

AI & Machine Learning for Biologists

Classical ML • Deep Learning • Protein AI • Drug Discovery • XAI

Program Overview

This 21-day curriculum equips biologists with the practical AI and machine learning skills reshaping modern life sciences. Starting from zero programming assumption and building progressively through classical machine learning, deep learning, large language models for proteins, generative AI, and explainability, every concept is anchored to real biological questions: Can we predict cancer from expression data? Can AI fold a protein we have never seen before? Can we design new drug molecules computationally?

Foundations of AI & ML for Biology

What is AI/ML? A Biologist's Map

- AI vs Machine Learning vs Deep Learning — definitions with biological analogies
- Supervised, unsupervised, and reinforcement learning with biology examples
- Where ML is transforming biology: drug discovery, protein folding, genomics
- Python ecosystem: NumPy, Pandas, scikit-learn, TensorFlow, PyTorch overview

- ✓ Distinguish AI, ML, and deep learning
- ✓ Map ML paradigms to biological problems
- ✓ Set up Python ML environment (Colab)

Google Colab, NumPy, Pandas, scikit-learn

Survey: map 5 real biological problems to the correct ML paradigm

Data Preparation & Biological Feature Engineering

- Loading biological tabular data with Pandas (CSV, Excel, TSV)
- Handling missing values in clinical/genomic datasets
- Feature scaling: StandardScaler, MinMaxScaler for expression data
- Encoding categorical data: one-hot encoding of amino acids, nucleotides

- ✓ Clean and preprocess biological datasets
- ✓ Scale and normalise expression features
- ✓ Encode categorical biological variables numerically

Pandas, scikit-learn preprocessing

Preprocess a gene expression CSV: impute, scale, encode, export ready-to-use matrix

<i>Exploratory Data Analysis & Visualisation</i>	<ul style="list-style-type: none"> • Descriptive statistics on biological datasets • Correlation heatmaps of gene expression matrices (seaborn) • PCA scatter plots for sample clustering • Violin plots, box plots for comparing conditions with matplotlib/seaborn 	<ul style="list-style-type: none"> ✓ Perform EDA on biological datasets ✓ Interpret correlation patterns in expression data ✓ Visualise sample groupings with PCA 	Pandas, seaborn, matplotlib, scikit-learn PCA	EDA report on a cancer gene expression dataset: correlation heatmap + PCA plot
<i>Supervised Learning I — Classification</i>	<ul style="list-style-type: none"> • Train/test split, cross-validation, overfitting explained with biology • Logistic Regression for binary disease classification • K-Nearest Neighbours for cell type prediction • Evaluation: accuracy, precision, recall, F1, confusion matrix 	<ul style="list-style-type: none"> ✓ Split data and apply logistic regression ✓ Evaluate classifiers with biological context ✓ Avoid overfitting with cross-validation 	scikit-learn: LogisticRegression, KNeighborsClassifier, metrics	Classify tumour vs normal samples using gene expression features
<i>Supervised Learning II — Regression & Trees</i>	<ul style="list-style-type: none"> • Linear and Ridge Regression for predicting protein expression levels • Decision Trees for interpretable biological rules • Random Forest: ensemble method for biomarker discovery • Feature importance scores from Random Forest 	<ul style="list-style-type: none"> ✓ Apply regression to predict biological quantities ✓ Build and interpret decision trees ✓ Extract feature importances as biomarker candidates 	scikit-learn: LinearRegression, DecisionTreeClassifier, RandomForestClassifier	Predict drug IC50 from molecular features using Random Forest; rank top features
<i>Unsupervised Learning — Clustering & Dimensionality Reduction</i>	<ul style="list-style-type: none"> • K-Means clustering for gene expression pattern discovery • Hierarchical clustering with dendrograms (scipy) • UMAP and t-SNE for high-dimensional biological data • Silhouette score for optimal cluster selection 	<ul style="list-style-type: none"> ✓ Apply K-Means and hierarchical clustering to omics data ✓ Visualise high-dimensional data with UMAP/t-SNE ✓ Select optimal number of clusters biologically 	scikit-learn, scipy, umap-learn, seaborn	Cluster RNA-seq samples with K-Means; visualise with UMAP and label clusters
<i>Model Evaluation, Tuning & Pipelines</i>	<ul style="list-style-type: none"> • GridSearchCV and RandomizedSearchCV for hyperparameter tuning • ROC-AUC, precision-recall curves for imbalanced biological datasets 	<ul style="list-style-type: none"> ✓ Tune hyperparameters systematically ✓ Evaluate models on imbalanced biological data ✓ Build end-to-end reproducible ML pipelines 	scikit-learn: Pipeline, GridSearchCV, roc_curve, joblib	Build a tuned pipeline to classify cancer subtypes; plot ROC-AUC curve

