

MLOps & Monitoring Plan

*** **4.11.1 مقدمه**

اگرچه در حال حاضر مبتنی بر قوانین ساده `nobatnou.ir` سیستم پیش‌بینی زمان انتظار در ارتقاء خواهد یافت. این سند (ML) است، اما به تدریج به مدل‌های یادگیری ماشین (rule-based) و **مانیتورینگ مستمر** را تعریف (در محیط عملیاتی ML اداره مدل‌های) **MLOps** چارچوب می‌کند تا اطمینان حاصل شود که مدل‌ها قابل اعتماد، قابل نگهداری و فاقد ریسک‌های عملیاتی هستند.

*** **4.11.2 اهداف MLOps**

- مدل در محیط تولید همواره در دسترس و پایدار باشد (**Reliability): قابلیت اطمینان**.
- هر نسخه از مدل و نتایج آن قابل بازتولید باشد (**Reproducibility): قابلیت بازتولید**.
- فرآیندهای آموزش و استنتاج با رشد داده‌ها مقیاس پذیرند (**Scalability): مقیاس‌پذیری**.
- به طور مستمر پایش شود drift عملکرد مدل، کیفیت داده و (**Monitoring): مانیتورینگ**.
- کاهش مداخله دستی در چرخه حیات مدل (**Automation): خودکارسازی**.

