

Using Orange Data Mining and Predictive Analytics for Analytical Workflows in Academic Research

Raabia Riaz^{1*}, Muhammad Areeb Chatni², Tanzeel ur Rehman³, Mehsam Bin Tahir⁴, Sherbaz Khan⁵

¹Lead Analyst, Nauf Networks, United Kingdom

²Kingston University, London

³University of Gloucestershire, England

⁴Middlesex University, London

⁵Jinnah University for Women, Karachi, Pakistan

Abstract: This section examines methodological transparency and interpretability in computational social science through the use of Orange Data Mining. As analytical workflows are often distributed across multiple software environments, key preprocessing and modelling decisions can become difficult to trace, limiting reproducibility and interpretability. The section argues that transparent analytical design is essential, particularly in fields where explanation is as important as prediction. Orange is presented as a visual, open source platform that makes each stage of analysis explicit. By constructing workflows through connected modules, the platform allows researchers to inspect data transformations, parameter settings and evaluation procedures within a single environment. The argument is demonstrated using the Zoo dataset included in Orange. Through a structured sequence of exploratory analysis, dimensionality reduction, interpretable classification and cross validated evaluation, the section illustrates a complete analytical cycle. The example shows how visual workflow design can support methodological clarity and accountable research practice.

Keywords: Data Mining, AI, Orange, Academic analysis.

1. Introduction

Since the early 2000s, the digital ecosystems have become integral avenues to social, political and institutional life. The online public debate increasingly unfolds through digital platforms; behavioural and survey data are now routinely stored in structured digital formats. The researchers across the social sciences, humanities and applied policy fields regularly engage with datasets of categorical indicators, coded variables, behavioural attributes.

The availability of this kind of data has created new opportunities beyond basic measurement. At the same time, it has increased concerns about transparency, interpretability, and reproducibility. Analytical work is often spread across different software tools. Data might be collected in one system, cleaned or transformed in another, analysed in a third, and visualised in a fourth. At each stage, decisions are made about filtering, coding, weighting, and validation. These choices affect which

patterns appear and how the results are understood. When these steps are hidden in scripts or spread across different platforms, it becomes difficult to clearly trace and repeat the full analytical process.

In computational social science and digital humanities this fragmentation has renewed concern about methodological clarity. Lazer et al. (2020) note that computational methods offer strong analytical capacity but demand procedural transparency. Kapoor and Narayanan (2023) show that data leakage and undocumented preprocessing weaken reproducibility in machine learning research. When intermediate transformations are not documented, evaluation depends on technical skill rather than methodological reporting.

Interpretability is a further issue. Machine learning models can achieve high predictive accuracy, yet complex ensemble models and deep neural networks often operate as black boxes. In fields such as education public health and social policy, explanation is necessary. Researchers must explain why a model produces a given output and how it relates to theory. Rudin (2019) argues that interpretable models should be prioritised where explanation is required.

This section addresses these issues through a single analytical environment, Orange Data Mining. Orange is an open-source visual analytics platform in which analysis is built using connected modules called widgets. The workflow is constructed visually on a canvas. Each transformation parameter and evaluation step remains visible and open to inspection.

The argument is direct. Orange provides a transparent framework for exploratory and interpretable analysis of structured datasets. Its contribution lies in making analytical reasoning explicit. By integrating exploration modelling inspection and validation within one environment, it supports structured workflow design and methodological accountability.

The section uses the Zoo dataset included in Orange. The dataset contains 101 animal instances described by Boolean attributes such as hair feathers eggs milk and aquatic behaviour,

*Corresponding author: princydude@gmail.com

alongside a categorical class variable representing animal type. The dataset is used as a controlled example to demonstrate method without domain complexity.

The analysis follows a defined sequence: data familiarisation descriptive analysis dimensionality reduction interpretable classification and cross validated evaluation. These stages are presented in six figures and together represent a complete analytical cycle. Although the dataset is simple, the workflow generalises to applied research contexts including education media studies public health and political science.

The section adopts an exploratory and descriptive modelling approach rather than confirmatory hypothesis testing or predictive deployment. Exploratory modelling identifies patterns and structural relationships not specified in advance and produces insights that require theoretical interpretation. This position aligns with Grimmer and Stewart (2013), who describe computational analysis as a tool for structuring inquiry rather than replacing substantive reasoning.

Within this scope, Orange provides a clear environment for analytical work. Its visual structure makes variable designation preprocessing decisions and model evaluation explicit. Modelling is therefore presented as a sequence of reasoned methodological decisions rather than an automated process.

The following sections demonstrate this workflow in practice.

