# Production Stack Guide for Lead Scoring Model

## 1. Hosting

Recommended Platform: GCP (Google Cloud Platform)

GCP offers a comprehensive set of tools and services for deploying and managing machine learning models. It provides robust infrastructure, scalable services, and seamless integration with various machine learning workflows. Also, I have more knowledge with GCP.

**Deployment Strategy:**

1. Containerization with Docker:
   - Advantages: Consistency across different environments, ease of scaling, and isolation of dependencies.
   - Implementation: Use Docker to containerize FastAPI applications for example, including the model and preprocessing pipeline.

2. Orchestration with Google Kubernetes Engine (GKE):
   - Advantages: Managed Kubernetes service that simplifies deployment, management, and scaling of containerized applications.
   - Implementation: Deploy the Docker container to GKE for automated scaling, monitoring, and management.

3. Storage with Google Cloud Storage:
   - Advantages: Reliable, scalable, and secure storage for model artifacts, input data, and logs.
   - Implementation: Store the trained model, scaler, and any other necessary files in a Google Cloud Storage bucket, which can be accessed by the GKE pods.

4. Networking with Google VPC and Load Balancer:
   - Advantages: Ensures secure communication, scalability, and availability.
   - Implementation: Use Google VPC to isolate the network environment and Google Cloud Load Balancer to distribute incoming requests across multiple GKE pods.

## 2. Monitoring

Tools and Strategies for Continuous Monitoring

**Monitoring Components:**

1. Application Monitoring:
   - Tools: Google Cloud Monitoring

- Metrics to Monitor: API request counts, response times, error rates, CPU and memory usage.
   - Implementation: Set up Google Cloud Monitoring to collect and track metrics, log files, and set up alerts.

2. Model Performance Monitoring:
   - Tools: MLFlow
   - Metrics to Monitor: Prediction accuracy, latency, drift in data distribution, model confidence scores.
   - Implementation: Use MLFlow to log model performance metrics and track them over time to detect changes in performance.

3. Alerting:
   - Tools: Google Cloud Monitoring, PagerDuty
   - Implementation: Configure Cloud Monitoring to set up alerts and integrate with PagerDuty for incident management.

## 3. Performance Measurement

Measuring Effectiveness and Efficiency of the Model

**Key Metrics:**

### 1. Accuracy Metrics

**Examples:**
   • Precision
   • Recall
   • F1-Score
   • ROC-AUC
   • Accuracy

Before selecting the metric/metrics, it's crucial to select the right accuracy metrics based on the specific needs and impact on our business.

**Precision:**
   • **Use Case:** Precision is particularly important when the cost of false positives is high. For example, if false positives (incorrectly predicting a lead as a potential customer) result in wasted sales resources and efforts, then maintaining a high precision would be crucial.
   • **Tracking:** Monitor precision if false positives have a significant negative impact on your business.

**Recall:**
- **Use Case:** Recall is critical when missing true positives (false negatives) is costly. For instance, if failing to identify a potential lead results in significant missed revenue opportunities, then recall should be prioritized.
- **Tracking:** Focus on recall if false negatives have a more severe negative impact on your business.

**F1-Score:**
- **Use Case:** The F1-Score is useful when both false positives and false negatives carry significant costs, and a balance between precision and recall is desired.
- **Tracking:** Use the F1-Score when both types of errors (false positives and false negatives) are equally important and need to be minimized.

**ROC-AUC (Receiver Operating Characteristic - Area Under Curve):**
- **Use Case:** This metric is useful for evaluating the overall discriminative power of the model, especially in cases where you need to balance sensitivity (true positive rate) and specificity (true negative rate).
- **Tracking:** Use ROC-AUC for a comprehensive view of model performance across all classification thresholds.

**Accuracy:**
- **Use Case:** Accuracy is a useful metric when the classes are balanced and you want to get a general sense of the model's performance.
- **Tracking:** Track accuracy to understand the general performance of the model, especially when the dataset has a balanced class distribution.

**Stakeholder Discussions**

Before finalizing which metrics to track, it's essential to have detailed discussions with stakeholders to understand the business context and the specific impacts of false positives and false negatives.

2. Operational Metrics:
   - Examples: API response time, throughput, system uptime
   - These metrics measure the efficiency and reliability of the deployed system.

3. Business Metrics:
   - Examples: Conversion rates of scored leads, ROI from leads
   - These metrics link model performance to business outcomes, validating its impact on revenue and growth.

## Implementation Strategy

1. Automated Evaluation:
   - Implementation: Schedule regular evaluation runs using Cloud Functions or Cloud Run to compute performance metrics on recent data and compare them to historical benchmarks.

2. Drift Detection:
   - Implementation: Monitor changes in input data distribution and model predictions. Use MLFlow to log and analyze these metrics for detecting data and model drift.

3. Periodic Retraining:
   - Implementation: Automate the retraining process using Cloud AI Platform Pipelines or Google Cloud Functions. Retrain the model on new data periodically to maintain its accuracy and relevance.

## Possible Implementation Steps

1. Containerization and Deployment on GCP

Containerization:
- Dockerize your FastAPI application.

Deployment:
- Deploy the Docker container to GKE.

2. Setting Up Monitoring with Google Cloud Monitoring

Application Monitoring:
- Set up Google Cloud Monitoring to track metrics and logs.
- Create dashboards and set up alerts for key metrics.

3. Tracking Model Performance with MLFlow

MLFlow Setup:
- Set up an MLFlow server on GCP, possibly using Google Cloud Run or a VM instance.
- Configure model training script to log metrics to MLFlow.