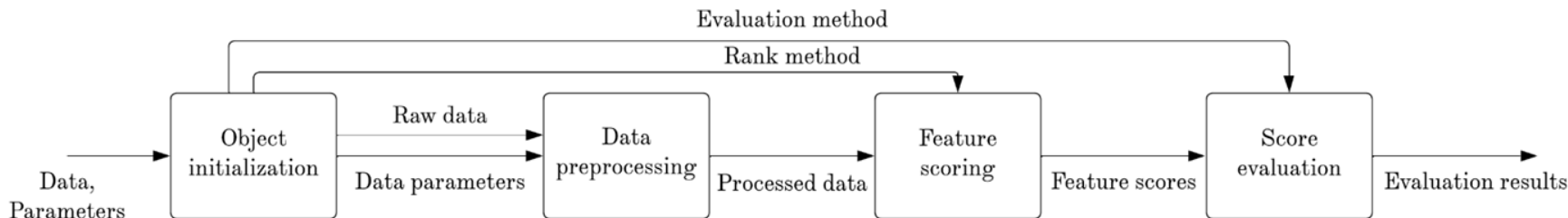# Feature selection
# on large and sparse datasets

Leon Hvastja, Amadej Pavšič

**Advisors:** Jure Demšar, Blaž Mramor, Blaž Škrlj
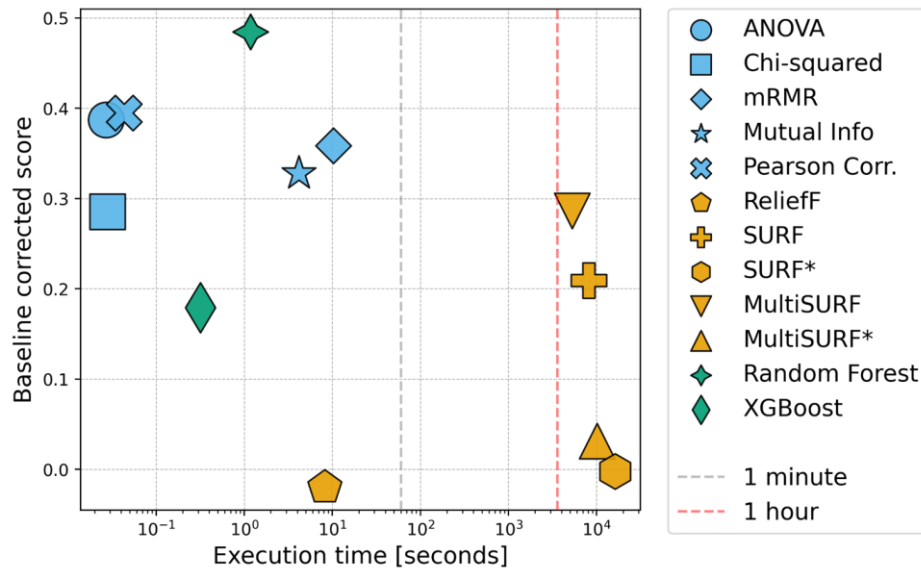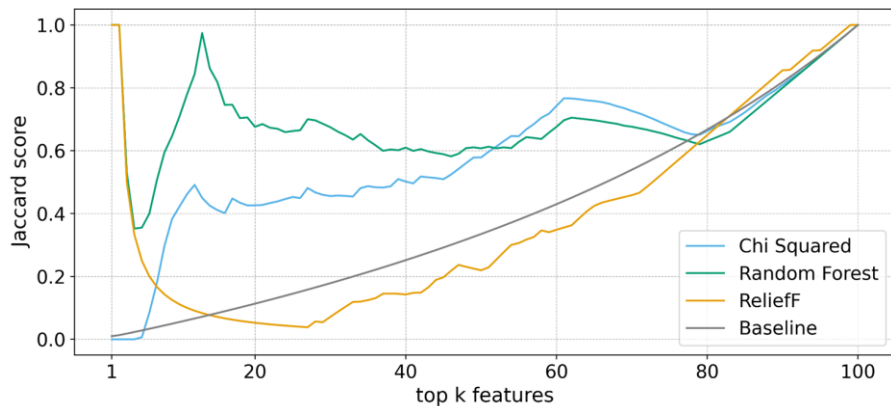
# Problem & dataset