A. Analytical Workflow as Research Logic

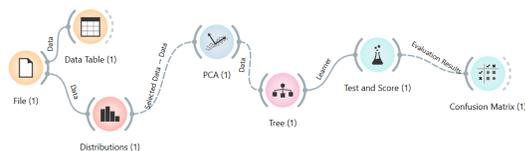


Fig. 1. Analytical workflow in Orange

File → Data Table → Distributions → PCA → Tree → Test & Score → Confusion Matrix

Although this figure appears technical, it represents research logic rather than software procedure. Each stage corresponds to a conceptual step in empirical reasoning. The workflow is not a sequence of buttons; it is a structured analytical argument.

The first stage establishes data integrity. Before modelling can occur, the researcher must understand the structure of the dataset. This includes confirming the number of observations, verifying variable types and designating the correct target variable.

The second stage situates the dataset descriptively. Distributional characteristics shape expectations about classification performance and interpretation.

The third stage introduces exploratory structural analysis. Dimensionality reduction techniques such as Principal Component Analysis allow researchers to visualise patterns of similarity across instances without imposing categorical decisions.

The fourth stage formalises those patterns through supervised classification. The decision tree translates structural relationships into explicit rule-based splits.

The final stages evaluate generalisation. Cross-validation tests model stability across subsets of the data. The confusion matrix provides diagnostic insight into specific patterns of misclassification.

This sequence reflects a disciplined progression from observation to interpretation. It avoids premature modelling and emphasises validation before inference. By representing this sequence visually, Orange externalises methodological reasoning and reduces the likelihood that analytical steps remain implicit.

2. Data Familiarisation and Structural Verification

The workflow begins by loading the Zoo dataset using the File widget and connecting it to the Data Table widget.

| | type | name | hair | feathers | eggs | milk | airborne | aquatic | predator |
|----|---------|----------|------|----------|------|------|----------|---------|----------|
| 63 | reptile | pilviper | 0 | 0 | 1 | 0 | 0 | 0 | 1 |
| 77 | reptile | seasnake | 0 | 0 | 0 | 0 | 0 | 1 | 1 |
| 81 | reptile | slowworm | 0 | 0 | 1 | 0 | 0 | 0 | 1 |
| 91 | reptile | tortoise | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| 92 | reptile | tuatara | 0 | 0 | 1 | 0 | 0 | 0 | 1 |
| 1 | mammal | sardark | 1 | 0 | 0 | 1 | 0 | 0 | 1 |
| 2 | mammal | antelope | 1 | 0 | 0 | 1 | 0 | 0 | 0 |
| 4 | mammal | bear | 1 | 0 | 0 | 1 | 0 | 0 | 1 |
| 5 | mammal | boar | 1 | 0 | 0 | 1 | 0 | 0 | 1 |
| 6 | mammal | buffalo | 1 | 0 | 0 | 1 | 0 | 0 | 0 |
| 7 | mammal | bull | 1 | 0 | 0 | 1 | 0 | 0 | 0 |
| 10 | mammal | cow | 1 | 0 | 0 | 1 | 0 | 0 | 0 |
| 11 | mammal | cheetah | 1 | 0 | 0 | 1 | 0 | 0 | 1 |
| 18 | mammal | deer | 1 | 0 | 0 | 1 | 0 | 0 | 0 |
| 20 | mammal | dolphin | 0 | 0 | 0 | 1 | 0 | 1 | 1 |
| 22 | mammal | elephant | 1 | 0 | 0 | 1 | 0 | 0 | 0 |
| 28 | mammal | frigate | 1 | 0 | 0 | 1 | 1 | 0 | 0 |
| 29 | mammal | giraffe | 1 | 0 | 0 | 1 | 0 | 0 | 0 |
| 30 | mammal | girl | 1 | 0 | 0 | 1 | 0 | 0 | 1 |
| 32 | mammal | goat | 1 | 0 | 0 | 1 | 0 | 0 | 0 |
| 33 | mammal | gorilla | 1 | 0 | 0 | 1 | 0 | 0 | 0 |
| 36 | mammal | hamster | 1 | 0 | 0 | 1 | 0 | 0 | 0 |
| 37 | mammal | hare | 1 | 0 | 0 | 1 | 0 | 0 | 0 |

Fig. 2. Data table view of the zoo dataset

The Zoo dataset contains 101 animal instances described by Boolean variables indicating attributes such as hair, feathers, eggs, milk, airborne ability and aquatic behaviour. A categorical variable labelled “type” identifies the animal class.

The Data Table stage serves a methodological function that extends beyond simple viewing. The researcher must verify that:

- The dataset contains the expected number of observations.
- Predictor variables are correctly typed as categorical or Boolean.
- The class variable “type” is explicitly designated as the target.

In Orange, the target variable is clearly marked. This designation is not procedural detail. In classification analysis, the target variable determines what the model attempts to predict. Misidentifying the target fundamentally alters the analytical task. By making variable roles explicit on the interface, Orange encourages deliberate specification rather than implicit assumption.

