

Interview Task: High-Level Architecture for Predictive Health Platform with Personalized Supplement Protocols

Task Description

You are tasked with sketching a high-level architecture for a consumer healthcare platform that:

1. Predicts health risks by integrating microbiome data, genomic data, and blood biomarker data. How might this differ for cardiovascular vs digestive disorders?
2. Recommends personalized supplement protocols based on these predictions to support individual health optimization.

The focus is on data integration, predictive modeling, and personalized recommendations.

- **Data Integration:**
 - **Data Sources:**
 - **Microbiome Data:** Taxonomic profiles, abundance tables, and diversity metrics derived from 16S rRNA or shotgun metagenomics sequencing.
 - **Genomic Data:** Variant call data (VCF), polygenic risk scores, and annotated genetic variants influencing nutrient metabolism and disease susceptibility.
 - **Blood Biomarkers:** Routine clinical markers (vitamins, lipids, inflammatory markers) measured via lab tests.

- **Integration Pipeline:**
 - **Data Ingestion & Standardization:**

Ingest raw microbiome (OTU/ASV tables, taxonomic classifications), genomic (variant reports, PRS files), and biomarker data (CSV, database tables) into a unified data lake.

- **Feature Engineering:**

Convert each data type into a structured, analytics-ready feature set. For microbiome data, extract key taxa counts and diversity indices. For genomic data, incorporate risk alleles or polygenic scores. For biomarkers, normalize and standardize measurements (e.g., adjust vitamin levels to reference ranges).

- **Unified Feature Store:**

Combine these engineered features—microbial diversity metrics, genomic risk factors, and biomarker readings—into a single integrated table. This forms the input for the predictive models.

Predictive Modeling:

- **Multi-Modal ML Models:**

- Use ensemble machine learning or neural networks that can handle heterogeneous inputs (e.g., microbial abundance vectors, numeric genomic scores, continuous biomarker values).
- Train models to produce condition-specific risk-scores. For example:
 - **Cardiovascular Disorders:** Place more emphasis on **lipid profiles, inflammation markers (CRP), and genomic variants related to heart disease**, while also incorporating gut-derived markers associated with systemic metabolic health.
 - **Digestive Disorders:** Focus more on **gut microbial composition and diversity, genetic predispositions affecting nutrient absorption or gut barrier integrity**, and biomarkers indicating digestive inflammation or nutrient malabsorption.
- **Model Validation & Continuous Improvement:**
 - Validate predictive accuracy using cross-validation and test sets.
 - Continuously retrain models as new data arrives, ensuring that predictions remain up-to-date with evolving population-level insights and improved genomic annotations.

Personalized Recommendations:

- **Recommendation Engine:**
 - **Mapping Risk Scores to Supplement Needs:**

Once a user's risk for cardiovascular or digestive issues is predicted, map these scores to a recommended supplement protocol. For cardiovascular risks, suggest supplements known to improve lipid profiles or provide anti-inflammatory benefits (e.g., omega-3 fatty acids). For digestive issues, focus on probiotics, prebiotics, or targeted micronutrients supporting gut health.

- **Genomic and Microbiome-Informed Choices:**

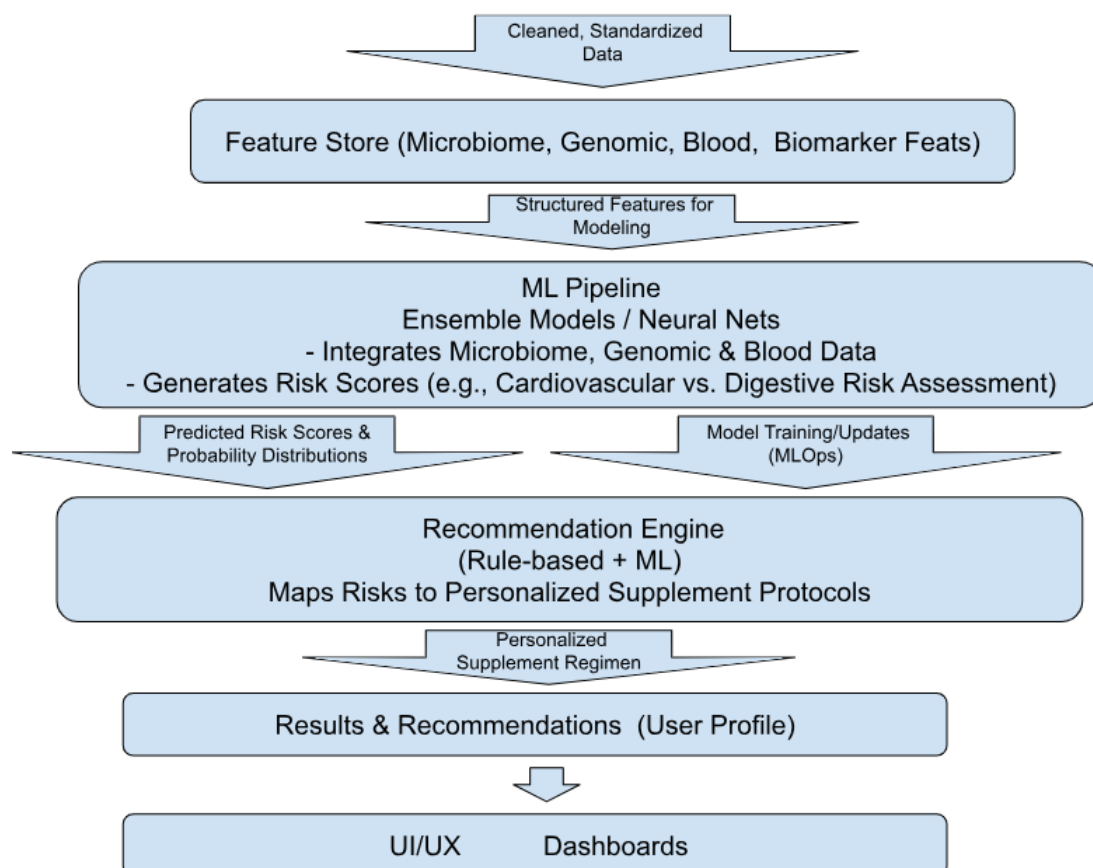
Integrate the user's genomic background (e.g., genes affecting nutrient absorption) and microbiome composition to refine the supplement list. If genomic data suggests poor B-vitamin utilization, include appropriate forms of these vitamins. If the microbiome analysis shows low abundance of beneficial microbes, suggest specific probiotic strains.

