# Integrating inconsistent data in a probabilistic model
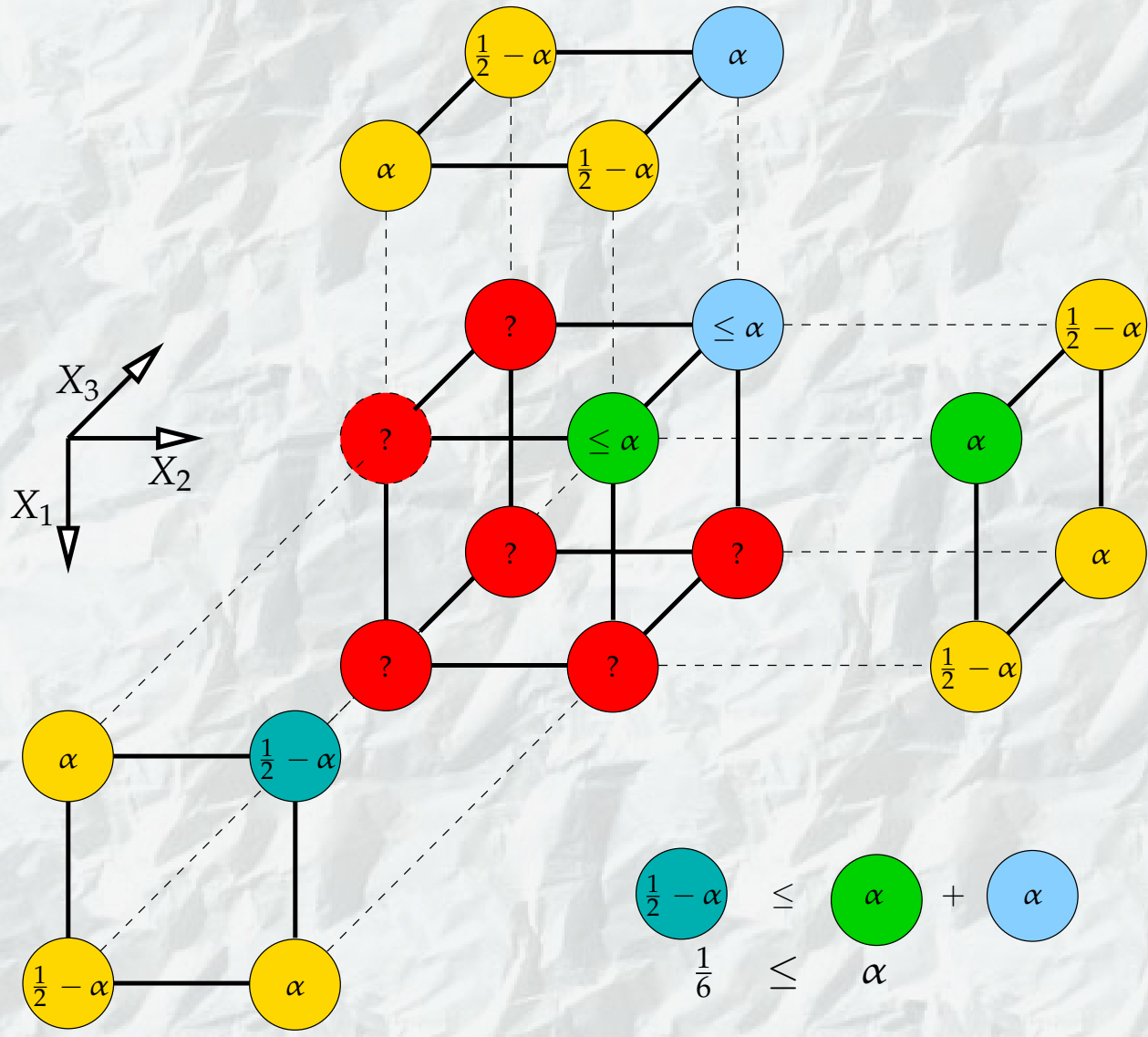
**Jiří Vomlel**

This presentation is available at
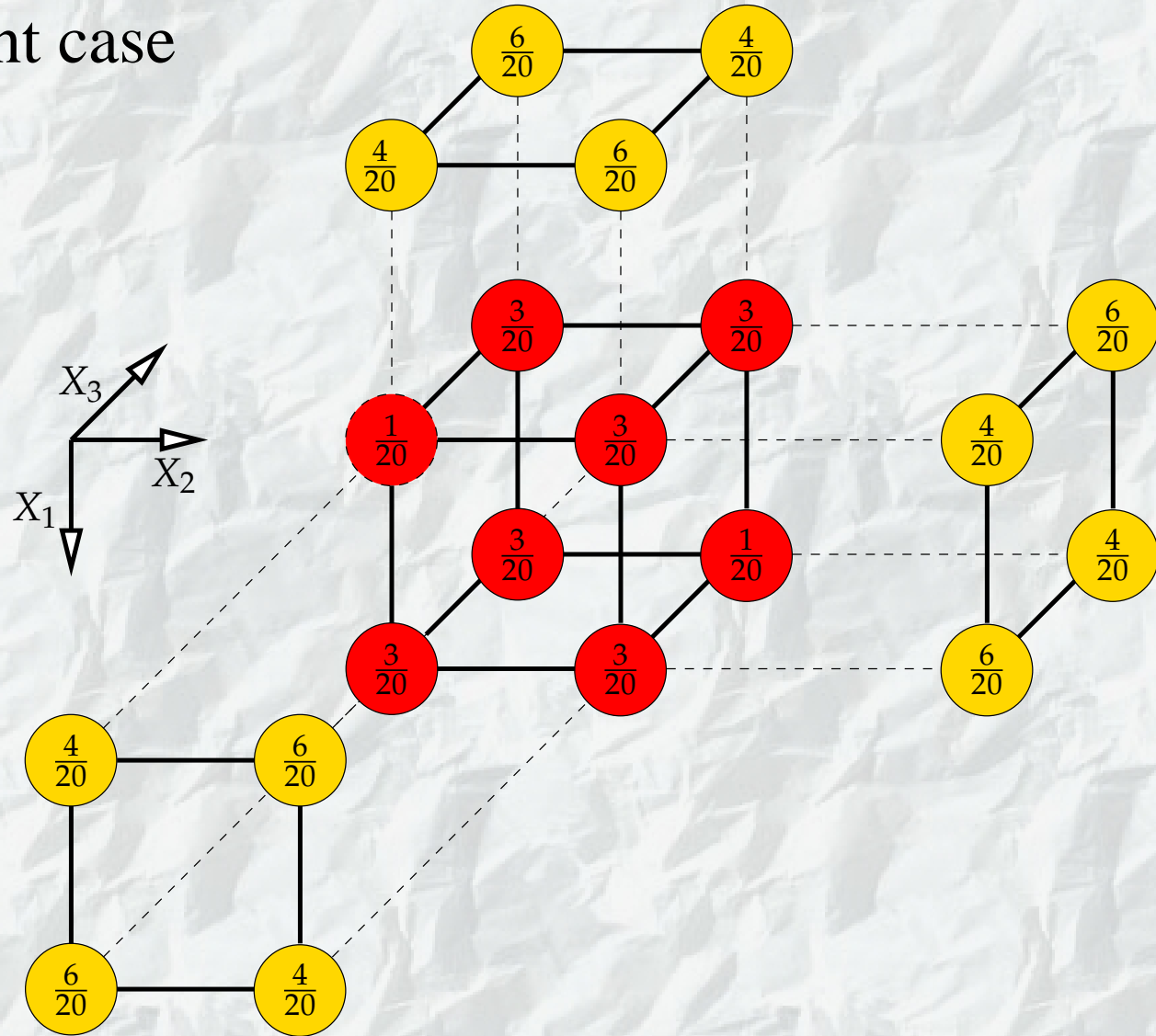http://www.utia.cas.cz/vomlel/

# Knowledge integration

- Discrete random variables $X_i$ indexed by natural numbers from $V = \{1, \ldots, n\} \subset \mathbb{N}$.

- Low-dimensional probability distributions $P_j, j = 1, \ldots, k$ defined on variables $\{X_\ell\}_{\ell \in E_j}, E_j \subseteq V$.

- *Knowledge integration* is the process of building a joint probability distribution $Q(X_1, \ldots, X_n)$ from a set of low-dimensional probability distributions $\mathcal{P} = \{P_1, \ldots, P_k\}$.

