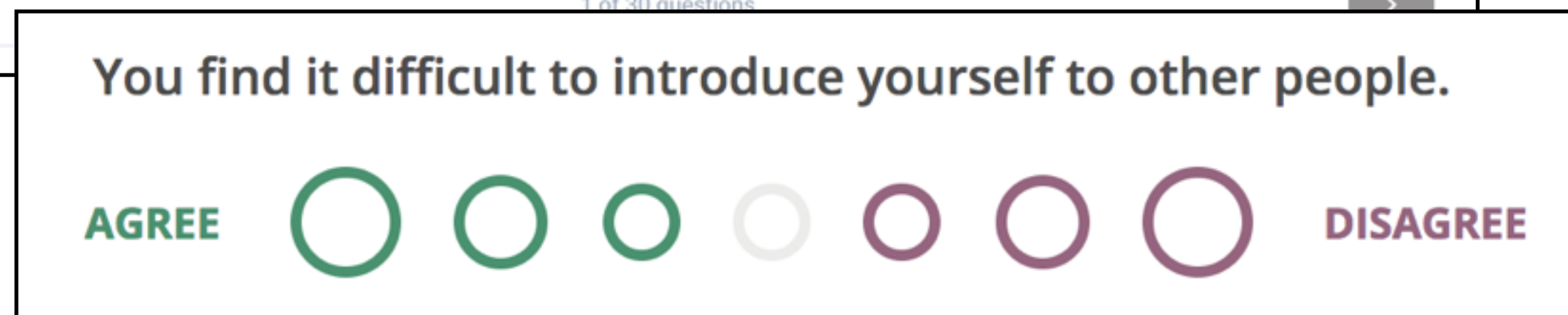
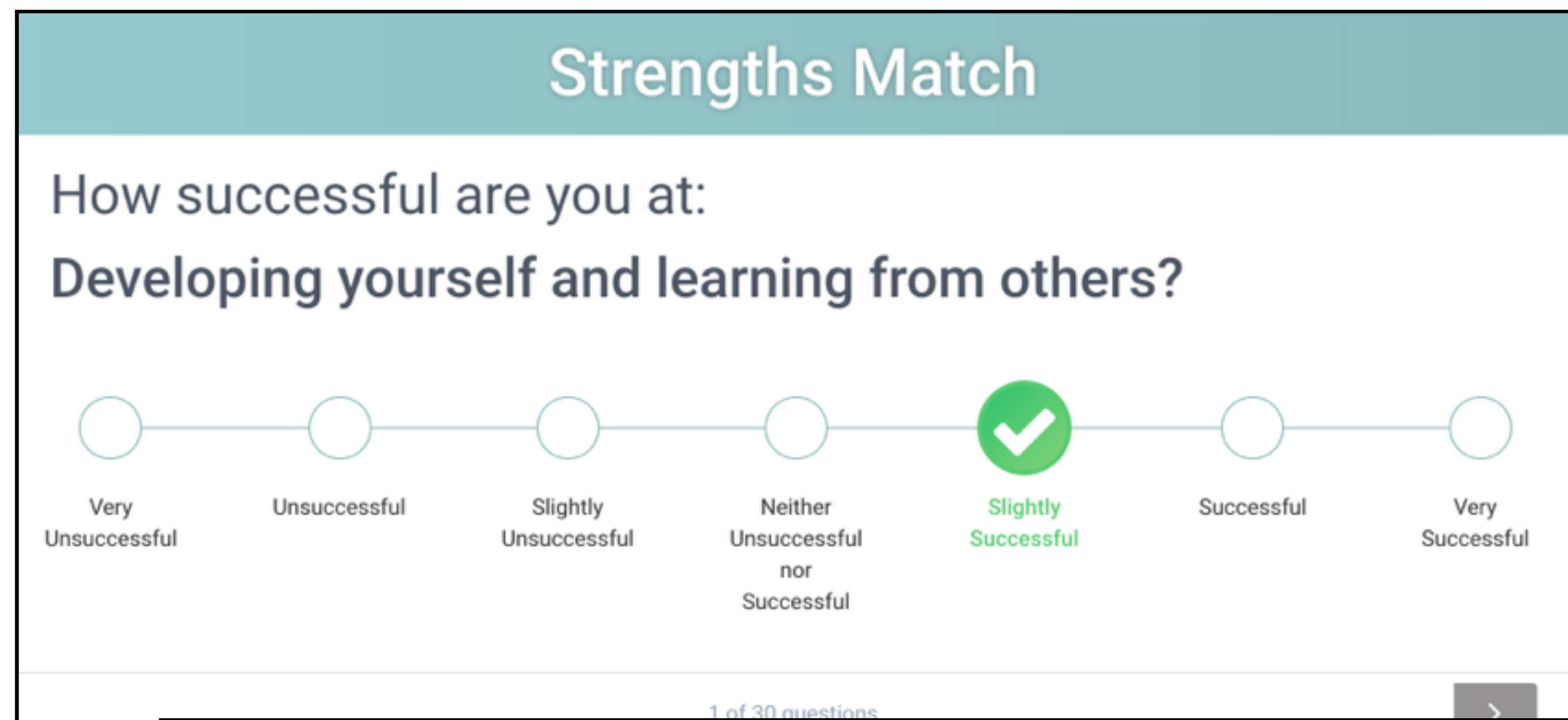

SHORTING NUMBER OF QUESTIONS IN LONG PSYCHOLOGICAL QUESTIONNAIRES

© Capp & Co Ltd, 2017

THE CONTEXT



- Psychological tests are useful for companies
- Discrete type of response
- Paper-based / online-based

* <https://www.jobmi.com/>

** <https://www.16personalities.com/>

THE PROBLEM

I

Are the questions capturing what we want to capture?

II

Are there redundancy among questions such that we can reduce the size of the test?

THE PROBLEM

I

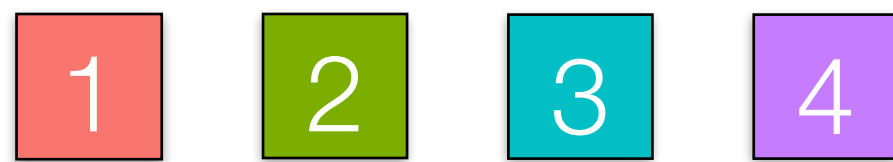
Are the questions capturing what we want to capture?

II

Are there redundancy among questions such that we can reduce the size of the test?

THE PROBLEM

- Dataset: collection of user responses (~30000)
- In our case the test has 90 questions with 4 possible answers:

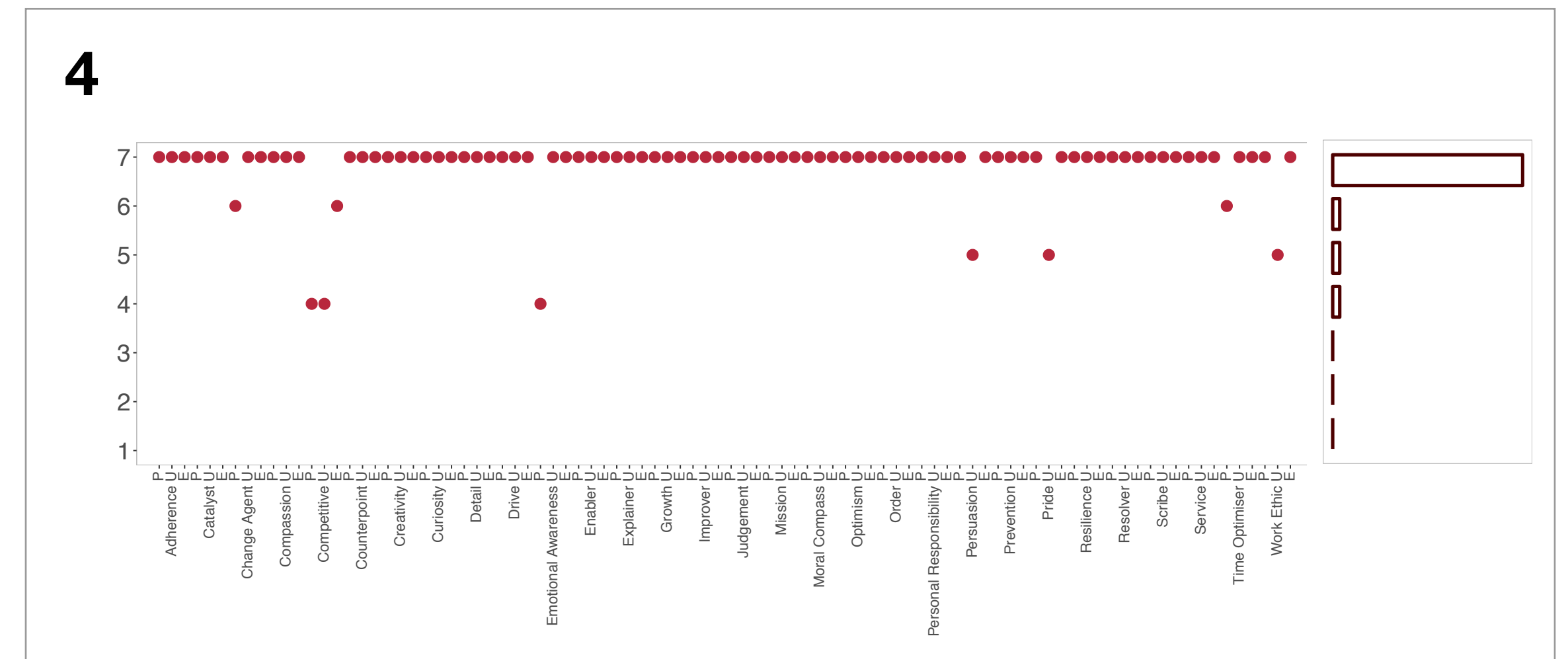
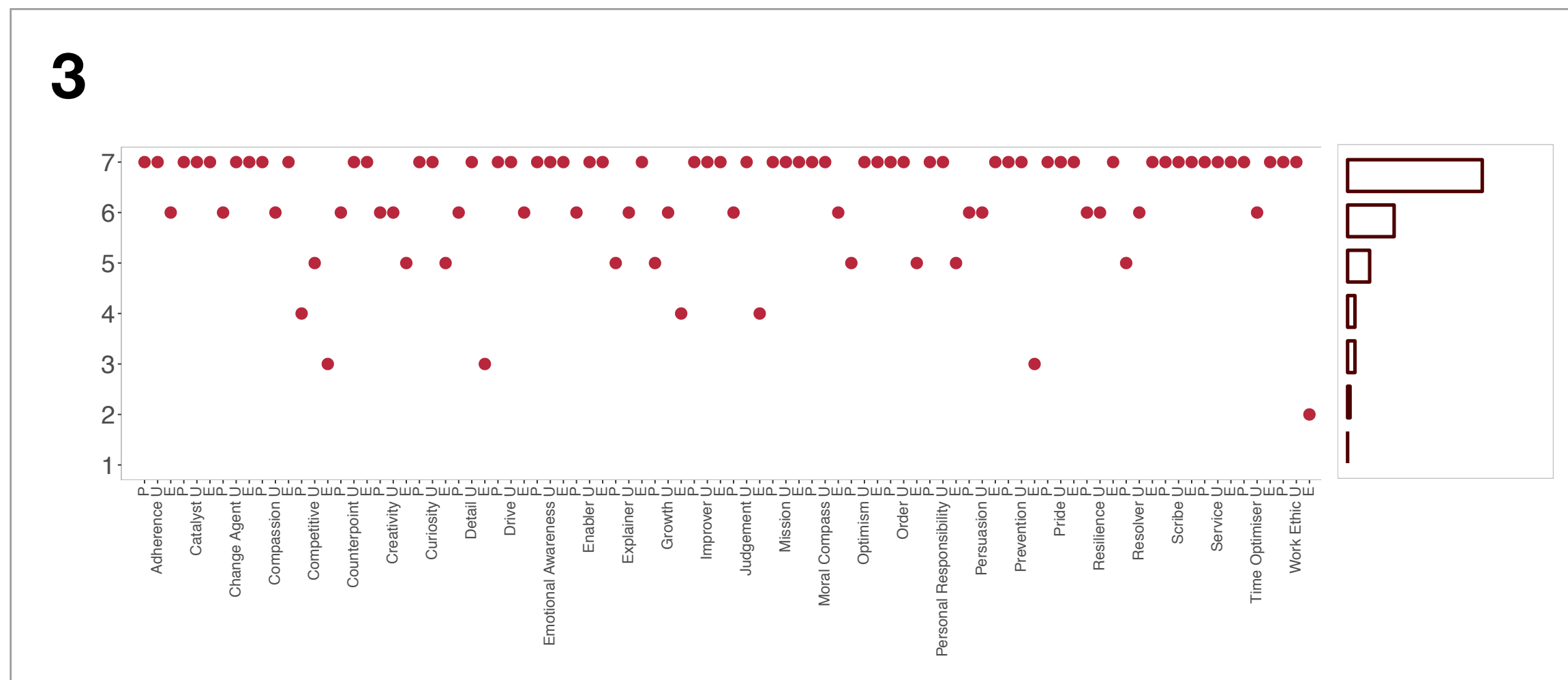
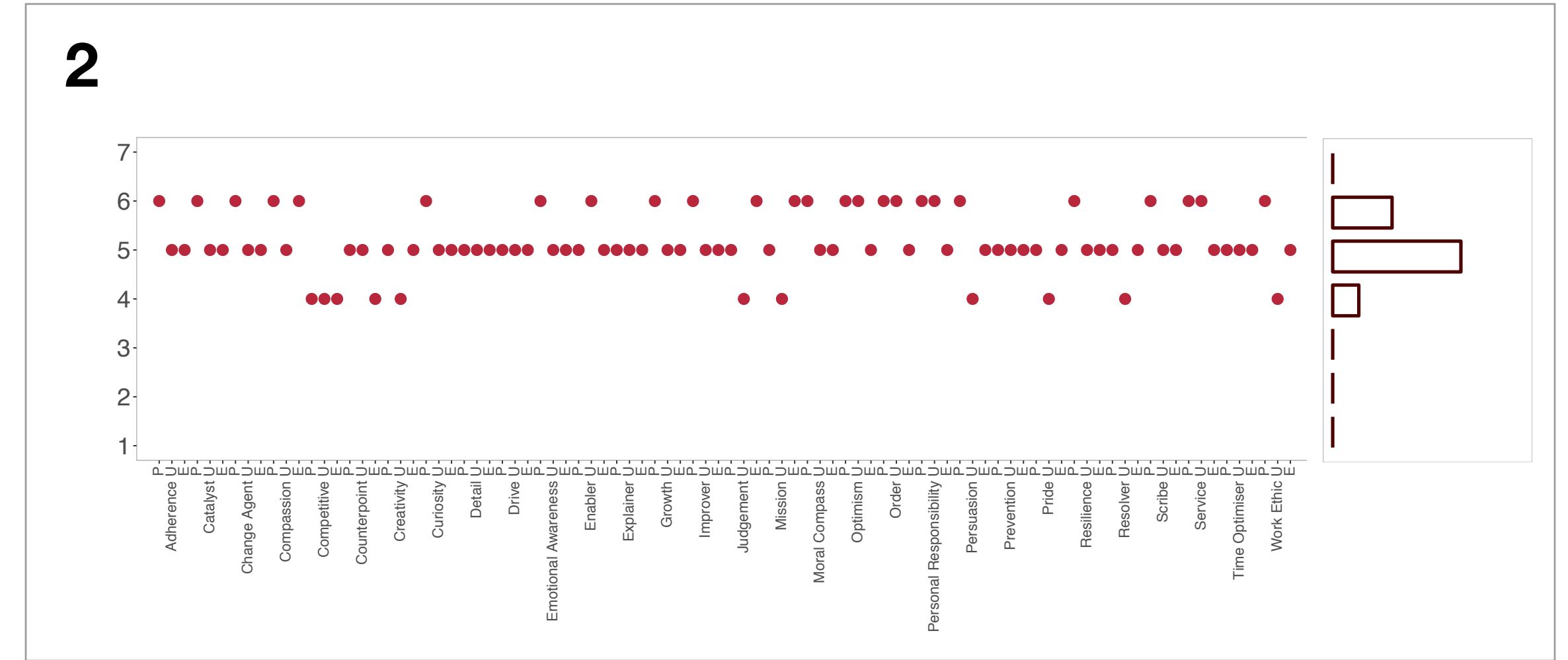
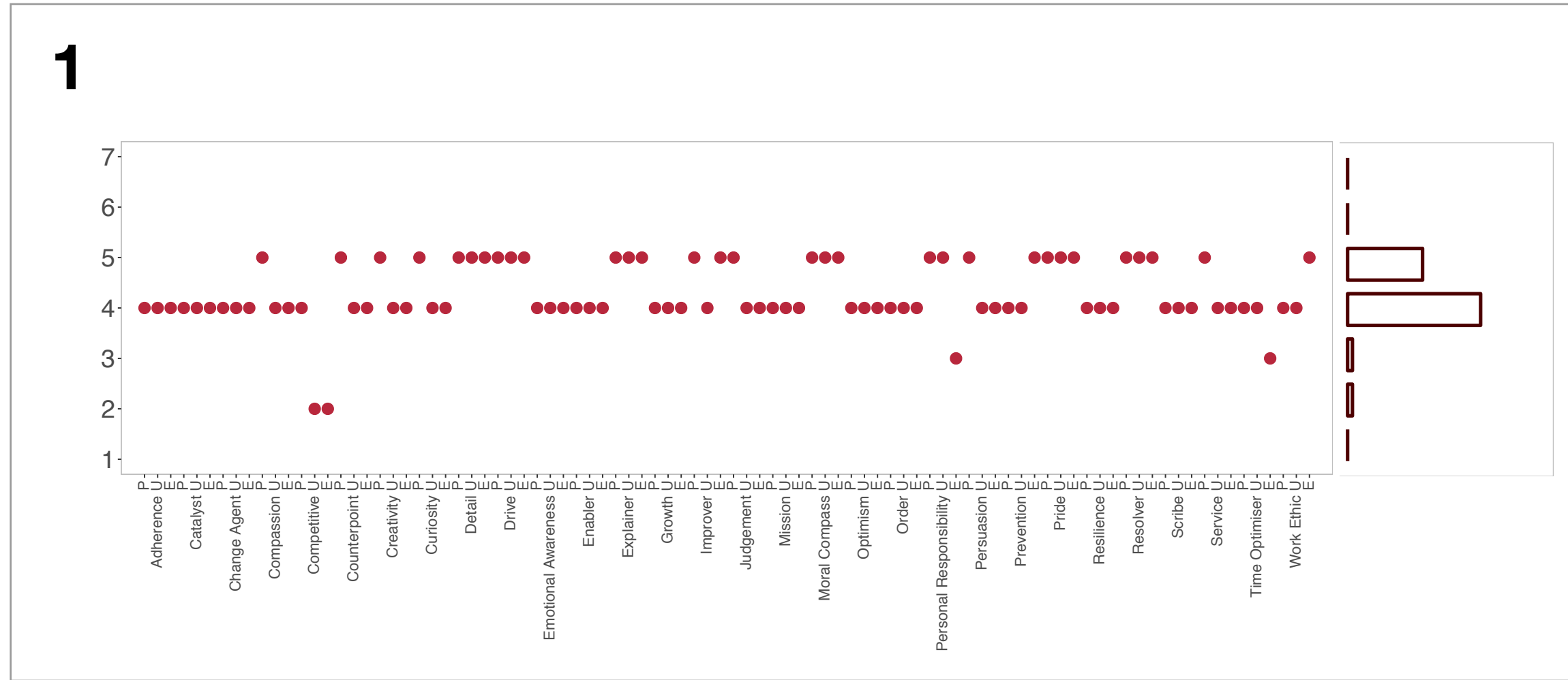


