

Study Guide: Training Models by Using AutoML

Subject: Google Cloud Professional Machine Learning Engineer Exam **Topics:**

- BigQuery ML Model Selection
- Feature Engineering & The TRANSFORM Clause
- Model Evaluation & Interpretability
- Time-Series Analysis (ARIMA_PLUS)
- Unsupervised Learning & Dimensionality Reduction

Summary

This study guide focuses on Vertex AI AutoML, Google Cloud's managed service for training high-quality custom machine learning models with minimal ML expertise. It covers the complete workflow from data preparation and feature engineering to model training, evaluation, and debugging across multiple data types. The guide emphasizes practical techniques for preparing datasets, selecting appropriate AutoML model types, leveraging managed feature stores, and implementing responsible AI practices including privacy-preserving techniques and bias mitigation.

Key Concepts

- **AutoML Workflow Philosophy:** AutoML democratizes ML by automating architecture search, hyperparameter tuning, and model selection. However, success still depends heavily on quality data preparation, appropriate feature selection, and proper training data labeling. AutoML handles the 'how' of model building, but practitioners must still define the 'what' (business objective) and 'why' (feature relevance).
- **Data Preparation for AutoML:** The most critical phase that determines model success. Key activities include feature selection (identifying relevant features that have predictive power while removing redundant or noisy features), data labeling (for supervised learning, high-quality labels are essential), and data splitting (AutoML automatically handles train/validation/test splits).
- **Tabular Workflows on AutoML:** AutoML Tables excels at structured data problems and automatically handles feature type detection, missing value imputation strategies, categorical encoding, feature scaling and normalization, feature crossing for interaction effects, and ensemble model training.
- **Training Custom Models by Data Type:** AutoML Tables for structured data, AutoML Vision for image classification and object detection, AutoML Natural Language for text classification and entity extraction, and AutoML Video for video classification and action recognition.

- **Forecasting Models with AutoML:** AutoML Forecasting provides specialized capabilities for time-series prediction, automatically detects and models trend, seasonality, and holiday effects, handles multiple time series with hierarchical relationships, and provides prediction intervals.
- **Feature Store Integration:** Vertex AI Feature Store provides centralized feature management, stores and serves features with low-latency online serving, maintains feature lineage and versioning, and prevents training-serving skew by using identical feature computation.

- **Model Evaluation and Interpretability:** AutoML provides comprehensive evaluation metrics including precision, recall, F1-score, AUC-ROC for classification; MAE, RMSE, R² for regression; and MAPE, quantile loss for forecasting. Feature importance scores help understand which features drive predictions.
- **Debugging and Optimization Strategies:** Address insufficient training data by collecting more diverse examples, handle class imbalance with class weights or oversampling, leverage domain-specific feature engineering, and adjust training budget for architecture exploration.
- **Responsible AI in AutoML:** Vertex AI includes Explainable AI for feature attributions, model monitoring to detect drift, fairness indicators to evaluate performance across demographic slices, and support for privacy preservation techniques.
- **Training-Serving Skew Prevention:** Use Feature Store to ensure identical feature computation in training and serving, monitor feature distributions, implement feature validation schemas, and set up automated retraining pipelines.

Vocabulary List

- **AutoML:** A managed service in Vertex AI that automates the process of applying machine learning to real-world problems, handling model architecture search, hyperparameter tuning, and model selection.
- **Feature Engineering:** The process of creating, transforming, and selecting features from raw data to improve model performance. AutoML automates many tasks but still benefits from domain-specific feature creation.
- **Training Budget:** The computational resources (measured in node hours) allocated for training. Higher budgets allow exploration of more model architectures but with diminishing returns.
- **Feature Store:** A centralized repository for storing, managing, and serving ML features with versioning, lineage tracking, and low-latency online serving.
- **Active Learning:** A machine learning approach where the algorithm identifies the most informative unlabeled examples to be labeled next, reducing total labeling effort.
- **Neural Architecture Search (NAS):** An automated technique used by AutoML to discover optimal neural network architectures for a given task.
- **Ensemble Learning:** A technique where multiple models are combined to produce better predictions than any individual model. AutoML Tables automatically creates ensembles.
- **Data Drift:** Changes in the statistical distribution of input features over time in production, which can degrade model performance.

- **Explainable AI (XAI):** Techniques for making ML model predictions interpretable and understandable to humans through feature attributions.
- **Training-Serving Skew:** A critical ML production issue where the data distribution or feature computation differs between training and production serving.

Key Questions

- 1. You're building a classification model with AutoML Tables and notice a large gap between training accuracy (95%) and validation accuracy (75%). What are the three most likely causes, and what specific actions would you take to address each?**
- 2. A retail company wants to forecast daily sales for 500 store locations using AutoML Forecasting. They have 2 years of historical data and want to incorporate planned promotional events. How should they structure their dataset, and what AutoML Forecasting features should they leverage?**
- 3. You're deploying an AutoML Vision model for medical image classification. The model performs well overall (90% accuracy) but shows significantly lower performance (60% accuracy) for images from one particular hospital. What steps would you take to diagnose this issue?**
- 4. Explain the difference between using AutoML Tables directly on raw CSV data versus building a Feature Store pipeline first. In what scenarios would the additional complexity be justified?**
- 5. You're training an AutoML Natural Language model for sentiment analysis and need to label 100,000 customer reviews with a limited labeling budget. Describe how you would use Vertex AI's managed labeling service with active learning to minimize costs.**