# Integrating AI to assess Students Well-being

**新 民 中 學**
**XINMIN SECONDARY SCHOOL**
壁韌不拔, 殷勤為众. Steadfast in Spirit, Diligent in Service

**Nguyen Le Hoang, Dinh Hoang Nam, Tan Yong Xiang**

## Problems

- Hard to recognise stress for the students
- Hard to evaluate students' stress level
- Many signs of stress are quantitative, it is harder to recognise the severity of stress
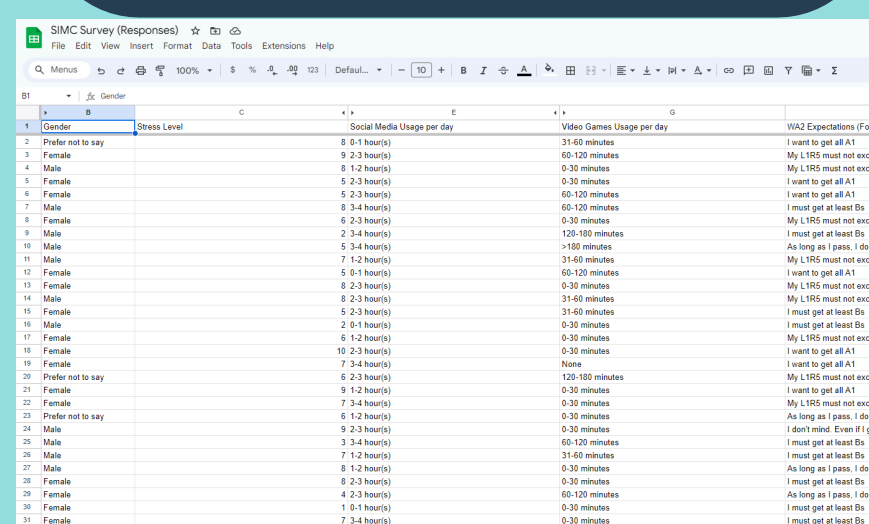
## How Singapore students responded

**On test anxiety**
"I feel very anxious even if I am well prepared for a test."
**76%**
agreed or strongly agreed with the statement compared with the OECD average of **55%**

**On achievement motivation**
"I want to be one of the best students in my class."
**82%**
said that they did, compared with the OECD average of **60%**

**On being bullied**
During the past 12 months, how often have you had the following experiences in school?
**18.3%**
said they were made fun of at least a few times a month, compared with the OECD average of **10.9%**

"Other students left me out of things on purpose."
**11.9%**
said they were left out on purpose, compared with the OECD average of **7.2%**

Source: PROGRAMME FOR INTERNATIONAL STUDENT ASSESSMENT STUDENTS' WELL-BEING STUDY 2015

## Methodology

## Solutions

- Our aim: Using AI to identify student's well-being throughout their academic year based on their daily schedule
- Prompt: Build an AI model that can predict the student's stress level based on the 4 features we gathered: Sleeping time, Gaming time, Outdoor Activities time and Social Media Usage

## Overview:

- Collect data from students
- Pre-processing data
- Split data
- Create a model to predict the stress level based on given features using Random Forest Classifier with GridSearchCV
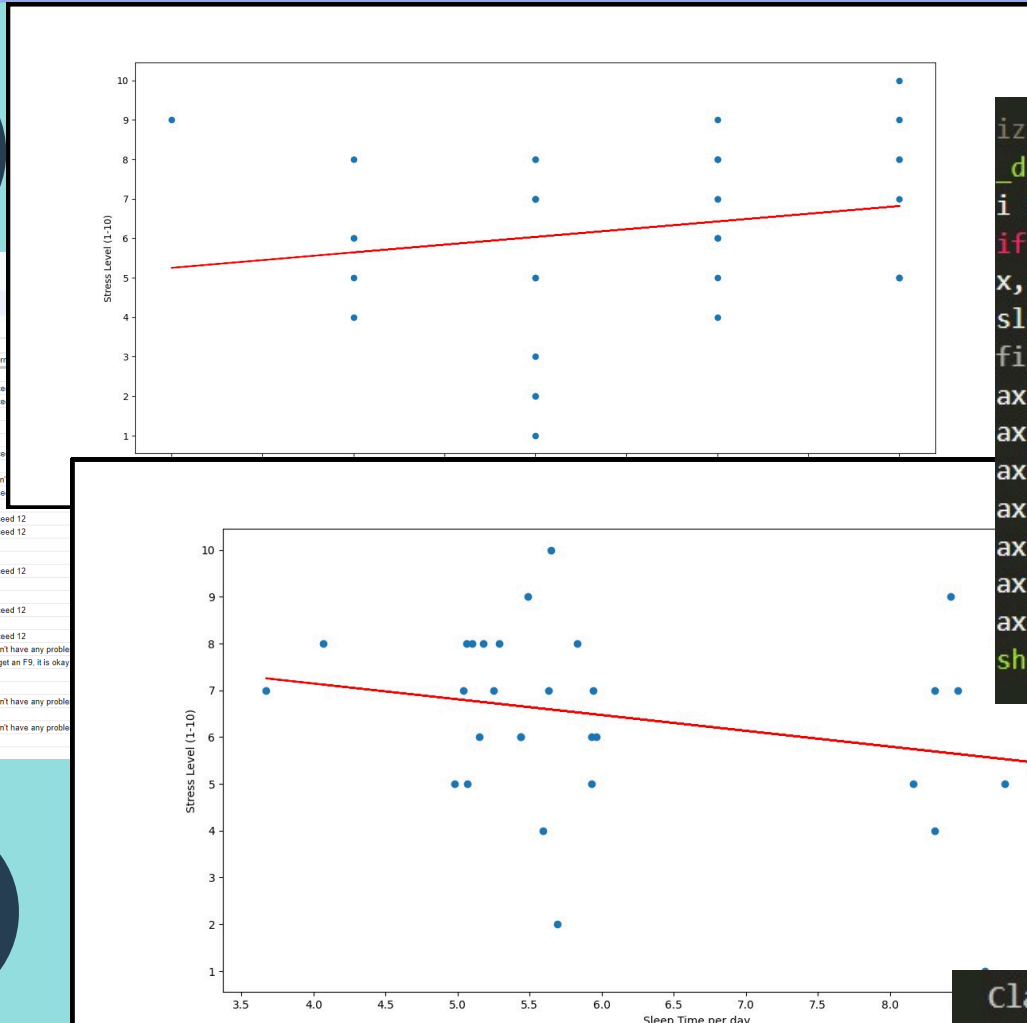- Train the model