*** **4.11.3 چرخه حیات مدل (Model Lifecycle)**

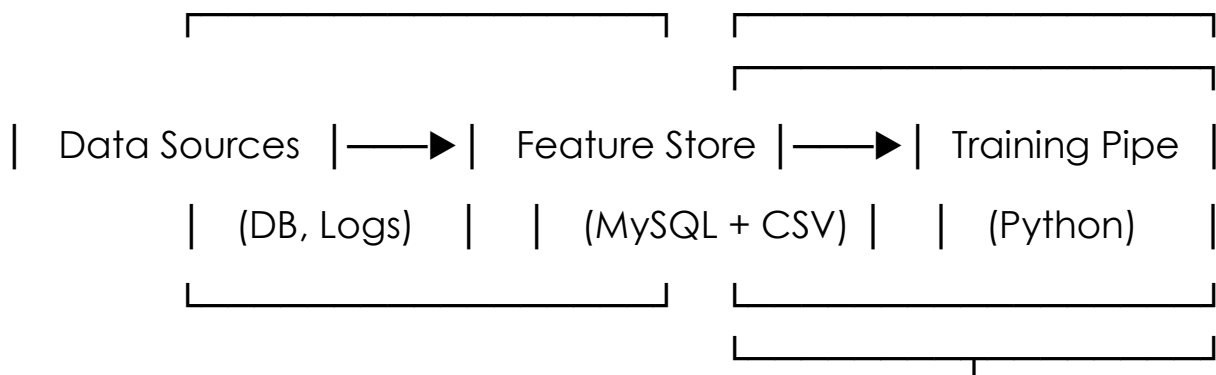
:مراحل چرخه حیات مدل در نوبت نو

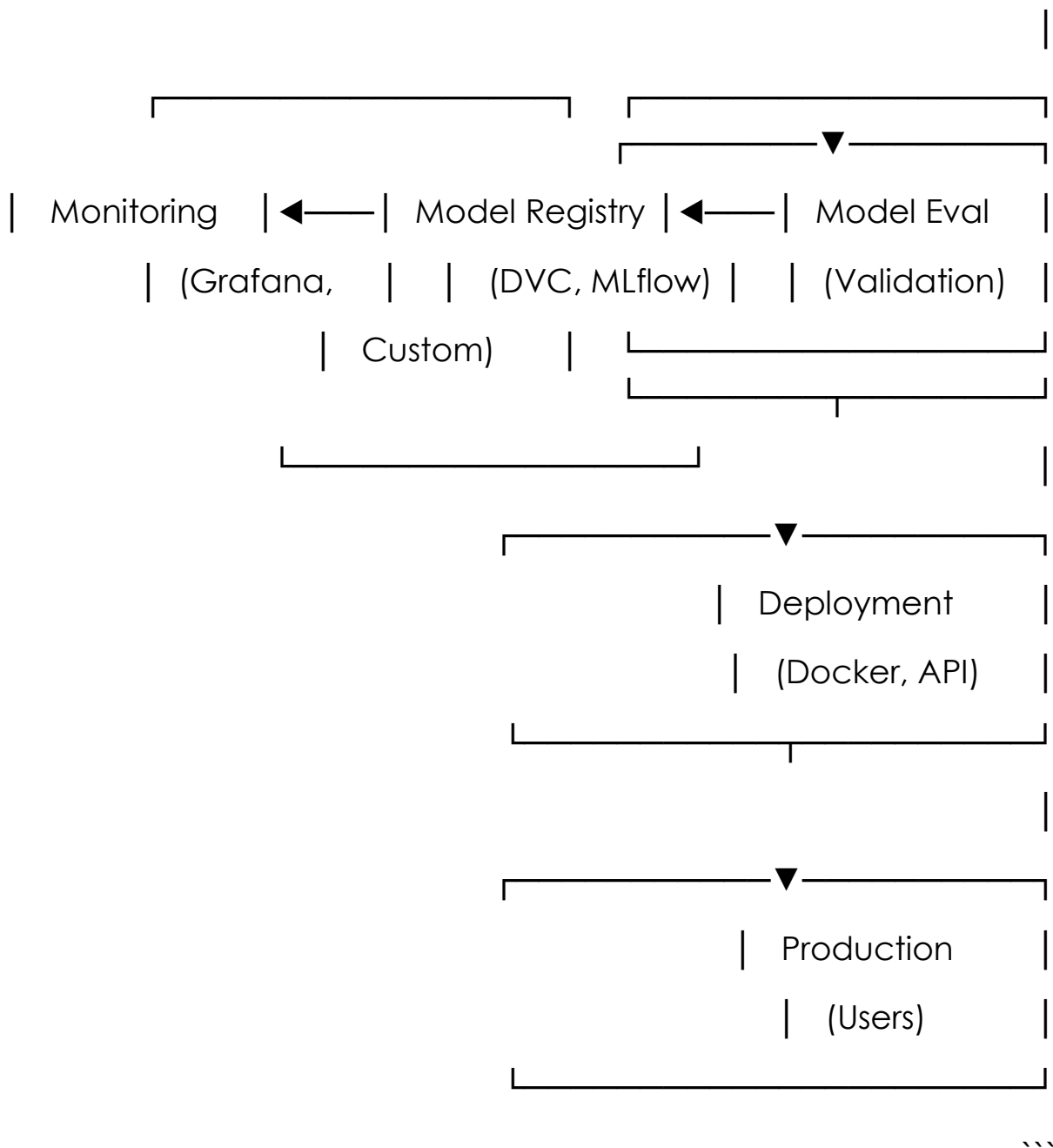
1. **جمع‌آوری داده (Data Collection)**
2. **آماده‌سازی داده (Data Preparation)**
3. **آموزش مدل (Model Training)**
4. **ارزیابی (Evaluation)**
5. **تست استقرار (Deployment Testing)**
6. **استقرار (Deployment)**
7. **مانیتورینگ و پایش (Monitoring & Observability)**
8. **بازآموزی (Retraining)**
9. **نسخه‌بندی و بایگانی (Versioning & Archival)**

4.11.4 معماری MLOps (MLOps Architecture)

ساختار پیشنهادی (بر اساس مقیاس فعلی) ###

...





اجزای کلیدی

- **Feature Store**: مثلاً ذخیره ویژگی‌های از پیش محاسبه‌شده
(`avg_wait_time_last_week`).

- **Model Registry**: **DVC** یا **MLflow** ثبت و نسخه‌بندی مدل‌ها با استفاده از.

- **GitHub Actions** اسکریپت‌های خودکار آموزش با **Training Pipeline**.
- **Docker** درون **FastAPI** سرویس مدل با **Serving API**.
- **Monitoring Dashboard** نمایش متریک‌های مدل و داده.

Model Monitoring (مانیتورینگ مدل 4.11.5)

Performance Monitoring (الف) مانیتورینگ عملکرد

متریک‌های کلیدی تحت پایش

متریک	هدف	آستانه هشدار	ابزار اندازه‌گیری
دقت پیش‌بینی	$< 20\%$ دقیقه	محاسبه روزانه بر	MAE (Mean Absolute Error)
درصد پیش‌بینی‌های درون محدوده	$> 50\%$ (محدوده $\pm 20\%$ دقیقه)	محاسبه	
APM (New Relic/Sentry)	p95) ثانیه > 1	کارایی	Latency زمان پاسخگویی مدل
Nginx	نرخ درخواست‌های ناموفق	در دسترس بودن	$< 5\%$ لاگ‌های

