

**هدف:** پیش‌بینی کیفیت محیط داخلی

**دیتاست:** FullDataset\_CombinedSimple.csv

ما مدل شخصی برای هر فرد نمی‌سازیم، بلکه یک مدل عمومی داریم که روی توالی‌های زمانی همه افراد آموزش می‌بیند.

**نوع مسئله:**

ابتدا regression، تولید خروجی اعشاری و محاسبه  $R^2$

سپس تبدیل خروجی به عدد صحیح بین ۱ تا ۵ با رند کردن یا clip به بازه [1,5] و محاسبه F1، Accuracy

### Preprocessing & Feature Engineering

ستون هایی که برای تحلیل استفاده می‌شوند وارد مدل نمی‌شوند:

برچسب زمانی برای مرتب سازی و محاسبه ساعت روز (Hour)، خستگی زمانی (Time Fatigue):

TimeVote,

Time\_Troom, Time\_RH, Time\_CO2, Time\_EA, Time\_Ttrend, Time\_Sound, Time\_VOC, Time\_Lighting,

برای ساخت فیلد location:

LocationFront, LocationBack, LocationLeft, LocationRight, LocationMiddle,

ID, Student, Gender, Age, CaseStudy

هم مقدار فعلی هم مقادیر قبلی به عنوان lag: (مقدار پویا سمت در طول زمان دیتاست تغییر می‌کند)

Moment, CLO\_Simple, Activity, Troom, RH, CO2, Sound, Lighting, Ttrend, Trm, WindowsClosed,

فقط به عنوان lag

Thermal\_Warm, Thermal\_Cold, Thermal\_Draft, Thermal\_Other,  
IAQ\_Stuffy, IAQ\_Odour, IAQ\_Dry, IAQ\_Humid, IAQ\_Other,  
AC\_Classroom, AC\_Hallway, AC\_Device, AC\_Outside, AC\_Other, AC\_Ventilation,  
Vis\_Blackboard, Vis\_Underexposure, Vis\_Overexposure, Vis\_Glare, Vis\_Other,

ThermalSatisfaction, IAQSatisfaction, AcousticSatisfaction, VisualSatisfaction, IEQSatisfaction

ستون هایی که باید پیش‌بینی شوند: دارای مقدار ۱ برای کمترین تا ۵ برای پیشترین رضایت

ستون هایی که کنار گذاشته می‌شوند:

دارای correlation زیاد:

EA, CLO\_Detailed,

اطلاعات تکراری از Student و CaseStudy می‌دهد:

Subgroup,

حجم کم داده‌ها:

VOC, GroupObs, AF

WindowMajority, ActivityMajority,

### مرحله ۱: پاکسازی و اصلاح مقادیر (Data Cleaning)

**اصلاح Age \***: اگر Age خالی است، بر اساس میانگین سن در آن CaseStudy خاص (مثلاً PrimarySchool) پر شود.

(در دیتاست برای تمام سطرها با Age بدون مقدار، اگر CaseStudy PrimarySchool باشد student باشد TRUE باشد) برای student باشد

### مدیریت مقادیر مفقود (Missing Values):

برای ستون‌های عددی: به جای حذف سطر، خانه‌های خالی با مقدار ۰- پر شوند (مگر اینکه مقدار پیش فرض دیگری پیشنهاد شود) تا مدل بفهمد این داده در آن لحظه ثبت نشده است.

برای ستون‌های متغیر: با کلمه Unknown جایگزین شوند.

### عددی کردی ستون‌های متغیر:

CLO\_Simple

0.5:Lightly

0.7:Neutral

1:Warmly

برای سطرهای خالی: عدد میانگین لیاس بقیه افراد در همان روز

.Moment

مقدار ۱ برای Start: (شروع کلاس؛ زمانی که فرد تازه وارد شده و هنوز با محیط تطبیق پیدا نکرده است).

مقدار ۲ برای End: (پایان کلاس؛ اوج خستگی و پیشترین زمان صرف شده در محیط).

