

Uni-Variate Statistics

Sourav SEN GUPTA
Lecturer, SCSE, NTU



Sample
COLLECTION



Practical
MOTIVATION

Data
PREPARATION



Problem
FORMULATION

Exploratory
ANALYSIS



Statistical
DESCRIPTION

Analytic
VISUALIZATION



Pattern
RECOGNITION

Algorithmic
OPTIMIZATION



Machine
LEARNING

Information
PRESSENTATION



Statistical
INFERENCE

Ethical
CONSIDERATION



Intelligent
DECISION

Data Science Uni-Variate Statistics

Exploratory Analysis

What are the Variables in the Data?
How to characterize the Variables?
How to find relation between them?

**How to intelligently
explore acquired Data?**



Data Science

The Pokemon Dataset

#	Name	Type 1	Type 2	Total	HP	Attack	Defense	Sp. Atk	Sp. Def	Speed	Generation	Legendary
430	Honchkrow	Dark	Flying	505	100	125	52	105	52	71	4	False
338	Solrock	Rock	Psychic	440	70	95	85	55	65	70	3	False
32	Nidoran♂	Poison	NaN	273	46	57	40	40	40	50	1	False
442	Spiritomb	Ghost	Dark	485	50	92	108	92	108	35	4	False
480	Uxie	Psychic	NaN	580	75	75	130	75	130	95	4	True
536	Palpitoad	Water	Ground	384	75	65	55	65	55	69	5	False
360	Wynaut	Psychic	NaN	260	95	23	48	23	48	23	3	False
478	Froslass	Ice	Ghost	480	70	80	70	80	70	110	4	False
76	Golem	Rock	Ground	495	80	120	130	55	65	45	1	False
177	Natu	Psychic	Flying	320	40	50	45	70	45	70	2	False

Source : Kaggle Datasets | **Pokemon with stats** by Alberto Barradas | <https://www.kaggle.com/abcsds/pokemon>

45, 60, 80, 80, 39, 58, 78, 78, 78, 44, 59, 79, 79,
45, 50, 60, 40, 45, 65, 65, 40, 63, 83, 83, 30, 55,
40, 65, 35, 60, 35, 60, 50, 75, 55, 70, 90, 46, 61,
81, 70, 95, 38, 73, 115, 140, 40, 75, 45, 60, 75, 35,
60, 60, 70, 10, 35, 40, 65, 50, 80, 40, 65, 55, 90,
40, 65, 90, 25, 40, 55, 55, 70, 80, 90, 50, 65, 80,
40, 80, 40, 55, 80, 50, 65, 90, 95, 95, 25, 50, 52,
35, 60, 65, 90, 80, 105, 30, 50, 30, 45, 60, 60, 35,
60, 85, 30, 55, 40, 60, 60, 95, 50, 60, 50, 50, 90,
40, 65, 80, 105, 250, 65, 105, 105, 30, 55, 45, 80, 30,
60, 40, 70, 65, 65, 65, 65, 75, 20, 95, 95, 130,
48, 55, 130, 65, 65, 65, 35, 70, 30, 60, 80, 80, 160,
90, 90, 41, 61, 91, 106, 106, 106, 100, 45, 60, 80,
39, 58, 78, 50, 65, 85, 35, 85, 60, 100, 40, 55, 40,
70, 85, 75, 125, 20, 50, 90, 35, 55, 40, 65, 55, 70,
90, 90, 75, 70, 100, 70, 90, 35, 55, 75, 55, 30, 75,
65, 55, 95, 65, 95, 60, 95, 60, 48, 190, 70, 50, 75,
100, 65, 75, 75, 60, 90, 65, 70, 70, 20, 80, 80, 55,
60, 90, 40, 50, 50, 100, 55, 35, 75, 45, 65, 65, 45,
75, 75, 75, 90, 90, 85, 73, 55, 35, 50, 45, 45, 45,
95, 255, 90, 115, 100, 50, 70, 100, 100, 106, 106, 100, 40,
50, 70, 70, 45, 60, 80, 80, 50, 70, 100, 100, 35, 70,
38, 78, 45, 50, 60, 50, 60, 40, 60, 80, 40, 70, 90,
40, 60, 40, 60, 28, 38, 68, 68, 40, 70, 60, 60, 60,
80, 150, 31, 61, 1, 64, 84, 104, 72, 144, 50, 30, 50,
70, 50, 50, 50, 50, 60, 70, 70, 30, 60, 60, 40,
70, 70, 60, 60, 65, 65, 50, 70, 100, 45, 70, 70, 130,
170, 60, 70, 70, 70, 60, 80, 60, 45, 50, 80, 50, 70,
45, 75, 75, 73, 73, 70, 70, 50, 110, 43, 63, 40, 60,
66, 86, 45, 75, 20, 95, 70, 60, 44, 64, 64, 20, 40,
99, 65, 65, 65, 95, 50, 80, 80, 70, 90, 110, 35, 55,
55, 100, 43, 45, 65, 95, 95, 40, 60, 80, 80, 80, 80,
80, 80, 80, 80, 100, 100, 100, 105, 105, 100, 50,

