

Chapter 2: TỔNG QUAN TIỀN XỬ LÝ DỮ LIỆU



Ex1: Điểm thi THPT Quốc Gia 2016

- Cho tập tin Diemthi_thpt_quocgia_2016.xlsx chứa bộ dữ liệu điểm thi THPT Quốc Gia năm 2016 của gần 35.000 thí sinh.
- Đọc dữ liệu. Xem thông tin dữ liệu:

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 34826 entries, 0 to 34825
Data columns (total 6 columns):
SOBAODANH      34826 non-null object
HO_TEN         34826 non-null object
NGAY_SINH      34826 non-null object
TEN_CUMTHI     34826 non-null object
GIOI_TINH      34826 non-null object
DIEM_THI       34826 non-null object
dtypes: object(6)
memory usage: 1.6+ MB
```

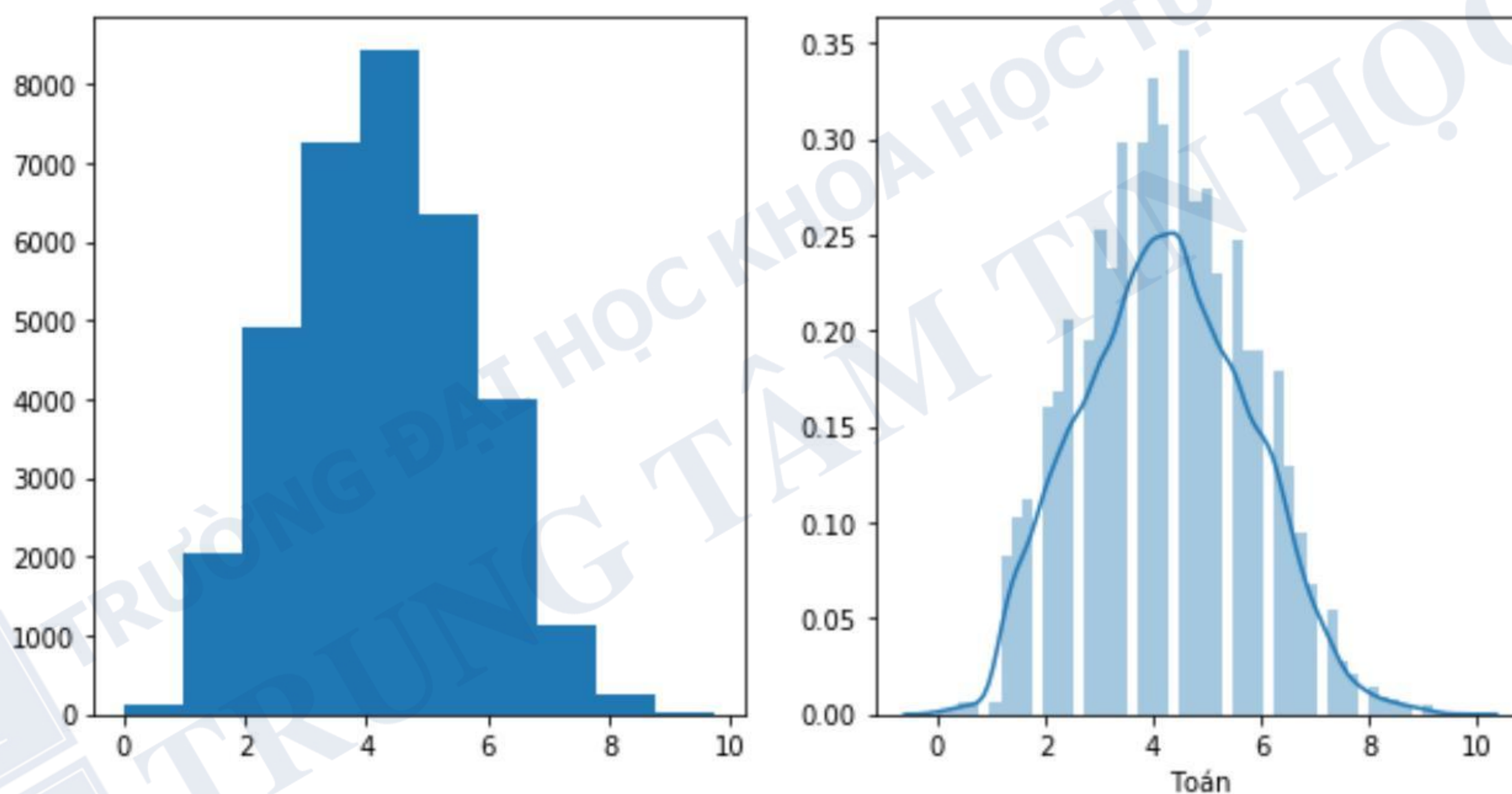
- Biết danh sách các môn thi là: "Toán", "Ngữ văn", "Địa lí", "Tiếng Anh", "Sinh học", "Vật lí", "Hóa học", "Lịch sử". Một thí sinh chỉ thi các môn bắt buộc chung còn các môn tự chọn có thể khác nhau.
- Với dữ liệu hiện tại, cột DIEM_THI là chuỗi chứa điểm thi của tất cả các môn mà một thí sinh thi:

| | SOBAODANH | HO_TEN | NGAY_SINH | TEN_CUMTHI | GIOI_TINH | DIEM_THI |
|---|-----------|---------------|------------|-------------------|-----------|--|
| 0 | 018000001 | DƯƠNG VIỆT AN | 12/03/1998 | Sở GDĐT Bắc Giang | Nam | Toán: 2.00 Ngữ văn: 5.50 Lịch sử: 3.00 |
| 1 | 018000002 | ĐỖ VĂN AN | 09/12/1998 | Sở GDĐT Bắc Giang | Nam | Toán: 5.50 Ngữ văn: 5.25 Địa lí: 5.50 |
| 2 | 018000003 | ĐỖ XUÂN AN | 12/08/1997 | Sở GDĐT Bắc Giang | Nam | Toán: 4.50 Ngữ văn: 5.50 Địa lí: 3.75 |
| 3 | 018000004 | ĐẶNG PHÚC AN | 19/03/1998 | Sở GDĐT Bắc Giang | Nữ | Toán: 3.00 Ngữ văn: 6.00 Địa lí: 5.50 |
| 4 | 018000005 | ĐẶNG VĂN AN | 25/10/1998 | Sở GDĐT Bắc Giang | Nam | Toán: 2.25 Ngữ văn: 4.75 Địa lí: 5.25 |