In applied research contexts, this stage corresponds to verifying coding schemes in survey data, confirming that categorical labels are consistent across observations or checking for missing values. Errors at this stage propagate throughout the workflow. Careful structural verification is therefore foundational rather than preliminary.

Data familiarisation also reinforces interpretive awareness. Seeing the dataset in tabular form reminds the researcher that

modelling operates on concrete observations. Even when later stages abstract these observations into components or splits, the empirical basis remains visible.

3. Descriptive Context and Distributional Awareness

Following structural verification, the workflow moves to descriptive contextualisation using the Distributions widget.

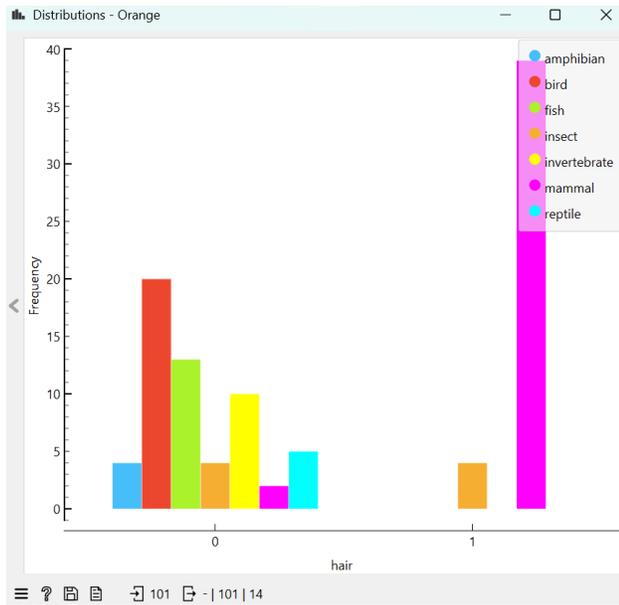


Fig. 3. Distribution of animal classes

The bar chart reveals how many instances belong to each animal category. Distributional awareness plays a crucial role in classification analysis.

First, it contextualises evaluation metrics. If one class overwhelmingly dominates the dataset, a model predicting that class exclusively could achieve deceptively high accuracy. This phenomenon, known as class imbalance, can distort interpretation of performance metrics. Accuracy alone is therefore insufficient without knowledge of distribution.

Second, distribution shapes theoretical interpretation. Rare categories may carry substantive significance. Misclassification of rare categories may be more analytically consequential than misclassification of common ones.

In the Zoo dataset, the distribution is uneven but not severely imbalanced. This informs expectations regarding classification results and reduces the likelihood that accuracy will be artificially inflated by dominance of a single class.

Descriptive analysis is not merely preparatory. It situates modelling within empirical context and establishes the structural conditions under which classification will operate.

4. Transition to Structural Exploration

With structural verification and descriptive grounding established, the workflow proceeds to structural exploration. At this stage, the objective is not prediction but pattern identification. The researcher seeks to understand how observations relate to one another in multidimensional space before imposing categorical boundaries.

The next section introduces Principal Component Analysis as a method for visualising similarity structure.

A. Structural Exploration Through Principal Component Analysis

After establishing structural integrity and descriptive context, the workflow proceeds to exploratory analysis through Principal Component Analysis (PCA).

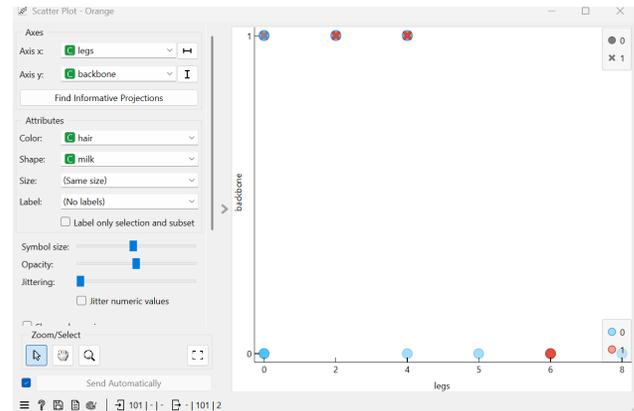


Figure 4. PCA projection of the zoo dataset

The Zoo dataset contains multiple Boolean attributes describing each animal. Conceptually, each animal can be represented as a point in multidimensional feature space, where each dimension corresponds to a variable such as hair, feathers, eggs or milk. In this space, similarity between animals is defined by shared attributes.

However, direct visualisation of high-dimensional space is not possible. PCA addresses this limitation by transforming the original variables into a new set of orthogonal components that capture maximal variance in descending order (Jolliffe & Cadima, 2016). The first principal component identifies the direction along which the data vary most strongly. The second principal component captures the next greatest variance independent of the first. Together, these components provide a reduced representation of the dataset that preserves as much structural information as possible in fewer dimensions.

