

# MEDALLION ARCHITECTURE

# DIABETES RISK ANALYSIS

Built a production-level data pipeline on Databricks using PySpark

# MEDALLION ARCHITECTURE

- Read and save raw data
- Change column names,
  - Handle missing values,
  - Handle outliers,
- Calculate BMI and add its column,
- Save the data as a table
- Develop custom features in the Gold Layer to capture complex health-risk patterns that raw data alone cannot represent



**BRONZE**

**SILVER**

**GOLD**

## GOLD LAYER

# FEATURE ENGINEERING HIGHLIGHTS

### BMI\_GROUP

Grouped BMI into 'Obese', 'Overweight', and 'Normal'. This transforms non-linear physiological risks into interpretable segments for the model and business reporting

### IS\_OVERWORKED

Flagged users working \$>40\$ hours/week. This acts as a proxy for stress-related metabolic risks and lifestyle imbalance.

### WORK\_METABOLIC\_BENEFIT

Quantified the "sweet spot" of physical labor. I assigned positive weights (\$+1.5\$) to moderate intensity but applied a penalty (\$-0.5\$) to high-intensity work to account for potential stress-related health trade-offs.

### LEISURE\_VITALITY\_SCORE

Weighted high-intensity leisure activities (\$2.0\$) more heavily than moderate ones (\$1.0\$). This emphasizes the quality of physical recovery outside of working hours as a key preventative factor.

# AD-HOC (EDA)

While high-intensity at work often failed to show benefits due to occupational stress, high-intensity in leisure leads to pure metabolic gain. It effectively lowers blood sugar because it is voluntary and enjoyable.

This protective effect exists independent of BMI. This means that even without significant weight loss, the physiological act of engaging in vigorous fun directly improves your insulin response.

If you have limited time, "vigorous leisure" (High Intensity) offers the highest "Return on Investment" for diabetes prevention. One day of intense play is more valuable than one day of moderate effort.

1

2

3

**"High Intensity is Joy, Not Stress"**

**The "Active Fun" Shield**

**Efficiency Comparison (ROI)**