	<ul style="list-style-type: none"> • scikit-learn Pipelines: preprocessing + model in one object • Saving and loading models with joblib for reproducibility 			
--	--	--	--	--

Deep Learning for Biological Sequences & Images

<i>Neural Networks & Deep Learning Foundations</i>	<ul style="list-style-type: none"> • Perceptron, activation functions (ReLU, sigmoid, softmax) — biological neuron analogy • Feedforward networks: layers, weights, backpropagation • Building an MLP with Keras Sequential API • Loss functions and optimisers: binary crossentropy, Adam 	<ul style="list-style-type: none"> ✓ Build and train a multi-layer perceptron in Keras ✓ Understand forward and backward propagation ✓ Apply appropriate loss functions for classification 	TensorFlow, Keras, NumPy	Train a neural network to classify cancer type from gene expression profiles
<i>Sequence Encoding & Embeddings for Genomics</i>	<ul style="list-style-type: none"> • One-hot encoding of DNA and protein sequences • k-mer frequency vectors for sequence representation • Word2Vec-style embeddings for amino acids (ProtVec concept) • Embedding layers in Keras for biological sequences 	<ul style="list-style-type: none"> ✓ Encode DNA and protein sequences for ML input ✓ Generate k-mer feature vectors ✓ Use embedding layers for sequence data 	NumPy, TensorFlow/Keras, Biopython	Encode 500 DNA sequences as one-hot and k-mer matrices; compare classification accuracy
<i>Convolutional Neural Networks for Sequence & Image Data</i>	<ul style="list-style-type: none"> • 1D CNNs for motif detection in DNA sequences • Conv1D, MaxPooling1D, Flatten, Dense layers in Keras • 2D CNNs for microscopy image classification • Transfer learning intro: using pre-trained image models for cell images 	<ul style="list-style-type: none"> ✓ Build 1D CNNs for sequence motif detection ✓ Apply 2D CNNs to biological image data ✓ Understand transfer learning for small datasets 	TensorFlow/Keras: Conv1D, Conv2D, MaxPooling	Train a 1D CNN to detect splice sites in DNA sequences
<i>Recurrent Neural Networks & LSTMs for Sequences</i>	<ul style="list-style-type: none"> • RNNs for sequential biological data: time series gene expression • LSTM and GRU layers to capture long-range sequence dependencies • Protein secondary structure prediction as a sequence labelling task • Bidirectional LSTMs for reading sequences in both directions 	<ul style="list-style-type: none"> ✓ Build LSTM models for biological sequence tasks ✓ Apply bidirectional LSTMs to protein sequences ✓ Model temporal gene expression with RNNs 	TensorFlow/Keras: LSTM, GRU, Bidirectional	Predict protein secondary structure (helix/sheet/coil) from amino acid sequence

<i>Transformers & Protein Language Models</i>	<ul style="list-style-type: none"> • Attention mechanism — biological sequence context analogy • Transformer architecture overview (BERT, GPT concept) • ESM-2 (protein language model) for protein representation • ProtTrans, ESMFold: embedding proteins and predicting properties 	<ul style="list-style-type: none"> ✓ Understand attention and transformer architecture ✓ Use pre-trained protein language models (ESM-2) ✓ Extract protein embeddings for downstream ML tasks 	HuggingFace transformers, fair-esm, ESM-2	Extract ESM-2 embeddings for 50 proteins; cluster by function using embeddings
<i>AlphaFold2 & Protein Structure Prediction</i>	<ul style="list-style-type: none"> • AlphaFold2 architecture: multiple sequence alignment + structure module • Running AlphaFold2 via ColabFold (Google Colab interface) • Interpreting pLDDT confidence scores and PAE plots • Visualising predicted structures with py3Dmol and nglview 	<ul style="list-style-type: none"> ✓ Run AlphaFold2 predictions via ColabFold ✓ Interpret confidence metrics (pLDDT, PAE) ✓ Visualise and compare predicted structures 	ColabFold, py3Dmol, nglview, Bio.PDB	Predict and visualise the structure of a novel protein; annotate confidence regions
<i>Generative Models — VAEs & GANs in Biology</i>	<ul style="list-style-type: none"> • Variational Autoencoders (VAE) for drug molecule generation • Generative Adversarial Networks (GAN) concept for synthetic biology data • Latent space exploration: interpolating between molecule properties • SMILES representation of molecules; RDKit for chemistry 	<ul style="list-style-type: none"> ✓ Understand VAE architecture and latent spaces ✓ Generate new molecule representations with VAEs ✓ Interpret latent space biology 	TensorFlow/Keras, RDKit, DeepChem	Train a VAE on SMILES drug molecules; sample new molecule candidates from latent space