- **Biomarker-Guided Adjustments:** Use blood biomarkers to address current deficiencies directly. Low vitamin D prompts vitamin D supplementation, while elevated inflammatory markers may lead to recommending anti-inflammatory compounds or antioxidants.
- **Dynamic, Iterative Refinement:**

- As users recheck their biomarkers or update their genomic and microbiome data, the system recalculates risk and updates recommendations accordingly.
- Incorporate user feedback and adherence data, refining future supplement protocols to improve personalization and health outcomes.

In essence, the platform brings together diverse biological datasets, runs predictive models tailored to the health domain of interest, and then translates those predictive insights into individualized supplement recommendations. *Cardiovascular-focused predictions* rely on integration of *lipid-related genes*, microbiome-linked metabolites, and inflammatory biomarkers, while digestive-focused predictions emphasize gut microbial composition, genetic variants impacting nutrient absorption, and GI-specific biomarkers. This cohesive integration and predictive modeling enable the creation of highly personalized, data-driven supplement protocols for health optimization. Below is an example solution outlining a high-level architecture diagram.

High-Level Architecture Diagram (Conceptual)



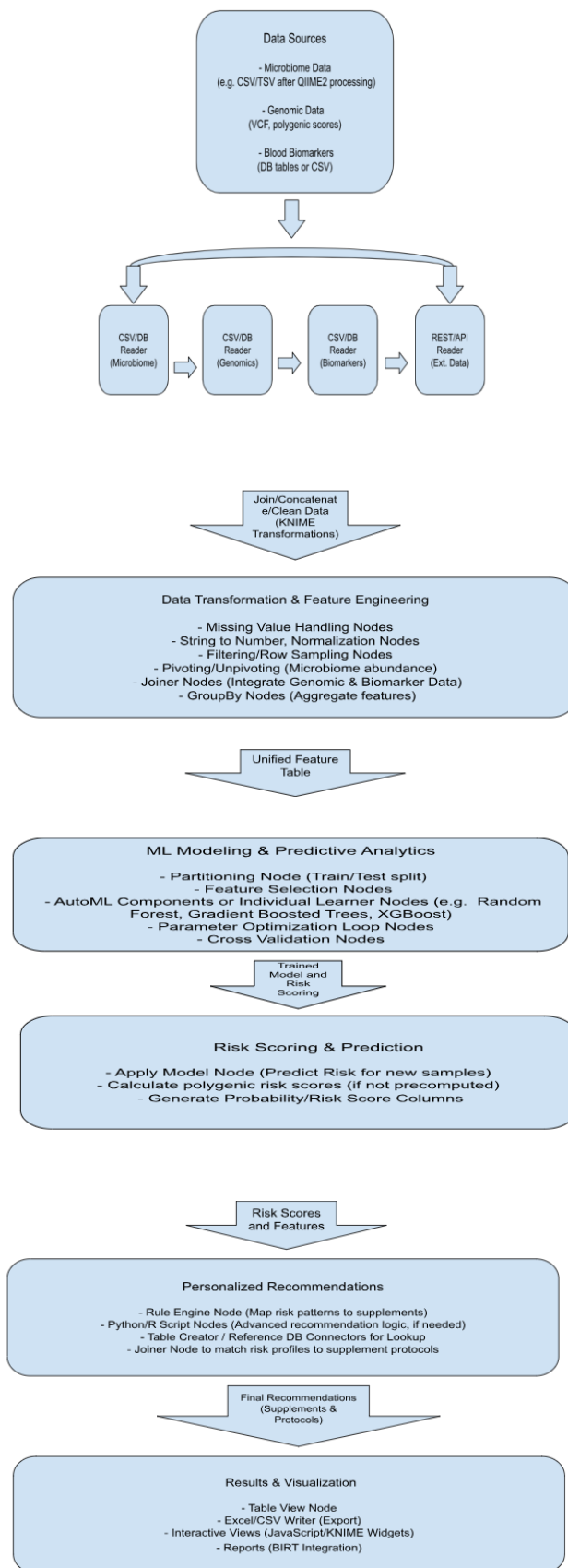
Introduction to Knime Analytics platform

As your AI-based startup explores effective **data analytics and machine learning platforms**, both KNIME and Databricks might be considered. KNIME offers a user-friendly, low-code environment that is excellent for rapidly prototyping models and integrating diverse data sources at an early stage. Meanwhile, Databricks provides a scalable, cloud-based infrastructure suitable for handling large datasets and complex workflows. At this point, I recommend starting with KNIME, as it allows quick experimentation, easy onboarding, and cost-effective analytics that can help your team iterate rapidly and efficiently.

In this project, you can utilize the KNIME Analytics Platform to design and orchestrate end-to-end workflows for integrating microbiome, genomic, and blood biomarker data. KNIME provides a user-friendly interface, extensive node library, and seamless extensibility, making it well-suited for data preprocessing, feature engineering, and predictive modeling in a startup environment. As the business grows and the complexity of data operations increases, investing in a KNIME Business Hub could offer substantial benefits: secure team collaboration, centralized management of workflows and data, scalable computational resources, and easier integration with enterprise IT systems. This approach ensures a robust, agile foundation for rapidly developing and iterating on data-driven healthcare solutions.

