**Exploratory data analysis**

Tell me (analyst/user) something interesting about the data

|   | Louise | Emma | Marie | Olivia | Nora | *Points scored?* |
|---|--------|------|-------|--------|------|------------------|
| 1 | *Plays* | P | P | *Rests* | P | $+$ |
| 2 | P | P |   | P | P | $-$ |
| 3 | P | P | P |   |   | $+$ |
| 4 | P |   |   | P |   | $-$ |
| 5 |   | P |   | P | P | $-$ |

**Discovering descriptions of coherent *interesting* data regions**

| Pattern $p$ | $\lvert p \rvert$ | $\lvert p^+ \rvert$ | $\lvert p^- \rvert$ | |
|---|---|---|---|---|
| {Louise, Emma} | 3 | 2 | 1 | Frequent |
| {Marie} | 2 | 2 | 0 | Discriminative (only $+$) |
| {Louise, Emma, Marie} | 2 | 2 | 0 | Discriminative & long |

**Problem statement**

| Pattern $p$ | $\|p\|$ | $\|p^+\|$ | $\|p^-\|$ | |
|---|---|---|---|---|
| {Louise, Emma} | 3 | 2 | 1 | Frequent |
| {Marie} | 2 | 2 | 0 | Discriminative (only $+$) |
| {Louise, Emma, Marie} | 2 | 2 | 0 | Discriminative & long |

**Given** dataset $\mathcal{D}$, constraints $\mathcal{C}$, and/or quality measure $\varphi$
**Mine** top patterns according to $\varphi$

**Problem statement**

| Pattern $p$ | $\lvert p \rvert$ | $\lvert p^+ \rvert$ | $\lvert p^- \rvert$ | |
|---|---|---|---|---|
| $\{\text{Louise}, \text{Emma}\}$ | 3 | 2 | 1 | Frequent |
| $\{\text{Marie}\}$ | 2 | 2 | 0 | Discriminative (only $+$) |
| $\{\text{Louise}, \text{Emma}, \text{Marie}\}$ | 2 | 2 | 0 | Discriminative & long |

**Given** dataset $\mathcal{D}$, constraints $\mathcal{C}$, and/or quality measure $\varphi$
**Mine** top patterns according to $\varphi$

**Problem statement**

| Pattern $p$ | $|p|$ | $|p^+|$ | $|p^-|$ | |
|---|---|---|---|---|
| {Louise, Emma} | 3 | 2 | 1 | Frequent |
| {Marie} | 2 | 2 | 0 | Discriminative (only $+$) |
| {Louise, Emma, Marie} | 2 | 2 | 0 | Discriminative & long |

**Given** dataset $\mathcal{D}$, constraints $\mathcal{C}$, and/or quality measure $\varphi$
**Mine** top patterns according to $\varphi$

*Subjectivity* in exploratory data analysis

| Pattern $p$ | $|p|$ | $|p^+|$ | $|p^-|$ | |
|---|---|---|---|---|
| {Louise, Emma} | 3 | 2 | 1 | Opponent |
| {Marie} | 2 | 2 | 0 | Coach |
| {Louise, Emma, Marie} | 2 | 2 | 0 | Journalist |

1. What patterns are interesting depends on the user
2. Non-experts can't tune constraints $\mathcal{C}$ or quality measure $\varphi$

# Direct user involvement is essential

**The user ranks patterns by their *subjective* interestingness**

{Louise, Emma}, {Louise, Emma, Marie}, {Marie}

$$\Downarrow$$

**1** {Louise, Emma, Marie}

**2** {Marie}

**3** {Louise, Emma}

Dzyuba, van Leeuwen, Nijssen, De Raedt. (2014) *Interactive learning of pattern rankings*, IJAIT

# Learn: Pattern ranking/scoring function

**Instance of *preference learning***

$$\varphi_{\text{user}}(p) = A + \frac{1-A}{1 + e^{-\vec{w} \cdot \vec{p}}}$$

$\vec{p}$ Pattern features (e.g., items, frequency, length…)

$\vec{w}$ Which features make a pattern interesting to the user
Learned from ordered feedback with *stochastic coordinate descent*

$A$ Technical parameter, see the paper

| | |
|---|---|
| **Interesting** | according to the current (approximate) $\varphi_{\text{user}}$ |
| **Compact** | so that feedback is easy to provide |
| **Diverse** | to ensure *exploration* necessary for learning |
| **Quick to obtain** | e.g., in an *anytime* manner |