- Và như vậy thì chúng ta sẽ không phân tích được điểm thi của thí sinh. Do đó, việc đầu tiên là phải tiền xử lý dữ liệu. Từ dữ liệu trong cột DIEM_THI, hãy tạo ra các cột tương ứng với danh sách các môn thi nói trên và đưa điểm của thí sinh từ chuỗi vào các cột, môn nào thí sinh không thi thì sẽ để NaN, như sau:

| | SOBAODANH | HO_TEN | NGAY_SINH | TEN_CUMTHI | GIOI_TINH | DIEM_THI | Toán | Ngữ văn | Địa lí | Tiếng Anh | Sinh học | Vật lí | Hóa học | Lịch sử |
|---|-----------|---------------|------------|-------------------|-----------|--|------|---------|--------|-----------|----------|--------|---------|---------|
| 0 | 018000001 | DƯƠNG VIỆT AN | 12/03/1998 | Sở GDĐT Bắc Giang | Nam | Toán: 2.00 Ngữ văn: 5.50 Lịch sử: 3.00 | 2.00 | 5.50 | 5.00 | NaN | NaN | NaN | NaN | 3.0 |
| 1 | 018000002 | ĐỖ VĂN AN | 09/12/1998 | Sở GDĐT Bắc Giang | Nam | Toán: 5.50 Ngữ văn: 5.25 Địa lí: 5.50 | 5.50 | 5.25 | 5.50 | 3.68 | NaN | NaN | NaN | NaN |
| 2 | 018000003 | ĐỖ XUÂN AN | 12/08/1997 | Sở GDĐT Bắc Giang | Nam | Toán: 4.50 Ngữ văn: 5.50 Địa lí: 3.75 | 4.50 | 5.50 | 3.75 | 2.25 | NaN | NaN | NaN | NaN |
| 3 | 018000004 | ĐẶNG PHÚC AN | 19/03/1998 | Sở GDĐT Bắc Giang | Nữ | Toán: 3.00 Ngữ văn: 6.00 Địa lí: 5.50 | 3.00 | 6.00 | 5.50 | 1.50 | NaN | NaN | NaN | NaN |
| 4 | 018000005 | ĐẶNG VĂN AN | 25/10/1998 | Sở GDĐT Bắc Giang | Nam | Toán: 2.25 Ngữ văn: 4.75 Địa lí: 5.25 | 2.25 | 4.75 | 5.25 | 2.00 | NaN | NaN | NaN | NaN |

- Hãy vẽ biểu đồ phân phối tần suất điểm thi, mỗi điểm thi là một biểu đồ, nhận xét trên từng biểu đồ: các thống kê mô tả, phân phối chuẩn hay nghiêng? Đường cong cao hơn hay thấp hơn phân phối chuẩn...
- Ví dụ: Môn “Toán”



- Nhận xét: ...
- Lưu dữ liệu điểm thi sau khi đã chuẩn hóa để sử dụng.