When plotted along the first two principal components, animals cluster according to shared traits. Mammals group together because they share attributes such as milk production and hair. Birds cluster because of feathers and egg-laying. Reptiles and amphibians appear in proximity where attributes overlap.

This projection does not impose categories. It does not predict class membership. Rather, it reveals structure inherent in the attribute space. PCA reorganises variance to make similarity patterns visible.

Methodologically, this stage serves several purposes.

First, it provides evidence regarding separability. If distinct clusters are visible in PCA space, subsequent classification has a structural basis. If clusters overlap extensively, classification may be unstable or dependent on minor distinctions.

Second, PCA encourages reflection on dimensional reduction. Transforming multiple observed variables into fewer components involves abstraction. The components are linear

combinations of original variables. They do not correspond directly to observable features. They are statistical constructs summarising variation. Interpretation therefore requires theoretical reasoning rather than mechanical reading of axes.

In survey research, PCA is often used to identify latent dimensions such as political ideology or socioeconomic status. However, such dimensions do not exist independently of the variables used to construct them. They are analytical artefacts that gain meaning only through interpretation. The same logic applies here. While mammals cluster in PCA space, the principal components themselves are not biological categories. They are mathematical transformations.

Third, PCA reinforces the distinction between exploratory and confirmatory modelling. At this stage, no predictive claims are made. The objective is to understand structural relationships before imposing classification rules. Exploratory analysis functions as diagnostic groundwork.

In academic research more broadly, PCA is valuable not because it provides definitive conclusions but because it structures inquiry. It reveals whether variance is concentrated along identifiable axes. It suggests potential groupings. It indicates whether feature sets meaningfully differentiate instances.

Within the Orange environment, the PCA widget makes these transformations visible. Loadings, variance explained and projections can be inspected directly. The researcher can examine how much variance is captured by each component and adjust the number of components displayed. This visibility reinforces methodological transparency. Dimensional reduction is not hidden within code; it is presented as an explicit transformation.

Importantly, PCA also reveals the limits of reduction. If the first two components explain only a modest proportion of total variance, the projection may obscure structure present in higher dimensions. The researcher must interpret the projection cautiously.

In the Zoo example, the PCA projection reveals relatively clear separation among major animal classes. This suggests that the attribute set meaningfully differentiates types. The exploratory evidence provides justification for proceeding to supervised classification.

At this stage, the workflow has achieved three outcomes:

- Structural integrity has been verified.
- Distributional context has been established.
- Similarity structure has been explored without imposing categorical decisions.

The next step transitions from exploratory abstraction to formal classification.

B. Interpretable Classification Through Decision Trees

Following exploratory structural analysis, the workflow proceeds to supervised classification using a decision tree model.

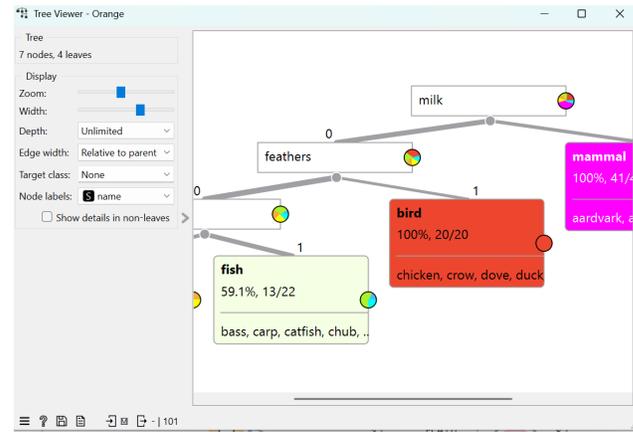


Fig. 5. Decision Tree predicting animal class

Unlike PCA, which reorganises variance without reference to a target variable, supervised classification explicitly models the relationship between predictors and a designated outcome. In this case, the outcome is the animal “type” variable.

A decision tree operates by recursively partitioning the dataset into increasingly homogeneous subsets. At each node, the algorithm selects the predictor that best reduces class impurity according to a splitting criterion, commonly Gini impurity or information gain. Impurity measures the degree to which multiple classes are mixed within a subset. A pure node contains instances belonging to a single class. The objective of each split is to increase homogeneity.

In the Zoo dataset, the decision tree identifies variables such as milk and feathers as primary discriminators. For example, a split on the milk variable separates mammals from non-mammals. Subsequent splits further partition the remaining instances based on additional attributes.

What distinguishes decision trees from more complex models is their interpretability. Each split corresponds to an explicit rule. A path from the root node to a leaf node can be read as a sequence of conditional statements. For instance, if milk equals one, then classify as mammal. If milk equals zero and feathers equals one, then classify as bird.