**Interesting**        according to the current (approximate) $\varphi_{\text{user}}$

**Compact**            so that feedback is easy to provide

**Diverse**            to ensure *exploration* necessary for learning

**Quick to obtain**    e.g., in an *anytime* manner

**Given** dataset $\mathcal{D}$, constraints $\mathcal{C}$, and quality measure $\varphi$
**Sample** patterns proportional to $\varphi_{\text{user}}$

**Efficient black-box pattern sampler with strong performance guarantees**

To sample one pattern (at least):

1. Generate an *implicit* random partitioning of all patterns

2. Enumerate all patterns in a random partition

3. Generate a perfect sample from this partition

Dzyuba, van Leeuwen, De Raedt. (2017) *Flexible constrained sampling with guarantees for pattern mining*, DMKD

PCA representation: similar patterns are close to each other

# Illustration: random partition ("cell")

Patterns are different from each other ⇒ Good *exploration*



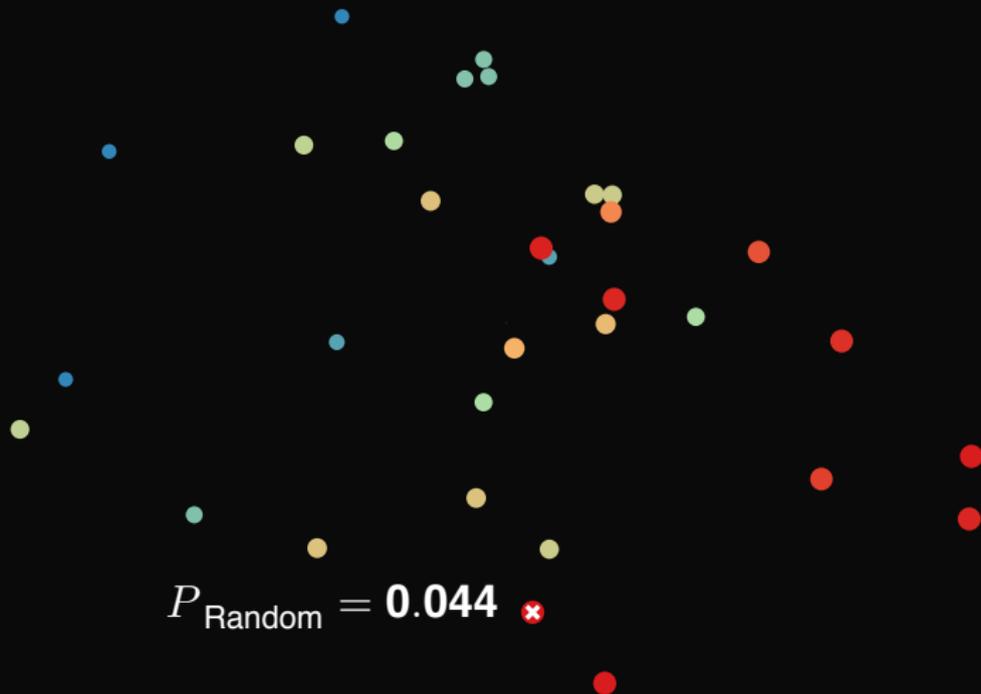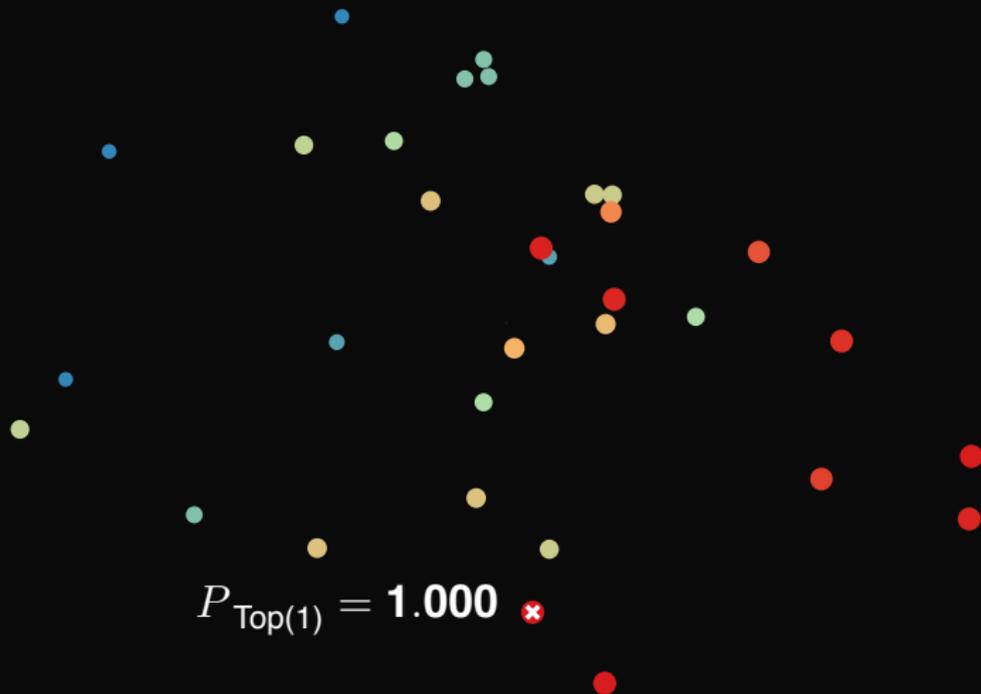$$P_{\mathsf{Random}} = \textbf{0.044}$$