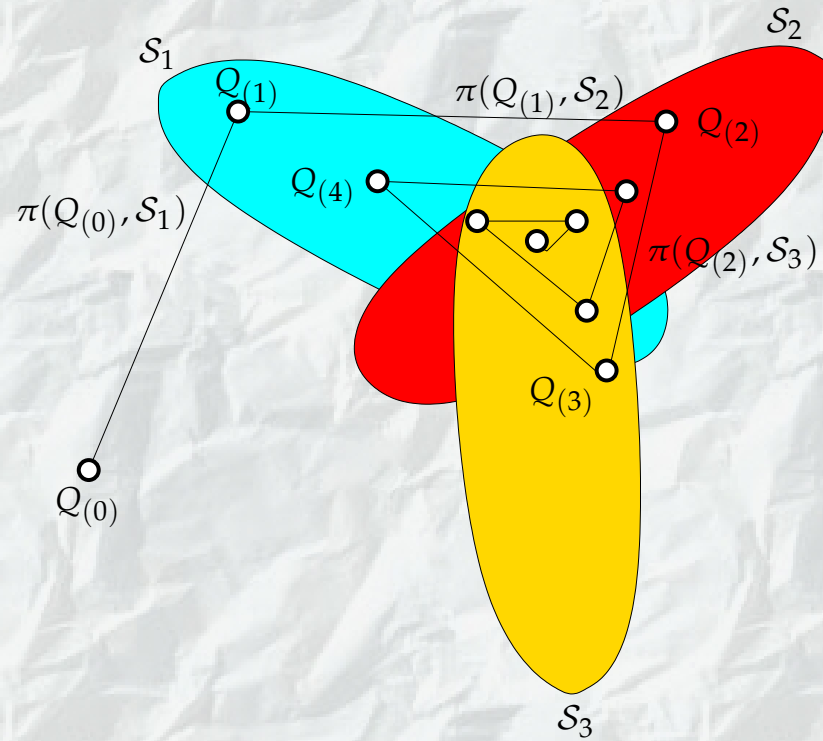Consistent case

$$\alpha = \frac{4}{20}$$

- input set $\mathcal{P} = \{P_1, \ldots, P_k\}$

- set of all distributions having $P_j$ as its marginal
$$\mathcal{S}_j = \{Q : Q^{E_j} = P_j\}$$

- set of all distributions having $\{P_1, \ldots, P_k\}$ as its marginals
$$\mathcal{S} = \cap_{j=1}^{k} \mathcal{S}_j$$

- *I*-projection of $Q_0$ to $\mathcal{S}$
$$\pi(Q_0, \mathcal{S}) = \arg \min_{Q \in \mathcal{S}} I(Q \parallel Q_0)$$

- Kullback-Leibler divergence

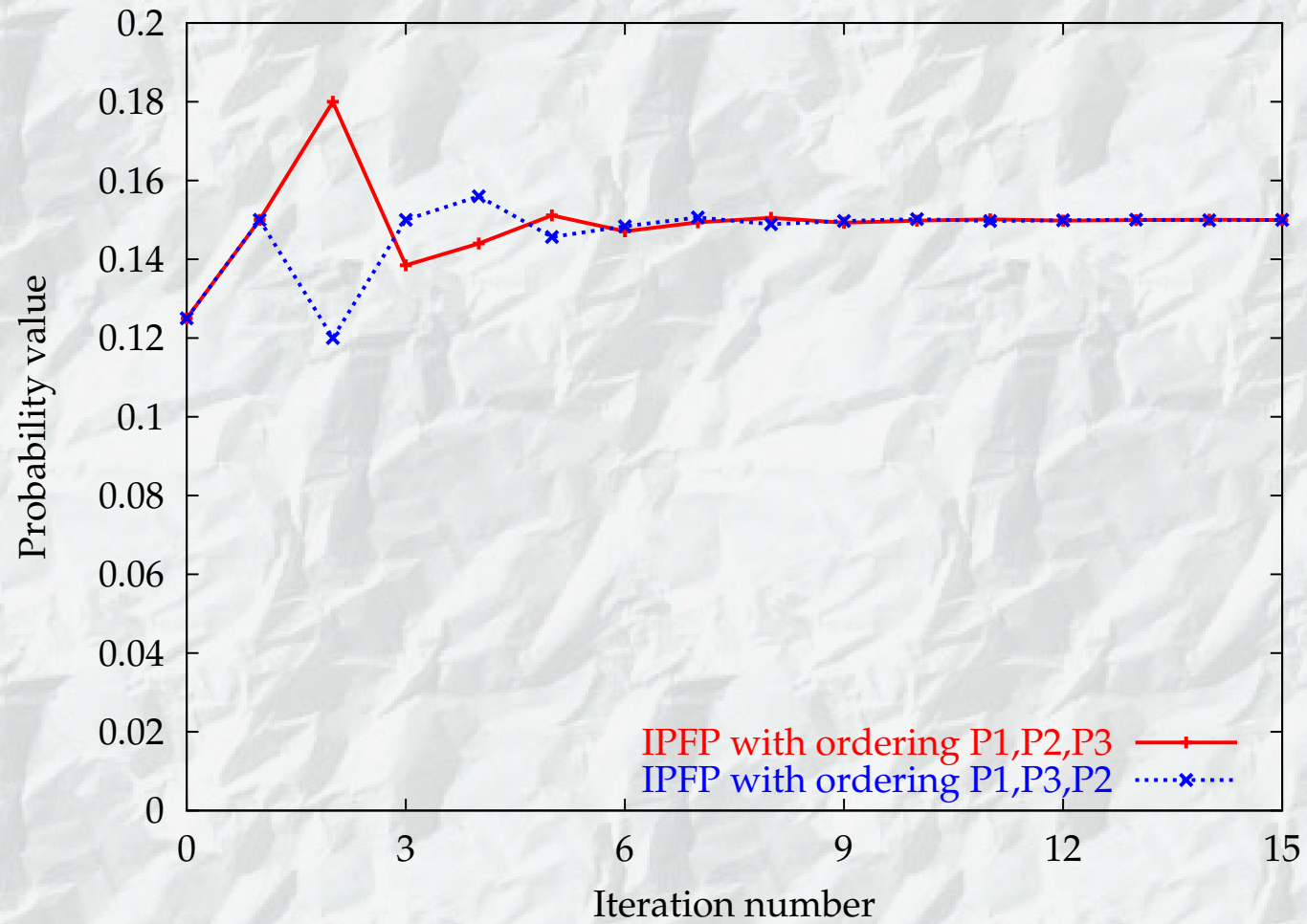$$I(P \parallel Q) \quad = \quad \sum_x P(X = x) \log \frac{P(X = x)}{Q(X = x)}$$