Below is an example of how you might conceptualize a KNIME-based workflow, from data ingestion to final recommendations, along with a high-level architecture diagram. This scenario assumes that the raw sequencing, genomic variant calling, and biomarker assays have already been performed externally, and that the data you're ingesting into KNIME is in a format suitable for analysis (e.g., CSV files, TSV files, database tables, or results from external pipelines accessed via REST APIs).

High-Level Architecture Diagram (KNIME-Focused)



Step-by-Step Workflow Description

1. Data Ingestion:

- **Microbiome Data:** Use a **CSV Reader** node (if the data is in CSV format) or a **File Reader** node for more complex formats. The data typically includes microbial relative abundances or OTU (Operational Taxonomic Unit)/ASV tables from external processing (e.g., QIIME2).
- **Genomic Data:** If variant calls and polygenic risk scores are precomputed and stored in CSV, VCF, or a database, use a **CSV Reader**, **Database Connector + DB Table Selector**, or **VCF Reader (community extension)**.
- **Blood Biomarkers:** Connect to a database (e.g., **Database Connector** and **DB Reader** nodes) or read from CSV. For example, you might have a table with patient ID and measured vitamin levels, cholesterol, CRP, etc.

2. Data Preprocessing & Integration:

- **Quality Checks & Cleaning:** Use **Missing Value** nodes to handle missing data, **Column Filter** to remove irrelevant columns, and **String to Number** nodes to ensure correct data types.
- **Feature Engineering:** Aggregate microbial features by taxonomic level (if needed) using **GroupBy** nodes. Normalize numeric data with **Normalizer** nodes.
- **Integration:** Join all datasets by a common patient/sample ID using the **Joiner** node to create a unified feature table that includes microbiome metrics, genomic variants or scores, and blood biomarkers.

3. Feature Store & Transformation:

- After integration, you have a single table with one row per subject/sample and columns representing microbial features (e.g., diversity indices, abundance of key taxa), genomic risk scores, and numeric biomarker readings. This serves as your “feature store.”
- Further refine features: **Column Filter** (remove low-variance features), **Math Formula** (derive new features), **One Hot Encoding** (if categorical genomic variants exist).

4. Modeling & Prediction:

- **Partitioning Node:** Split data into training and test sets.
- **Machine Learning Nodes:** Use a **Gradient Boosted Trees Learner**, **Random Forest Learner**, or integrate Python/R scripts via **Python Script** or **R Predictor** for advanced models (e.g., neural networks or ensemble approaches).

- **Cross Validation Loop:** Validate model performance and select the best model.
- **Scorer Node:** Assess metrics (AUC, Accuracy, F1-score).

5. Risk Assessment:

- **Apply Model Node:** Once the best model is chosen, apply it to new/prediction data to generate risk scores for conditions such as cardiovascular or digestive disorders.
- **Output:** A table with patient ID and predicted risk scores/labels.

6. Recommendation Engine:

- **Rule Engine Node:** Define rules that map certain risk profiles and biomarker deficiencies to recommended supplements. For instance, if vitamin D level < threshold + high genomic risk for cardiovascular disease → recommend vitamin D + Omega-3 supplements.
- For more complex logic, use **Python Script** nodes to implement a more sophisticated recommendation logic that references a supplemental knowledge base.
- **Table Creator / Reference Lookup:** Store a mapping of risk profiles to supplement products and protocols in a reference table, then join this table with the risk predictions to generate final recommendations.

7. Output & Visualization:

- **Excel Writer or CSV Writer:** Export final recommendations.
- **KNIME Interactive Views** or JavaScript nodes: Create dashboards to visualize microbial diversity, genomic risk scores distribution, and recommended supplements.
- **Report Designer (BIRT):** Generate professional reports summarizing patient risk profiles and suggested supplement regimes.

Considerations:

- **Scalability:** KNIME Server or Executors can be used to scale and schedule workflows, enabling orchestration of regular data updates and predictions.
- **Data Security:** Use database connections with secure credentials and consider KNIME's support for encrypted workflows and secure connections to ensure sensitive genomic and health data is protected.

- **Extensibility:** Python and R integration nodes allow for custom advanced analytics, machine learning models, and integration with specialized bioinformatics libraries.

This KNIME-centric approach provides a visual, modular solution. Each step (ingestion, integration, modeling, recommendation) corresponds to a segment of the KNIME workflow that can be adjusted or extended as new data sources or modeling techniques become available.

Example workflow and results

We can utilize the KNIME Analytics Platform to design and orchestrate end-to-end workflows for integrating microbiome, genomic, and blood biomarker data. KNIME provides a user-friendly interface, extensive node library, and seamless extensibility, making it well-suited for data preprocessing, feature engineering, and predictive modeling in a startup environment. As the business grows and the complexity of data operations increases, investing in a KNIME Business Hub could offer substantial benefits: secure team collaboration, centralized management of workflows and data, scalable computational resources, and easier integration with enterprise IT systems and automation. This approach ensures a robust, agile foundation for rapidly developing and iterating on data-driven healthcare solutions.

KNIME allows seamless integration with various programming languages and external tools, providing flexibility for advanced analytics and custom processing. Common integrations include:

- **Python:**
 - Through the **Python Script** nodes, you can leverage Python libraries (e.g., pandas, scikit-learn, TensorFlow) directly within the KNIME environment.
- **R:**

Using **R Scripting** nodes, incorporate popular R packages (e.g., dplyr, ggplot2, caret) to perform statistical tests, build machine learning models, or create custom visualizations.