- Drawbacks?

Var	Var 4	Var 5	Var 6	Var 7	Var 8	Var 9	Var 10	Var 11	Var 12	Var 13	Var 14	Var 15	Var 16	Var 17	Var 18	Var 19	Var 20
1	4	4	3	3	1	3	1	3	3	3	4	1	3	1	3	1	...
1	1	4	2	2	3	3	2	3	1	1	3	4	3	3	3	2	...
1	1	1	4	3	4	3	3	3	3	2	3	1	3	3	1	1	...
2	3	2	1	1	3	1	1	1	3	1	1	1	2	1	3	1	...
1	1	2	3	1	3	3	2	1	4	1	1	2	1	3	1	3	...
1	2	1	1	3	2	2	2	2	2	3	2	1	1	3	1	1	...
2	3	4	1	1	1	3	3	4	3	3	3	4	2	2	3	2	...
1	3	2	1	1	3	3	3	2	3	3	3	1	3	1	1	1	...
2	2	1	1	1	1	3	1	3	3	1	2	1	1	3	3	3	...
4	3	4	3	2	3	3	2	1	3	3	1	1	3	2	3	3	...
2	2	3	3	3	1	3	3	3	3	4	1	1	3	1	1	2	...
3	1	1	3	1	3	1	2	3	3	2	3	1	4	4	1	4	...
3	3	4	3	1	4	3	3	3	3	3	3	2	4	2	2	1	...
2	1	3	3	3	3	3	3	2	3	3	1	1	3	2	3	1	...
1	4	1	2	2	3	3	1	1	3	4	1	1	1	3	1	4	...

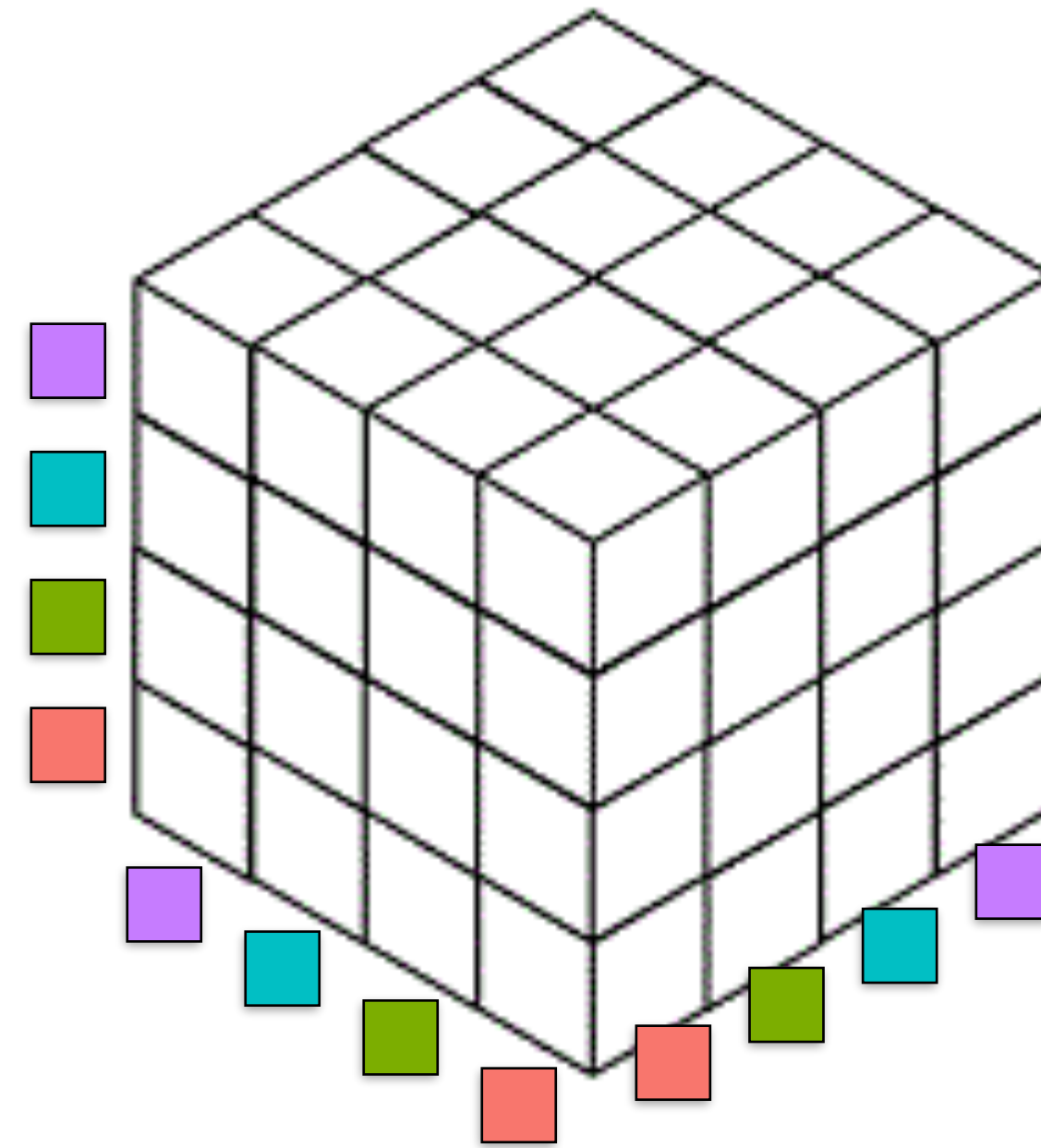
Drawback: users perceive the scale in different way

Drawback: users tend to choose high values



THE PROBLEM

- Dataset can be mapped into $\{1, 2, 3, 4\}^{90}$
- How does it look like?



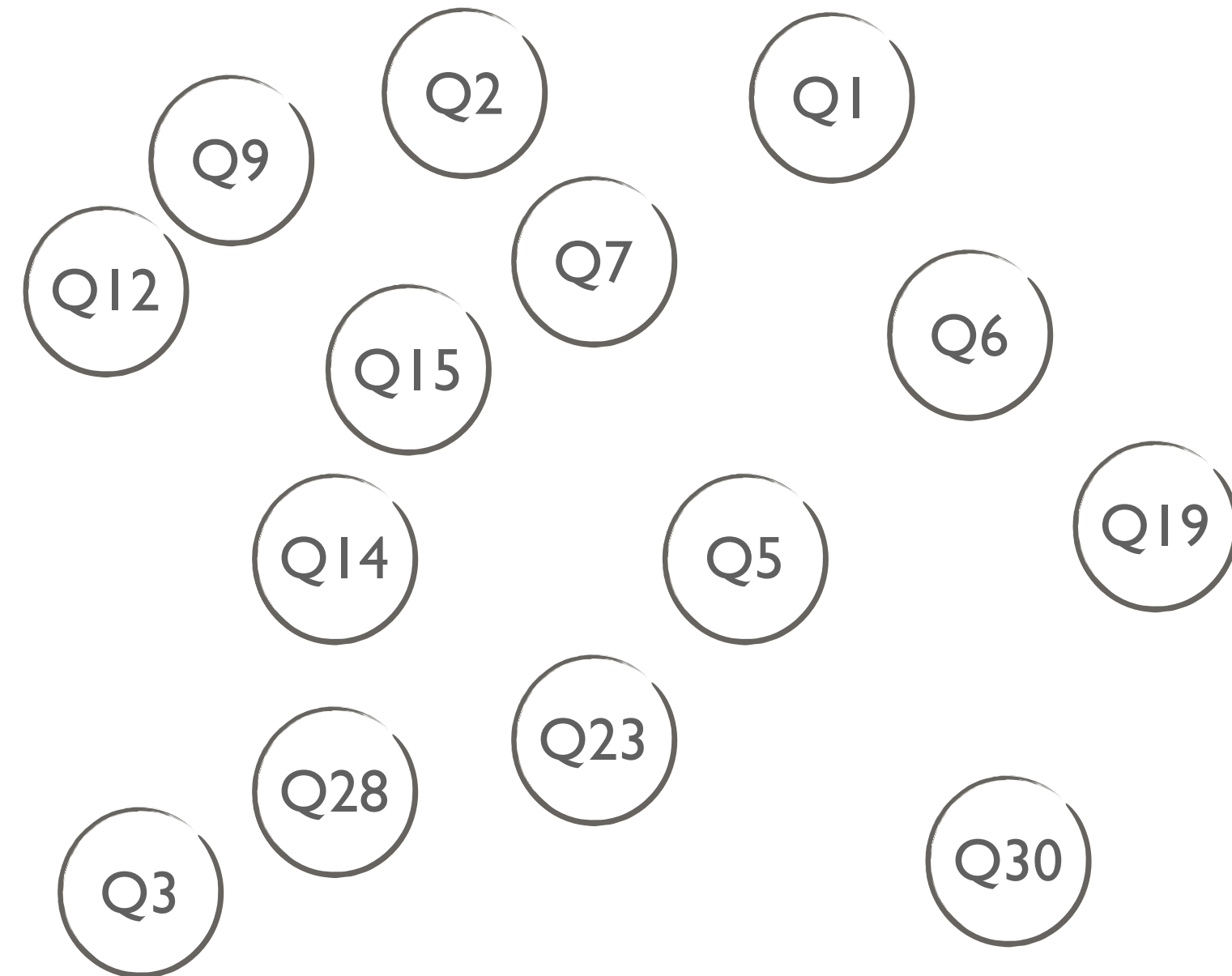




Drawback: sparse data.

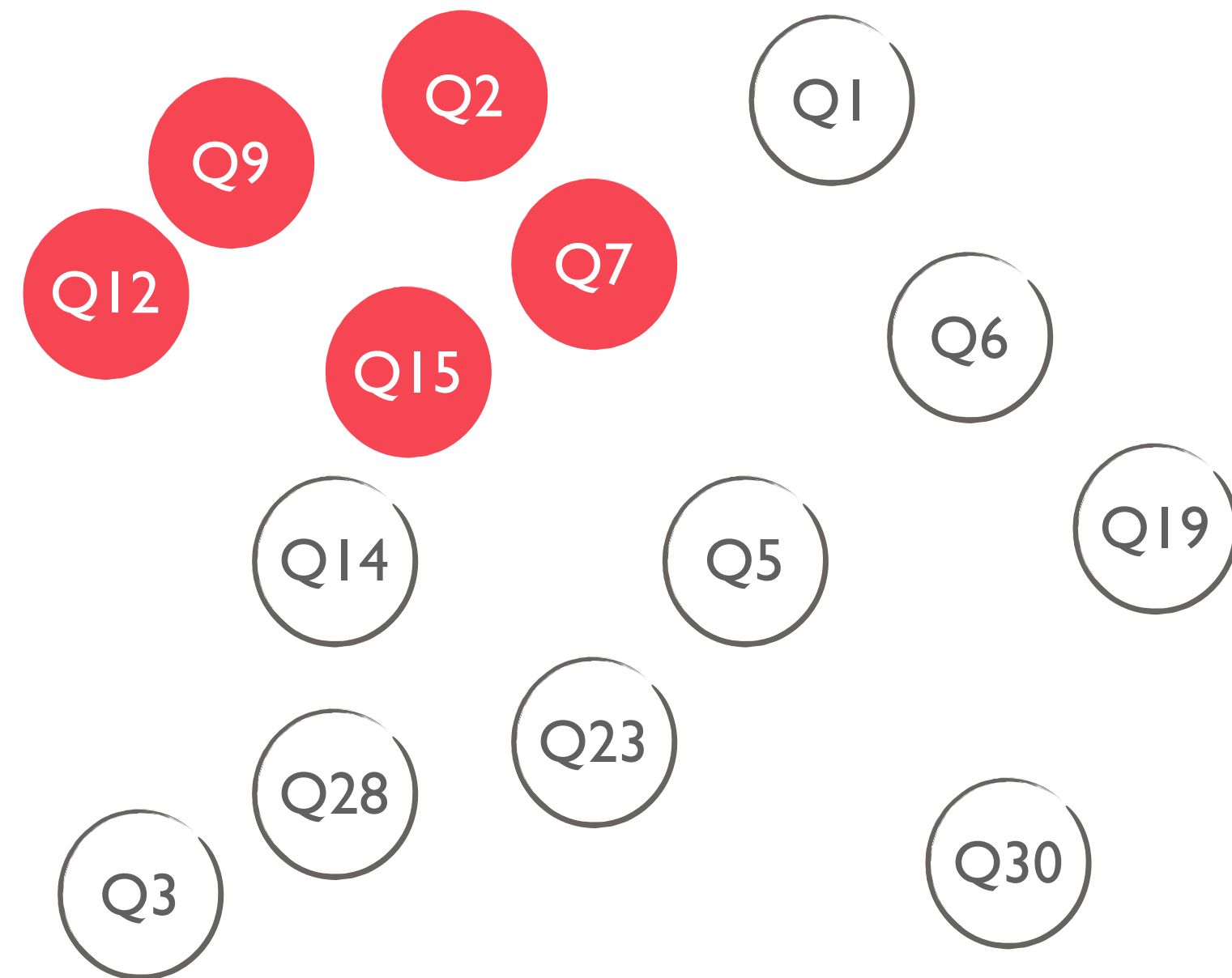
Drawback: data concentrated in few cells

HOW TO SOLVE IT



- Divide the problem:

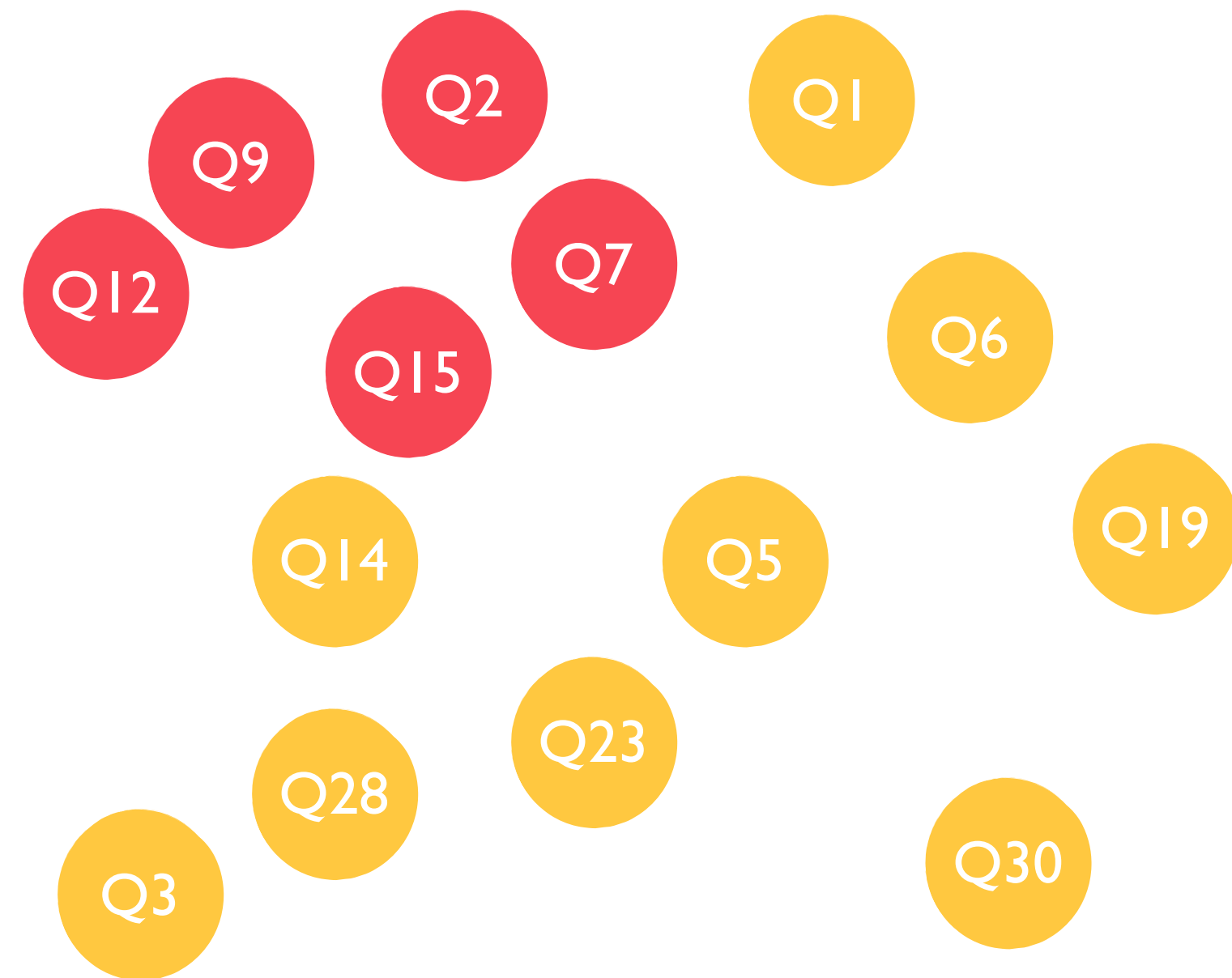
HOW TO SOLVE IT



- Divide the problem:
- I. Find a set of predictors and a set of questions to be predicted

$$Q = P \dot{U} S$$

HOW TO SOLVE IT



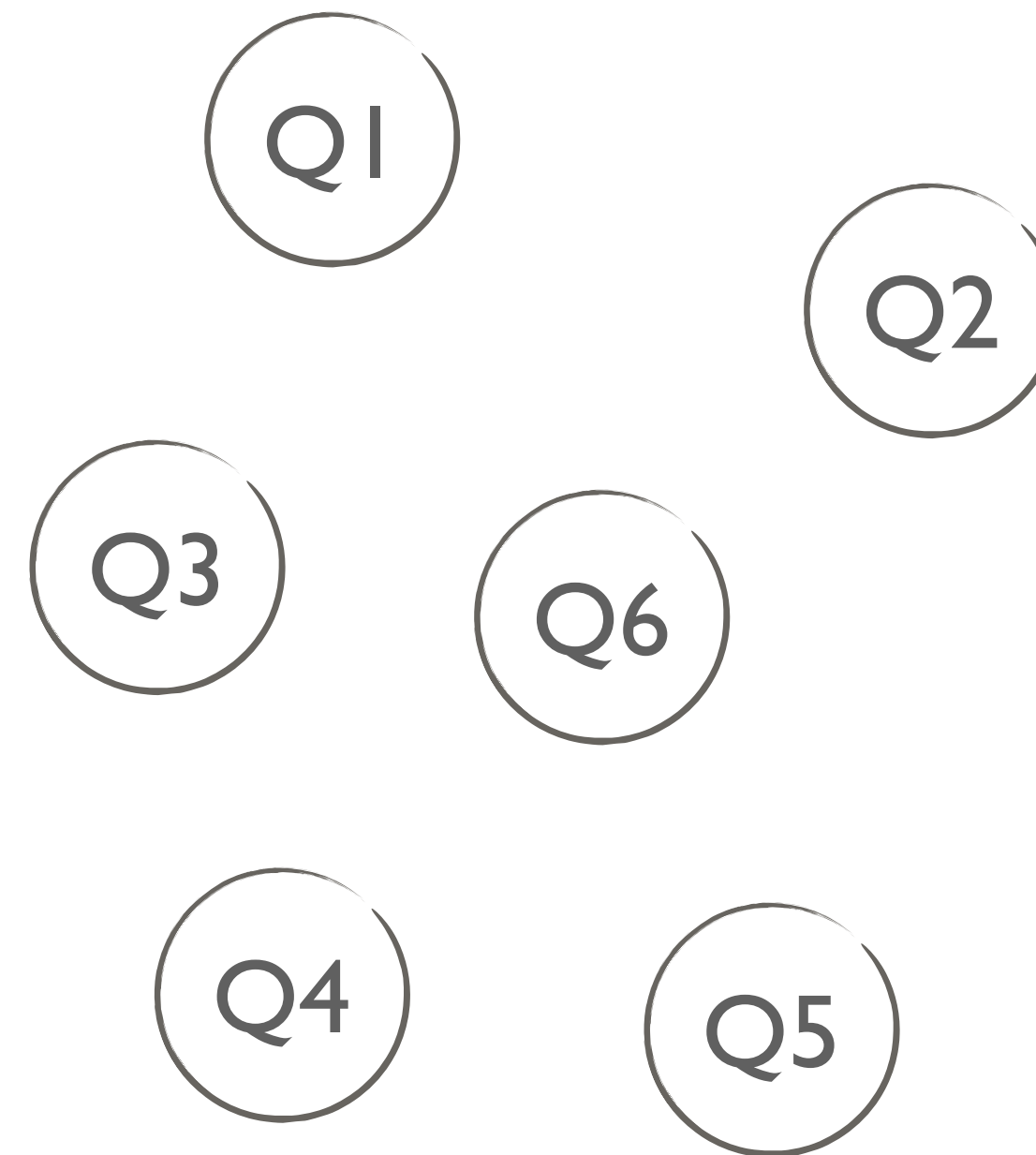
- Divide the problem:
- 1. Find a set of predictors and a set of questions to be predicted

$$Q = P \dot{U} S$$

- 2. Predict P using S

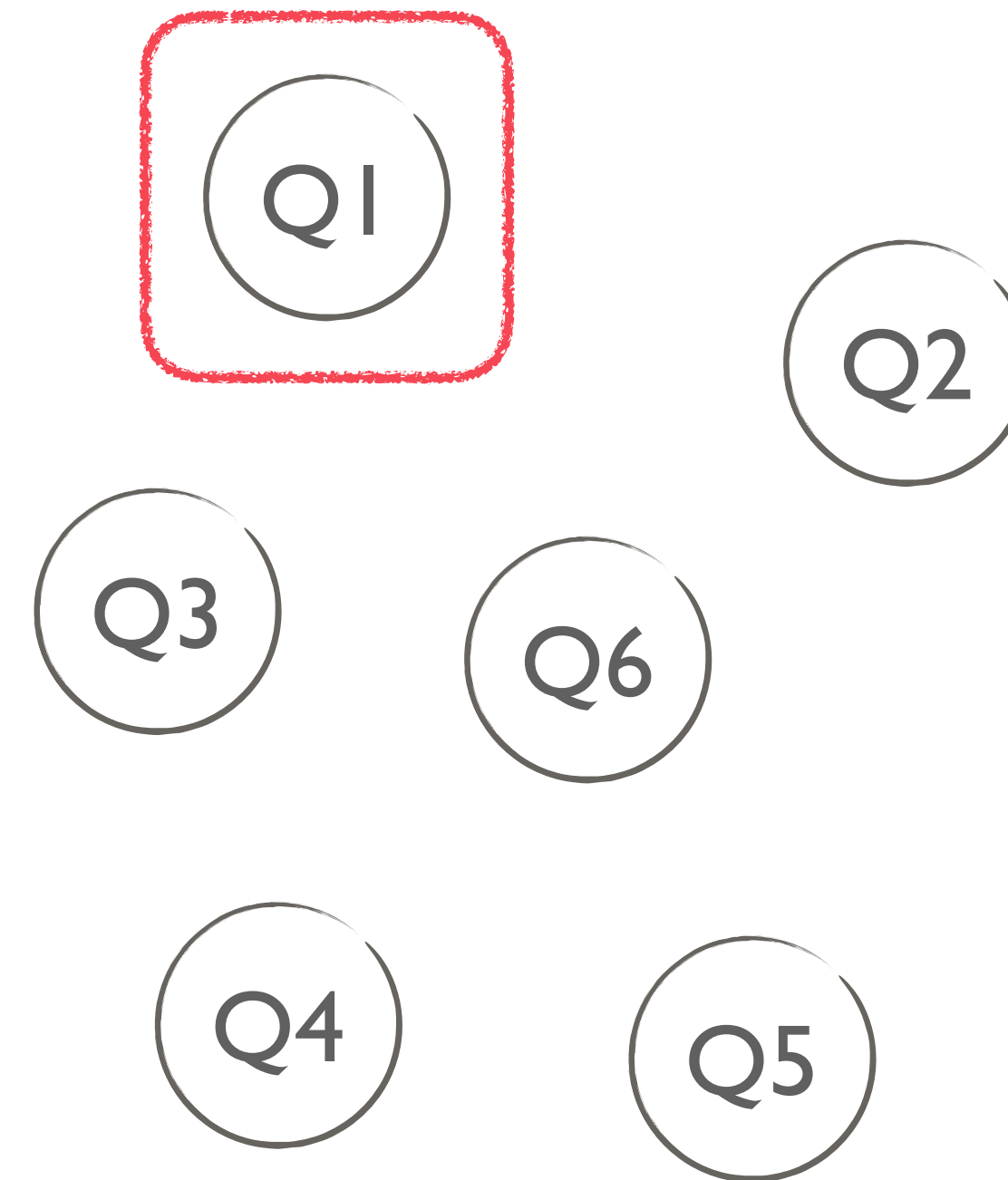
I. FEATURE SELECTION

- First part is a *feature selection* problem
- Ideally, find P and S automatically
- In reality, divide the problem



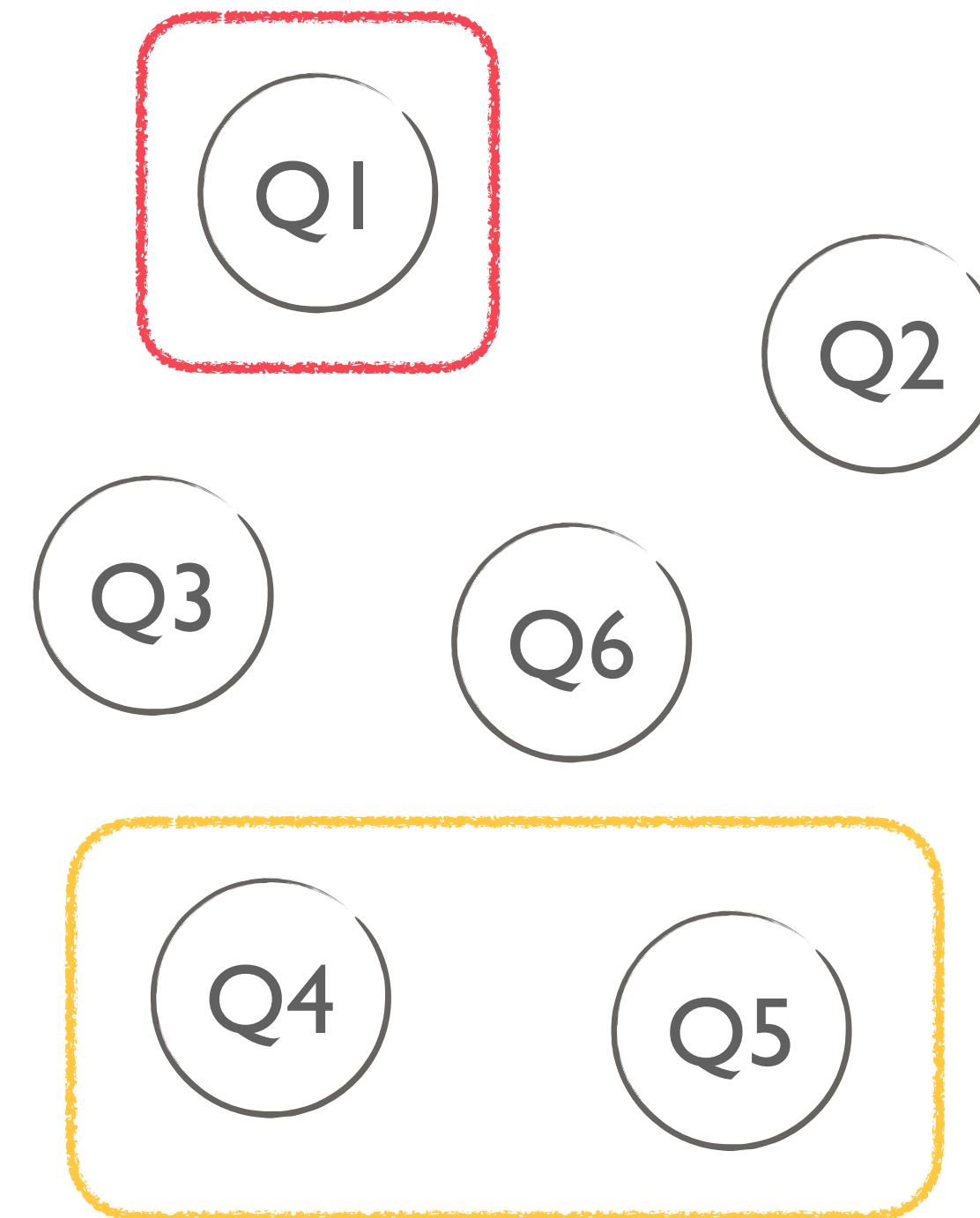
I. FEATURE SELECTION

- First part is a *feature selection* problem
- Ideally, find P and S automatically
- In reality, divide the problem
- **Fix one question** and find the best subset of predictors



I. FEATURE SELECTION

- First part is a *feature selection* problem
- Ideally, find P and S automatically
- In reality, divide the problem
- **Fix one question** and find the **best subset** of predictors



I. FEATURE SELECTION

	Filter methods
Description	Intrinsic properties of data
Advantages	Computationally simple and fast
Disadvantages	Ignore interaction with the classifier
Examples used in our problem	Correlation-based Mutual Information

