

Model Type	Best For	Key Features	Medical Applications	Strengths	Limitations
Logistic Regression	Binary classification, risk scoring	Interpretable coefficients, probabilistic output, linear decision boundary	Disease prediction (diabetes, cancer), mortality risk, readmission prediction	• Highly interpretable (feature importance) Fast to train and deploy Probabilistic output variance	• Limited to linear relationships Struggles with correlated features Underperforms with complex data
Decision Tree	Rule-based classification, feature importance	Tree structure, interpretable decision paths, handles mixed data types	Patient triage, diagnostic rule systems, treatment pathways	• Intuitive visualization Handles nonlinear relationships No feature scaling needed Captures feature interactions	• Prone to overfitting Unstable (high variance) Biased toward dominant classes Struggles with small datasets
Random Forest	Ensemble classification/regression, robust prediction	Multiple randomized trees, bagging, out- of-bag error	Medical imaging diagnosis, genomics, clinical outcome prediction	• Reduces overfitting Handles missing data feature importance Robust to outliers	• Less interpretable than single trees Computationally intensive Biased toward numerical features
Gradient Boosting	Structured/tabular data, high-performance prediction	Sequential boosting of weak learners, gradient optimization	Predicting treatment outcomes, length of stay, readmission risk	• State-of-the-art performance on tabular data Handles mixed data types Robust to outliers Feature importance	• Risk of overfitting Slower training Memory intensive interpretable

Model Type	Best For	Key Features	Medical Applications	Strengths	Limitations
SVM	Margin-based classification, small-to- medium datasets	Maximizes margin, kernel trick for non- linearity	Histopathology classification, protein structure prediction	• Effective in high-dimensional spaces Memory efficient Versatile through kernels Robust to overfitting in high dimensions	• Slow for large datasets Sensitive to parameter tuning box with non-linear kernels brobabilistic by default
k-NN	Instance-based learning, similarity search	No training phase, distance-based predictions	Similar patient profile lookup, rare disease diagnosis	• Zero training time br>• Intuitive approach • Naturally handles multi-class • Works with any distance metric	• Slow prediction time br>• Curse of dimensionality Sensitive to feature 
Naive Bayes	Text classification, high- dimensional sparse data	Probabilistic, assumes feature independence	Medical text classification, symptom- based diagnosis	• Fast training and prediction Works well with small data Handles high-dimensional data to irrelevant features	• "Naive" independence assumption Poor calibration of probabilities Sensitive to feature engineering

### Feature Processing Techniques

Technique	Purpose	Methodology	Medical Applications	Advantages	Limitations	Imple
PCA	Dimensionality reduction, feature extraction	Orthogonal projection, variance maximization	Gene expression analysis, medical image compression	• Reduces overfitting Handles multicollinearity Speeds up training visualization	• Loses interpretability Sensitive to scaling br>• May lose important information Linear transformation only	pca PCA(
t-SNE	Dimensionality reduction for visualization	Probability- based mapping preserving local structure	Visualizing high- dimensional biomedical data, clustering patients	• Preserves local structure Effective for visualization Handles non-linear relationships	• Computationally expensive Non- deterministic Not suitable for very large datasets Cannot project new samples	TSNE
UMAP	Fast non-linear dimensionality reduction	Manifold learning and topological data analysis	Single-cell genomics visualization, medical imaging analysis	• Faster than t- SNE SNE Preserves global structure Supports supervised version 	• Complex parameter tuning tuning intuitive than PCA br>• Stochastic results	(uma <sub>l</sub>
Feature Selection	Identifying relevant features	Statistical tests, model- based importance	Biomarker discovery, genomic analysis	• Improves model performance Reduces overfitting Faster training Better interpretability	• Risk of removing useful information computationally intensive br>• May require domain expertise	Sele k=10