- **Java and Groovy:**

KNIME's **Java Snippet** and **Groovy Script** nodes allow for embedding custom Java or Groovy code, giving you direct access to Java-based libraries and tools.

- **Weka and H2O.ai:**

KNIME provides integration nodes for **Weka** and **H2O.ai**, enabling access to a wide range of machine learning algorithms, automated machine learning, and large-scale modeling techniques.

- **Deep Learning Frameworks:**

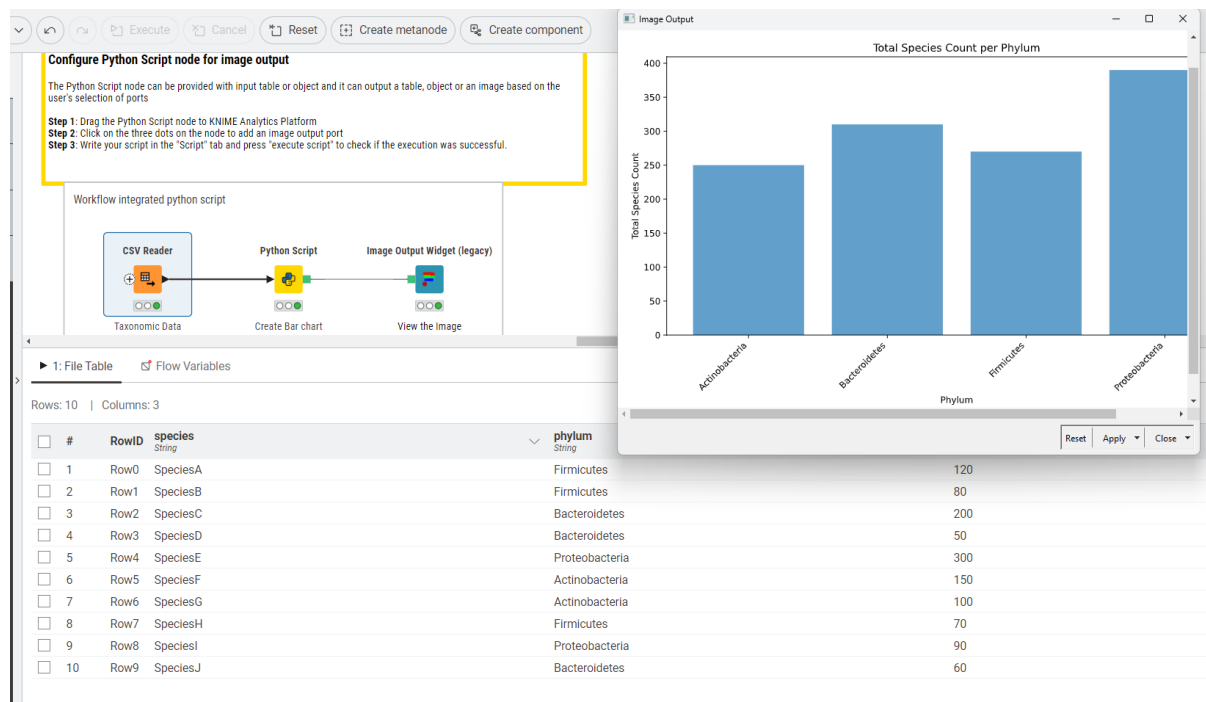
With dedicated extensions, KNIME can interface with **TensorFlow**, **Keras**, and **ONNX**, allowing you to train and deploy deep learning models as part of your KNIME workflows.

- **External Command Line Tools (CLI):**

By using the **External Tool** node, KNIME can run shell commands or invoke bioinformatics pipelines (e.g., QIIME2, SAMtools) and other domain-specific software. This lets you integrate specialized tools into your KNIME pipeline without leaving the platform.

In essence, KNIME's extensibility through scripting nodes and dedicated integration extensions makes it a versatile hub for combining multiple analytic and computational environments under a single, visual workflow framework.

Integrate python script: The first example shows knime workflow that integrate python script and visualize the results from task 2.

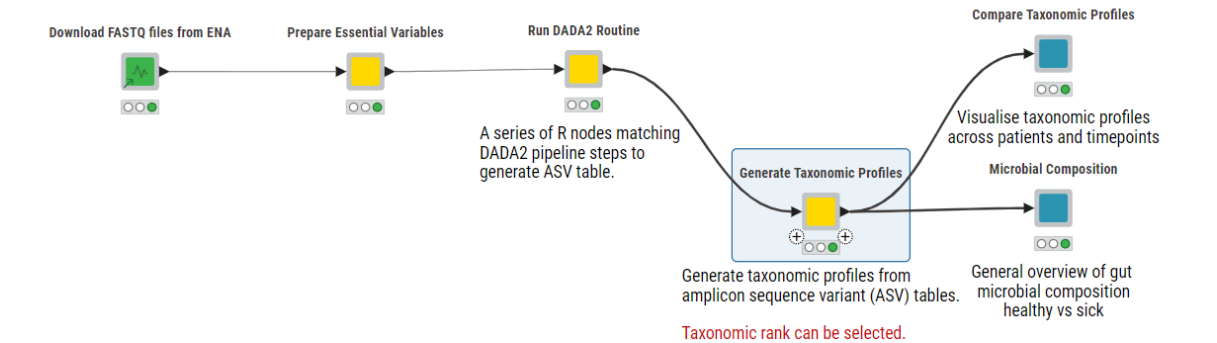


Microbiome Analysis

Gut Microbiome Analysis of Inflammatory bowel disease (IBD) Patients (16S)

This workflow downloads environmental DNA in the form of 16S amplicon sequences from ENA and analyzes it to create and compare taxonomic profiles of a microbial community.