مقدار ۱.۵ برای Blank: (چون نمی‌دانیم دقیقاً کی بوده، آن را حد وسط قرار می‌دهیم تا مدل فرض کند جایی بین شروع و پایان بوده است).

### مرحله ۲: مهندسی ویژگی‌ها (Feature Engineering):

ایجاد مختصات مکانی (Location): یک مختصات دکارتی (Y, Z) که هر مقدار آن نماینده یک جهت است، با اولویت‌بندی جهات اصلی و استفاده از Middle به عنوان جایگزین.

x=0

y=0

if front=true then y=1

if back=true then y=3

if Left=true then x=1

if Right=true then x=3

if x=0 then if middle=true then x=2

if y=0 then if middle=true then y=2

محاسبه امتیاز حساسیت (User Sensitivity Score): جمع جبری تیک‌های زده در ستون‌های Thermal\_Cold (مثلاً Thermal\_Specific\_Issues) برای هر سطر (نشان‌دهنده سخت‌گیری فرد).

نرخ تولید گرمای بدن (Metabolic Rate): تبدیل متون Activity به اعداد معادل بر اساس استانداردهای آسایش حرارتی (ASHRAE)

مقادیر پیشنهادی برای فعالیت (Activity) به Met:

مقادیر ۱۰۰ برای Test: (کمترین فعالیت فیزیکی، تمرکز ذهنی بالا در حالت نشسته).

مقدار ۱.۱ برای Passive course: (نشستن معمولی و گوش دادن در کلاس).

مقدار ۱.۲ برای Other: (مقدار متوسط فعالیت‌های پیش‌بینی نشده).

مقدار ۱.۳ برای Not applicable: (خنثی‌سازی داده‌های نامشخص).

مقدار ۱.۴ برای Blank (خالی): (جایگزینی برای جلوگیری از خطای محسوساتی).

مقدار ۱.۵ برای Presentation: (ایستان و صحبت کردن که ارزی بیشتری می‌برد).

مقدار ۱.۶ برای Active course: (بیشترین فعالیت؛ شامل کارهای گروهی و جابه‌جانی).

استخراج ویژگی‌های زمانی:

استخراج ساعت از TimeVote: (فقط مقدار فعلی به مدل ارسال می‌شود)

تفاوت زمان فعلی با اولین حضور فرد در آن روز: (فقط مقدار فعلی به مدل ارسال می‌شود)

تفاوت زمان فعلی با سطر قبلی (بدون در نظر گرفتن ID): (فقط مقدار فعلی به مدل ارسال می‌شود)

یافتن تعداد lag به عنوان hyperparameter مدل

یافتن تعداد lag مناسب:

حلقه به ازای هر از لیست ۳, ۵, ۱۰, ۱۵, ۲۰

مرحله ۳ مرتب‌سازی و ساخت Lag‌ها (Vectorization)

مرحله ۴ بازگرداندن نظم و تقسیم‌بندی

اجرای پایپلاین (Train + TimeSeriesSplit + GridSearch)

Test

ارزیابی

{

شروع حلقه

مرحله ۳: مرتب‌سازی و ساخت Lag‌ها (Vectorization)

مرتب‌سازی اول: کل دیتا است به صورت صعودی بر اساس ID و سپس TimeVote مرتب شود. (این کار باعث می‌شود سوابق هر فرد پشت سر هم قرار بگیرد).

ساخت ستون‌های Lag (از ۱ تا n):

برای تمام سوابق (سنتسورها، مکان، فعالیت، پوشش، زمان، رضایت‌ها و دلایل نارضایتی)، مقدار لحظه فعلی (t) حذف شود و فقط Lag‌ها باقی بمانند تا نشت داده (Data Leakage) رخ ندهد.

مرحله ۴: بازگرداندن نظم و تقسیم‌بندی

مرتب‌سازی دوم: حالا که تمام سوابق در هر سطح ترتیب شده، دیتا است را دوباره فقط بر اساس TimeVote مرتب کنید. (این کار برای حفظ واقع‌گرایی مدل در زمان آموزش حیاتی است).

