



ROYAL UNIVERSITY OF
PHNOM PENH

DIABETES PREDICTION USING MACHINE LEARNING

LECTURE: KHIM CHAMREOUN
GROUP TWO

Phork Norak, Phon Lihour, Vay Mithona, Sim Liheng

December 16 2025

Bussiness Understanding

The Challenge: A Chronic Metabolic Disorder

Diabetes mellitus is a chronic metabolic disorder characterized by the body's impaired ability to utilize blood sugar (glucose) effectively. It poses a significant challenge to global public health, necessitating better tools for early diagnosis and management.



Type 1 Diabetes

An autoimmune response, often manifesting in childhood, leading to an insulin deficiency. Caused by the destruction of insulin-producing beta cells.



Type 2 Diabetes

The most prevalent type, arising from insufficient insulin secretion or insulin resistance. Associated with family history, obesity, and inactivity.



Gestational Diabetes (GDM)

A temporary type developing during pregnancy due to hormonal changes. While it resolves post-childbirth, it increases the mother's future risk of Type 2 diabetes.



Predictive Modeling for Early Diagnosis

Many machine learning models have been built to assist doctors in the diagnosis of diabetes, leveraging patient data to identify patterns and predict outcomes. These models serve as powerful decision-support tools.



The Precedent

The PIMA Indian Diabetes Dataset

A foundational dataset that has been instrumental in developing and benchmarking diabetes prediction models. Your work will build upon this legacy.



The Advancement

Our Mission

We will use a more recent, larger dataset to build a new predictive model and deploy it as a functional web application.

Data Understanding

The Taipei Medical Center Dataset

Based on the 2023 study by Chou et al., J. Pers. Med., your work utilizes a contemporary and relevant dataset.

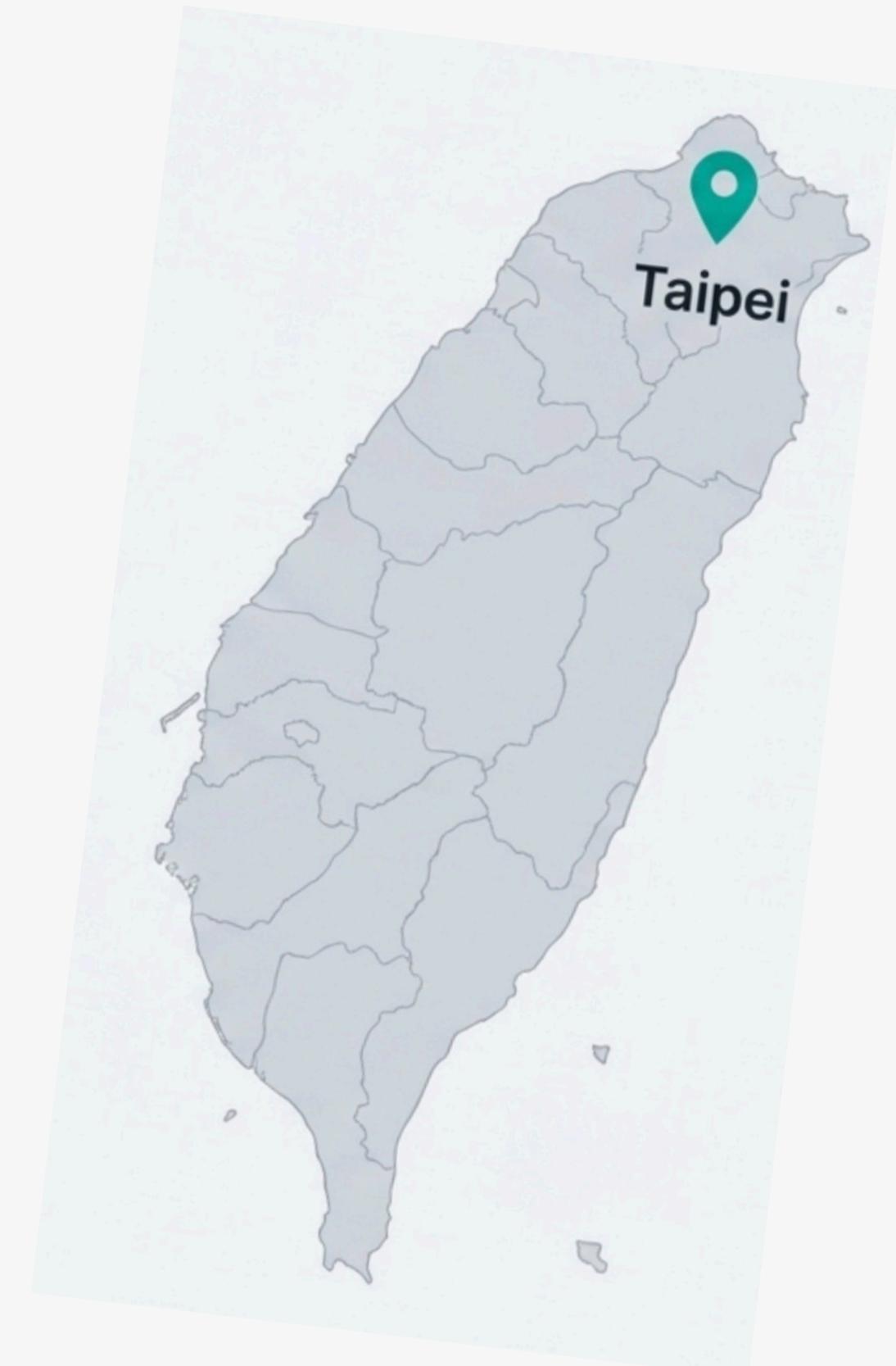
Name: 'TAIPEI diabetes. csv'

