# Testing Metrics: Classic AI vs Generative AI Models

## Classification & Prediction Metrics

| Metric | Classic AI | Generative AI | What It Measures | How It Works |
|---|---|---|---|---|
| **Accuracy** | ✅ Core metric for classification tasks | ⚠️ Limited use, mainly for classification fine-tuning | Percentage of correct predictions | (True Positives + True Negatives) / Total Predictions |
| **Precision** | ✅ Essential for imbalanced datasets | ⚠️ Used in specific evaluation scenarios | Proportion of positive predictions that are correct | True Positives / (True Positives + False Positives) |
| **Recall (Sensitivity)** | ✅ Critical for detecting positive cases | ⚠️ Applied to content generation evaluation | Proportion of actual positives correctly identified | True Positives / (True Positives + False Negatives) |
| **F1-Score** | ✅ Balances precision and recall | ⚠️ Sometimes used for classification tasks | Harmonic mean of precision and recall | 2 × (Precision × Recall) / (Precision + Recall) |
| **AUC-ROC** | ✅ Standard for binary classification | ❌ Not applicable | Area under receiver operating characteristic curve | Plots True Positive Rate vs False Positive Rate |
| **Mean Squared Error (MSE)** | ✅ Standard for regression | ❌ Not directly applicable | Average squared differences between predicted and actual values | $\Sigma(y\_actual - y\_predicted)^2 / n$ |

## Language & Content Quality Metrics

| Metric | Classic AI | Generative AI | What It Measures | How It Works |
|---|---|---|---|---|
| **BLEU Score** | ❌ Not applicable | ✅ Translation and text generation quality | N-gram overlap between generated and reference text | Geometric mean of modified n-gram precisions with brevity penalty |
| **ROUGE Score** | ❌ Not applicable | ✅ Summarization quality | Overlap of n-grams, word sequences, and word pairs | Compares generated text to reference summaries using recall-based metrics |
| **METEOR** | ❌ Not applicable | ✅ Machine translation evaluation | Alignment between generated and reference text | Uses exact, stem, synonym, and paraphrase matches with precision/recall |
| **BERTScore** | ❌ Not applicable | ✅ Semantic similarity of generated text | Contextual embeddings similarity | Computes cosine similarity between BERT embeddings of generated and reference text |
| **Perplexity** | ❌ Not applicable | ✅ Language model fluency | How well model predicts text sequences | $2^{(-1/N \times \Sigma \log_2 P(word\_i}$ |

## Generative AI Specific Metrics

| Metric | Classic AI | Generative AI | What It Measures | How It Works |
|---|---|---|---|---|
| **Inception Score (IS)** | ❌ Not applicable | ✅ Image generation quality | Quality and diversity of generated images | Measures KL divergence between conditional and marginal label distributions |
| **Fréchet Inception Distance (FID)** | ❌ Not applicable | ✅ Image generation quality | Statistical distance between real and generated images | Computes Fréchet distance between feature distributions from Inception network |
| **CLIP Score** | ❌ Not applicable | ✅ Text-to-image alignment | Semantic similarity between text and generated images | Uses CLIP model to measure cosine similarity between text and image embeddings |
| **Human Evaluation** | ⚠️ Occasionally used | ✅ Essential for quality assessment | Subjective quality ratings | Human annotators rate outputs on relevance, fluency, creativity, factuality |
| **Diversity Metrics** | ❌ Not applicable | ✅ Output variety | Uniqueness across generated samples | Self-BLEU (lower is more diverse), distinct n-grams, semantic diversity |

## Robustness & Safety Metrics

| Metric | Classic AI | Generative AI | What It Measures | How It Works |
|---|---|---|---|---|
| **Adversarial Robustness** | ✅ Important for security | ✅ Critical for safe deployment | Resistance to malicious inputs | Tests model performance under adversarial attacks and input perturbations |
| **Bias Detection** | ✅ Fairness evaluation | ✅ Essential for responsible AI | Unfair treatment across demographic groups | Statistical parity, equalized odds, demographic parity across protected attributes |
| **Hallucination Rate** | ❌ Not applicable | ✅ Critical for factual accuracy | Frequency of generating false information | Automated fact-checking against knowledge bases or human annotation |
| **Toxicity Detection** | ⚠️ For content moderation models | ✅ Essential for text generation | Harmful or inappropriate content generation | Uses toxicity classifiers (Perspective API, custom models) to score outputs |
| **Calibration** | ✅ Confidence accuracy | ✅ Uncertainty quantification | Alignment between confidence and actual performance | Reliability diagrams, Expected Calibration Error (ECE) |

## Performance & Efficiency Metrics

| Metric | Classic AI | Generative AI | What It Measures | How It Works |
|---|---|---|---|---|
| **Inference Speed** | ✅ Important for real-time applications | ✅ Critical for user experience | Time to generate predictions/outputs | Measures latency from input to output completion |
| **Memory Usage** | ✅ Resource optimization | ✅ Deployment feasibility | RAM consumption during inference | Monitors peak memory usage during model execution |
| **Energy Consumption** | ⚠️ Growing concern | ✅ Major sustainability metric | Power usage during training/inference | Measures GPU/CPU power draw, carbon footprint calculations |
| **Scalability** | ✅ System design metric | ✅ Production deployment | Performance under increased load | Throughput vs latency trade-offs, concurrent user handling |

## Specialized Evaluation Approaches

| Metric | Classic AI | Generative AI | What It Measures | How It Works |
|---|---|---|---|---|
| **Cross-Validation** | ✅ Standard practice | ⚠️ Limited by computational cost | Model generalization | k-fold validation, stratified sampling, time-series splits |
| **A/B Testing** | ✅ Production evaluation | ✅ User preference measurement | Real-world performance comparison | Randomized controlled trials with different model versions |
| **Reinforcement Learning from Human Feedback (RLHF)** | ❌ Not applicable | ✅ Alignment optimization | Human preference learning | Trains reward models from human comparisons, optimizes policy with PPO |
| **Constitutional AI Evaluation** | ❌ Not applicable | ✅ Ethical behavior assessment | Adherence to principles and values | Tests model responses against predefined constitutional principles |
| **Red Team Testing** | ⚠️ Security testing | ✅ Safety evaluation | Identification of harmful capabilities | Systematic attempts to elicit dangerous or inappropriate outputs |

## Key Differences Summary

**Classic AI** focuses on:

- Statistical accuracy and error metrics
- Performance on specific, well-defined tasks
- Quantitative evaluation with ground truth
- Computational efficiency

**Generative AI** emphasizes:

- Content quality and human preference
- Safety and alignment considerations
- Subjective evaluation methods
- Emergent capabilities assessment

## Legend

- ✅ Primarily used and essential
- ⚠️ Sometimes used or limited application

- ❌ Not applicable or rarely used

- ❌ Not applicable or rarely used