# Iterative Proportional Fitting Procedure (IPFP)

## Deming & Stephan, 1940



$$\pi(Q, \mathcal{S}_j) \;=\; Q\frac{P_j}{Q^{E_j}}$$
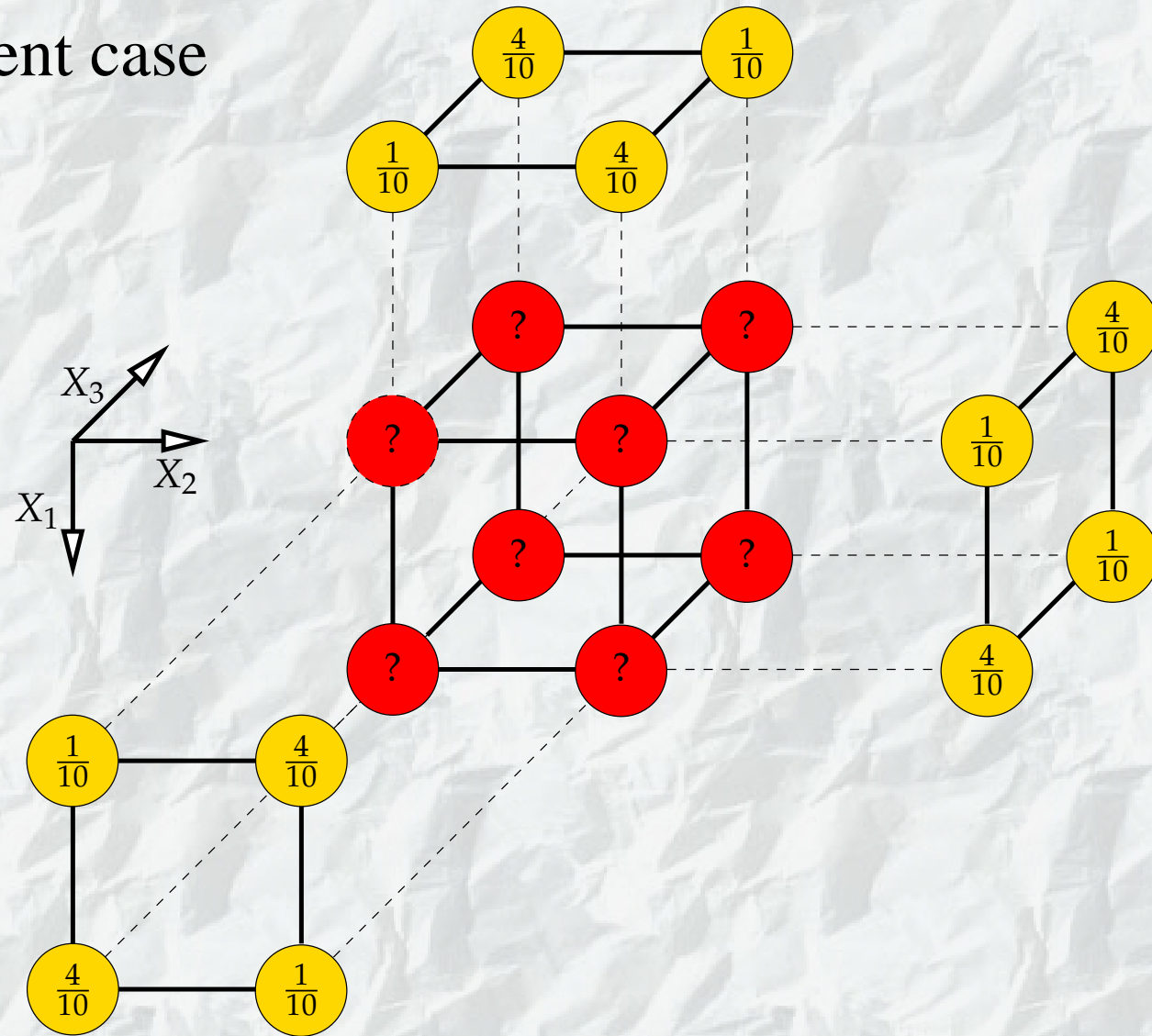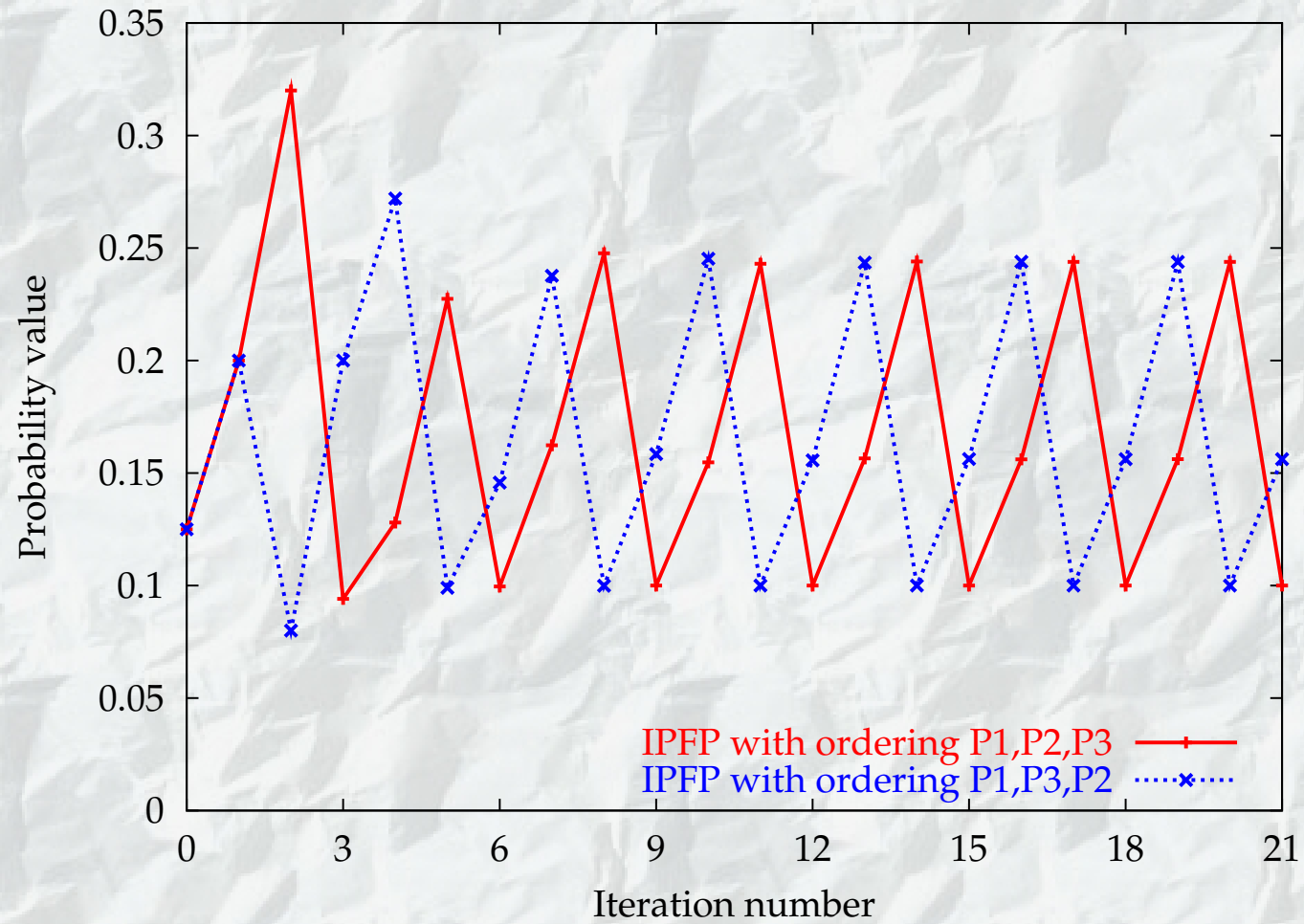
Inconsistent case