This structure allows researchers to trace classification logic directly. In academic research contexts, this traceability is methodologically significant. Models are not merely evaluated for performance; they are examined for explanatory alignment with substantive reasoning. Rudin (2019) argues that in settings where explanation is central, interpretable models should be preferred over opaque alternatives.

In the Zoo example, the tree’s reliance on milk and feathers aligns with intuitive biological distinctions. This alignment reinforces confidence that the model is capturing meaningful structure rather than artefactual correlations.

However, interpretability does not eliminate limitations. A decision tree reflects patterns present in the dataset. It does not establish causality. The split on milk does not explain why mammals produce milk; it merely uses the variable as a discriminative feature.

Moreover, decision trees can become overly complex. If allowed to grow without constraint, a tree may create numerous splits that capture idiosyncratic variations specific to the

training data. This leads to overfitting.

C. Overfitting and the Need for Validation

Overfitting occurs when a model learns noise or accidental structure present in the training data rather than generalisable patterns. An overfitted model may achieve near-perfect accuracy on the training set but perform poorly on new data.

Decision trees are particularly susceptible to overfitting because they can continue splitting until each leaf node contains only a small number of observations. While this increases purity on training data, it reduces generalisability.

Addressing overfitting requires explicit validation procedures. Performance must be evaluated not only on the data used to train the model but also on unseen subsets.

In this workflow, validation is implemented using ten-fold cross-validation within the Test & Score widget.

D. Ten-Fold Cross-Validation

Ten-fold cross-validation partitions the dataset into ten approximately equal subsets, or folds. The model is trained on nine folds and tested on the remaining fold. This process repeats ten times so that each fold serves once as the test set. Performance metrics are then averaged across all iterations.

This procedure approximates the model's expected performance on new data drawn from the same distribution. It reduces reliance on a single arbitrary train-test split and mitigates the variance associated with small datasets.

Ten folds represent a balance between bias and variance in performance estimation. Using too few folds increases variability in estimates. Using too many folds increases computational cost and may offer limited additional stability in moderate-sized datasets.

In academic research, cross-validation serves as a safeguard against overly optimistic interpretation. It reinforces the principle that modelling should not be evaluated solely on training performance.

E. Diagnostic Evaluation Through the Confusion Matrix

Following cross-validation, results are examined using the Confusion Matrix widget.

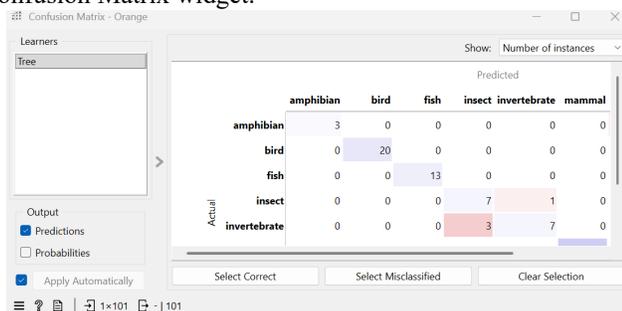


Fig. 6. Confusion matrix from ten-fold cross-validation

The confusion matrix provides a detailed account of classification outcomes. Rows correspond to actual classes. Columns correspond to predicted classes. Values along the diagonal represent correct predictions. Off-diagonal values represent misclassifications.

While overall accuracy summarises performance in a single

metric, the confusion matrix reveals patterns of error. For example, if reptiles are frequently misclassified as amphibians, this may indicate similarity in attribute representation. Such patterns can guide further investigation.

In the Zoo dataset, most predictions fall along the diagonal, indicating stable generalisation across folds. However, the confusion matrix remains essential because it reveals which classes are more difficult to distinguish.

In applied contexts, diagnostic granularity is critical. In education research, systematic misclassification of underperforming students could obscure structural inequalities. In public health, misclassifying medium-risk patients as low-risk could have practical consequences. The confusion matrix enables researchers to identify such asymmetries.

F. Integrating Modelling and Interpretation

At this stage, the workflow has moved from structural exploration to formal classification and then to validation and diagnostic evaluation. Importantly, these stages are interconnected.

The PCA projection provided preliminary evidence of separability. The decision tree formalised this separability into rule-based partitions. Cross-validation tested whether these partitions generalise beyond the training set. The confusion matrix identified specific patterns of error.

Each stage informs interpretation. If PCA had revealed substantial overlap among classes, the researcher might anticipate lower classification performance. If cross-validation had revealed unstable performance across folds, the researcher might reconsider tree depth or feature relevance.

Orange's visual architecture supports this integration. Adjustments to the decision tree automatically propagate to validation results. The workflow remains cohesive rather than fragmented across environments.