DADA2 R package required



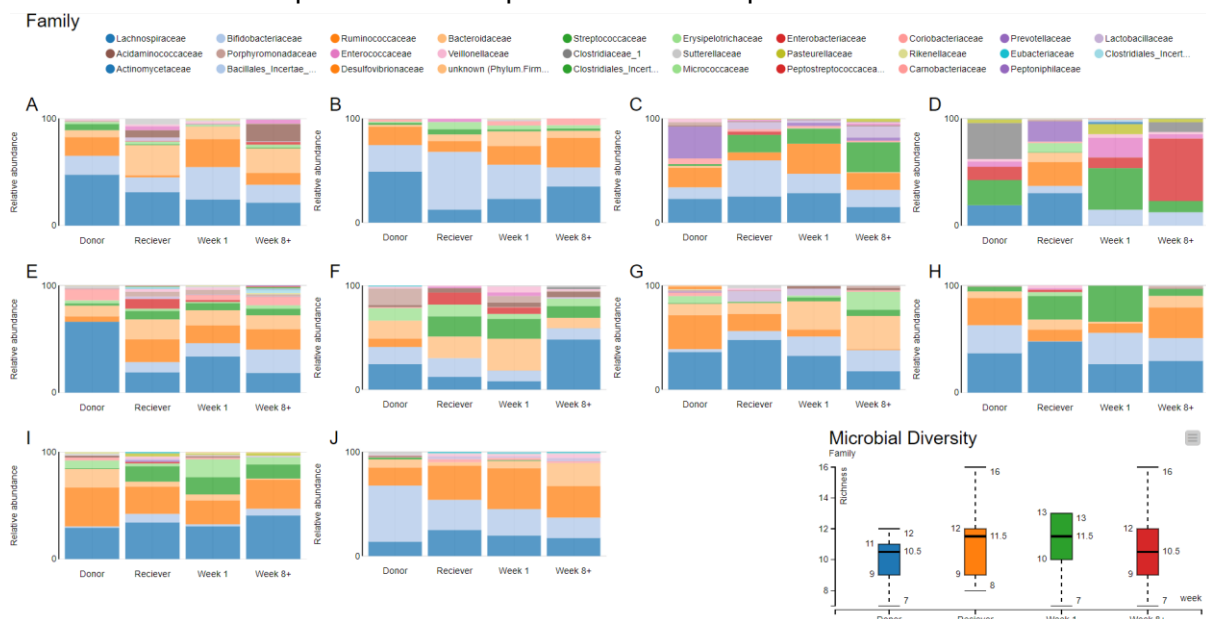
onomic profile ☒ Flow Variables

Columns: 47

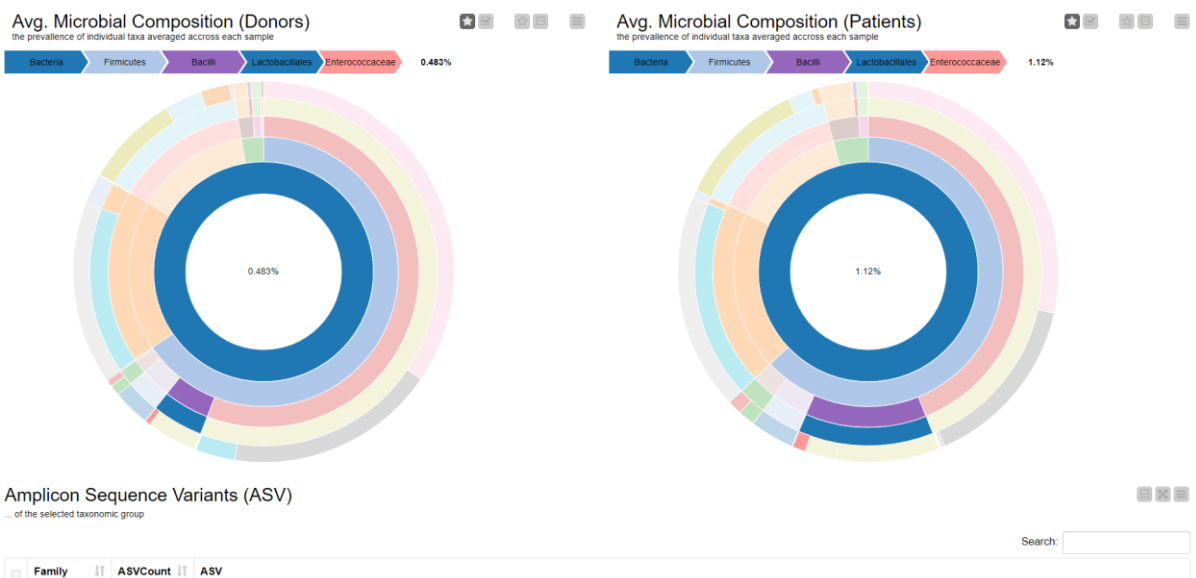
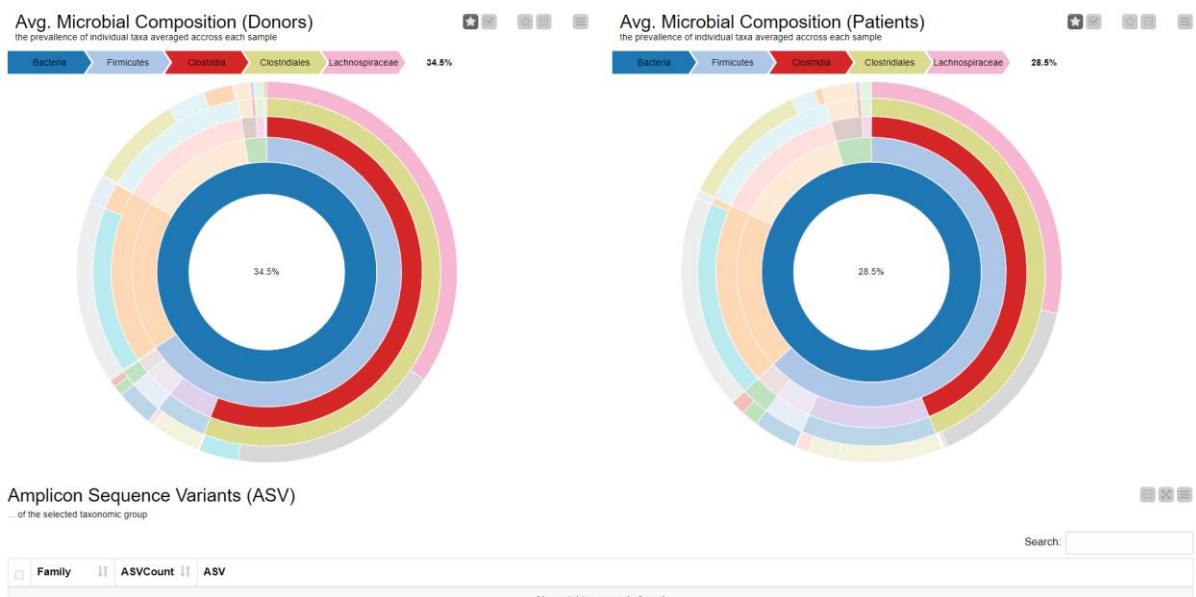
Table ☒ Statistics ☐