Applied AI in Bioinformatics & Drug

<i>Graph Neural Networks for Molecular Biology</i>	<ul style="list-style-type: none"> • Graphs as molecules: atoms as nodes, bonds as edges • Message passing neural networks (MPNN) concept • PyTorch Geometric and DGL-LifeSci for molecular GNNs • Predicting molecular toxicity and bioactivity with GNNs 	<ul style="list-style-type: none"> ✓ Represent molecules as graphs for GNN input ✓ Build and train a simple molecular GNN ✓ Predict molecular properties from graph structure 	PyTorch Geometric, DGL-LifeSci, RDKit, DeepChem	Predict ESOL aqueous solubility for 100 drug molecules using a GNN
<i>AI for Drug Discovery & Virtual Screening</i>	<ul style="list-style-type: none"> • QSAR modelling: molecular descriptors as ML features 	<ul style="list-style-type: none"> ✓ Build QSAR models for bioactivity prediction 	DeepChem, RDKit, scikit-learn, mordred	QSAR model to predict hERG cardiotoxicity from molecular fingerprints

	<ul style="list-style-type: none"> • DeepChem for molecular property prediction (ADMET) • Docking score prediction vs physics-based docking • Active learning for hit identification in virtual screening 	<ul style="list-style-type: none"> ✓ Use DeepChem for ADMET property prediction ✓ Apply active learning to prioritise compounds 		
<i>Deep Learning for Genomics & Variant Effect Prediction</i>	<ul style="list-style-type: none"> • Predicting gene expression from DNA sequence (Enformer concept) • Variant effect scoring with CNN-based models • DeepSEA / Basset architecture for regulatory genomics • Interpreting models: saliency maps and in-silico mutagenesis 	<ul style="list-style-type: none"> ✓ Apply CNNs to regulatory sequence prediction ✓ Score variant effects on gene expression ✓ Use saliency maps to interpret model decisions 	TensorFlow/Keras, Biopython, Captum (PyTorch)	Train a CNN to predict enhancer vs non-enhancer sequences; compute saliency maps
<i>AI for Medical Imaging & Digital Pathology</i>	<ul style="list-style-type: none"> • Histopathology image classification with CNNs • Transfer learning with ResNet50 / EfficientNet for cell images • Data augmentation strategies for limited biomedical image datasets • Grad-CAM visualisation for model interpretability in pathology 	<ul style="list-style-type: none"> ✓ Fine-tune a pre-trained CNN on pathology images ✓ Apply augmentation to overcome limited data ✓ Visualise CNN decisions with Grad-CAM 	TensorFlow/Keras, torchvision, OpenCV, tf-explain	Fine-tune EfficientNet on H&E histology images to classify tumour grade
<i>Single-Cell AI — scRNA-seq & Cell Type Annotation</i>	<ul style="list-style-type: none"> • scVI: variational autoencoder for scRNA-seq integration • Scanpy + scANVI for automated cell type annotation • SCGAN and CellGAN for synthetic cell generation • Trajectory inference with PAGA and Monocle3 concepts 	<ul style="list-style-type: none"> ✓ Use scVI for batch-corrected scRNA-seq integration ✓ Automate cell type annotation with scANVI ✓ Infer differentiation trajectories from scRNA-seq 	scvi-tools, Scanpy, anndata, PyTorch	Integrate two scRNA-seq batches with scVI; annotate cell clusters automatically
<i>Explainable AI (XAI) & Model Interpretability in Biology</i>	<ul style="list-style-type: none"> • SHAP values for feature importance in biological ML models • LIME for local explanations of individual predictions • Attention weight visualisation in transformer-based protein models • Biological validation of ML-identified features and biomarkers 	<ul style="list-style-type: none"> ✓ Apply SHAP to explain gene importance in classifiers ✓ Use LIME for individual patient prediction explanations ✓ Validate ML biomarkers against biological literature 	SHAP, LIME, Captum, matplotlib	Apply SHAP to Random Forest cancer classifier; identify and validate top biomarker genes

<p><i>Capstone — End-to-End AI Biology Pipeline</i></p>	<ul style="list-style-type: none"> • Problem framing: biological question to ML pipeline design • Full pipeline: data ingestion → preprocessing → model training → evaluation → interpretation • Model card and reproducibility: documenting assumptions, data, metrics • Presenting AI biology results to mixed scientific audiences 	<ul style="list-style-type: none"> ✓ Design and execute a complete ML biology pipeline ✓ Document model decisions and limitations transparently ✓ Present findings to both computational and wet-lab audiences 	<p>All prior tools integrated; Colab, MLflow (optional)</p>	<p>Full AI pipeline on a chosen dataset: sequence, omics, or imaging with written report</p>
---	---	---	---	--

Ready to Build Real-World AI Solutions in Biology?

Join the next generation of researchers combining Biology + AI + Innovation.

APPLY NOW

Visit Website : www.omniedgesci.com