5. Methodological Framing: Exploratory Modelling and Knowledge Production

The analytical sequence demonstrated above is exploratory rather than confirmatory. This distinction is methodologically significant.

Confirmatory research begins with predefined hypotheses and evaluates them through formal statistical testing. The objective is to determine whether evidence supports or rejects specified claims. Exploratory modelling operates differently. It seeks to identify patterns, structures and relationships that may not be fully specified at the outset of inquiry. Both approaches are legitimate within academic research, but they generate different forms of knowledge.

Exploratory computational modelling does not replace theory. It interacts with theory. Grimmer and Stewart (2013), in their foundational discussion of automated text analysis, argue that computational methods reveal structure at scale, but interpretation remains a theoretical act. The analyst must decide which patterns matter and how they connect to substantive questions. Exploratory modelling therefore functions as structured inductive reasoning rather than automated discovery.

Within the Orange workflow, exploratory reasoning occurs

at multiple stages.

The inspection of class distribution reveals the empirical shape of the dataset. It situates modelling expectations.

The PCA projection reduces dimensional complexity to identify similarity patterns without imposing categorical boundaries. It makes structural relationships visible before prediction is attempted.

Even the decision tree participates in exploratory reasoning. Although it is supervised, it reveals which features most strongly differentiate classes. The tree highlights structural regularities that may not have been anticipated.

Principal Component Analysis illustrates the epistemological implications of computational abstraction. PCA constructs new axes of representation by transforming observed variables into orthogonal components that maximise variance (Jolliffe & Cadima, 2016). These components do not exist independently of the data. They are mathematical summaries of variation. In survey research, PCA is often used to identify latent attitudinal dimensions. Yet such dimensions are analytical constructs. They gain meaning only through interpretation.

Similarly, classification models impose boundaries within multidimensional space. A decision tree discretises complexity by splitting feature space into rule-based segments. This simplification is analytically useful because it makes prediction tractable. However, simplification is also reduction. Models illuminate structure but do not exhaust it.

Exploratory modelling therefore produces provisional knowledge. It generates structured representations that must be situated within broader theoretical frameworks. Orange's visual workflow reinforces this awareness. Each transformation is visible. Each modelling decision can be inspected. The analyst is reminded that computational analysis consists of sequential choices rather than automatic revelation.

A. Interpretability, Transparency and the Trade-Off with Complexity

The increasing sophistication of machine learning techniques has intensified debates about interpretability. Complex ensemble methods and deep neural networks often achieve high predictive accuracy but operate as opaque systems. Their internal logic is distributed across numerous parameters that resist straightforward explanation.

In many academic research contexts, particularly those involving education, public health or policy analysis, opacity raises both ethical and epistemological concerns (Rudin, 2019). Researchers must be able to explain why a model produces a particular classification. Without such explanation, results cannot be meaningfully integrated into theoretical frameworks.

Decision trees represent a form of interpretable modelling. Each split corresponds to an explicit rule. In the Zoo example, the split on milk corresponds directly to biological classification. In applied contexts, splits on variables such as education level, income or media source type can be examined in relation to theoretical expectations.

The trade-off between interpretability and predictive performance is well recognised. Highly flexible models can capture subtle nonlinear patterns but may sacrifice

transparency. Simpler models may achieve slightly lower accuracy while providing clearer insight into structure. In exploratory academic research, interpretability often carries greater value than marginal performance gains.

Transparency also extends beyond model structure. It includes documentation of preprocessing decisions, variable selection and validation procedures. Mitchell *et al.* (2019) argue for systematic model reporting to promote accountability. While their discussion focuses on machine learning ethics, the underlying principle applies broadly. Research credibility depends on the traceability of analytical decisions.

Orange's architecture supports this transparency. The workflow canvas externalises the analytical sequence. Rather than embedding parameter settings in scripts, the interface displays them. Readers and collaborators can trace how data move from inspection to modelling to evaluation.

Importantly, interpretability does not guarantee correctness. An interpretable model may still reflect biased data or flawed assumptions. Transparency facilitates scrutiny but does not eliminate responsibility. Responsible research requires both interpretability and critical reflection.

B. Reflective Analytical Practice and Iteration

Empirical research rarely unfolds in a strictly linear fashion. Analytical workflows are iterative.

Initial exploratory steps may reveal unexpected patterns that prompt revision of feature sets. Cross-validation results may indicate overfitting, requiring adjustment of model complexity. Confusion matrices may reveal systematic misclassification that suggests the need for additional variables.

In the Zoo example, if the PCA projection had shown substantial overlap between classes, the researcher might question whether the selected attributes adequately capture distinguishing features. If cross-validation had produced unstable results across folds, the tree's depth or splitting criteria might be reconsidered.

Overfitting in particular invites iteration. A tree that perfectly classifies training data but performs poorly in cross-validation likely captures noise rather than stable structure. Reducing tree depth or limiting minimum samples per leaf may improve generalisability.