$$\alpha = \frac{1}{10}$$

Let $r = \sqrt{-3\alpha^2 + 2\alpha}$, $\beta = 0.5(1 - \alpha - r)$, and $\gamma = 0.5(-\alpha + r)$.

### The limit cycle for the ordering $\textcolor{red}{P_1, P_2, P_3}$

| $x$ | 000 | 001 | $\textcolor{red}{010}$ | 011 | 100 | 101 | 110 | 111 |
|---|---|---|---|---|---|---|---|---|
| $lim_{n\to\infty} Q_{3n+1}(x)$ | 0 | $\alpha$ | $\textcolor{red}{\gamma}$ | $\beta$ | $\beta$ | $\gamma$ | $\alpha$ | 0 |
| $lim_{n\to\infty} Q_{3n+2}(x)$ | 0 | $\gamma$ | $\textcolor{red}{\beta}$ | $\alpha$ | $\alpha$ | $\beta$ | $\gamma$ | 0 |
| $lim_{n\to\infty} Q_{3n+3}(x)$ | 0 | $\beta$ | $\textcolor{red}{\alpha}$ | $\gamma$ | $\gamma$ | $\alpha$ | $\beta$ | 0 |
| arithm. average | 0 | $\frac{1}{6}$ | $\frac{1}{6}$ | $\frac{1}{6}$ | $\frac{1}{6}$ | $\frac{1}{6}$ | $\frac{1}{6}$ | 0 |