هشدارها

از 20 دقیقه عبور کرد → **هشدار سطح 2** (ایمیل به تیم داده) MAE اگر -

از ۳۰ دقیقه عبور کرد → **هشدار سطح ۱** (تماس با مدیر فنی) MAE اگر -

** (انحراف) Drift (ب مانیتورینگ) ****

نحوه تشخیص	فرکانس	آستانه هشدار	Drift نوع
-----	-----	-----	-----
بین هفته (`hour_of_day` مانند) مقایسه توزیع ویژگی‌های ورودی **Data Drift**	KL Divergence > 0.2 جاری و گذشته هفتگی		
هفته جاری با میانگین ۴ هفته گذشته هفتگی افزایش MAE مقایسه **Concept Drift**	MAE > ۱۰٪		
ماهانه (`actual_wait` تغییر توزیع زمان واقعی انتظار **Label Drift**	Kolmogorov-Smirnov test p-value < 0.05		

ابزار تشخیص

- Python: `alibi-detect`, `evidently` کتابخانه‌های -

- Dashboard: **Grafana** محاسبه‌شده

** (Infrastructure Monitoring) (ج مانیتورینگ زیرساخت) ****

Netdata** (با CPU، RAM، I/O منابع سرور **).

API endpoint، سلامت Docker سرویس مدل: وضعیت کانتینر ** -

دیتابیس ویژگی‌ها: ** اتصالات، حجم داده ** -

*** **4.11.6 (Model Versioning) نسخه‌بندی و مدیریت مدل

*** سیاست نسخه‌بندی

- مدیریت می‌شوند `MAJOR.MINOR.PATCH` نسخه‌ها** به صورت** -
- تغییرات اساسی در معماری مدل یا ویژگی‌ها: `MAJOR` -
- بهبود دقت یا افزودن ویژگی جدید: `MINOR` -
- رفع باگ یا بهینه‌سازی جزئی: `PATCH` -

*** ذخیره‌سازی

- | مؤلفه | محل ذخیره | ابزار |
- | ----- | ----- | ----- |
- | Git | (tagging با) GitHub | **کد آموزش** |
- | MLflow Tracking | MLflow در Metadata + مدل آموزش‌دیده** | مدل** |
- | DVC | DVC + (مثلاً MinIO) Object Storage | **دیتاست‌ها** |
- | Git | در مخزن کد YAML هایپرپارامترها** | فایل** |

*** نمونه ساختار فایل‌ها

...

models/

└─ v1.0.0/

| └─ model.pkl

4. (مجازی A/B testing) **مقایسه با مدل فعلی**.

5. Model Registry اگر بهتر بود: **ثبت در**.

6. staging تست استقرار **در محیط**.

7. انتخاب برای استقرار **توسط مدیر مدل**.

Airflow یا **GitHub Actions** خودکارسازی: **با**.

*** (Pre-deployment Testing) تست قبل از انتشار 4.11.8 ***

**** **مراحل تست**

| مرحله | هدف | روش |

| ----- | ----- | ----- |

| Python روی کد `pytest` | تست توابع آماده‌سازی داده و آموزش | **Unit Testing** |

| داده → مدل → end-to-end pipeline تست | **Integration Testing** |

| روی داده‌های ساختگی pipeline پیش‌بینی | اجرای

| (مجازی) ** | مقایسه مدل جدید با مدل فعلی روی داده‌های تاریخی | محاسبه A/B Testing ** |

| متریک‌ها برای هر دو

| (۱٪) تست مدل جدید روی بخش کوچکی از ترافیک واقعی | **Canary Deployment** |

| routing با load balancer |

| `k6` یا `locust` | سرویس مدل API تست بار بر روی | **Stress Testing** |

*** (Acceptance Criteria) معیارهای پذیرش ***

- کمتر ** از مدل فعلی داشته باشد (یا حداکثر ۵٪ بدتر) MAE ** مدل جدید باید -
- هیچ ویژگی حیاتی ** نباید حذف شده باشد ** -
- زمان استنتاج ** باید زیر ۵۰۰ میلی ثانیه باشد ** -
- تست‌های یکپارچگی ** همه باید عبور ** کنند -

*** 4.11.9 Rollback (Deployment & Rollback) استقرار و ***

*** استراتژی استقرار ***

Blue-Green Deployment: ** -

- v1.0.0) مدل فعلی: ** (Blue) محیط ** آبی -
- v1.1.0) مدل جدید: ** (Green) محیط ** سبز -
- پس از تست، ترافیک به تدریج به سبز منتقل می‌شود -