I. FEATURE SELECTION

Description

Advantages

Disadvantages

Examples used
in our problem

Filter methods

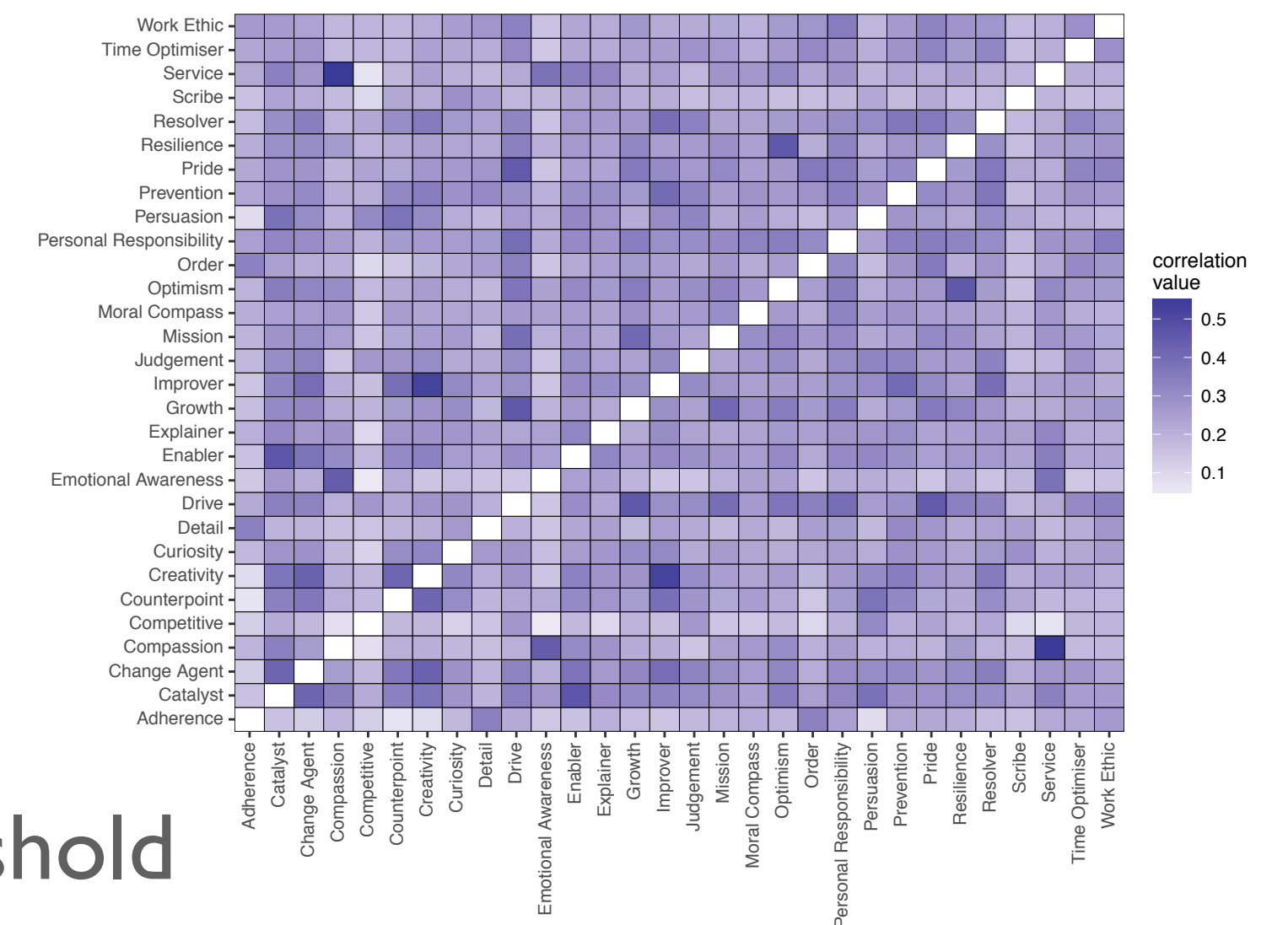
Intrinsic properties of data

Computationally simple and fast

Ignore interaction with the classifier

Correlation-based → Select a threshold

Mutual Information



I. FEATURE SELECTION

Filter methods

Description

Intrinsic properties of data

Advantages

Computationally simple and fast

Disadvantages

Ignore interaction with the classifier

Examples used
in our problem

Correlation-based

Mutual Information →

- 1) Initialization: Set $F \leftarrow$ “initial set of n features”; $S \leftarrow$ “empty set.”
- 2) Computation of the MI with the output class: For each $f_i \in F$, compute $I(C; f_i)$.
- 3) Selection of the first feature: Find the feature f_i that maximizes $I(C; f_i)$; set $F \leftarrow F \setminus \{f_i\}$; set $S \leftarrow \{f_i\}$.
- 4) Greedy selection: Repeat until $|S| = k$.
 - a) Computation of the MI between variables: For all pairs (f_i, f_s) with $f_i \in F$ and $f_s \in S$, compute $I(f_i; f_s)$, if it is not yet available.
 - b) Selection of the next feature: Choose the feature $f_i \in F$ that maximizes

$$I(C; f_i) - \beta \sum_{f_s \in S} I(f_s; f_i)$$

Set $F \leftarrow F \setminus \{f_i\}$; set $S \leftarrow \{f_i\}$.

- 5) Output the set S containing the selected features.

I. FEATURE SELECTION

Description

Advantages

Disadvantages

Examples used
in our problem

Embedded methods

The search of methods is built into the classifier

Include interaction with the classifier

Classifier dependent selection

Random forest

GLM using regularisation

I. FEATURE SELECTION

Description

Advantages

Disadvantages

Examples used
in our problem

Embedded methods

The search of methods is built into the classifier

Include interaction with the classifier

Classifier dependent selection

Feature Importance ← **Random forest**

GLM using regularisation

I. FEATURE SELECTION

Description

Advantages

Disadvantages

Examples used
in our problem

Embedded methods

The search of methods is built into the classifier

Include interaction with the classifier

Classifier dependent selection

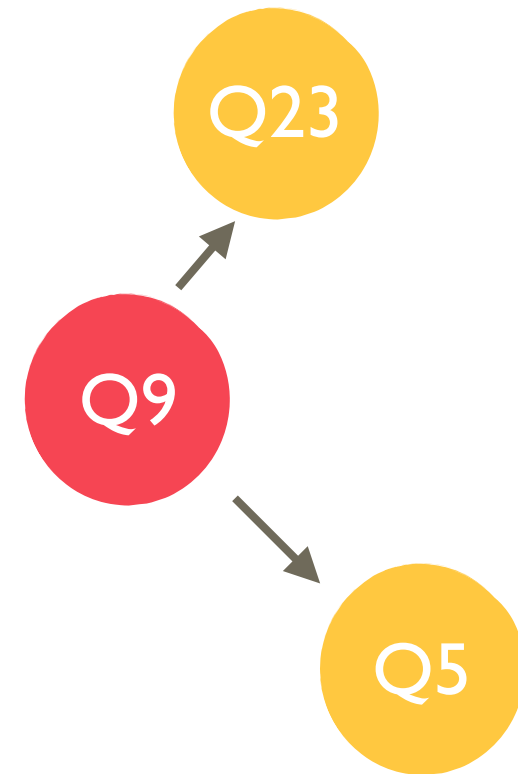
Random forest

GLM using regularisation

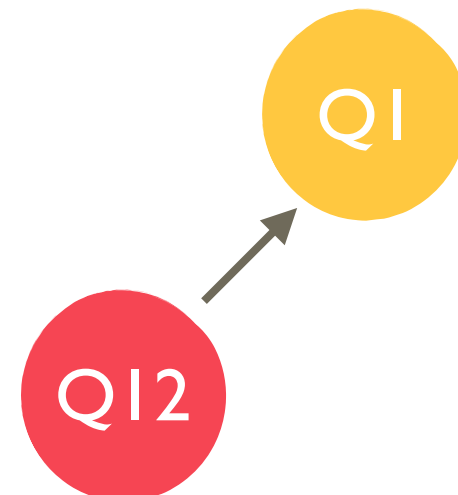
$$\min_{w \in \mathbb{R}^p} \frac{1}{n} \|\hat{X}w - \hat{Y}\|^2 + \lambda(\alpha \|w\|_1 + (1 - \alpha) \|w\|_2^2), \alpha \in [0, 1]$$

I. BUILDING A GRAPH

- Feature selection methods are applied

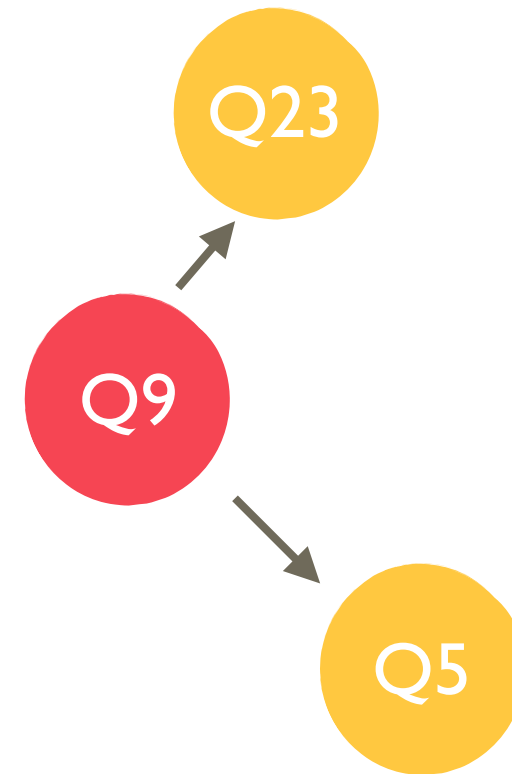


- For each question we have a set of predictors

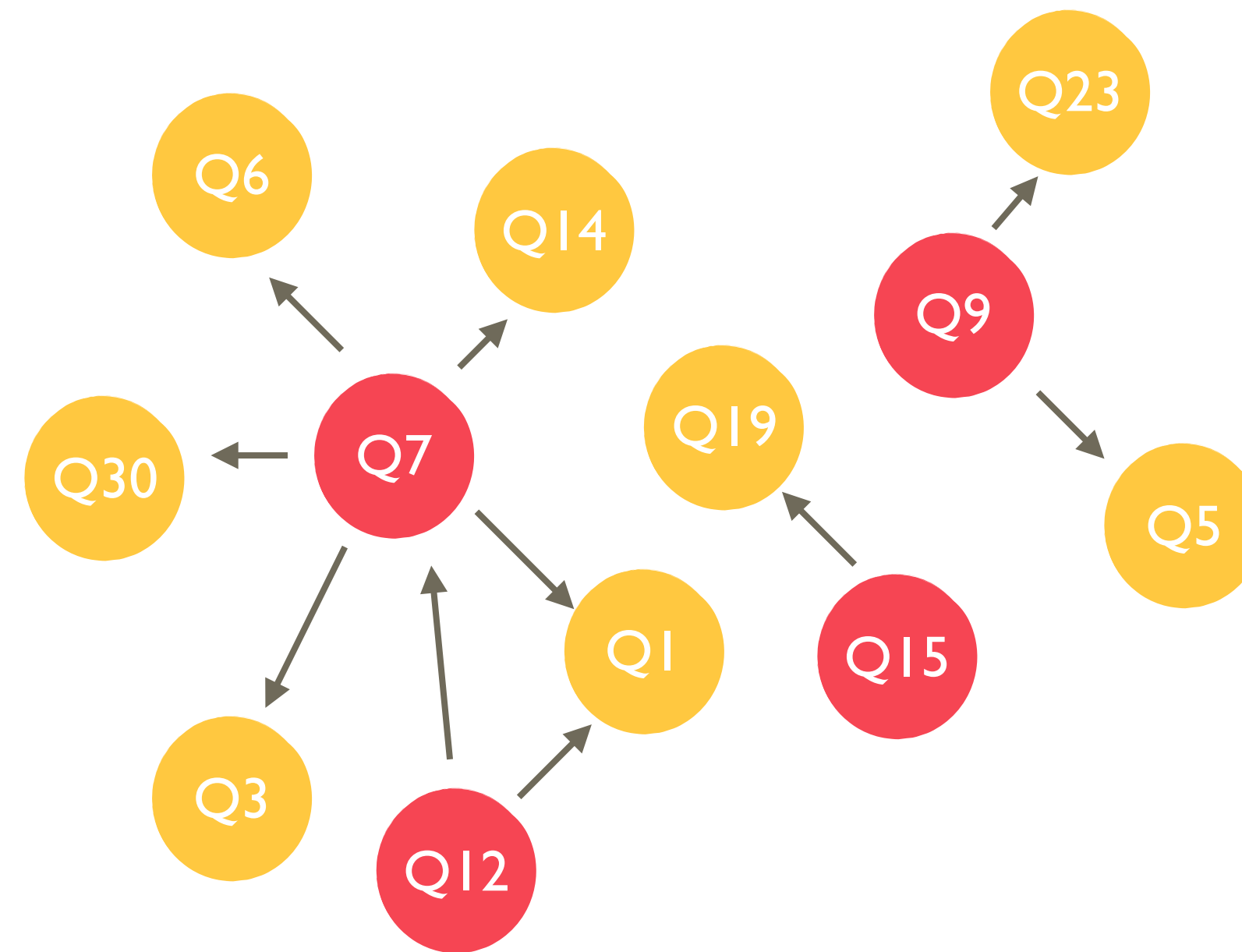
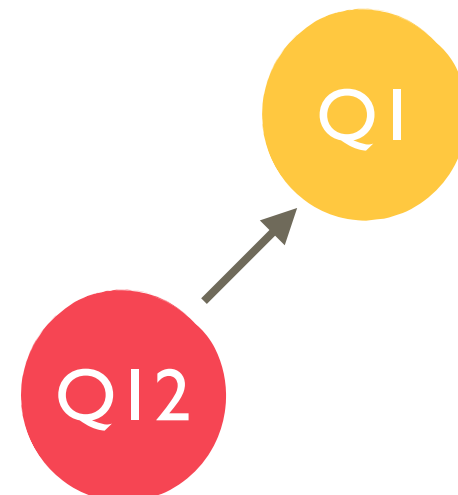


I. BUILDING A GRAPH

- Feature selection methods are applied



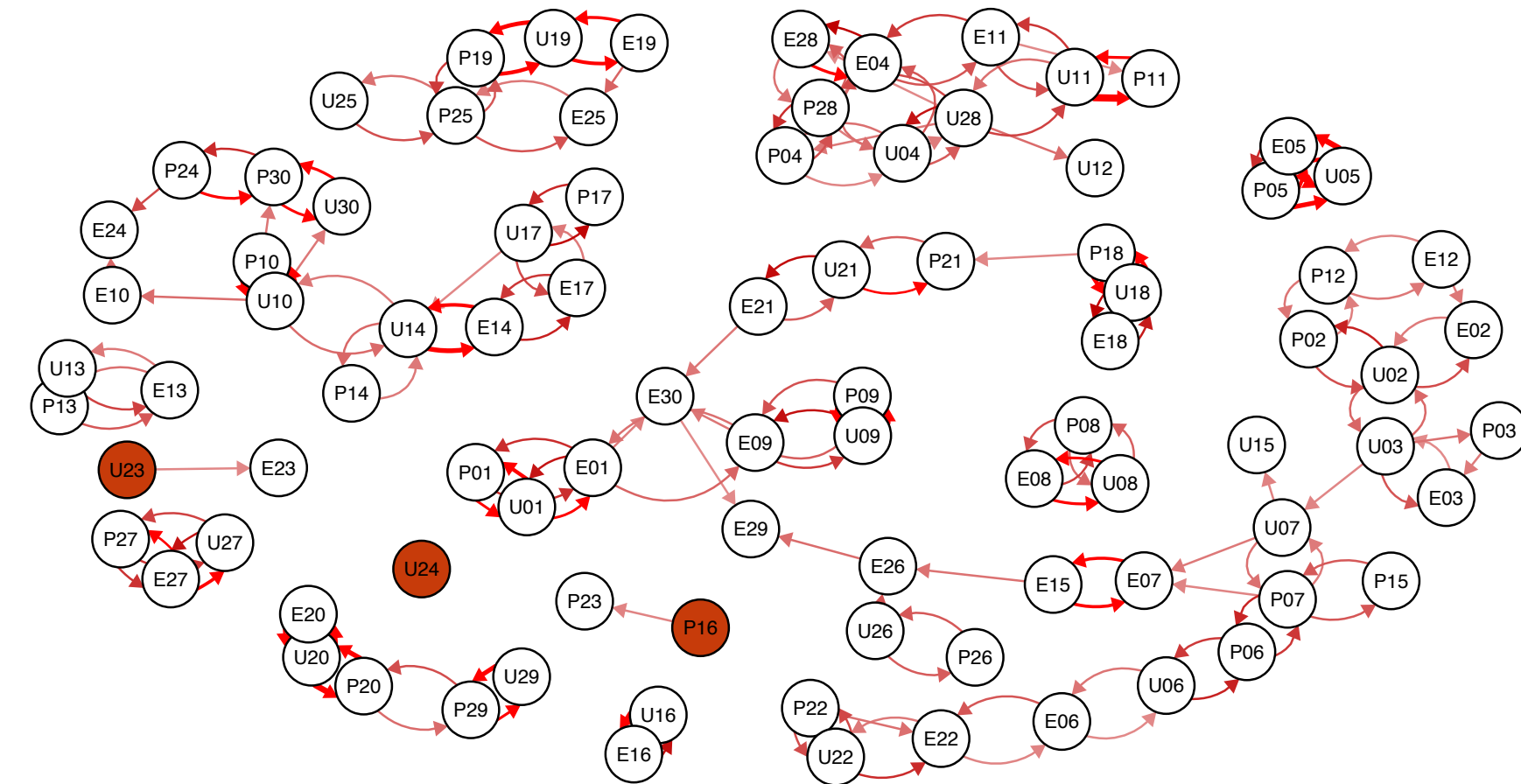
- For each question we have a set of predictors