تقسیم‌بندی (Time-Series Split) ۲۰-۸۰

۸۰ درصد داده‌های قدیمی‌تر (بر اساس زمان) برای آموزش (Train).

۲۰ درصد داده‌های جدیدتر برای تست (Test).

نکته: دقت کنید که چون داده‌ها را بر اساس زمان تقسیم می‌کنید، مدل باید بتواند آینده را بر اساس گذشته پیش‌بینی کند.

یافتن بقیه hyperparameter TimeSeriesSplit, Grid Search مدل

تعداد درخت: بهینه‌سازی: بهترین راه برای کردن عدد دقیق، استفاده از روش TimeSeriesSplit و Grid Search است

عمق درخت: بهینه‌سازی: بهترین راه برای پیدا کردن عدد دقیق، استفاده از روش TimeSeriesSplit و Grid Search است

پیدا کردن بهترین عمق (Max Depth) و تعداد درخت با استفاده از TimeSeriesSplit برای حفظ نظم زمان.

split بر حسب زمان و غیرتصادفی،

گذشته → train

(test) → آینده

بهترین تعداد درخت: از ۱۰۰ درخت شروع کرده تعداد را زیاد می‌کند تا زمانی که تأثیری بر دقت نداشته باشد

Mدل با الگوریتم Train Random Forest و در مرحله بعد XGBoost سپس neural network سپس SVM

هدف: آموزش ۵ مدل مجزا برای ۵ هدف (Target)

Train و Grid Search یک فرآیند ترکیبی هستند

در نهایت شما فقط یک «شیء مدل» دارید که هم هوشمندترین عمق را انتخاب کرده و هم آموزش دیده است

تست

ارزیابی مدل

الف) میانگین مجدد خطای Mean Squared Error - MSE: خطای جرمهدار (برای کاهش خطاهای بزرگ).

این معیار شناس می‌دهد که پیش‌بینی‌های مدل چقدر از واقعیت فاصله دارند. هرچه این عدد به صفر نزدیک‌تر باشد، مدل بهتر است.

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

ب) دقت (Accuracy): درصد کلی تشخیص‌های درست (دسته‌بندی)

اگر بخواهیم مدل در چند درصد موارد دقیقاً همان عدد داشت آموز را درست حدس زده است:

$$Accuracy = \frac{\text{Correct Predictions}}{\text{Total Predictions}} \times 100$$

ج)  $score R^2$  (ضریب تعیین): درصد موفقیت مدل در توضیح رفتار محیط.  
این عدد بین ۰ تا ۱ است. اگر ۱ باشد یعنی مدل شما ۱۰۰٪ رفتار دانشآموزان را فهمیده است

$$R^2 = 1 - \frac{\sum(y_i - \hat{y}_i)^2}{\sum(y_i - \bar{y})^2}$$

د) MAE (Mean Absolute Error): میانگین خطای واقعی (به واحد نمره).  
برخلاف MSE که خطا را به توان ۲ می‌رساند، MAE دقیقاً به شما می‌گوید مدل به طور متوسط چند نمره اشتیاه می‌کند.

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

ه) F1-Score (برای ماتریس درهم‌ریختگی): تعادل دقت در کلاس‌های مختلف.  
اگر بخواهید بگویید مدل چقدر در تشخیص افراد «ناراضی» موفق بوده، این معیار تعادل بین «دقت» و «صحت» را نشان می‌دهد.

$$F1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

#### Visualization

۱. نمودار پیش‌بینی در برابر واقعیت (Actual vs Predicted Plot)

۲. نمودار اهمیت ویژگی‌ها (Feature Importance Plot)

۳. نمودار باقیمانده‌ها (Residuals Plot)

۴. نمودار سری زمانی پیش‌بینی (Time-Series Forecast Plot)

۵. ماتریس درهم‌ریختگی (Confusion Matrix) - اگر داده‌ها را طبقه‌بندی کنیم

پایان حلقه