### The limit cycle for the ordering $\textcolor{blue}{P_1, P_3, P_2}$

| $x$ | 000 | 001 | $\textcolor{blue}{010}$ | 011 | 100 | 101 | 110 | 111 |
|---|---|---|---|---|---|---|---|---|
| $lim_{n\to\infty} Q_{3n+1}(x)$ | 0 | $\alpha$ | $\textcolor{blue}{\beta}$ | $\gamma$ | $\gamma$ | $\beta$ | $\alpha$ | 0 |
| $lim_{n\to\infty} Q_{3n+2}(x)$ | 0 | $\gamma$ | $\textcolor{blue}{\alpha}$ | $\beta$ | $\beta$ | $\alpha$ | $\gamma$ | 0 |
| $lim_{n\to\infty} Q_{3n+3}(x)$ | 0 | $\beta$ | $\textcolor{blue}{\gamma}$ | $\alpha$ | $\alpha$ | $\gamma$ | $\beta$ | 0 |
| arithm. average | 0 | $\frac{1}{6}$ | $\frac{1}{6}$ | $\frac{1}{6}$ | $\frac{1}{6}$ | $\frac{1}{6}$ | $\frac{1}{6}$ | 0 |

For $\alpha = 0.1$ we get $\beta \doteq 0.244$ and $\gamma \doteq 0.156$.

# Inconsistent input set $\mathcal{P} = \{P_1, \ldots, P_k\}$

It means that $\mathcal{S} = \cap_{j=1}^{k} \mathcal{S}_j = \emptyset$.

$Q(X_1, \ldots, X_n)$ is required to:

- **minimize a distance aggregate** with respect to $\mathcal{P}$:

$$\sum_{P_j \in \mathcal{P}} w_j \cdot d(P_j \parallel Q^{E_j})$$

- **factorize** with respect to $\mathcal{E} = \{E_1, \ldots, E_k\}$:
  there exist potentials $\psi_{E_i} : \mathbb{X}^{E_i} \mapsto \mathbb{R}, i = 1, 2, \ldots, k$ such that for all $x \in \mathbb{X}$

$$Q(x) = \prod_{E_i \in \mathcal{E}} \psi_{E_i}(x^{E_i}) \ .$$

# Distance

- measured by the Kullback-Leibler divergence

$$d(P_j \parallel Q^{E_j}) \;\; = \;\; (P_j \parallel Q^{E_j}) \;=\; \sum_{x^{E_j}} P_j(x^{E_j}) \log \frac{P_j(x^{E_j})}{Q^{E_j}(x^{E_j})}$$
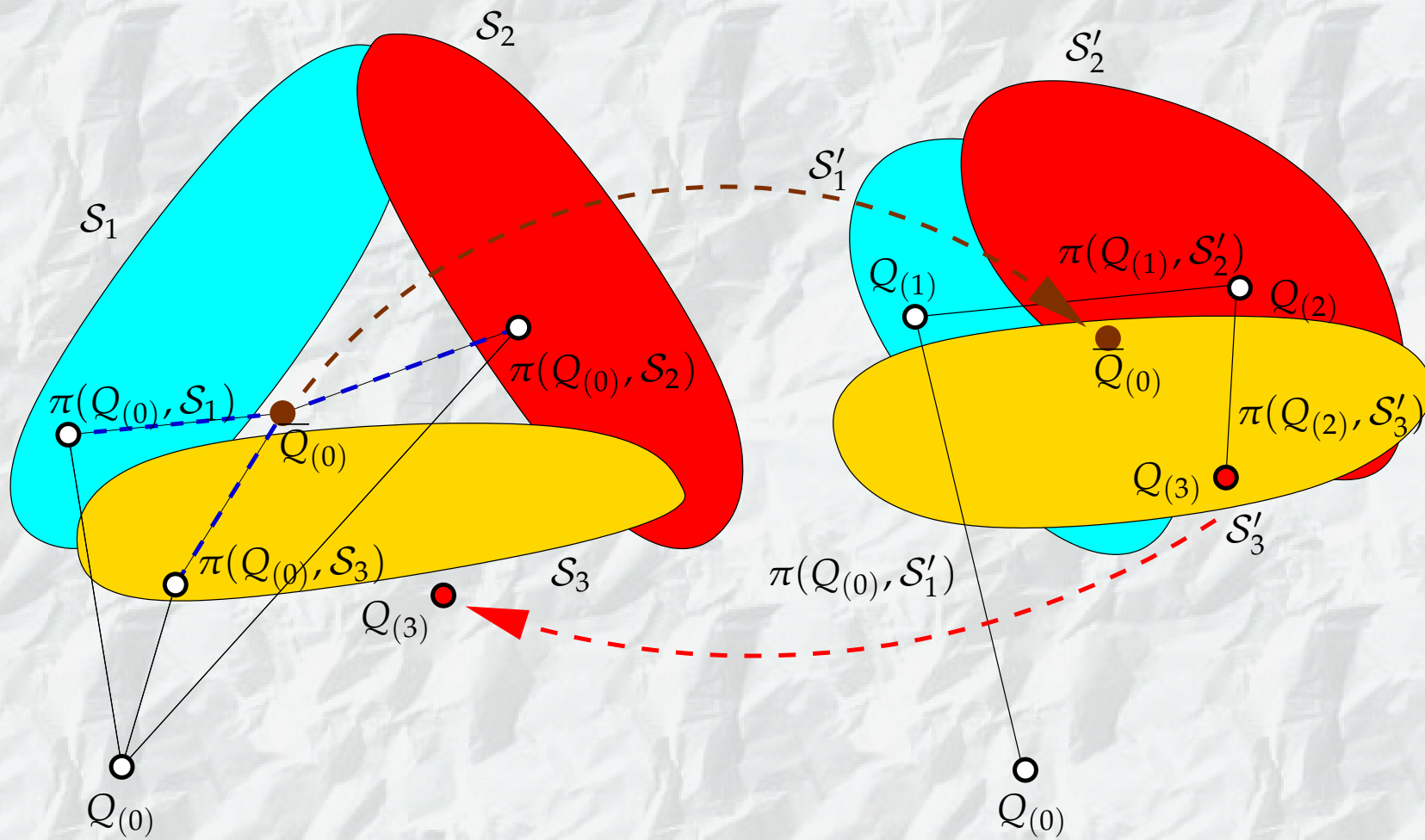
- measured by the total variance