Rollback روش ***

۱. شناسایی مشکل: ** از طریق هشدارها یا گزارش کاربران **
۲. مدیر مدل تأیید می‌کند Rollback: ** تصمیم **
۳. load balancer تغییر مسیر ترافیک: ** بازگشت به محیط آبی (مدل قبلی) از طریق **
۴. بررسی علت: ** تحلیل لاگ‌ها و داده‌ها برای شناسایی ریشه مشکل **

کمتر از ۱۵ دقیقه از شناسایی تا بازیابی **Rollback: زمان هدف**

4.11.10 نقش‌ها و مسئولیت‌ها (Roles & Responsibilities)*

نقش	مسئولیت‌ها	فرد/تیم
تأیید انتشار، نظارت بر عملکرد، تصمیم‌گیری	مدیر مدل (Model Owner)	**
مدیر فنی	Rollback	
آموزش مدل، نگهداری زیرساخت pipeline توسعه	ML (ML Engineer)	** مهندس
تیم داده	MLOps	
توسعه API یکپارچه‌سازی مدل با	مهندس نرم‌افزار (Software Engineer)	**
سرویس‌های همراه	تیم توسعه	
گزارش‌دهی عملکرد، پیشنهاد بازآموزی	Drift تحلیل	** (Data Analyst) تحلیلگر داده
تیم داده		
مانیتورینگ زیرساخت، هشدارها، پشتیبان‌گیری	تیم فنی	** (Ops) پشتیبان عملیاتی

4.11.11 ابزارهای پیشنهادی (Tool Stack)*

دسته | ابزار | توضیح

| ----- | ----- | ----- |
| کد و کنترل نسخه | Git, GitHub | **مدیریت کد و نسخه** |
| ردیابی آزمایش‌ها، نسخه‌بندی مدل و داده | MLflow, DVC | **مدیریت مدل** |
| خودکارسازی آموزش و | GitHub Actions, Apache Airflow | **پایپ‌لاین** |
| استقرار |
| نمایش متریک‌ها و | Grafana, Prometheus, Custom Scripts | **مانیتورینگ** |
| هشدار |
| و کانتینری‌سازی API | FastAPI, Docker, Nginx | **سرویس‌دهی مدل** |
| Drift تست کد، بار و | pytest, locust, evidently | **تست** |

****4.11.12 برنامه اجرایی (Implementation Roadmap)*****

**** (ماه ۳) MLOps فاز ۱: پایه*****

MLflow با Model Registry راه‌اندازی -

.(روزانه MAE) ایجاد مانیتورینگ ساده -

آموزش دستی pipeline راه‌اندازی -

**** فاز ۲: خودکارسازی (ماه ۶)*****

GitHub Actions خودکارسازی بازآموزی با -

Drift راه‌اندازی مانیتورینگ -

- Canary Deployment اجرای -

فاز ۳: تکمیل (۱۲ ماه) ###

- Blue-Green Deployment استقرار کامل -

- Grafana با Real-time مانیتورینگ -

- Feature Store یکپارچه‌سازی -

MLOps (MLOps Costs) هزینه‌های 4.11.13 ###

| مؤلفه | هزینه ماهانه تخمینی | توضیح |

| ----- | ----- | ----- |

| دلار | برای مانیتورینگ پیشرفته ۱۰ | **Grafana Cloud) ابزارهای مانیتورینگ** |

| Object Storage | ذخیره‌سازی مدل‌ها و داده‌ها** | ۵ دلار** |

| دلار | ۲ سرور مجازی ۲۰ | **Training/Staging) سرورهای اضافی** |

| دلار | حقوق و دستمزد ۵۰۰ | **نیمه‌وقت ML مهندس) نیروی انسانی** |

| مجموع** | ۵۳۵ دلار** | برای فاز اول** |

.توجه:** هزینه‌ها با مقیاس سیستم افزایش می‌یابد**

****نتیجه‌گیری 4.11.۱۴** ###**

به گونه‌ای طراحی شده است که **مدل‌ها را از `nobatnou.ir` و مانیتورینگ MLOps سیستم بمب ساعتی به دارایی پایدار ** تبدیل کند. با ترکیب **خودکارسازی، مانیتورینگ چندلایه و فرآیندهای کنترل‌شده انتشار **، اطمینان حاصل می‌شود که مدل‌ها نه تنها دقیق، بلکه **قابل اعتماد، قابل نگهداری و rule-واکنش‌گرا** به تغییرات هستند. این سند به عنوان نقشه راه تیم فنی برای حرکت از مدل ساده صنعتی عمل می‌کند ML به سیستم based