# Deep Learning Models

Model Type	Architecture	Best For	Medical Applications	Advantages	Limitations
Convolutional Neural Network (CNN)	Convolutional layers, pooling, fully connected layers	Image classification, object detection, segmentation	Medical imaging (X-ray, CT, MRI), histopathology, dermatology	• Translation invariance Parameter efficiency Hierarchical feature learning Transfer learning potential	• Requires large datasets Computationally intensive box nature Struggles with global context
Recurrent Neural Network (RNN)	Sequential processing with memory	Time-series analysis, sequence prediction	EHR temporal analysis, ECG/EEG interpretation, patient monitoring	• Handles variable- length inputs • Captures temporal dependencies • Suitable for monitoring data	• Vanishing/exploding gradients br>• Difficult to parallelize Limited context window Complex training
Transformer	Self- attention, positional encoding	NLP tasks, long sequence analysis	Clinical text analysis, genomic sequence analysis, multi- modal fusion	• Handles long- range dependencies • Parallelizable training • State-of-the-art NLP performance • Attention interpretability	• High computational requirements br>• Large datasets needed Quadratic complexity with sequence length
Generative Adversarial Network (GAN)	Generator + Discriminator networks	Data synthesis, augmentation, translation	Medical image synthesis, data augmentation for rare conditions	• Creates realistic synthetic data • Useful for imbalanced datasets • Can generate rare cases • Domain adaptation	• Training instability • Mode collapse Difficult to evaluate considerations

Model Type	Architecture	Best For	Medical Applications	Advantages	Limitations
Autoencoder	Encoder- decoder architecture	Dimensionality reduction, denoising, anomaly detection	Medical image denoising, anomaly detection, feature learning	<ul> <li>Unsupervised</li> <li>learning br&gt;</li> <li>Noise</li> <li>reduction </li> <li>Feature</li> <li>extraction </li> <li>Anomaly detection</li> </ul>	• Not task- specific • Can learn trivial encodings • Requires careful architecture design
Graph Neural Network (GNN)	Message passing on graph structures	Relational data, network analysis	Protein-protein interactions, brain connectivity, drug discovery	• Models relationships explicitly Handles irregular structures Inductive learning Combines structure and features	• Complex implementation • Scaling issues • Limited interpretability • Domain-specific adaptation

# Advanced Learning Paradigms

Paradigm	Methodology	Best For	Medical Applications	Advantages	Limitations
Self- Supervised Learning	Learning representations from unlabeled data via pretext tasks	Pretraining for limited labeled data	Medical image representation, clinical text understanding	• Leverages unlabeled data Better representations Transfer learning efficiency	• Task design complexity Computational overhead Domain adaptation challenges
Meta-Learning	Learning to learn across tasks/domains	Few-shot learning, domain adaptation	Rare disease diagnosis, cross-hospital generalization	• Quick adaptation to new tasks • Data efficiency • Domain generalization	• Complex implementation Computationally intensive Limited interpretability
Reinforcement Learning	Learning through environment interaction	Treatment optimization, adaptive trials	Personalized treatment planning, dosage optimization	• Sequential decision modeling optimizes long-term outcomes Handles uncertainty	• Sample inefficiency • Safety concerns • Reward design difficulty • Exploration challenges
Federated Learning	Distributed training without sharing raw data	Privacy- preserving collaborative learning	Multi-hospital collaboration, privacy-sensitive data	• Privacy preservation Regulatory compliance Access to diverse data	• Communication overhead • Non- IID data challenges • System heterogeneity
Few-Shot Learning	Learning from very limited examples	Rare disease diagnosis, new medical conditions	Rare anomaly detection, novel disease classification	• Minimal labeled data needed Generalizes from few examples Reduces annotation burden	• Limited accuracy vs. fully-supervised • Complex architecture design br>• Task similarity requirements