Take top patterns from the partition $\Rightarrow$ Emphasize *exploitation*



$$P_{\mathsf{Top(1)}} = \mathbf{1.000}$$

**More contributions & details in the paper**

| | |
|---|---|
| **Mine** | Sample with *Flexics* |
| | <span style="color:red">Modified "cell" sampling</span> |
| **Interact** | Ordered feedback |
| **Learn** | <span style="color:red">Learning $\varphi_{\text{user}}$ with SCD</span> |
| **Repeat** | Sample with updated $\varphi_{\text{user}}$ |

# Experimental evaluation

Frequency, surprisingness, discriminativity with $\chi^2$

|   |   | $\varphi_{\text{obj}}$ |
|---|---|---|
| 1 | {Louise, Emma, Marie} | $0.5$ |
| 2 | {Marie} | $0.3$ |
| 3 | {Louise, Emma} | $0.1$ |

**Values of $\varphi_{obj}$ are not known to the learner, only used to measure its performance**

|   |   | $\varphi_{obj}$ |
|---|---|---|
| 1 | {Louise, Emma, Marie} | 0.5 |
| 2 | {Marie} | 0.3 |
| 3 | {Louise, Emma} | 0.1 |

Regret w.r.t. $\max.\varphi = 1 - 0.5 \qquad\qquad\qquad = 0.5$

$\qquad\qquad\quad \text{avg}.\varphi = 1 - (0.5 + 0.3 + 0.1)/3 \quad = 0.3$

- ▶ Learn to sample from $\varphi_{obj}$ only from small orders

- ▶ 10 datasets; choose *min.frequency* so that there are 140 000+ frequent patterns

- ▶ 30 learning iterations, 5 patterns per iteration
  Regret $\in [0,\ 30]$, lower is better

**Modified cell sampling improves the performance**

|  |  | Regret w.r.t. | | |
|---|---|---|---|---|
|  |  | Avg.qual | Max.qual | Diversity |
| Cell | *Random* | $10.6 \pm 0.7$ | $1.9 \pm 0.6$ | **12.2 $\pm$ 0.6** |
| sampling | *Top(1)* | **5.1 $\pm$ 1.1** | **0.8 $\pm$ 0.5** | $13.7 \pm 1.0$ |

|  |  | Regret w.r.t. | | |
|---|---|---|---|---|
|  |  | Avg.qual | Max.qual | Diversity |
| Algorithm | LetSIP | **$3.1 \pm 0.8$** | **$0.3 \pm 0.2$** | **$13.1 \pm 0.8$** |
|  | APLe | $3.2 \pm 2.6$ | $2.6 \pm 2.4$ | – |
|  | *IPM* | $12.9 \pm 2.4$ | $5.1 \pm 2.3$ | $16.0 \pm 1.9$ |

# Mine, Interact, Learn, Repeat

▶ Pattern sampling with *Flexics* delivers compact diverse sets of high-quality patterns in an anytime fashion…

▶ …which helps balance exploration and exploitation in interactive mining with *LetSIP*

Thank you for your attention!

May I answer any questions?

0.034

0.044

0.039

0.034

0.038

0.036

0.041 ✕

0.032 ⊗

0.036

0.033