

Plan de Machine Learning

para el Hackathon

3 caminos prácticos para agregar ML al MVP (sin complicar).

Fecha: 10/02/2026

Objetivo: proponer 3 implementaciones de ML compatibles con lo que ya existe (backend + UI + datos) y con valor claro para el pitch.

- Mantener el MVP estable: ML como servicio opcional y trazable.
- Priorizar demostrabilidad: resultados visibles en la UI (Admin) y exportables (CSV).
- Permitir dividir el trabajo entre 5 integrantes de Data Science.

Repositorio: https://github.com/brodyandre/growth_equestre_hackathon_2026

Stack: Backend Node/Express · Postgres · UI Admin Streamlit · Docker Compose.

1) Lo que ya tenemos en el MVP

El proyecto ya funciona end-to-end: captura leads, calcula score y permite explorar socios/partners. La idea es acoplar ML sin romper el flujo actual.

UI Admin (Streamlit)	Backend (Node/Express)	Base de datos (Postgres)
Filtros · KPIs · Export CSV Acciones: score/handoff	Endpoints: /leads, /events, /partners, /handoff	Tablas: leads, events, partners, cnae_map

Salida esperada del ML: un nuevo endpoint (o ampliación del score existente) que devuelva predicciones/recomendaciones, y una sección visual en la UI para mostrarlo.

The screenshot shows the Streamlit-based Admin interface for the Growth Equestre (MVP) project. On the left, there's a sidebar with navigation links like 'Visão geral', 'Leads', 'CRM (Kanban)', 'Parceiros', 'Criar lead (demo)', and 'Roteiro de demo'. Below that is a 'Configurações' dropdown and a link to the active backend at 'http://localhost:3000/health'. The main area has a title 'Admin — Growth Equestre (MVP)' and a subtitle 'Acompanhe leads, priorize atendimento e explore parceiros por UF/Segmento — sem telas técnicas.' It features several cards with metrics: 'Leads (total) 14', 'Qualificados 1', 'Aquecendo 7', and 'Conversão p/ qualificado 7.1%'. Below this is a section titled 'Parceiros (diretório)' with a dropdown for 'UF' set to 'UF'. It shows four boxes: 'Total de parceiros 1003', 'Cavalos 251', 'Serviços 251', and 'Equipamentos 251'. A 'Resumo completo (tabela)' table follows, showing data for segments: Cavalos (251), Serviços (251), Equipamentos (251), and Eventos (250). At the bottom, there's a 'Como usar (bem direto)' section with a note about leads being leads, leads being leads, leads being leads, and a 'Deploy' button in the top right corner.

Figura: UI Admin (referencia).

2) Tres posibilidades de ML (vista rápida)

Las 3 propuestas usan datos que ya existen o que se pueden generar con poca instrumentación adicional.

Idea	Qué entrega	Datos mínimos	Dónde se integra
(A) Propensión de conversión de lead	Probabilidad de volverse “CALIFICADO” y sugerencia de acción	Leads + eventos del funnel + score actual	Backend: /leads/:id/score o /ml/lead-propensity UI: Leads / Visión general
(B) Recomendación de partners (matching)	Top-N partners por lead (ranking) + explicación simple	Lead (segmento/UF/ciudad) + partners (UF/segmento/CNAE)	Backend: /recommendations/partners?lead_id=... UI: Leads (detalle) / Partners
(C) Segmentación + anomalías	Clusters de leads/partners + alertas (outliers / duplicados)	Leads/partners con features numéricas/categóricas	Backend: /ml/segments y /ml/anomalies UI: KPIs, Kanban, listas

Recomendación: empezar por (A) o (B). Son fáciles de demostrar, generan números claros y se conectan bien con el flujo del producto.

3) (A) Propensión de conversión / calificación de lead

Objetivo: predecir, al momento de crear/actualizar un lead, la probabilidad de que se vuelva CALIFICADO y sugerir la próxima acción (priorizar vs nutrir).

Por qué encaja en el hackathon: convierte el scoring en “ML de verdad” y permite mostrar un KPI de mejora (ej.: mayor conversión).