- Together they form a graph

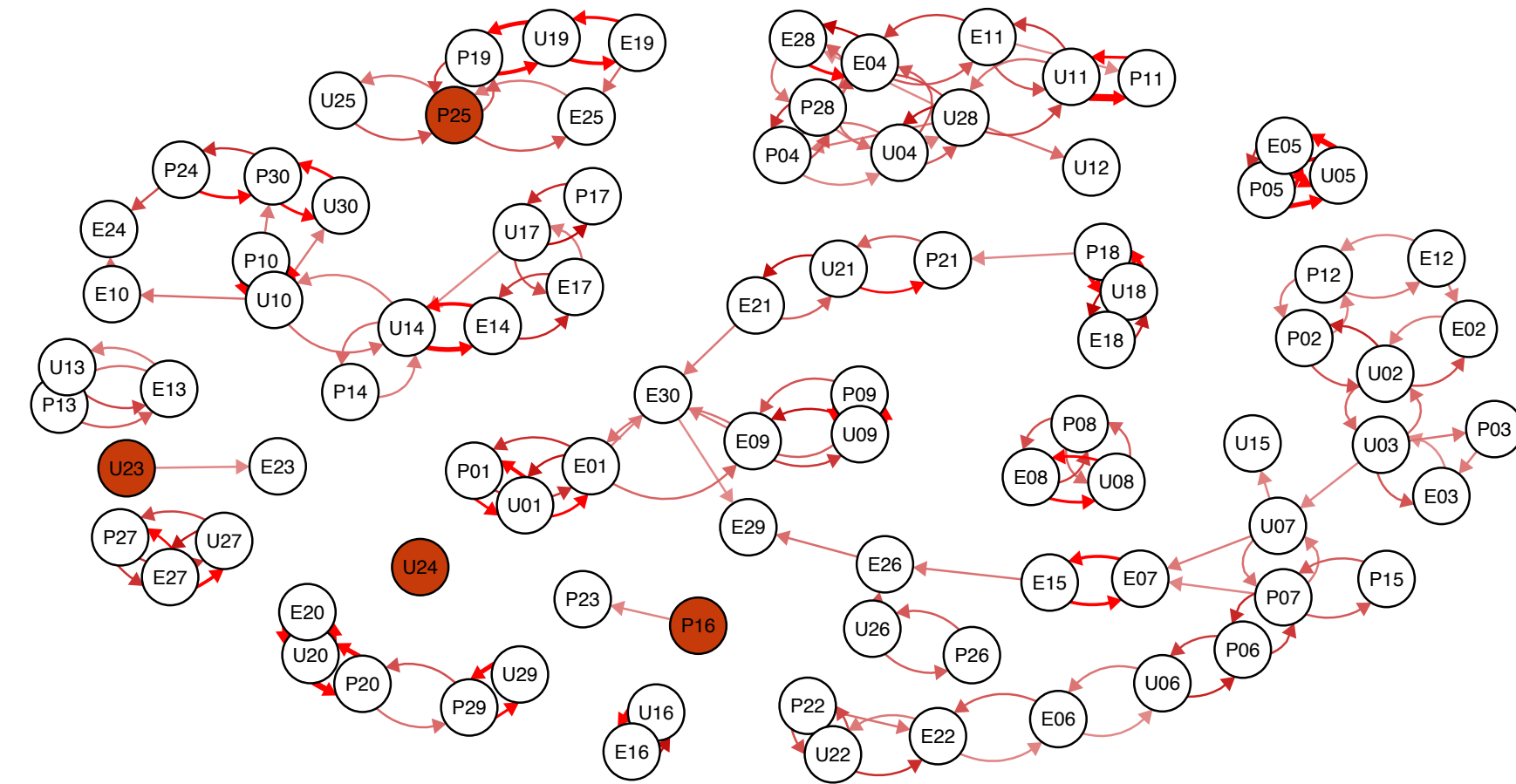
I. BUILDING A GRAPH

- Goal: find a 'minimal' subgraph with minimum number of predictors



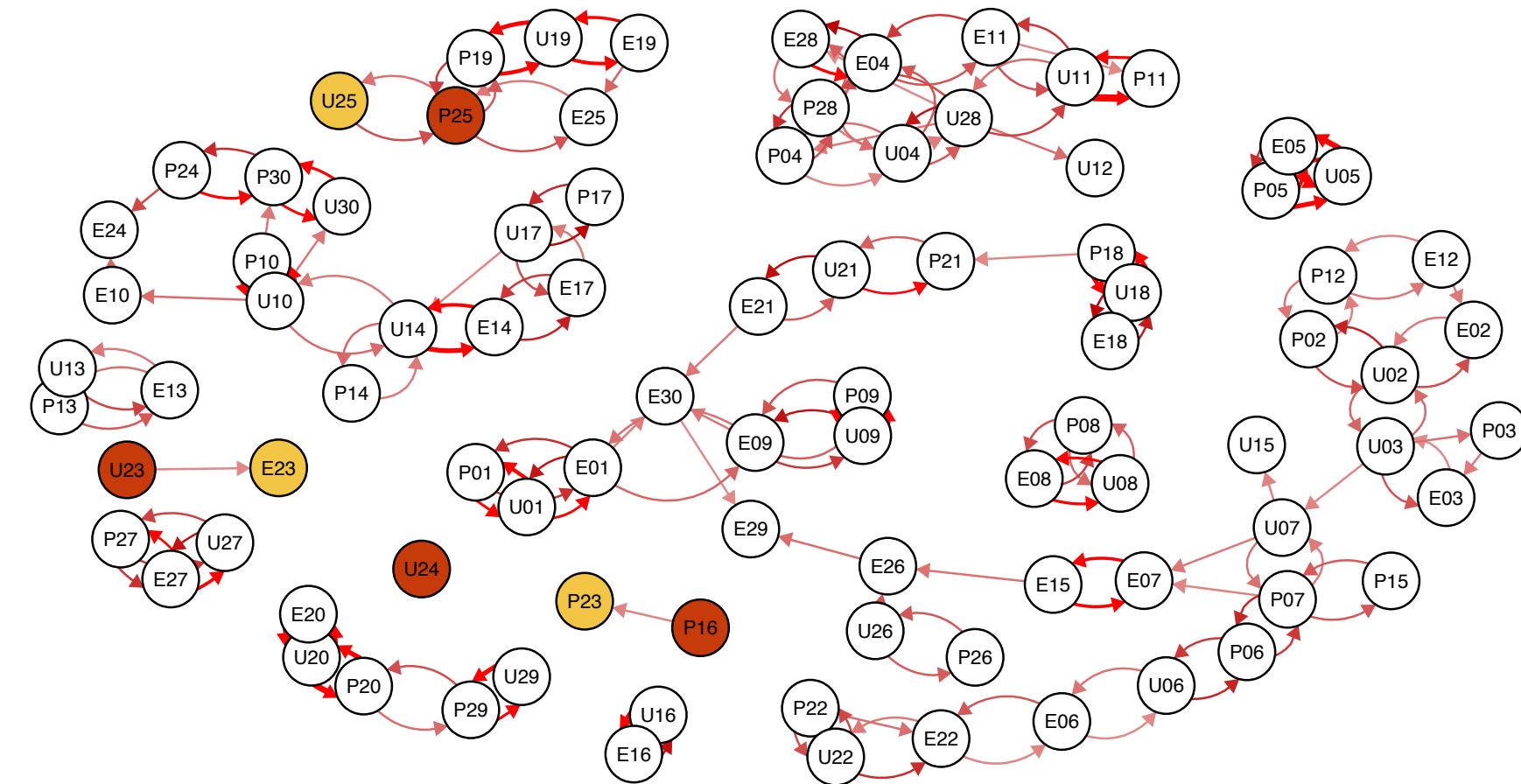
I. BUILDING A GRAPH

- Goal: find a 'minimal' subgraph with minimum number of predictors



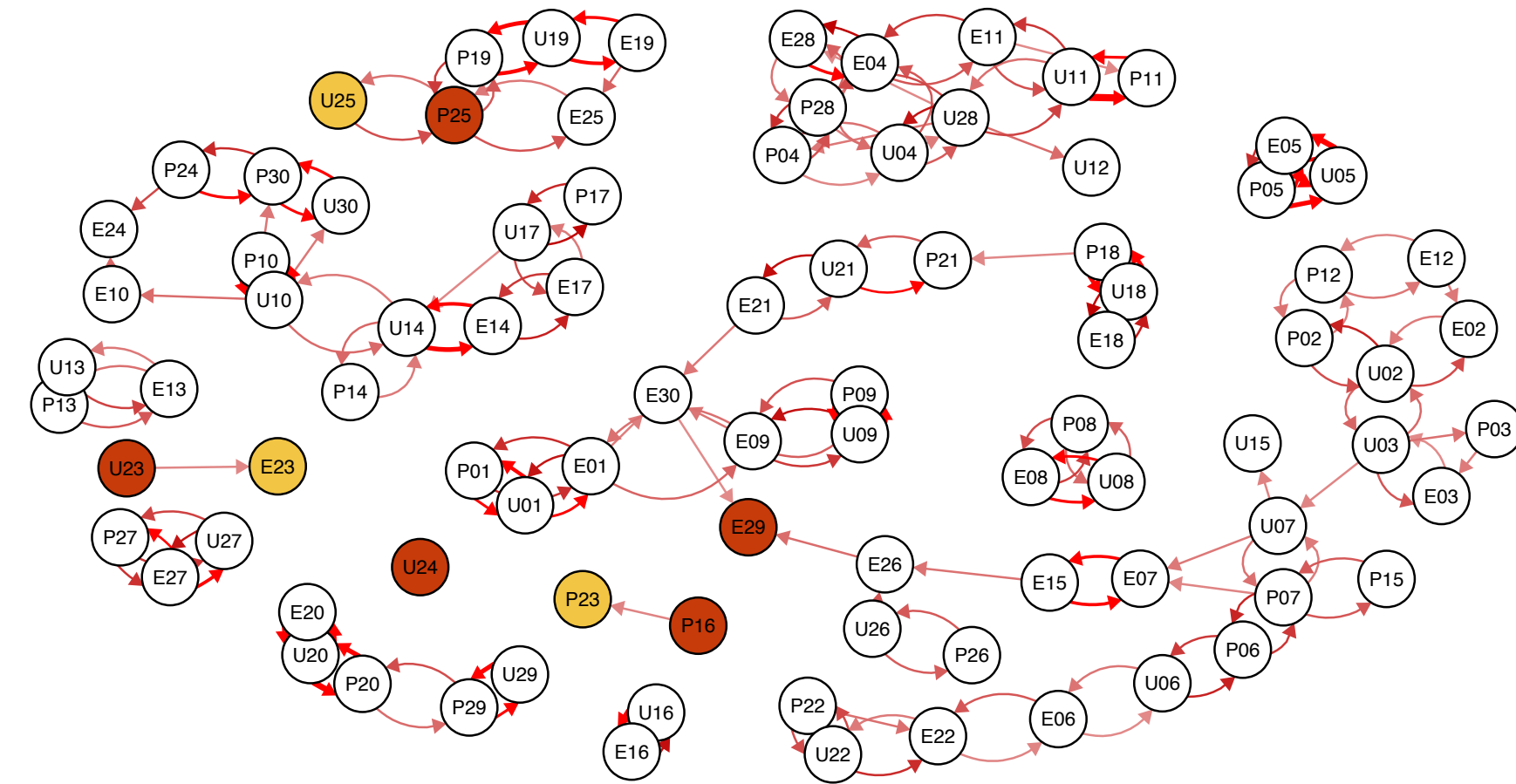
I. BUILDING A GRAPH

- Goal: find a 'minimal' subgraph with minimum number of predictors



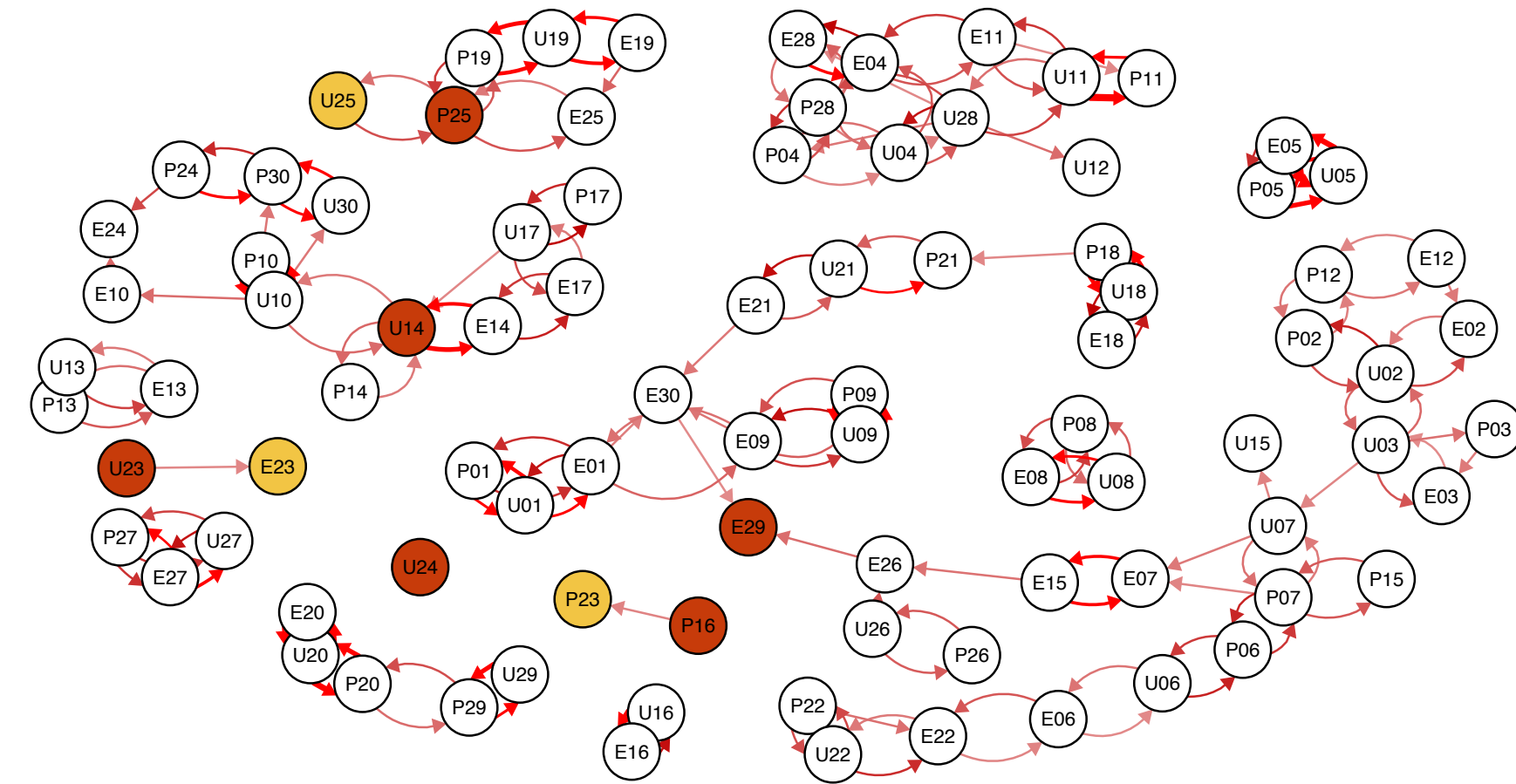
I. BUILDING A GRAPH

- Goal: find a 'minimal' subgraph with minimum number of predictors



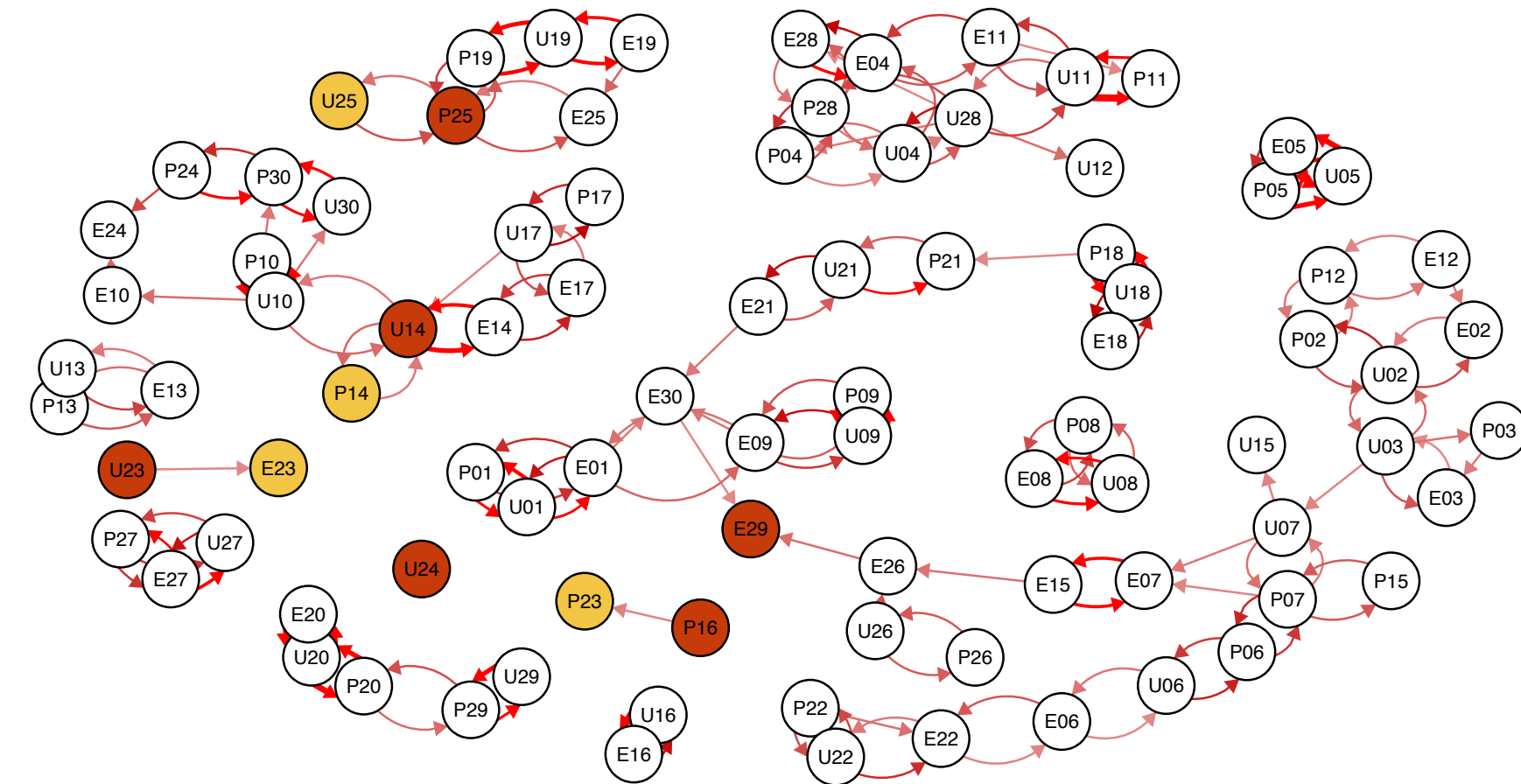
I. BUILDING A GRAPH

- Goal: find a 'minimal' subgraph with minimum number of predictors



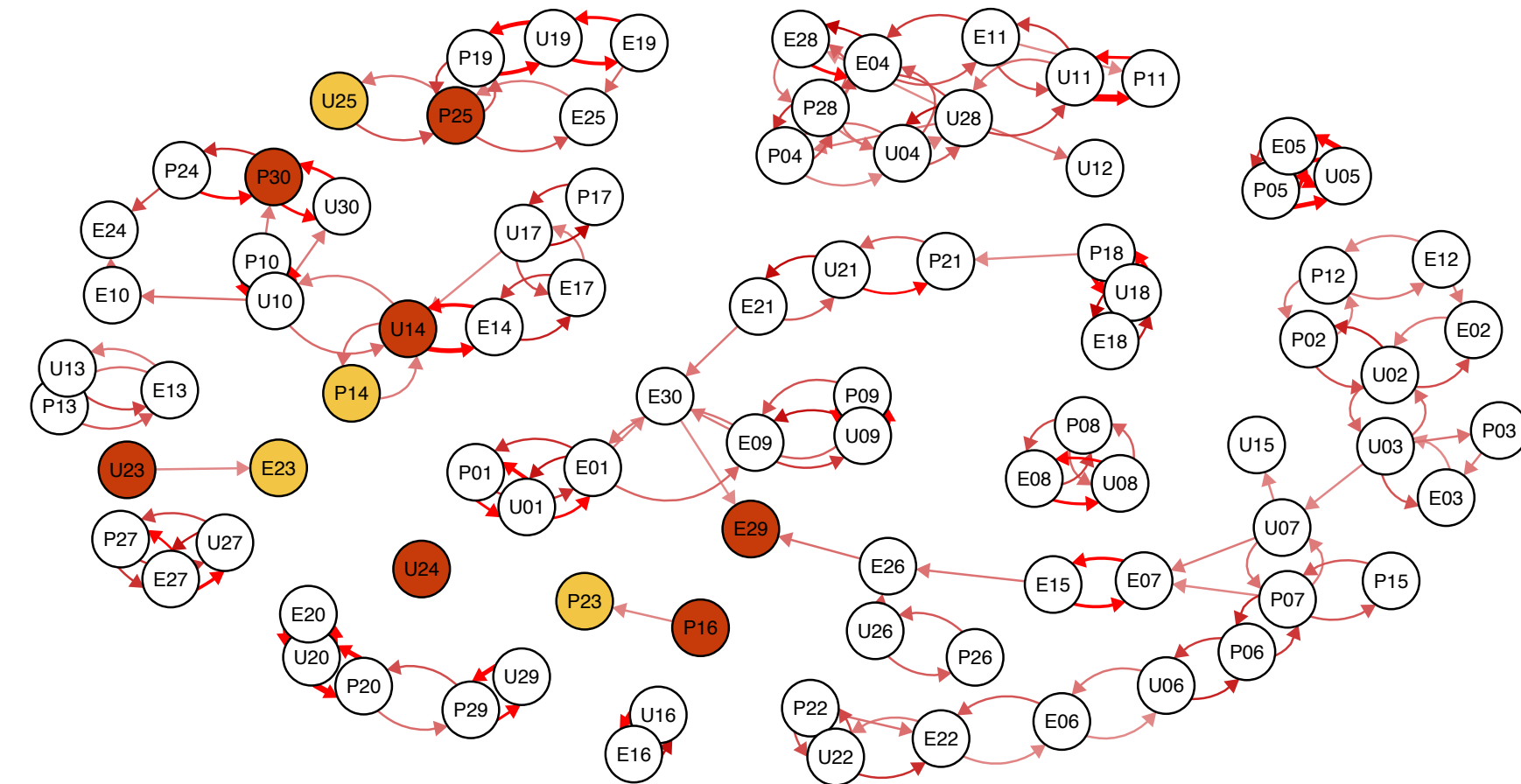
I. BUILDING A GRAPH

- Goal: find a 'minimal' subgraph with minimum number of predictors



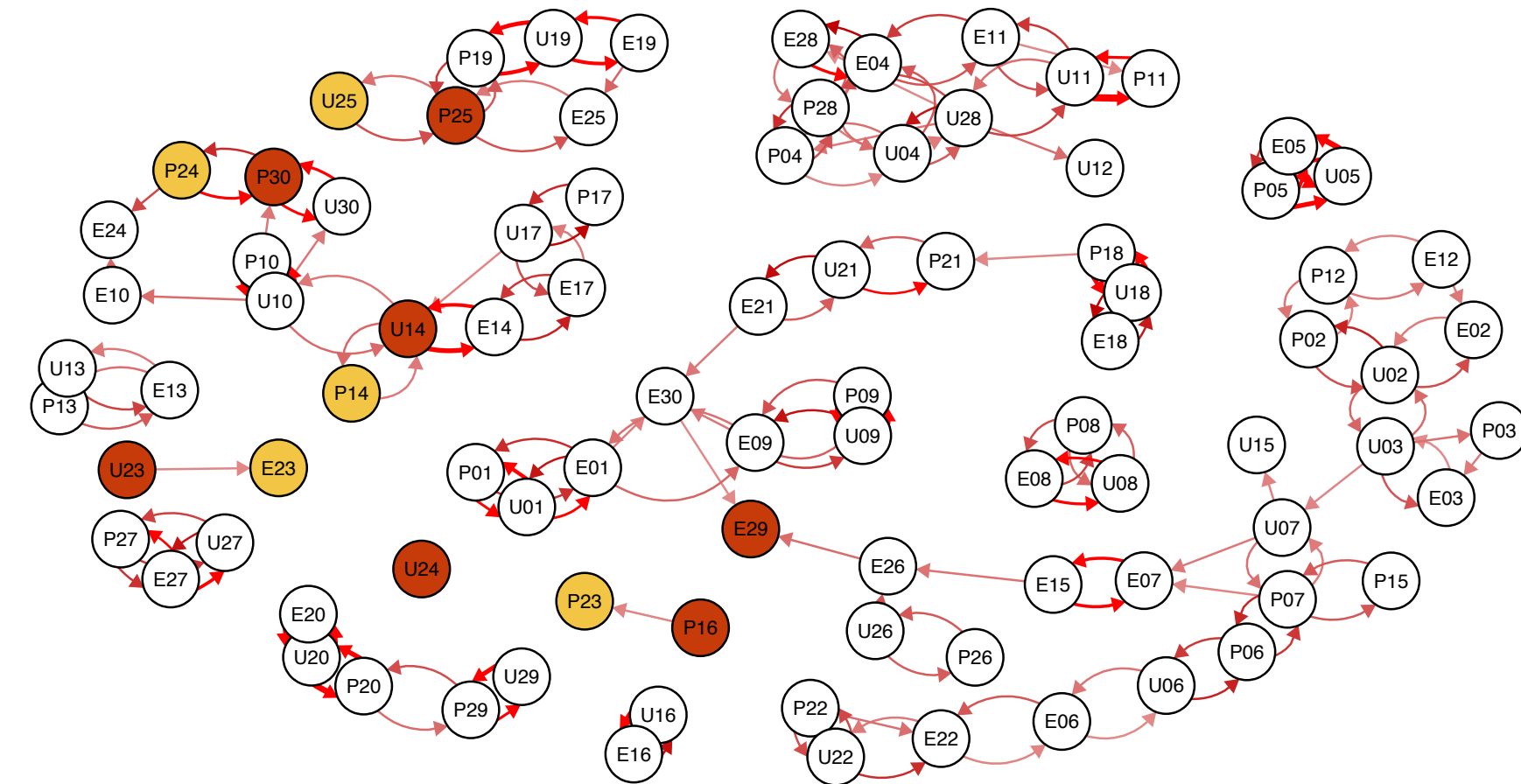
I. BUILDING A GRAPH

- Goal: find a 'minimal' subgraph with minimum number of predictors



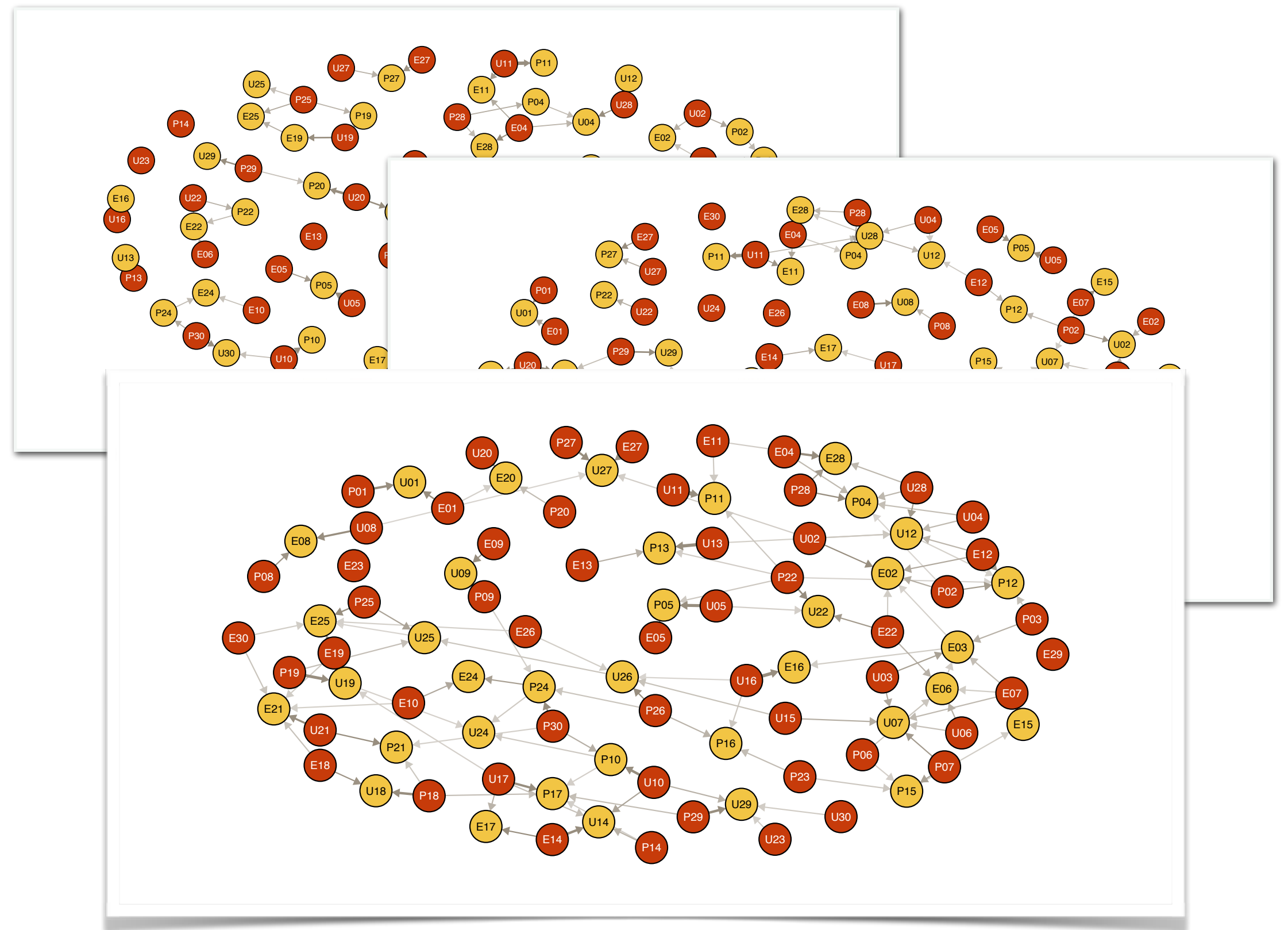
I. BUILDING A GRAPH

- Goal: find a 'minimal' subgraph with minimum number of predictors
- It can be formalised as a matrix problem



I. BUILDING A GRAPH

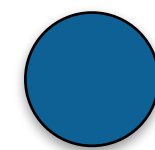