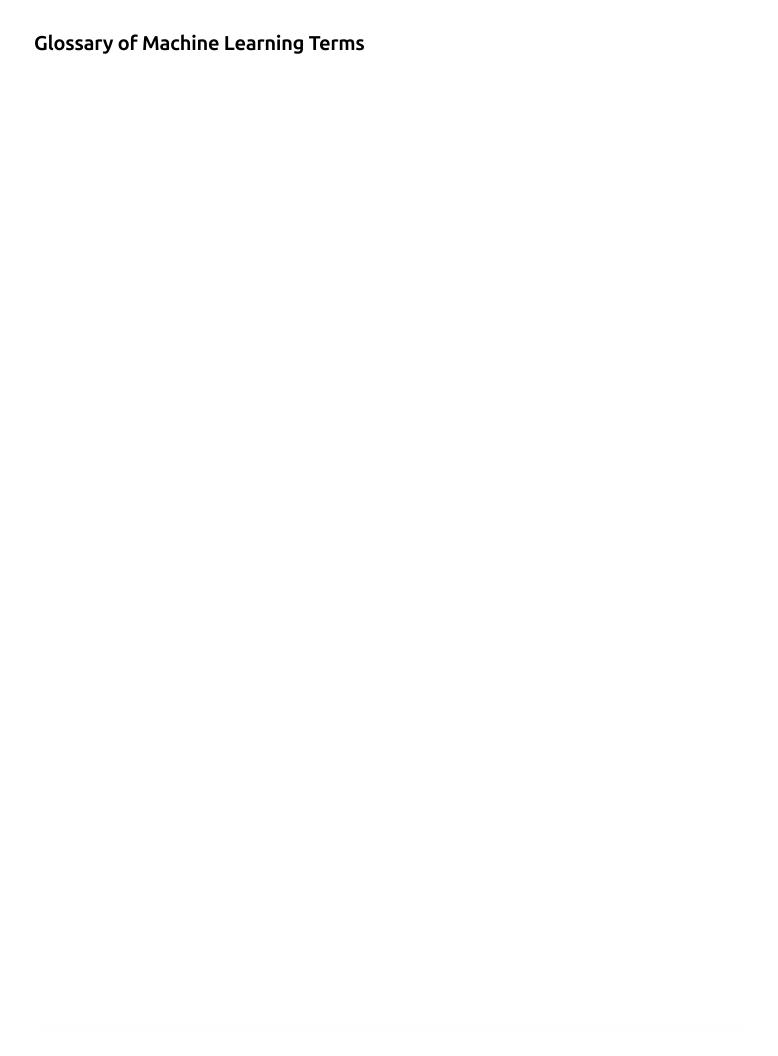
#### **Evaluation and Validation**

Method	Purpose	Implementation	Medical Considerations	Best Practices	Pitfalls	Tool
Cross- Validation	Robust performance estimation	k-fold, stratified, leave-one-out	• Patient-wise splitting Stratification by condition severity Temporal validation for longitudinal data	• Use stratified for imbalanced data Consider grouped CV for correlated samples Report variance across folds	• Data leakage Overfitting to validation data Temporal dependencies ignored	sciki tool split
ROC Analysis	Classification threshold optimization	ROC curve, AUC, precision-recall curve	• Sensitivity vs. specificity trade- off off br>• Different thresholds for screening vs. diagnosis Cost-sensitive evaluation	• Use PR curves for imbalanced data Consider clinical thresholds Report confidence intervals	• AUC insensitive to class imbalance Threshold optimization bias br>• Singlemetric focus	scikil matr
Statistical Testing	Model comparison, hypothesis validation	t-tests, McNemar's test, bootstrapping	• Power analysis for clinical significance • Multiple testing correction Confidence intervals for risk predictions	• Use appropriate tests for paired comparisons Bootstrap for robust intervals Account for dataset size	• p-hacking Misinterpretation of significance Ignoring effect size	scipy stats
Calibration Assessment	Reliability of probability estimates	Calibration curves, Brier score	• Critical for risk prediction br>• Impacts clinical decision thresholds Essential for shared decision making	• Use reliability diagrams br>• Apply calibration methods (Platt, isotonic) Separate calibration dataset	• Overlooking calibration Post-hoc adjustments limitations Dataset shift effects	sciki calib

#### Medical Domain-Specific Considerations

Aspect	Challenges	Solutions	Tools	Best Practices
				• Stratified
	Rare	• Class weighting •		sampling •
Class	conditions,	SMOTE/ADASYN •	imbalanced-learn,	Threshold
Imbalance	outcome	Focal loss •	WeightedRandomSampler	adjustment •
	imbalance	Ensemble methods		Consider prevalence-
				focused metrics
		• Multiple		
		imputation •		Analyze missingness
	Incomplete	Autoencoder		patterns •
Missing Data	medical	imputation •	missingpy, fancyimpute,	Document imputation
	records	Missingness	scikit-learn imputers	strategy •
		modeling • Domain-		Sensitivity analysis
		specific rules		
				• Involve clinicians in
	Clinical adoption requirements	SUAD L LINAT	SHAP, LIME, Captum, What-If Tool	interpretation •
		• SHAP values • LIME		Present explanations
Interpretability		explanations • Rule		alongside
		extraction 		predictions •
		Attention visualization		Balance accuracy vs.
				interpretability
		• Model versioning •		Performance
	ci I	Monitoring for		monitoring •
	Clinical	drift • Integration	MLflow, TensorFlow	Failsafe
Deployment	integration,	with EHR •	Serving, ONNX	mechanisms •
	regulatory	Regulatory		Clinical validation
		documentation		protocols
		Differential		Minimal sufficient
	LUDAA CODO	privacy • Federated	O	data • Formal
Data Privacy	HIPAA, GDPR	learning • Secure	OpenDP, TensorFlow	privacy
	compliance	enclaves • De-	Privacy, PySyft	guarantees •
		identification		Ethics board review
•	•		1	

