

ML Feature DSL Proposal

Overview

The ML Feature DSL is a domain-specific language that enables clear, reusable, and consistent definition of features used in machine learning models. It abstracts away complex SQL or Python transformation code, allowing teams to define feature logic once and generate compatible code across training, serving, monitoring, and reporting environments.

Goals

- Define machine learning features in a human-readable, declarative format
- Generate code for multiple targets (Python, SQL, API schemas, monitoring)
- Eliminate code duplication and inconsistency across teams
- Make feature logic auditable and shareable

DSL Syntax (Input)

Example:

```
feature: is_active = days_since_last_login < 30
```

```
feature: low_spender = total_spent < 100
```

```
feature: no_recent_support = support_tickets_last_90d == 0
```

Each feature: line declares a new machine learning feature. The syntax supports simple expressions with logical, comparison, and function operations.

Outputs Generated Automatically

1. Python (for training in Pandas)

```
df["is_active"] = df["days_since_last_login"] < 30  
df["low_spender"] = df["total_spent"] < 100  
df["no_recent_support"] = df["support_tickets_last_90d"] == 0
```

2. SQL (for data extraction or dashboards)

```
CASE WHEN days_since_last_login < 30 THEN TRUE ELSE FALSE END AS is_active,  
total_spent < 100 AS low_spender,  
support_tickets_last_90d = 0 AS no_recent_support
```

3. Monitoring Config (e.g. YAML)

```
features:  
  - name: is_active  
    drift_check: ks_test  
    threshold: 0.05  
  - name: low_spender  
    drift_check: js_divergence  
    threshold: 0.1
```

4. API Schema (JSON Input for Real-Time Prediction)

```
{  
  "input": {  
    "days_since_last_login": "number",  
    "total_spent": "number",  
    "support_tickets_last_90d": "number"  
  },  
  "features": ["is_active", "low_spender", "no_recent_support"]  
}
```

Supported DSL Input Types

Type	Example Syntax	Description
Comparison	<code>x > 5</code> , <code>y == 0</code>	Boolean conditions
Logical	<code>a AND b</code> , <code>x OR y</code>	Combine multiple conditions
Function	<code>hour(time) > 20</code> , <code>bucket(age, [18,30,50])</code>	Feature transformation functions
Arithmetic	<code>total / count > 5</code>	Derived numeric features

These inputs are parsed and transformed into executable code depending on the target environment.

Benefits

- **Single Source of Truth:** One DSL file generates all environments' code.
- **Team Collaboration:** Data scientists, engineers, and analysts share the same logic.
- **Auditability:** Easy to track, version, and explain.
- **Consistency:** Same features in training, production, and analytics.
- **Productivity:** Build new models faster with reusable feature libraries.

Use Cases

Churn Prediction

feature: `is_inactive = days_since_last_login > 45`

feature: `low_engagement = weekly_sessions < 2`

Fraud Detection

feature: `large_txn = transaction_amount > 1000`

feature: `night_txn = hour(transaction_time) > 22`

Recommendation Engine

feature: viewed_recently = days_since_last_view < 7

feature: repeat_buyer = total_orders > 2

Conclusion

This proposal introduces an ML Feature DSL that enables clean, reusable, and consistent feature engineering at scale. The simplicity of the syntax hides its powerful ability to unify ML workflows across environments. Adoption of this DSL can significantly reduce duplication, bugs, and inconsistencies in real-world ML projects.

DSL Example – Customer Lifetime Value Prediction

```
1  # User Activity
2  feature: is_recent_login = days_since_last_login < 14
3  feature: active_last_month = monthly_sessions_last_30d > 5
4  feature: multi_device_user = device_count > 1
5
6  # Purchase Behavior
7  feature: avg_order_value = total_spent / total_orders
8  feature: high_spender = total_spent > 1000
9  feature: frequent_buyer = total_orders > 10
10 feature: last_purchase_recent = days_since_last_purchase < 30
11
12 # Subscription & Retention
13 feature: subscription_active = subscription_status == "active"
14 feature: churn_risk = days_since_subscription_renewal > 45
15
16 # Support Interactions
17 feature: frequent_complainer = support_tickets_last_90d > 3
18 feature: unresolved_issues = open_support_tickets > 0
19
20 # Engagement
21 feature: viewed_products_last_week = product_views_last_7d > 10
22 feature: marketing_engaged = marketing_email_clicks_last_30d > 3
23
```

Python Output (Pandas)

```
1 df["is_recent_login"] = df["days_since_last_login"] < 14
2 df["active_last_month"] = df["monthly_sessions_last_30d"] > 5
3 df["multi_device_user"] = df["device_count"] > 1
4
5 df["avg_order_value"] = df["total_spent"] / df["total_orders"]
6 df["high_spender"] = df["total_spent"] > 1000
7 df["frequent_buyer"] = df["total_orders"] > 10
8 df["last_purchase_recent"] = df["days_since_last_purchase"] < 30
9
10 df["subscription_active"] = df["subscription_status"] == "active"
11 df["churn_risk"] = df["days_since_subscription_renewal"] > 45
12
13 df["frequent_complainer"] = df["support_tickets_last_90d"] > 3
14 df["unresolved_issues"] = df["open_support_tickets"] > 0
15
16 df["viewed_products_last_week"] = df["product_views_last_7d"] > 10
17 df["marketing_engaged"] = df["marketing_email_clicks_last_30d"] > 3
18
```

SQL Output

```
1  ✓ days_since_last_login < 14 AS is_recent_login,
2    monthly_sessions_last_30d > 5 AS active_last_month,
3    device_count > 1 AS multi_device_user,
4
5    total_spent / NULLIF(total_orders, 0) AS avg_order_value,
6    total_spent > 1000 AS high_spender,
7    total_orders > 10 AS frequent_buyer,
8    days_since_last_purchase < 30 AS last_purchase_recent,
9
10   subscription_status = 'active' AS subscription_active,
11   days_since_subscription_renewal > 45 AS churn_risk,
12
13   support_tickets_last_90d > 3 AS frequent_complainer,
14   open_support_tickets > 0 AS unresolved_issues,
15
16   product_views_last_7d > 10 AS viewed_products_last_week,
17   marketing_email_clicks_last_30d > 3 AS marketing_engaged
18
```

Monitoring Config (YAML)

```
1  ∨ features:
2  ∨   - name: is_recent_login
3      drift_check: ks_test
4      threshold: 0.05
5   - name: avg_order_value
6      drift_check: wasserstein_distance
7      threshold: 10
8   - name: churn_risk
9      drift_check: ks_test
10     threshold: 0.05
11   - name: marketing_engaged
12     drift_check: js_divergence
13     threshold: 0.1
14
```


API Input Schema (JSON)

```
1  {
2    "input": {
3      "days_since_last_login": "number",
4      "monthly_sessions_last_30d": "number",
5      "device_count": "number",
6      "total_spent": "number",
7      "total_orders": "number",
8      "days_since_last_purchase": "number",
9      "subscription_status": "string",
10     "days_since_subscription_renewal": "number",
11     "support_tickets_last_90d": "number",
12     "open_support_tickets": "number",
13     "product_views_last_7d": "number",
14     "marketing_email_clicks_last_30d": "number"
15   },
16   "features": [
17     "is_recent_login",
18     "active_last_month",
19     "multi_device_user",
20     "avg_order_value",
21     "high_spender",
22     "frequent_buyer",
23     "last_purchase_recent",
24     "subscription_active",
25     "churn_risk",
26     "frequent_complainer",
27     "unresolved_issues",
28     "viewed_products_last_week",
29     "marketing_engaged"
30   ]
31 }
32
```