- For each type of feature selection method we can build one graph



I. FEATURE SELECTION - REVIEW

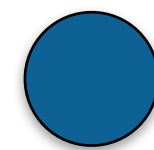
Graph construction

Method 1 (non-heuristic)



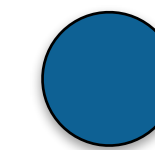
Choose model

1. Choose a model
2. Study the algorithm. Can it be adapted to the problem?
3. Set the parameters



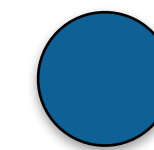
Separate data 70%-30%

- For each model:
4. Select training data (70%)
 5. Transform the data (1-7 or 0-1)
 6. Divide the problem in sub-problems



Feature selection

- For each sub-problem:
7. Perform cross-validation to avoid over-fitting
 8. Select features that minimise the error



Tree pruning

- Finally:
9. Construct a graph with all the questions
 10. Transform the graph into a tree
 11. Calculate the final number of questions to be included in test

Method 2



Heuristic graph construction

e.g. dropping "Performance" questions



Random graph construction

e.g. choosing 10 random questions

2. MULTI-CLASSIFICATION SUBPROBLEM

- Different types of models:
- Random model
- Generalised Linear models
- Random forest
- Support Vector Machines (linear kernel)
- Some basic Neural Networks