The confusion matrix functions as a diagnostic instrument within this iterative process. Persistent misclassification of certain categories may indicate structural similarity in feature space or insufficient representation of distinguishing variables.

Reflective practice therefore involves returning to earlier stages of the workflow. Exploration informs modelling. Evaluation informs revision. Interpretation remains provisional.

Orange supports this iterative reasoning because modifications propagate automatically through connected widgets. Adjusting the decision tree recalculates cross-validation results. The canvas becomes both an analytical environment and a record of experimentation.

C. Applied Academic Scenarios: Extending the Workflow Beyond the Zoo Example

The worked example using the Zoo dataset demonstrates a complete analytical cycle: inspection, contextualisation, exploration, modelling and validation. Although zoological classification is not a typical academic research objective, the structure of the workflow mirrors analytical reasoning across disciplines. The strength of the example lies not in its subject matter but in the clarity with which it illustrates methodological sequencing.

To demonstrate how the workflow generalises, it is useful to consider several applied research contexts.

In education research, scholars frequently analyse student-level datasets containing variables such as attendance rates, prior academic performance, socioeconomic background and engagement indicators. Suppose the objective is to classify students into performance categories such as high distinction, distinction, pass and fail. The workflow would begin with structural verification in the Data Table to ensure that grade categories are correctly encoded and that missing values are handled appropriately. Descriptive analysis would reveal whether certain categories dominate, shaping interpretation of classification accuracy.

PCA could then be used to explore whether behavioural indicators cluster along identifiable dimensions. Attendance and study hours might load strongly onto a first principal component that could be interpreted as engagement. If high-performing students cluster distinctly in projection space, this provides preliminary structural evidence. A decision tree might reveal that prior GPA serves as a primary split, followed by attendance or tutorial participation. Cross-validation would test the stability of these splits across subsets of students. The confusion matrix might reveal systematic misclassification between distinction and high distinction categories, indicating that available predictors insufficiently differentiate top-performing groups.

In media studies, researchers often work with coded datasets of news articles. Variables may include tone, framing, source type and issue category. Suppose the objective is to classify articles by outlet type. The workflow would again begin with inspection and distributional analysis to assess class balance. PCA might reveal clustering along axes corresponding to sensationalism or institutional authority. A decision tree could identify which coded features most strongly distinguish outlets. Cross-validation would assess stability. The confusion matrix might reveal that broadsheet and public broadcaster articles are frequently confused, suggesting overlapping framing strategies. Within the same environment, researchers could compare multiple classifiers or rank features to assess which variables contribute most strongly to discrimination.

In public health research, structured datasets frequently include physiological and behavioural indicators such as age, BMI, smoking status and exercise frequency. Suppose the objective is to classify patients into risk categories. After verifying structure and distribution, PCA might reveal clustering among physiological indicators. A decision tree might identify BMI as a primary split, followed by smoking

status. Cross-validation would test generalisability. If medium-risk patients are frequently misclassified as low-risk, this may suggest insufficient granularity in the available predictors. Because clinical contexts require explanation, the interpretability of decision trees becomes particularly important. Clinicians must understand why a classification is made, not merely that it is accurate.

In political science, survey datasets often include demographic and attitudinal variables. Suppose the objective is to classify respondents according to support for a policy initiative. PCA might reveal ideological clustering along axes derived from party identification and trust in government. A decision tree might identify party affiliation as the primary predictor, followed by education level. Cross-validation would test stability across subsets of respondents. Misclassification of undecided respondents may indicate ambiguity in predictor variables. Researchers could also use clustering to identify latent opinion segments independent of declared support categories.

Across these contexts, the domain changes but the analytical logic remains consistent. The workflow moves from inspection to contextualisation, from exploration to classification and from classification to validation and interpretation. The six figures in this section therefore represent a generalisable analytical template rather than a domain-specific procedure.

Figure 1 establishes the logic of structured reasoning. Fig. 2 verifies data integrity. Figure 3 situates modelling within distributional context. Figure 4 explores structural relationships without imposing categorical boundaries. Figure 5 formalises these relationships through interpretable classification. Figure 6 evaluates generalisation and error patterns.

The sequence reinforces methodological discipline. It prevents premature modelling and encourages validation before inference. This consistency is the central methodological takeaway.

D. Extending the Framework Within Orange

Although the section has focused on a specific workflow, Orange supports additional techniques that remain compatible with the same structured logic.

Feature selection can be used to reduce dimensionality prior to modelling, helping researchers identify which variables contribute most strongly to discrimination. Outlier detection can identify anomalous cases that warrant closer examination. Comparative model evaluation through the Test & Score widget allows researchers to assess performance across different algorithms while maintaining interpretability. Interactive visualisation tools support subgroup analysis and exploratory inspection of misclassified instances.