RowID	Kingdom String	Phylum String	Class String	Order String	ASVCount String	ASV List	DRR0659... Number (dou...	DRR0659... Number (dou...	DRR0659... Number (do
Lachnospiraceae	Bacteria	Firmicutes	Clostridia	Clostridiales	319	[AGAGTTTGATC	47.559	31.222	24.447
Bifidobacteriaceae	Bacteria	Actinobacteria	Actinobacteria	Bifidobacteriales	85	[AGAGTTTGATC	17.635	13.876	30.24
Ruminococcaceae	Bacteria	Firmicutes	Clostridia	Clostridiales	221	[AGAGTTTGATC	17.466	2.011	26.189

Visualise taxonomic profiles across patients and timepoints



General overview of gut microbial composition healthy vs sick.



Predictive modelling

Proposed Machine Learning Approaches:

When integrating diverse biological datasets—such as microbiome profiles, genomic risk scores, and blood biomarkers—it's often beneficial to experiment with multiple modeling approaches and then select the one that best suits the project's unique characteristics and performance criteria. A flexible strategy might include:

1. Ensemble Models (e.g., Gradient Boosted Trees, Random Forest):

Ensemble methods are well-established for handling heterogeneous feature sets. They naturally deal with different scales and missing values, while also providing feature importance metrics that can offer insights into which inputs (microbiome, genomic, or biomarker data) drive the predictions.

2. Neural Networks (e.g., Fully Connected Deep Learning, Multimodal Architectures):

Deep learning models can capture complex, nonlinear relationships and may excel when large amounts of data are available. If the project involves very large, rich datasets, neural networks can be extended into multimodal architectures that separately process microbial abundance vectors, genomic variants, and numeric biomarkers, and then combine those embeddings into a unified prediction.

3. Hybrid or Stacked Models:

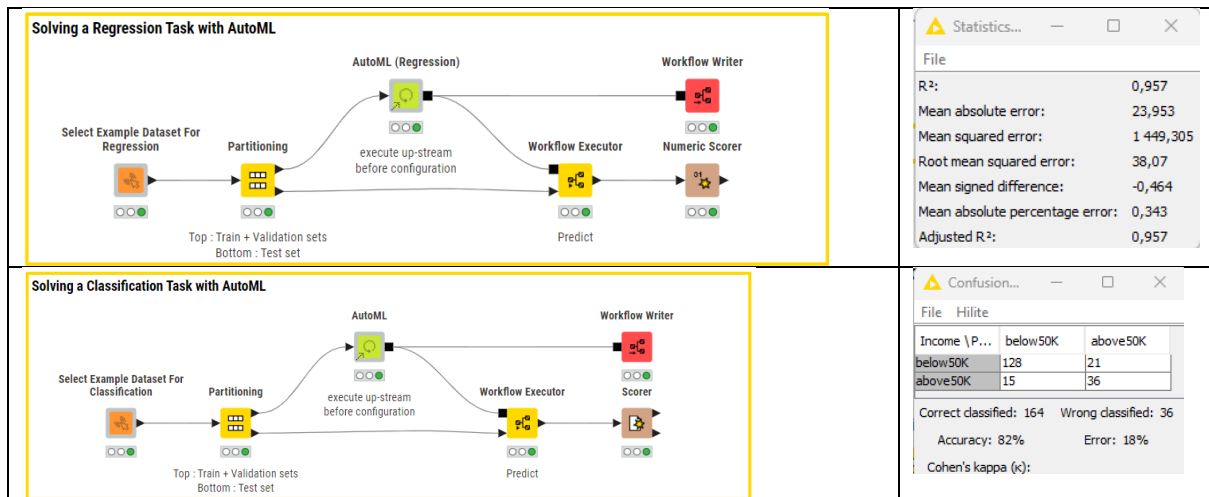
A combination of different algorithms can be tried in a stacked ensemble framework, where outputs of baseline models (e.g., linear models, random forests) feed into a higher-level model (e.g., gradient boosting or neural network). This often yields improved predictive performance by leveraging the strengths of multiple algorithms.