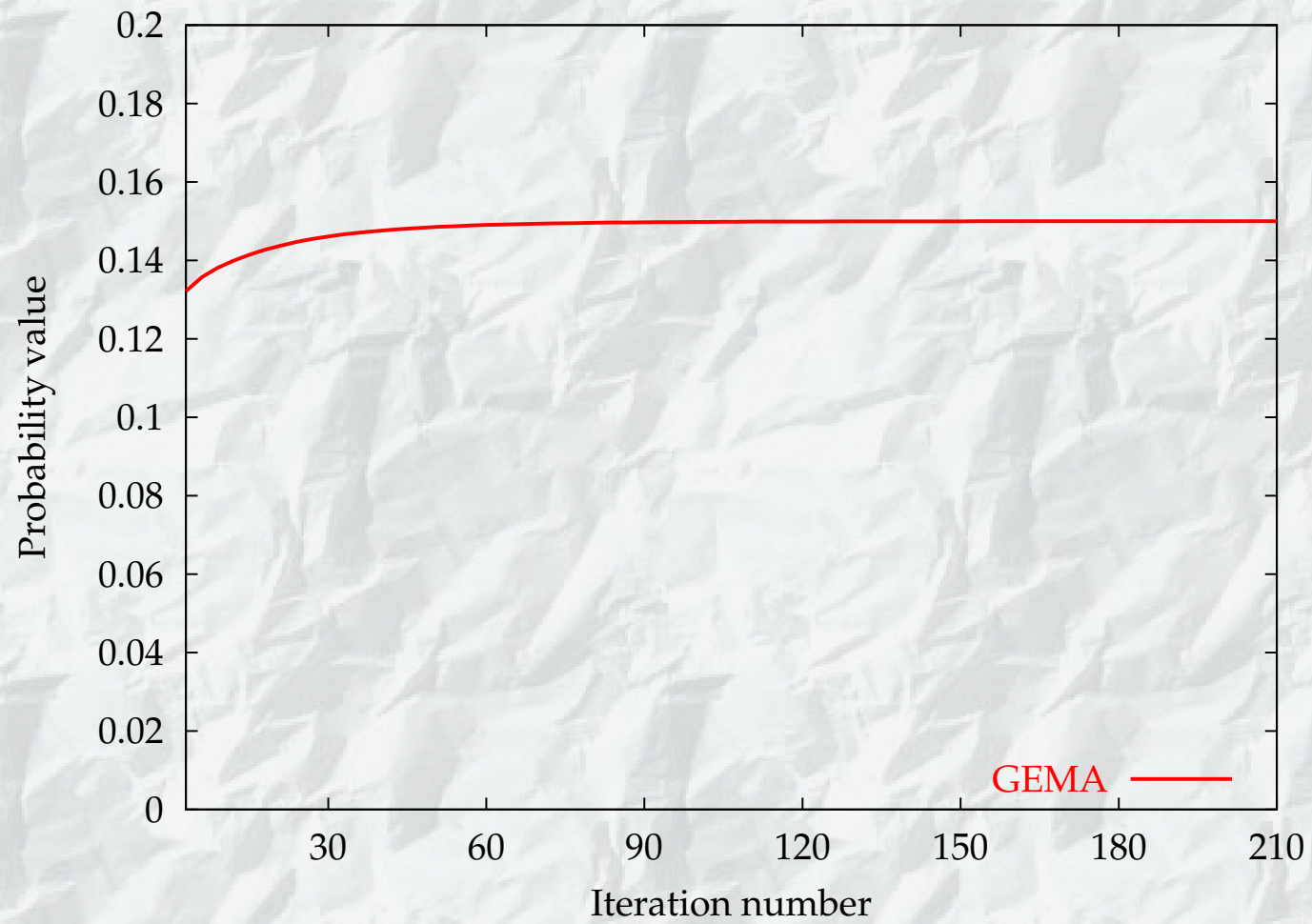$$d(P_j \parallel Q^{E_j}) \;\; = \;\; |P_j - Q^{E_j}| \;=\; \sum_{x^{E_j}} |P_j(x^{E_j}) - Q^{E_j}(x^{E_j})|$$