Pasos (end-to-end)

- 1 Construir dataset: unir leads + events (conteos por tipo) + atributos (UF, segmento, presupuesto, plazo).
- 2 Definir etiqueta: CALIFICADO=1, resto=0 (o multi-clase por estado).
- 3 Entrenar baseline: Logistic Regression / RandomForest / XGBoost (si disponible).
- 4 Evaluar: AUC/ROC, precision@k (top leads), matriz de confusión.
- 5 Exportar modelo: joblib/pickle + versionado.
- 6 Servir predicción: endpoint nuevo o dentro del score existente.
- 7 UI: mostrar “Probabilidad” + “Sugerencia” + 2-3 razones (features importantes).

Archivos/pastas sugeridos para trabajar

- /scoring o /ml: notebooks/scripts de entrenamiento (pandas, scikit-learn).
- backend/routes: agregar endpoint (ej.: /ml/lead-propensity) o ampliar /leads/:id/score.
- ui_admin (Streamlit): sección en Leads → Detalle del lead.
- data/: dataset versionado (CSV/Parquet) + artefacto del modelo.

The screenshot shows the 'Admin — Growth Equestre (MVP)' interface. On the left, there's a sidebar with navigation links like 'Visão geral', 'Leads', 'CRM (Kanban)', 'Parceiros', 'Criar lead (demo)', and 'Roteiro de demo'. Below it are 'Configurações' and 'Backend ativo: http://localhost:3000/health'. The main content area has a title 'Admin — Growth Equestre (MVP)' and a subtitle 'Acompanhe leads, priorize atendimento e explore parceiros por UF/segmento — sem telas técnicas.' It features a 'Roteiro de demo (para o pitch)' section with steps: 1) Gerar cenário completo em 1 clique (with a button 'Criar cenário de demo (3 leads)'), 2) O que mostrar (ordem recomendada) (with a list of 5 steps: Visão geral, Leads, Handoff, Parceiros, Export CSV), and 3) Checklist rápido (antes de apresentar) (with a bulleted list: Backend UP, UI OK, Criar cenário demo, Abrir Leads e filtrar por status/score, Mostrar motivos do score). There's also a 'Reset demo (limpar leads/parceiros demo)' button. In the top right corner, there are 'Deploy' and three-dot menu icons.

Figura: la UI ya contiene flujo de demo guiado; el modelo puede alimentar este guion.

4) (B) Recomendación de partners (matching lead ↔ partner)

Objetivo: dado un lead, devolver una lista de partners recomendados (Top-N) para facilitar el “handoff” y la prospección.

Cómo funciona (versión simple)

- Filtro duro: UF del lead (o UF vecina) + segmento compatible.
- Ranking: similitud de texto (TF-IDF) usando palabras-clave / CNAE / nombre (partner) + interés del lead.
- Explicación: mostrar 1-2 motivos (“mismo segmento”, “UF coincide”, “CNAE relacionado”).

Pasos (end-to-end)

- 1 Preparar base de partners: normalizar CNAE/segmento/UF y generar texto base por partner.
- 2 Entrenar o configurar: TF-IDF + cosine similarity (no necesita GPU).
- 3 Endpoint: /recommendations/partners?lead_id=... devuelve lista ordenada.
- 4 UI: en el detalle del lead, botón “Ver partners recomendados” + export CSV.
- 5 Validación rápida: comprobar si el Top-N tiene sentido para 10 leads de demo.

nome_fantasia	segmento	uf	cidade	prioridade	cnae	contato	razao_social	cnpj
Haras Exemplo	🐴 Cavalos	SP	Sao Paulo	1	0152-1/02		Haras Exemplo LTDA	12345678000190
Parceiro Demo 1	🔴 Serviços	MG	Uberlândia	1	9313-1/00		Parceiro Demo 1 LTDA	16862858863977
Parceiro Demo 10	🏡 Eventos	MG	Belo Horizonte	1	8230-0/01		Parceiro Demo 10 LTDA	13605362735169
Parceiro Demo 100	🐴 Cavalos	MG	Belo Horizonte	1	0152-1/02		Parceiro Demo 100 LTDA	16468704446293
Parceiro Demo 1000	🐴 Cavalos	MG	Belo Horizonte	1	0152-1/02		Parceiro Demo 1000 LTDA	1027214333738
Parceiro Demo 101	🔴 Serviços	GO	Anápolis	1	9313-1/00		Parceiro Demo 101 LTDA	10449795398677
Parceiro Demo 102	🏡 Eventos	SP	Ribeirão Preto	1	8230-0/01		Parceiro Demo 102 LTDA	10446585159571
Parceiro Demo 11	Equipamentos	GO	Anápolis	1	4647-0/01		Parceiro Demo 11 LTDA	14919813940032
Parceiro Demo 110	🏡 Eventos	GO	Goiânia	1	8230-0/01		Parceiro Demo 110 LTDA	16025883166039
Parceiro Demo 111	Equipamentos	SP	Campinas	1	4647-0/01		Parceiro Demo 111 LTDA	17363460136624