2. REVIEW

Model construction

- **Choose tree**
 - Non-heuristic tree
 - Heuristic tree
 - Random tree

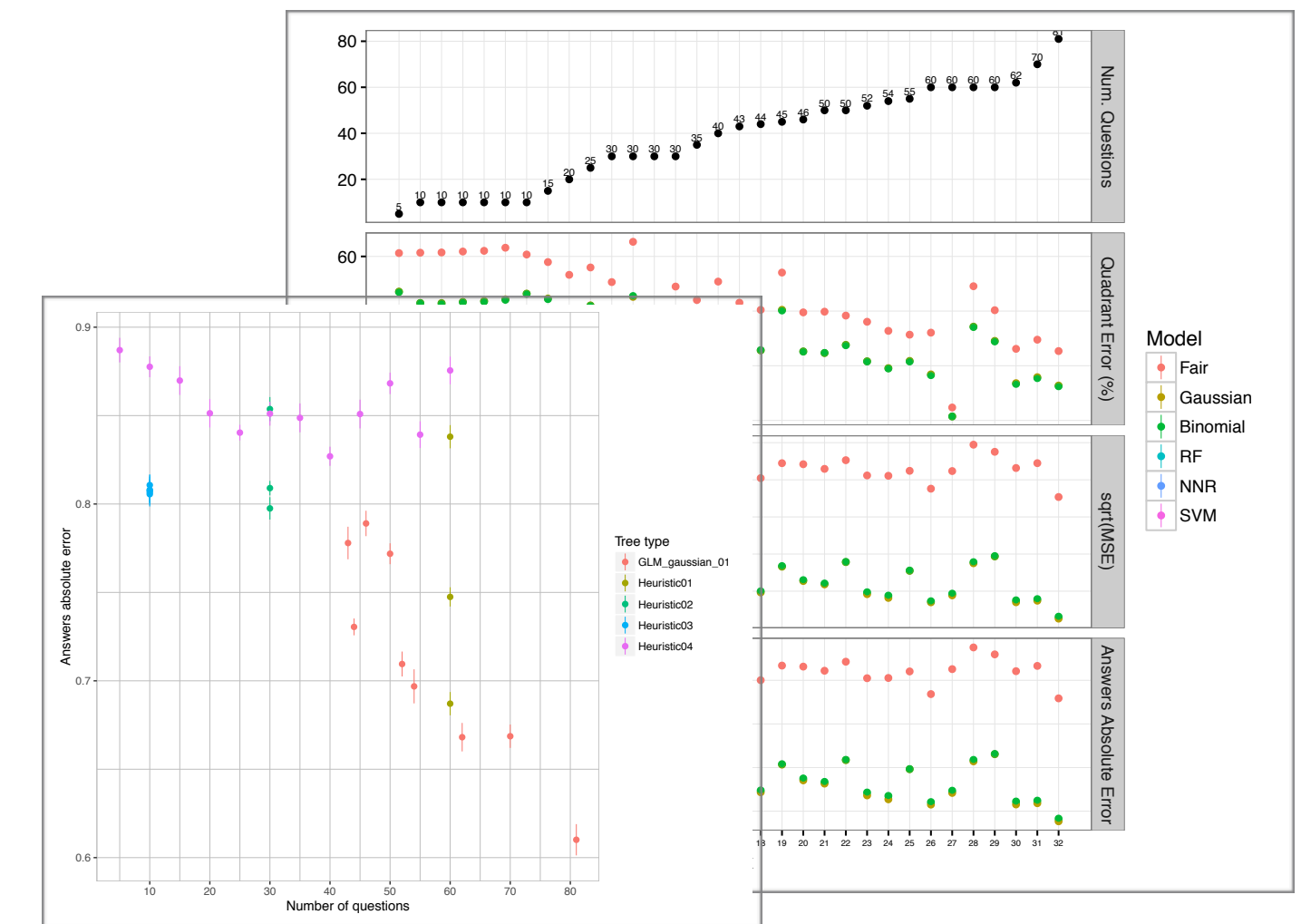
- **Choose model**
 - Fair model
 - Gaussian linear model
 - Binomial model
 - Random forest
 - Support v. machine

● Perform Monte Carlo Cross-Validation

- For each loop:
1. Simulate the results of the test
 - Ask N questions to the user
 - For each of the remaining 90-N questions:
 - train model (using 70%)
 - predict answer (using 30%)
 - calculate error

- Finally:
2. Calculate the average error:
 - *Answers absolute error*

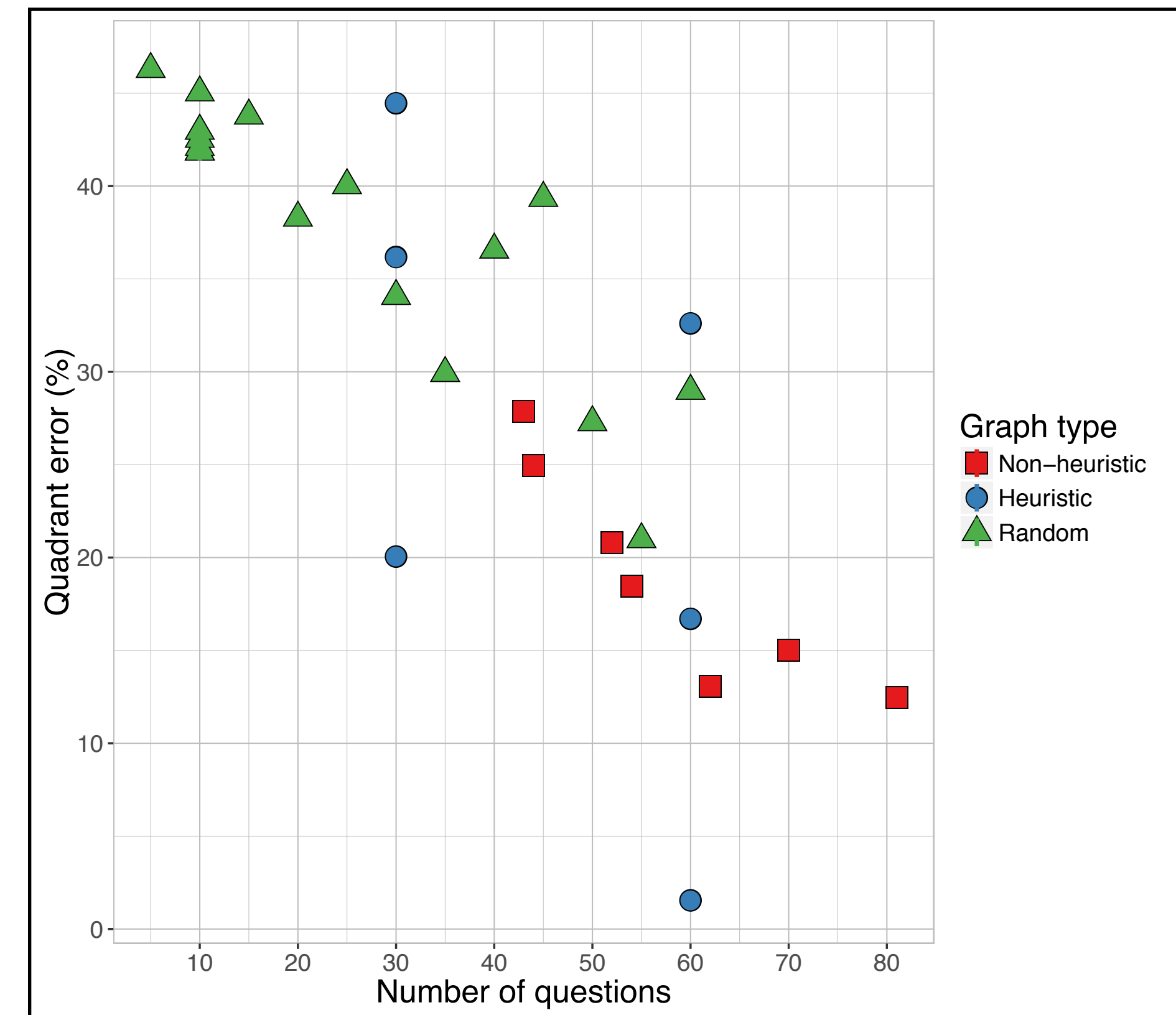
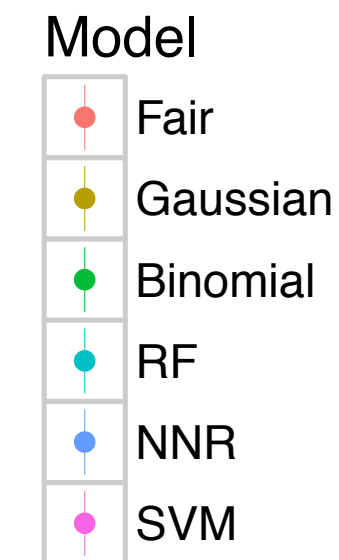
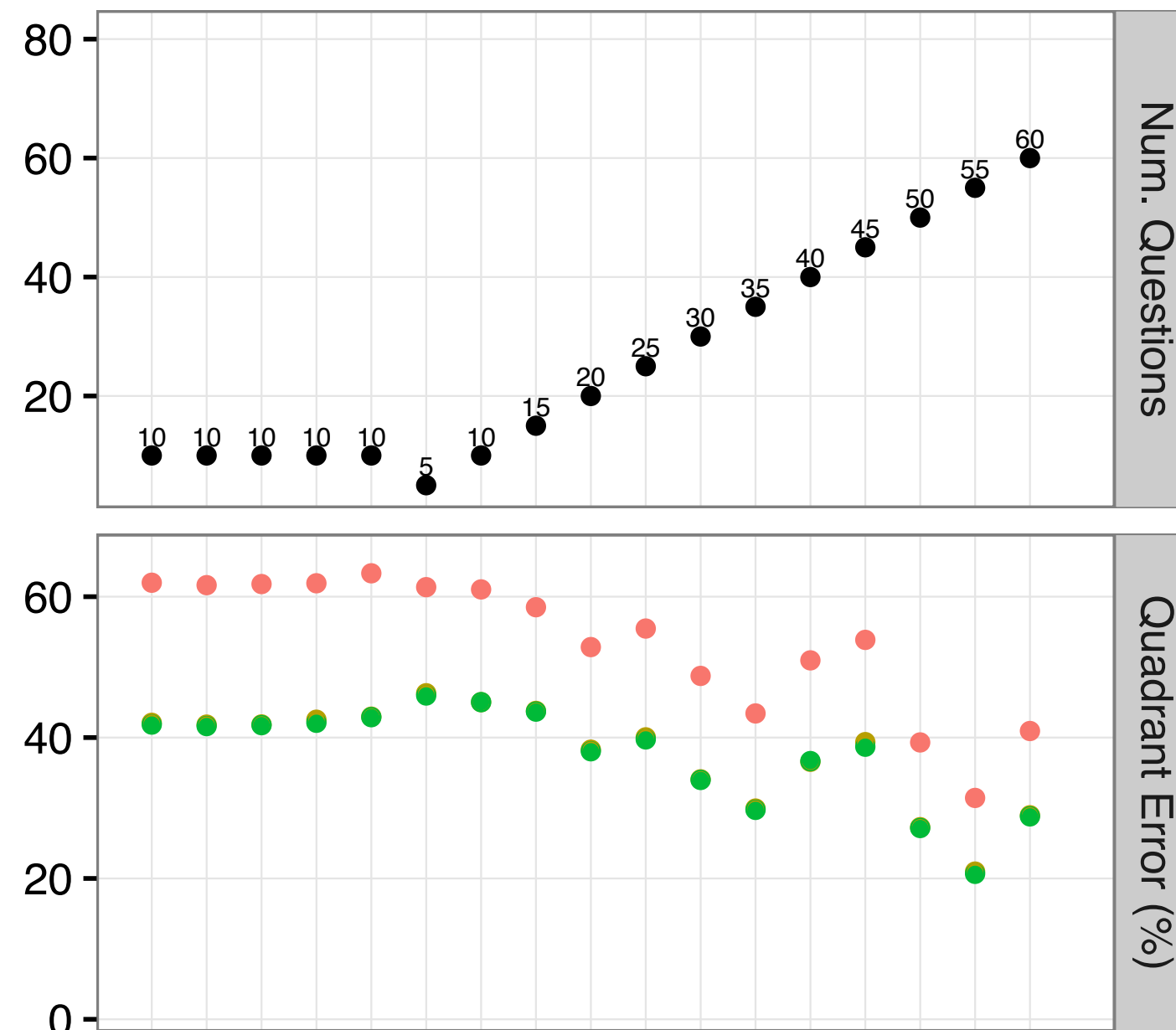
● Results



2. RESULTS

■ Error vs. number of predictors

- Why are they performing in a similar way?
- Data is too noisy?



2. RESULTS

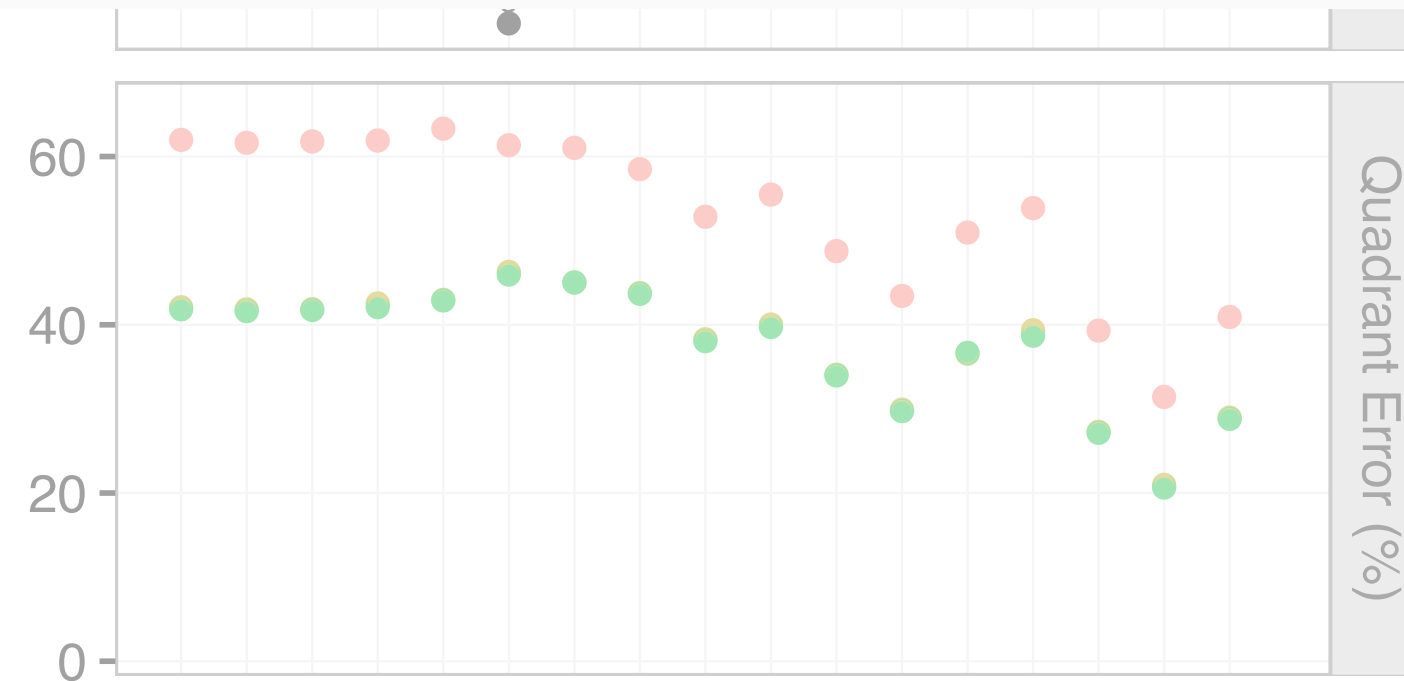
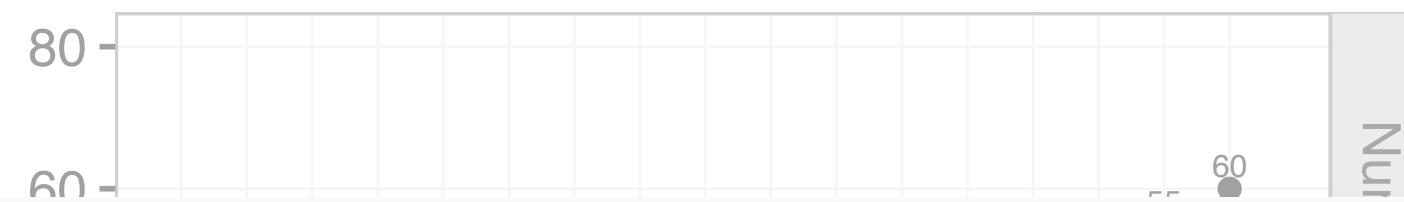
- Error vs. number