Subjects: 15,000 women

Age Range: 20 to 80 years

Source: Patients at Taipei Municipal Medical

Center Collection Period: 2018-2020 and
2021-2022



Understanding the Features: A Data Dictionary

Input Variables

**Pregnancies:**

Number of times pregnant

**PlasmaGlucose:**

Plasma glucose concentration (2 hrs post oral test)

**DiastolicBloodPressure:**

Diastolic blood pressure (mm Hg)

**TricepsThickness:**

Triceps skin fold thickness (mm)

**SerumInsulin:**

2-Hour serum insulin (mu U/ml)

**BMI:**

Body mass index (weight in kg / (height in m)²)

**DiabetesPedigree:**

A function scoring diabetes probability based on family history

**Age:**

Age in years

**Target Variable:**

Diabetic: 1 = Diabetes Diagnosed, 0 = No Diabetes

Result: Exploratory data analysis

```

df.duplicated().sum()

[4]
...
0

print("\nShowing the number of rows and columns of the dataset:")
df.shape

[5]
...
Showing the number of rows and columns of the dataset:

[6]
...
(15000, 10)

df.info()

[3]
...
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 15000 entries, 0 to 14999
Data columns (total 10 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   PatientID   15000 non-null   int64  
 1   Pregnancies  15000 non-null   int64  
 2   PlasmaGlucose 15000 non-null   int64  
 3   DiastolicBloodPressure 15000 non-null   int64  
 4   TricepsThickness 15000 non-null   int64  
 5   SerumInsulin  15000 non-null   int64  
 6   BMI          15000 non-null   float64 
 7   DiabetesPedigree 15000 non-null   float64 
 8   Age          15000 non-null   int64  
 9   Diabetic     15000 non-null   int64  
dtypes: float64(2), int64(8)
memory usage: 1.1 MB

print("\nShowing first 5 rows of the dataset:")
df.head()

[6]
...
Showing first 5 rows of the dataset:

[7]
...
PatientID Pregnancies PlasmaGlucose DiastolicBloodPressure TricepsThickness SerumInsulin    BMI  DiabetesPedigree Age  Diabetic
0   1354778        0       171             80            34         23  43.509726  1.213191  21    0
1   1147438        8       92              93            47         36  21.240576  0.158365  23    0
2   1640031        7       115             47            52         35  41.511523  0.079019  23    0
3   1883350        9       103             78            25         304 29.582192  1.282870  43    1
4   1424119        1       85              59            27         35  42.604536  0.549542  22    0

```