The objective is not to demonstrate every available widget but to illustrate how a single environment can accommodate flexible inquiry without fragmenting the analytical pipeline.

E. Pedagogical and Interdisciplinary Implications

Beyond its analytical functionality, Orange has pedagogical significance. Computational modelling can appear opaque to students and researchers without programming backgrounds.

Visual workflows reduce this barrier by externalising analytical logic.

In teaching contexts, instructors can demonstrate how descriptive statistics relate to modelling and how validation procedures function. Students observe the sequence of operations rather than encountering isolated fragments of code. This supports conceptual understanding of the relationship between exploration, classification and evaluation.

Interdisciplinary collaboration also benefits from visual workflows. Research teams often include members with varied methodological expertise. A shared visual canvas provides a common reference point for discussing analytical decisions. Conversation can focus on conceptual reasoning rather than technical syntax.

Visual tools do not eliminate the need for methodological rigour. Rather, they make methodological reasoning more accessible.

F. Limits and Responsible Use

No analytical environment is universally appropriate. Orange is particularly well suited to small-to-medium structured datasets and exploratory modelling. It is not optimised for large-scale neural network training or distributed computation.

Models reflect the structure of the data on which they are trained. If datasets are biased, incomplete or unrepresentative, model outputs will reflect those limitations. Interpretability facilitates scrutiny but does not guarantee validity. Ethical evaluation and theoretical reflection remain essential.

Recognising limits strengthens methodological contribution. By clarifying the scope within which Orange is most effective, researchers can deploy it appropriately.

6. Conclusion

This section has demonstrated how Orange Data Mining supports transparent, interpretable and reproducible analytical workflows in academic research. Through integration of descriptive analysis, dimensionality reduction, interpretable classification and cross-validation within a single visual environment, Orange aligns computational procedure with

scholarly norms.

The Zoo dataset serves as a methodological illustration. The analytical reasoning it demonstrates generalises across education research, media analysis, public health and political science.

The primary contribution of Orange lies not in computational novelty but in the visibility of analytical structure. By externalising each stage of reasoning, it reinforces disciplined research practice and supports interpretability in exploratory contexts.

References

- [1] D. Bamman, B. O'Connor, and N. A. Smith, "Learning latent personas of film characters," in *Proc. 52nd Annu. Meeting Assoc. Comput. Linguistics (ACL)*, Baltimore, MD, USA, 2014, pp. 352–361.
- [2] A. Bruns and J. Burgess, "Twitter hashtags from ad hoc to calculated publics," in *Hashtag Publics: The Power and Politics of Discursive Networks*, N. Rambukkana, Ed. New York, NY, USA: Peter Lang, 2015, pp. 13–28.
- [3] M. Graham, S. A. Hale, and D. Gaffney, *Digital Geographies: An Introduction*. London, U.K.: SAGE, 2021.
- [4] J. Grimmer and B. M. Stewart, "Text as data: The promise and pitfalls of automated content analysis methods for political texts," *Political Anal.*, vol. 21, no. 3, pp. 267–297, 2013.
- [5] S. A. Hale, G. Blank, and V. D. Alexander, "Live versus recorded social media data: Ethical considerations in computational research," *Big Data Soc.*, vol. 8, no. 1, 2021, Art. no. 2053951721996689.
- [6] I. T. Jolliffe and J. Cadima, "Principal component analysis: A review and recent developments," *Philos. Trans. Roy. Soc. A*, vol. 374, no. 2065, Art. no. 20150202, 2016.
- [7] S. Kapoor and A. Narayanan, "Leakage and reproducibility in machine learning-based science," *Patterns*, vol. 4, no. 9, Art. no. 100804, 2023.
- [8] D. Lazer *et al.*, "Computational social science: Obstacles and opportunities," *Science*, vol. 369, no. 6507, pp. 1060–1062, 2020.
- [9] M. Mitchell *et al.*, "Model cards for model reporting," in *Proc. Conf. Fairness, Accountability, and Transparency (FAT*)*, 2019, pp. 220–229.
- [10] C. Rudin, "Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead," *Nature Mach. Intell.*, vol. 1, pp. 206–215, 2019.
- [11] L. Sloan, A. Quan-Haase, and M. Thelwall, Eds., *The SAGE Handbook of Social Media Research Methods*, 2nd ed. London, U.K.: SAGE, 2020.
- [12] A. Törnberg and P. Törnberg, "How digital media drive affective polarisation," *Inf., Commun. Soc.*, vol. 25, no. 5, pp. 666–683, 2022.
- [13] M. Zappavigna, *Searchable Talk: The Linguistic Functions of Hashtags*. London, U.K.: Bloomsbury Academic, 2018.