Dataset	Domain	Size	Features	Tasks	Access	Notes
MIMIC-IV	Critical care	>40,000 patients	Vitals, labs, notes, medications, diagnoses	Mortality prediction, length of stay, phenotyping	PhysioNet (requires credentialing)	Deidentified ICU data from Beth Israel Deaconess
UK Biobank	Population health	>500,000 participants	Genetics, imaging, lifestyle, health outcomes	Disease risk prediction, biomarker discovery	Formal application process	Longitudinal study with extensive phenotyping
ChestX- ray14	Chest radiology	112,120 images	14 thoracic pathologies	Classification, localization, report generation	NIH open access	Frontal chest X- rays with disease labels
ADNI	Alzheimer's disease	>1,500 subjects	MRI, PET, genetic, clinical assessments	Progression prediction, biomarker identification	Application required	Longitudinal multi-modal Alzheimer's data
ISIC Archive	Dermatology	>150,000 images	Skin lesion images, metadata	Melanoma detection, lesion segmentation	Open access	International Skin Imaging Collaboration
TCGA	Cancer genomics	>11,000 patients	Genomic, proteomic, clinical data	Cancer subtyping, survival analysis, biomarker discovery	GDC Data Portal	Multi-omic data across 33 cancer types
PhysioNet Challenges	Various	Varies annually	Time series, signals, clinical data	ECG analysis, sepsis prediction, sleep staging	PhysioNet (some require credentialing)	Annual competitions with clinical relevance
APTOS	Ophthalmology	3,662 images	Retinal fundus images	Diabetic retinopathy grading	Kaggle	Asia Pacific Tele- Ophthalmology Society
•						<b>&gt;</b>



Term	Definition	Relevance to Medical AI
AUC-ROC	Area Under the Receiver Operating	Standard metric for diagnostic tests,
AUC-ROC	Characteristic - measures discrimination	accounts for all possible thresholds
Bias-Variance Tradeoff	Balance between underfitting (high bias)	Critical for generalizable clinical predictions
bids-variance fradeon	and overfitting (high variance)	across diverse populations
Confusion Matrix	Table showing TP, TN, FP, FN prediction	Maps to sensitivity/specificity in diagnostic
Coni asion Macrix	results	evaluation
Data Augmentation	Creating variations of training samples to	Helps with limited medical datasets,
Data Augmentation	increase diversity	especially imaging
Ensemble Learning	Combining multiple models for better	Improves robustness of clinical predictions
Ensemble Learning	performance	improves robustness of curriculty redictions
Feature Engineering	Creating new features from raw data	Critical for incorporating medical domain
reduce Engineering	Creating new reactives monitors acted	knowledge
Gradient Descent	Optimization algorithm for finding model	Fundamental training method for most ML
Greener	parameter minima	models
Hyperparameter	Model configuration not learned during	Requires careful tuning for medical
yperporeeee.	training	applications
Loss Function	Quantifies prediction error during	Should align with clinical objectives (e.g.,
	training	penalizing dangerous misses)
One-Hot Encoding	Converting categorical variables to binary	Common for medical categorical data like
	vectors	diagnoses, procedures
Overfitting	Model learns training data too well, fails	Risk with small or non-diverse medical
	to generalize	datasets
Precision/Recall	Metrics for positive prediction quality vs.	Precision: PPV in medical terms; Recall:
	completeness	sensitivity
Quadratic Weighted	Metric for ordinal classification	Standard for grading diseases with severity
Карра	agreement	levels (e.g., diabetic retinopathy)
Regularization	Techniques to prevent overfitting	Essential when features outnumber
		samples (common in genomics)
Sensitivity/Specificity	True positive rate vs. true negative rate	Fundamental diagnostic test evaluation
Эсполитену респисту	The positive race vs. a de freguence race	metrics
Transfer Learning	Using knowledge from one task to	Leverages general imaging features for
Transfer Learning	improve another	specific medical tasks
Validation Set	Data subset for tuning hyperparameters	Critical for unbiased model selection in
Validacion Sec	Data subsection turning hyperparameters	clinical applications
◀		<b>&gt;</b>