### Data Collection

### Model Training

### Preprocessing

```
# Preprocessing
encoder = OneHotEncoder(sparse_output=False)
one_hot_encoded = encoder.fit_transform(data[['Gender']])
one_hot_df = pd.DataFrame(one_hot_encoded, columns=encoder.get_feature_names_out(['Gender']))
data = pd.concat([data, one_hot_df], axis=1)
data = data.drop(['Gender'], axis=1)

columns_to_scale = ['Social Media Usage per day', 'Video Games Usage per day', 'WA2 Expectations', 'O
data[columns_to_scale] = StandardScaler().fit_transform(data[columns_to_scale])

pca_columns = ['Social Media Usage per day', 'Video Games Usage per day']
pca = PCA(n_components=1).fit_transform(data[pca_columns])
pca_df = pd.DataFrame(pca, columns=['Screen Time per day'])
data = pd.concat([data['Stress Level'], pca_df, data.loc[:, 'WA2 Expectations':]], axis=1)
print(data)
```

### Split data

```
# Input data
data = pd.read_csv('d:/Office files/SIMC/src/SIMC Survey Responses.csv')
data = data.drop(columns=['WA2 Expectations (Form)'])
data = data.fillna(0)

ref = {"Social Media Usage per day":{"0-1 hour(s)":(0,1), "1-2 hour(s)":(1,2), "2-3 hour(s)":(2,3), "3-4 hour(s)":(3,4)},
       "Video Games Usage per day":{"0-30 minutes":(0,0.5), "31-60 minutes":(0.5,1), "60-120 minutes":(1,2),"120-180 minutes"
       "Outdoor Time per day":{"<1 hour":(0,00,1), "1-2 hour(s)":(1,2), "2-4 hour(s)":(2,4)},
       "Sleep Time per day":{"<5 hours":(3,5),"5-6 hours":(5,6),"8-10 hours":(8,10)}}

for label, content in data.items():
    if label in ref:
        for i in range(len(content)):
            if content[i] in ref[label]:
                x, y = ref[label][content[i]]
                content[i] = round(uniform(x,y), 2)
print(data)
```

```
# GridSearchCV
def tuning():
    param_grid = {'n_estimators':[1, 1000],
                  'max_features':['sqrt','log2',None],
                  'criterion':['gini','entropy','log_loss']}
    model = RandomForestClassifier()
    CV_model = GridSearchCV(estimator=model, param_grid=param_grid, cv=5, scoring='accuracy', n_jobs=-1)
    CV_model.fit(X_train, y_train)
    print(CV_model.best_params_)

# Training and testing
model = RandomForestClassifier(criterion='gini', max_features='sqrt', n_estimators=1000, n_jobs=-1)
model.fit(X_train, y_train)
predicted = model.predict(X_test)

# Evaluation Metrics
print(f'Classification Report:\n{classification_report(y_test, predicted)}')
print(f'Confusion Matrix:\n{confusion_matrix(y_test, predicted)}')
```

```
Classification Report:
              precision    recall  f1-score   support

           1       1.00      1.00      1.00         2
           2       1.00      1.00      1.00         2
           3       1.00      1.00      1.00         1
           4       1.00      1.00      1.00         2
           5       1.00      1.00      1.00        11
           6       1.00      1.00      1.00         6
           7       1.00      1.00      1.00         7
           8       1.00      1.00      1.00         4
           9       1.00      1.00      1.00         4
          10       1.00      1.00      1.00         1

    accuracy                           1.00        40
   macro avg       1.00      1.00      1.00        40
weighted avg       1.00      1.00      1.00        40

Confusion Matrix:
[[ 2  0  0  0  0  0  0  0  0  0]
 [ 0  2  0  0  0  0  0  0  0  0]
 [ 0  0  1  0  0  0  0  0  0  0]
 [ 0  0  0  2  0  0  0  0  0  0]
 [ 0  0  0  0 11  0  0  0  0  0]
 [ 0  0  0  0  0  6  0  0  0  0]
 [ 0  0  0  0  0  0  7  0  0  0]
 [ 0  0  0  0  0  0  0  4  0  0]
 [ 0  0  0  0  0  0  0  0  4  0]
 [ 0  0  0  0  0  0  0  0  0  1]]
```

## Reflection

- The model was trained pretty well according to expectation, when it is able to predict the stress level based on given factor
- What we can improve is the breadth of the data set, that is we should collect more responses so that resembling artificial sample is not nescessary.
- What we learn from the project is to create a model that can perform our prompt – and we were able to somehow train it!

## References

- SIMC Workshop Resources
- Scikit-learn Resources
- The Strait Times, "Study says Singapore students suffer from high levels of anxiety"