Figura: pantalla de Partners; es el lugar natural para mostrar recomendaciones/filtrado inteligente.

5) (C) Segmentación y detección de anomalías (para insights)

Objetivo: agrupar leads/partners en perfiles (clusters) y detectar casos fuera del patrón (outliers), mejorando priorización y calidad de datos.

Casos de uso típicos

- Segmentos de leads por intención: alto presupuesto + plazo corto vs curiosos.
- Partners “estrella” por densidad de demanda (UF/segmento) y CNAE.
- Alertas de datos: leads duplicados, valores incoherentes, partners con CNAE raro para el segmento.

Pasos (end-to-end)

- 1 Definir features: one-hot de segmento/UF + presupuesto/plazo (ordinal) + conteos de eventos.
- 2 Clustering: KMeans (3–6 clusters interpretables).
- 3 Anomalías: IsolationForest o LocalOutlierFactor para marcar top outliers.
- 4 Endpoint(s): /ml/segments y /ml/anomalies.
- 5 UI: mostrar “Cluster” en la tabla y una lista “Alertas” para revisión.

The screenshot shows the Growth Equestre (MVP) application interface. On the left, there's a sidebar with navigation links like 'Visão geral', 'Leads', 'CRM (Kanban)', 'Parceiros', 'Criar lead (demo)', and 'Roteiro de demo'. Below that is a 'Configurações' section and a 'Backend ativo: http://localhost:3000/health' status message. The main area is a Kanban board with three columns: 'AQUELENDO', 'AQUECENDO', and 'AQUECENDO'. Each column has a dropdown menu with a 'Mover' button. The first column contains a card for 'Visitante Demo' with details: 'AQUELENDO - score=40', 'UF: MG', 'Cidades: Belo Horizonte/MG', 'Proxima ação: não definida'. The second and third columns contain cards for 'Lead Demo - Médio' with similar details. To the right of the board, there's a 'Parceiros sugeridos' section with a dropdown showing '1. Parceiro (SP) - EVENTOS - compat=90'. Below it is a 'Resumo do parceiro' section with fields: 'Nome: Parceiro', 'UF: SP', 'Cidade: São Paulo', 'Segmento: EVENTOS', 'CNAE: 8230-0/01', 'Contato: Sem contato', and 'Compatibilidade: 90'. There are buttons for 'Enviar para parceiro selecionado' and 'Quando? hoje 09:00'. At the bottom, there's a 'Notas' section with a 'Adicionar nota' button and a 'Histórico' button.

Figura: vista Kanban; la segmentación/anomalías puede alimentar badges y alertas.

6) Plan de trabajo para 5 integrantes de Data Science

Recomendación: elegir 1 opción como “principal” para el pitch y dejar otra como “extra” si hay tiempo.

Rol	Responsabilidad	Entregable
DS-1 (Data/ETL)	Dataset (leads+events+partners), limpieza, versionado	dataset_v1.csv + notebook ETL
DS-2 (Model A)	Propensión de lead (A): features, entrenamiento, métricas	model_lead.joblib + reporte
DS-3 (Model B)	Recomendación (B): TF-IDF, ranking, explicaciones	reco_partners.py + ejemplos
DS-4 (Integración)	Endpoints en backend + contratos JSON + tests básicos	routes /ml/* + docs
DS-5 (Producto/Demo)	UX en UI + guion demo + KPIs + validación manual	pantallas actualizadas + script demo

Checklist final (antes de presentar)

- Backend activo y endpoints de ML respondiendo (health + /ml/*).
- Modelos/artefactos versionados y reproducibles (script de entrenamiento).
- UI mostrando predicciones/recomendaciones sin términos técnicos.
- Datos de demo listos (3 leads) y un ejemplo de “antes vs después”.