4. AutoML Tools and Model Selection Pipelines:

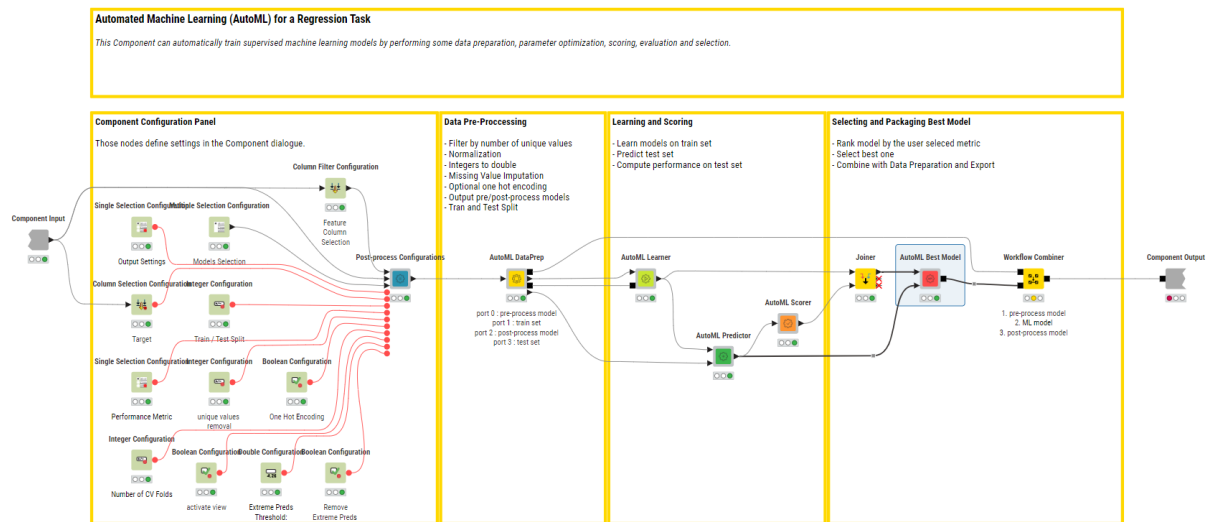
Using AutoML solutions or implementing a systematic model selection pipeline allows you to run multiple models simultaneously, comparing their performance on a validation set. This approach ensures that you don't rely on one particular technique prematurely and can dynamically choose the best-performing model for the given dataset and objectives.

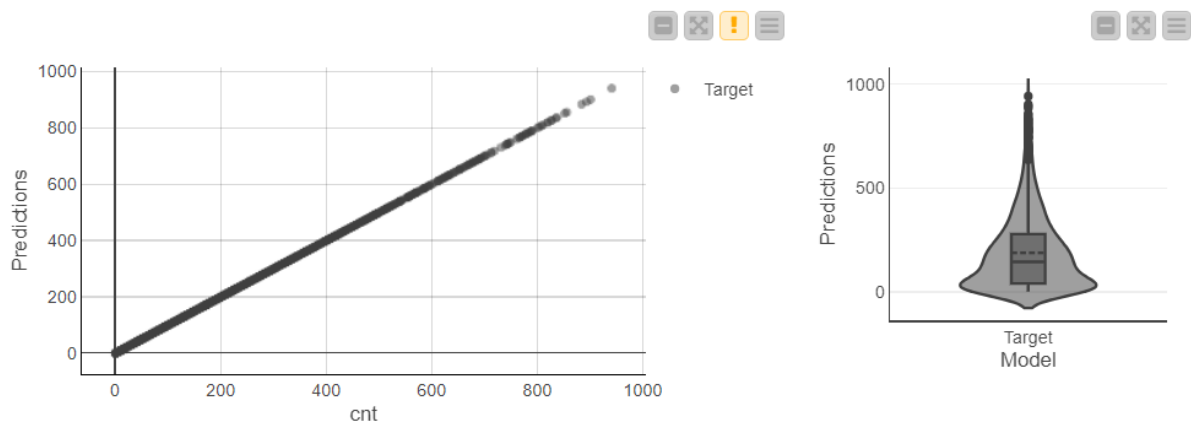
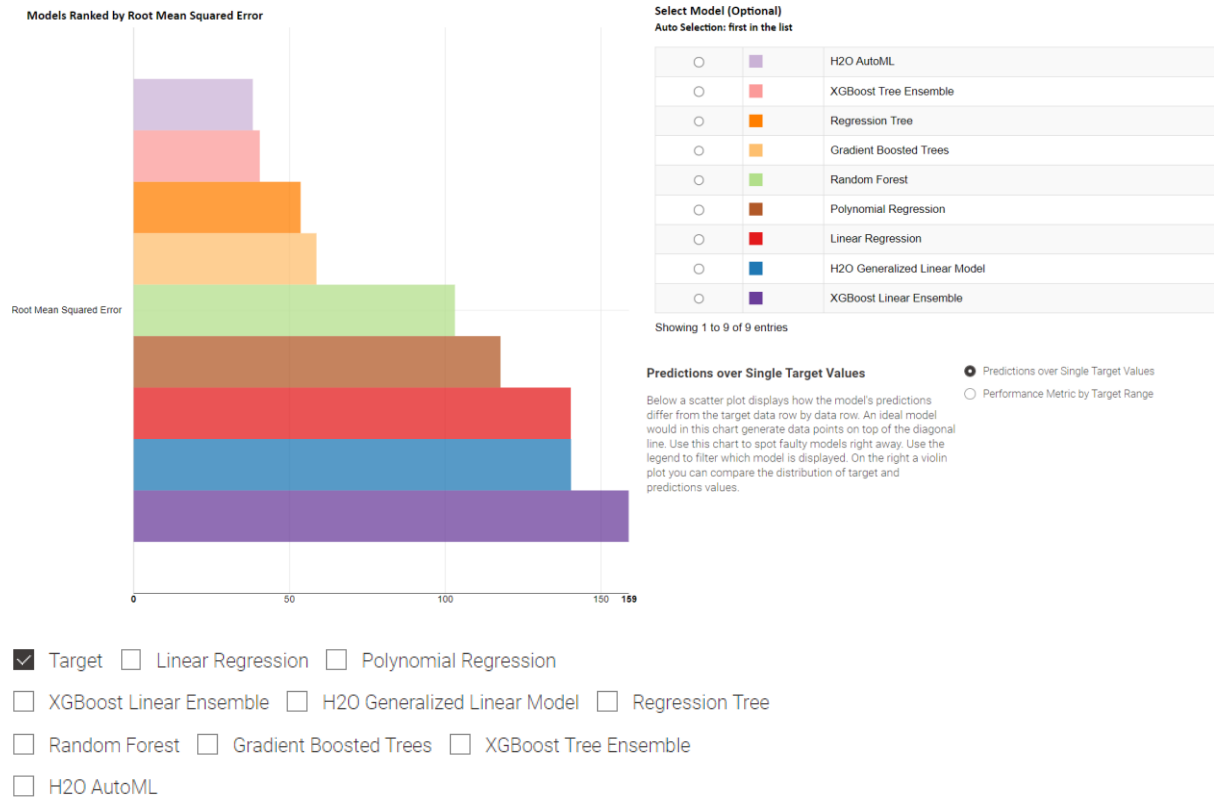
Demonstration AutoML with KNIME Workflows:

To illustrate this process, I've attached example KNIME workflows that incorporate various model learner nodes (Random Forest, Gradient Boosted Trees, Neural Network Learners) applied to a representative dataset. These workflows also include validation nodes, scoring nodes, and interactive views for comparing model performance metrics such as accuracy, ROC plot or F1 score. By examining the visualized outputs, one can see how different approaches perform and select the model that best aligns with the project's requirements.



In summary, I recommend starting with a variety of ML approaches—ensemble methods, neural networks, and hybrid/stacked models—and then systematically comparing their performance. This ensures an evidence-based selection of the best predictive model for your integrated omics dataset.





Following shows another AutoML example workflow and its output. The model predicts if customer churns from the telco company or not. The list of the chosen models for the dataset is as following.

1. Data Access

Connect to the pre-defined data source.

Insert database credentials to read from database instead.

User

Password

OR

Select example datasets to predict..

- ☒ churn of customers of telco company (customers.csv)
- ☐ ..income from census data (adult.csv)
- ☐ ..survival of individuals on Titanic (titanic.csv)
- ☐ ..species of flowers (iris.csv)
- ☐ ..quality of wines from chemical data (wine.csv)
- ☐ ..departure delay of flights from Chicago airport (flights.csv)

2. AutoML Settings

Define what you want to predict as target column and what columns should be used as input features. Finally select which models you would like to train.

Select Target

Churn

Select Models

[7 selected]

- ☒ Naive Bayes
- ☒ Logistic Regressi..
- ☒ Neural Network
- ☒ Gradient Booste..
- ☒ Decision Tree
- ☒ Random Forest
- ☒ XGBoost Trees
- ☐ Generalized Line..
- ☐ Deep Learning (K..
- ☐ H2O AutoML

Model Selection by:

F-measure

Select Input Columns

Excludes

Includes

- VMail Message
- Day Mins
- Eve Mins
- Night Mins
- Intl Mins
- CustServ Calls
- Day Calls
- Day Charge
- Eve Calls

Numeric Nominal Data Preview

Search:

Column	Minimum	Maximum
VMail Message	0	51
Day Mins	0	350.800
Eve Mins	0	363.700
Night Mins	23.200	395
Intl Mins	0	20
CustServ Calls	0	9
Day Calls	0	165
Day Charge	0	59.640
Eve Calls	0	170
Eve Charge	0	30.910
Night Calls	33	175
Night Charge	1.040	17.770
Intl Calls	0	20
Intl Charge	0	5.400
Area Code	408	510
Account Length	1	243

The result visually shows that XGBoost Trees model provides the best performance in this case to show the best performance based on F-measure of test set.

3. Results for Trained Model

Best Model Selected by F-measure :

XGBoost Trees

Next Actions

Controls for downloading the model (as a packaged workflow for KNIME Analytics Platform) and for deploying it (directly to KNIME Server).

[Download](#)

☒ **Deploy**

Statistics on XGBoost Trees

Stats are measured on new *Validation Data*

	churned (...)	not churn...
churned (...)	67	30
not churn...	3	567

Overall Accuracy

95.05%

MORE AUTOML MODELS

F-measure measured on *Test Data*

F-measure

Phone String	Account Number (inte..)	Churn String	Intl Plan Number (inte..)	VMail Plan Number (inte..)	State String
376-7145	78	not churned	0	0	IN
375-2975	55	churned	1	0	AL
376-8573	92	churned	1	0	ME
366-7360	129	not churned	0	1	RI
347-7898	18	not churned	0	0	MD
390-7328	161	not churned	1	0	FL
373-3251	144	not churned	0	0	AL
343-1965	75	churned	1	0	ME
378-8019	95	not churned	0	0	WV

By applying this approach, once the best-performing model is identified, it can be seamlessly deployed on the server and integrated into an automated production pipeline. This ensures that the selected model is consistently applied to new incoming data, maintaining a robust, real-time predictive service that can scale as the project and data grow.

Prediction heart disease

Example workflow for prediction using XGBoost Tree regression. Confusion matrix with more than 87% accuracy is shown in top-right corner.