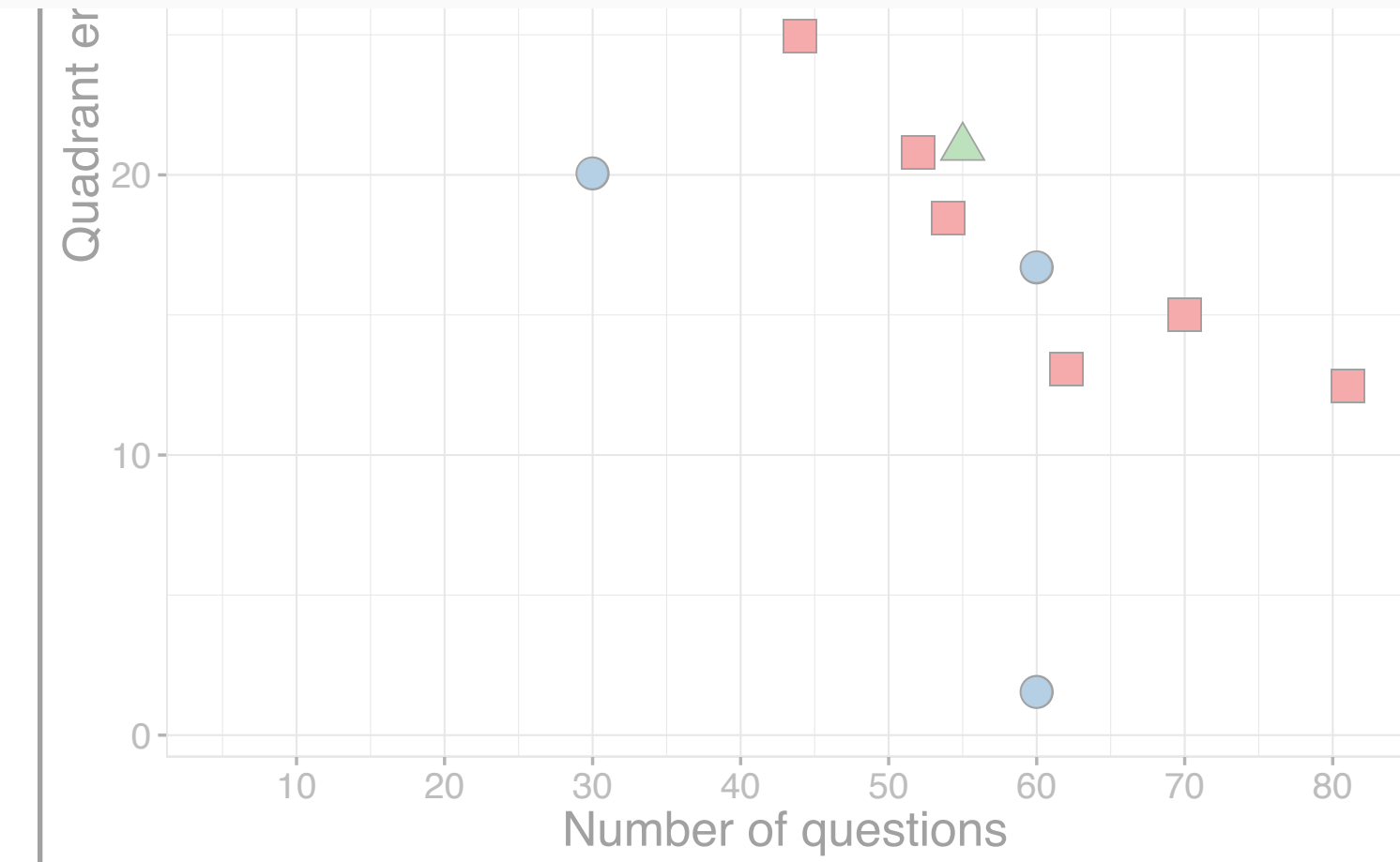
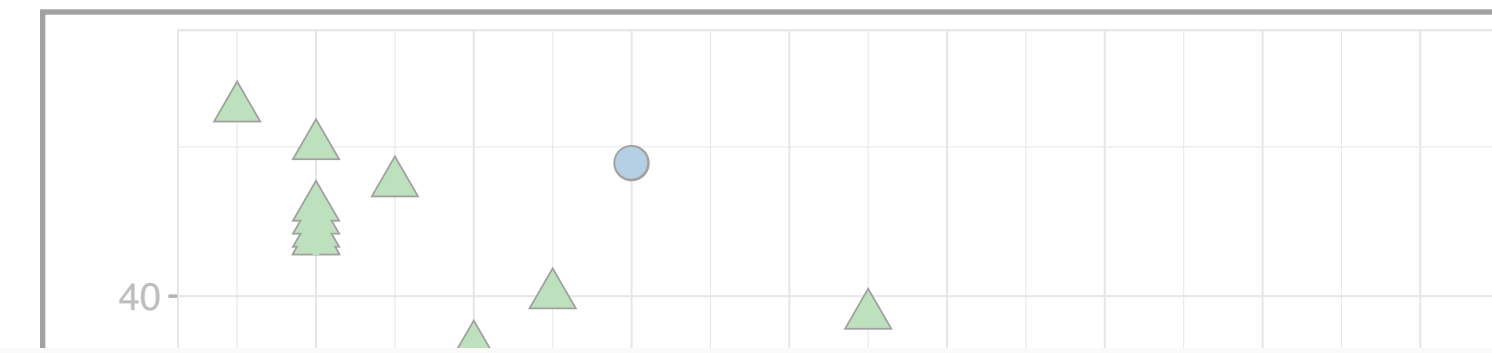
is my problem solvable?



- Why are they performing in a similar way?
- Data is too noisy?



- Binomial
- RF
- NNR
- SVM



- Graph type
- Non-heuristic
 - Heuristic
 - Random

OTHER APPROACHES

- Information theory approach (no ML but useful)
- We want to find **redundancy** in data
- What is the *most redundant* set of questions

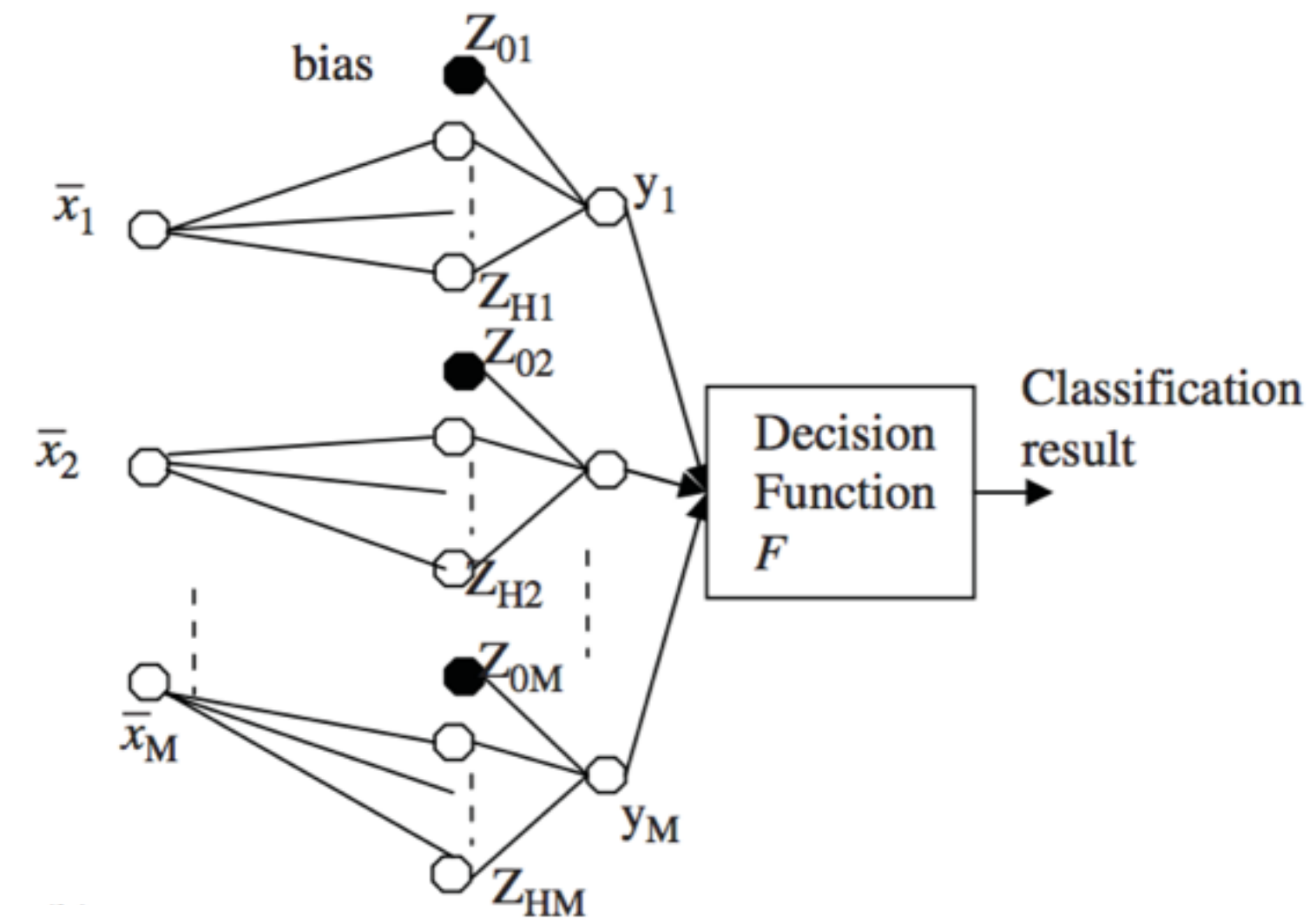
- Entropy!

$$H(X) = - \sum_{i=1}^n p(x_i) \log p(x_i).$$

- $H(X) = 0$. High redundancy
- $H(X)$ max. No redundancy

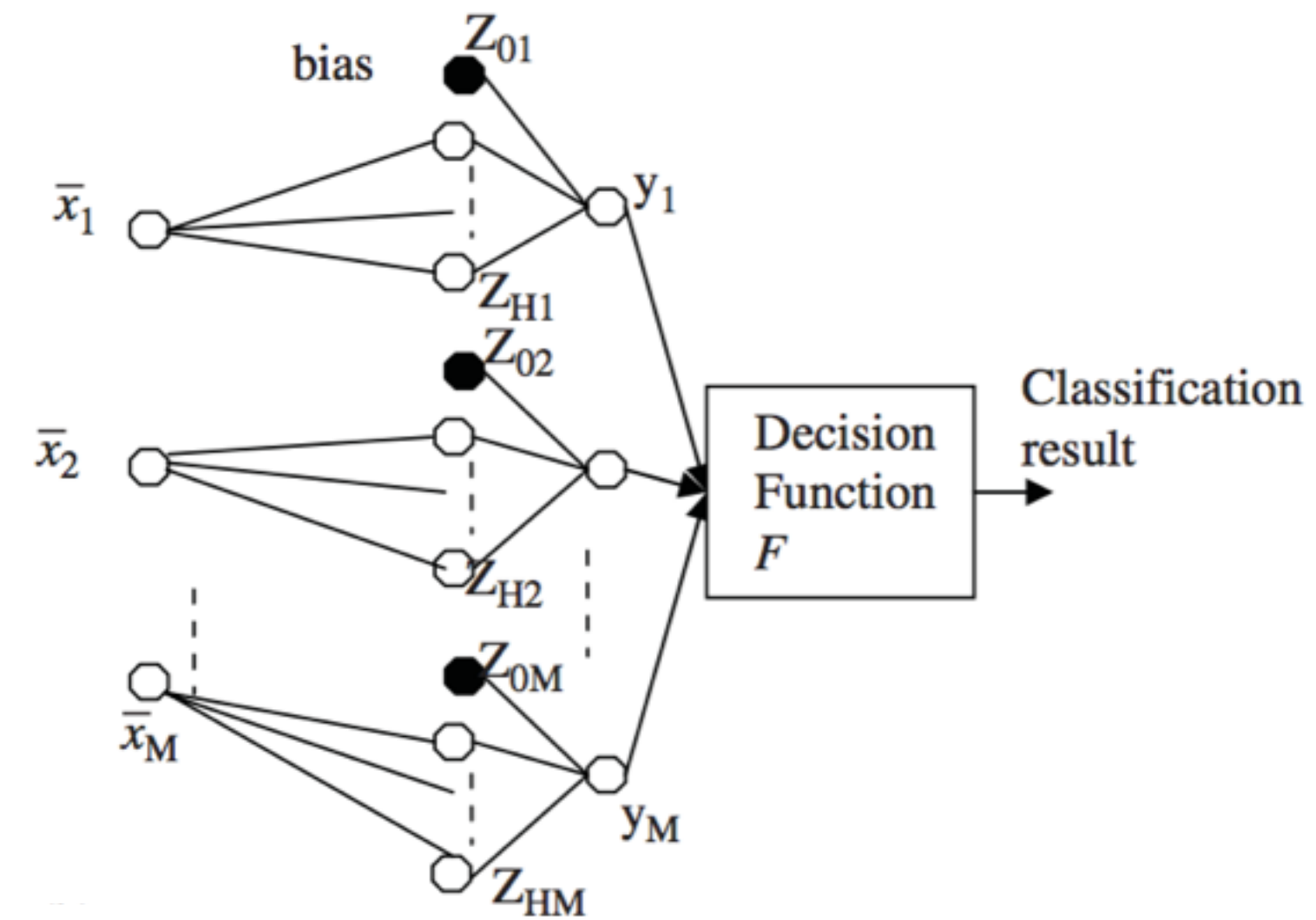
NEURAL NETWORKS

- Again, space $Q = P \cup S$
- Fix s in S
- Classification problem with 4 classes
- We want a 0/1 output
- I. 4 NN, one output node for each class
- II. 1 NN, four output nodes



NEURAL NETWORKS

- Again, space $Q = P \cup S$
- Fix s in S
- Classification problem with 4 classes
- We want a 0/1 output
- I. 4 NN, one output node for each class
- II. 1 NN, four output nodes



Problem with unbalanced classes:

zoom-in regions with lower numbers of points
(discrete space)

COMMENTS

The problem is still open

Although, there are other steps omitted here that have been useful

If it cannot be solved, how to prove it?

REVIEW OF THE PROBLEM

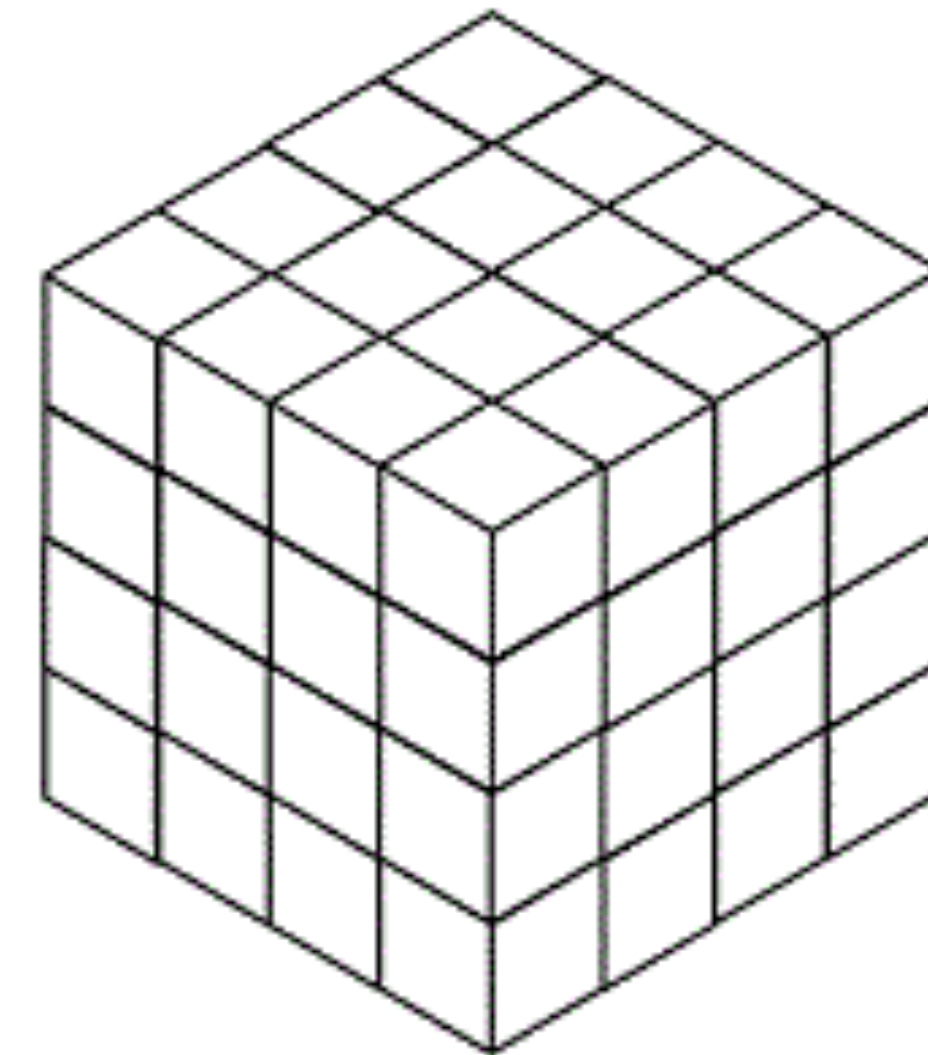
- Dataset: collection of user responses (~30000)
- In our case the test has 90 questions with 4 possible answers:

Drawback: sparse data.

Drawback: data concentrated in few cells

Drawback: users perceive the scale in different way

Drawback: users tend to choose high values



$$Q = PUS$$

CONCLUSIONS

Psychological test are relents in industry

Special type of data

How to extend traditional DM techniques to deal with these type of data?

SOME REFERENCES

- Many interesting ones...

[1] (Review FS) Saeys, Y. et al, A review of feature selection techniques in bioinformatics. *Bioinformatics*, 2007.

[2] (Normalised MI) Estévez, P. et al, Normalized Mutual Information Feature Selection. *IEEE Transactions On Neural Networks*, 2009.

[3] (Multiclass. NN)

[4] (Unbalanced data)

[5] (high-dimensional NN)

Thank you! Questions?