# IPFP properties in the inconsistent case

- "converges" to a limit cycle

- distributions in the limit cycle are different for different orderings of the input set

- in the example, the average of the distributions in the limit cycle does not depend on the ordering – but generally, it is not true

- in the example, the distributions in the limit cycles minimized the aggregate of the total variance – but generally it is not known

- there are also other distributions that minimize the aggregate of the total variance that are not computed with IPFP

- generally, the distributions in the limit cycles **do not** minimize the aggregate of the Kullback-Leibler divergence

- distributions computed within a finite number of iterations factorize with respect to $\mathcal{E} = \{E_1, \ldots, E_k\}$
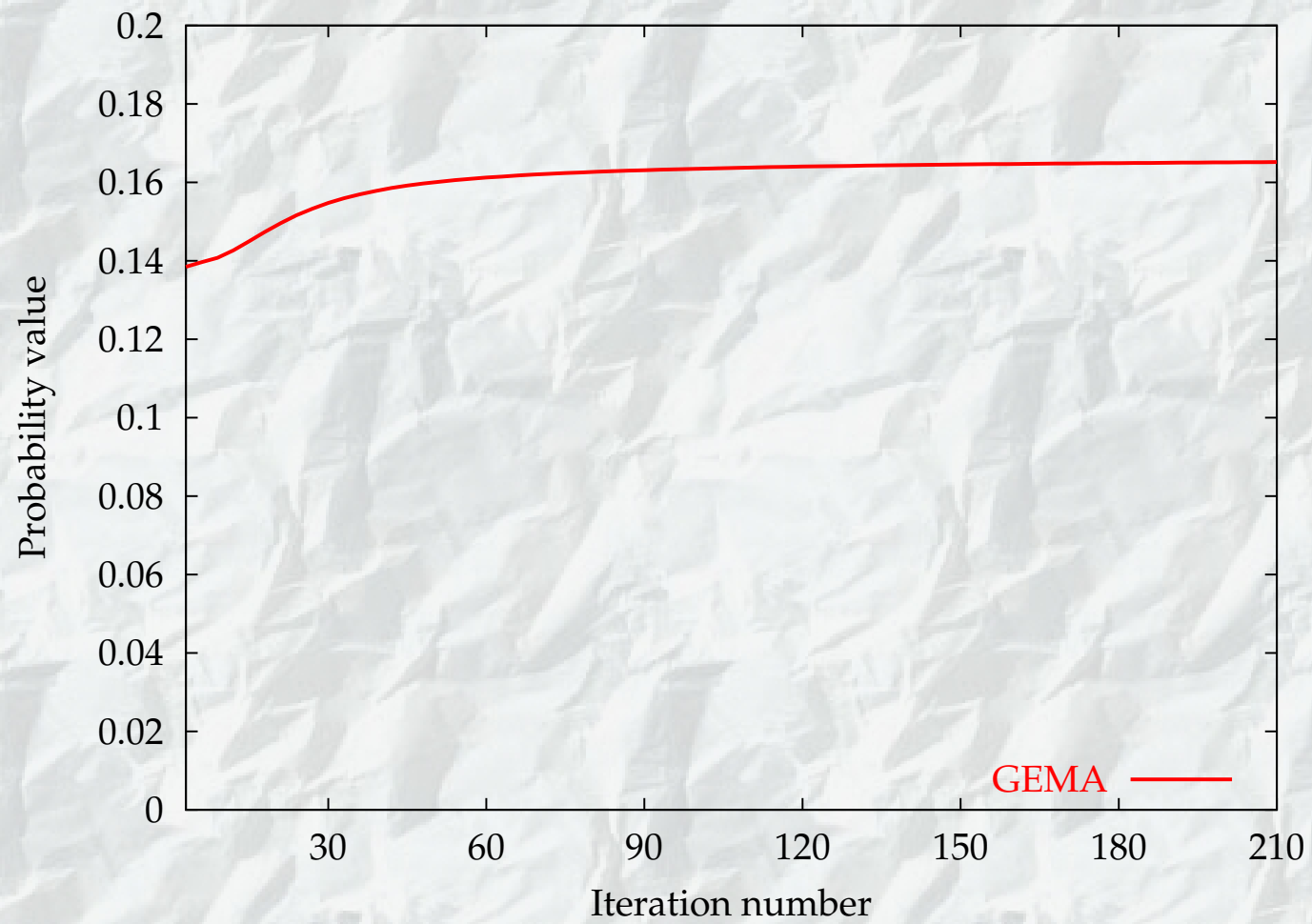
GEMA

GEMA on the inconsistent input set

# GEMA properties

- converges also in the inconsistent case

- the limit distribution satisfies the necessary condition for the local minima of the <span style="color:red">aggregate of the Kullback-Leibler divergence</span>

- the distributions computed within a finite number of iterations factorize with respect to $\mathcal{E} = \{E_1, \ldots, E_k\}$