- Outbrain, web recommendation platform
- Sparsity, high cardinality, uninformative features
- Efficient feature selection approach
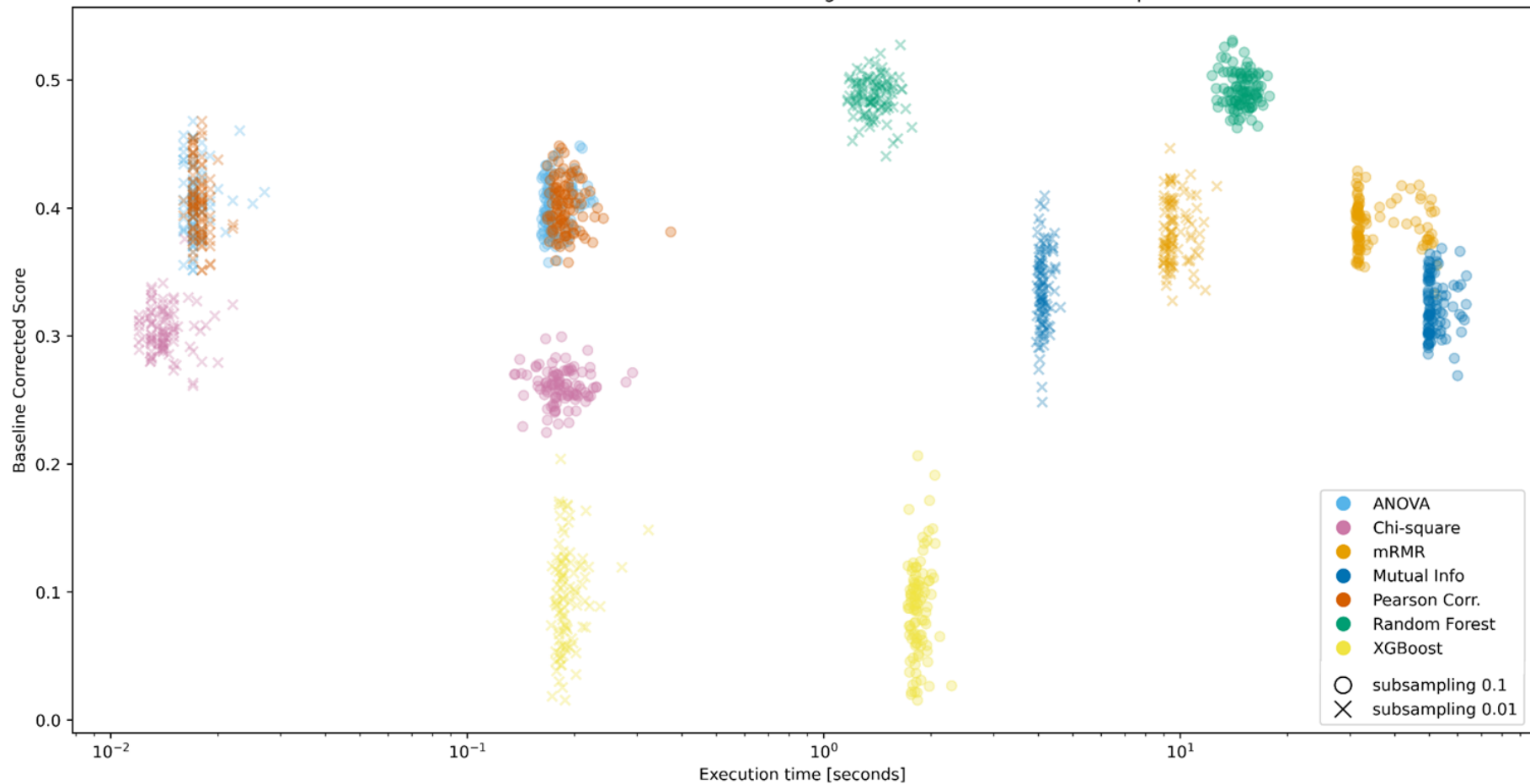- Ranking & evaluation pipeline

# Algorithms

- Filter and wrapper methods
- Using statistical properties
- Relief based approaches
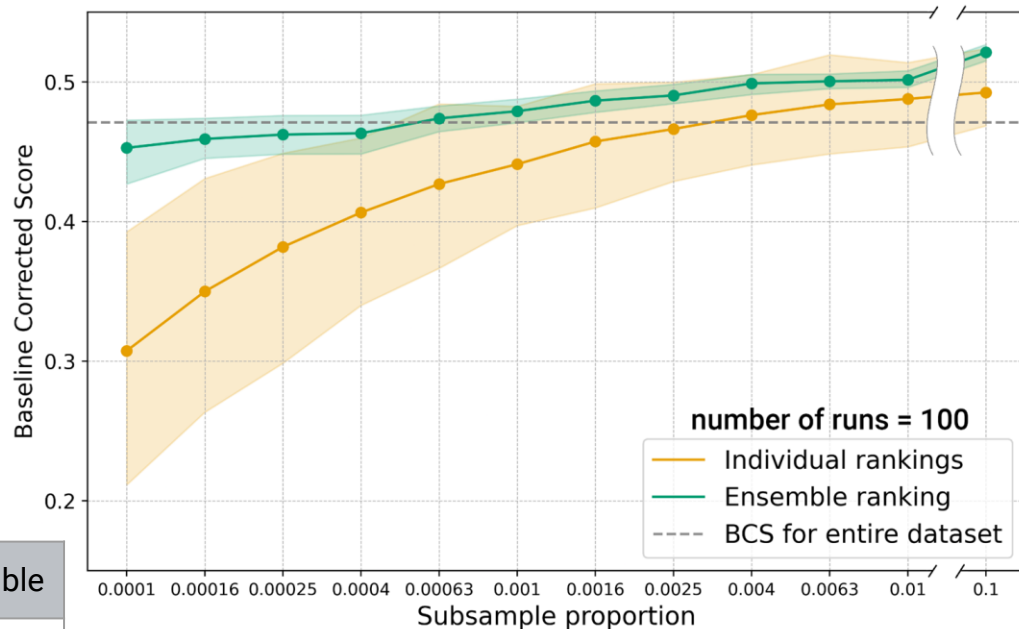- Baseline corrected score
- First eliminations

Performance over 100 runs for several algorithms with 1% and 10% samples.

# RF performance & ensemble

- Subsampling optimization
- Ensemble approach
- Borda count voting system
- Performance & uncertainty



| Feature | Run 1 | Run 2 | Run 3 | Borda | Ensemble |
|---------|-------|-------|-------|-------|----------|
| A | 1 | 2 | 1 | 4 | 1 |
| B | 2 | 1 | 3 | 6 | 2 |
| C | 3 | 3 | 2 | 8 | 3 |

# Conclusions

- Better scores
- Faster execution
- Less data
- Smaller uncertainty