Data Science

Uni-Variate Statistics

Numeric Uni-Variate Data

HP

Hit Point of Pokemon

Numeric Variable

800 Values in Total

Pertinent Questions

- How to describe the Data?
- How to analyze the Data?

45, 60, 80, 80, 39, 58, 78, 78, 78, 44, 59, 79, 79,
45, 50, 60, 40, 45, 65, 65, 40, 63, 83, 83, 30, 55,
40, 65, 35, 60, 35, 60, 50, 75, 55, 70, 90, 46, 61,
81, 70, 95, 38, 73, 115, 140, 40, 75, 45, 60, 75, 35,
60, 60, 70, 10, 35, 40, 65, 50, 80, 40, 65, 55, 90,
40, 65, 90, 25, 40, 55, 55, 70, 80, 90, 50, 65, 80,
40, 80, 40, 55, 80, 50, 65, 90, 95, 95, 25, 50, 52,
35, 60, 65, 90, 80, 105, 30, 50, 30, 45, 60, 60, 35,
60, 85, 30, 55, 40, 60, 60, 95, 50, 60, 50, 50, 90,
40, 65, 80, 105, 250, 65, 105, 105, 30, 55, 45, 80, 30,
60, 40, 70, 65, 65, 65, 65, 75, 20, 95, 95, 130,
48, 55, 130, 65, 65, 65, 35, 70, 30, 60, 80, 80, 160,
90, 90, 41, 61, 91, 106, 106, 106, 100, 45, 60, 80,
39, 58, 78, 50, 65, 85, 35, 85, 60, 100, 40, 55, 40,
70, 85, 75, 125, 20, 50, 90, 35, 55, 40, 65, 55, 70,
90, 90, 75, 70, 100, 70, 90, 35, 55, 75, 55, 30, 75,
65, 55, 95, 65, 95, 60, 95, 60, 48, 190, 70, 50, 75,
100, 65, 75, 75, 60, 90, 65, 70, 70, 20, 80, 80, 55,
60, 90, 40, 50, 50, 100, 55, 35, 75, 45, 65, 65, 45,
75, 75, 75, 90, 90, 85, 73, 55, 35, 50, 45, 45, 45,
95, 255, 90, 115, 100, 50, 70, 100, 100, 106, 106, 100, 40,
50, 70, 70, 45, 60, 80, 80, 50, 70, 100, 100, 35, 70,
38, 78, 45, 50, 60, 50, 60, 40, 60, 80, 40, 70, 90,
40, 60, 40, 60, 28, 38, 68, 68, 40, 70, 60, 60, 60,
80, 150, 31, 61, 1, 64, 84, 104, 72, 144, 50, 30, 50,
70, 50, 50, 50, 50, 60, 70, 70, 30, 60, 60, 40,
70, 70, 60, 60, 65, 65, 50, 70, 100, 45, 70, 70, 130,
170, 60, 70, 70, 70, 60, 80, 60, 45, 50, 80, 50, 70,
45, 75, 75, 73, 73, 70, 70, 50, 110, 43, 63, 40, 60,
66, 86, 45, 75, 20, 95, 70, 60, 44, 64, 64, 20, 40,
99, 65, 65, 65, 95, 50, 80, 80, 70, 90, 110, 35, 55,
55, 100, 43, 45, 65, 95, 95, 40, 60, 80, 80, 80, 80,
80, 80, 80, 80, 100, 100, 100, 105, 105, 100, 50,

Data Science

Uni-Variate Statistics

Basic Summary of the Data

HP

The Average Hit Point

Deviation from Average

Maximum and Minimum

Statistical Questions

- What is the Central Tendency?
- What is the Spread of the Data?



45, 60, 80, 80, 39, 58, 78, 78, 78, 44, 59, 79, 79,
45, 50, 60, 40, 45, 65, 65, 40, 63, 83, 83, 30, 55,
40, 65, 35, 60, 35, 60, 50, 75, 55, 70, 90, 46, 61,
81, 70, 95, 38, 73, 115, 140, 40, 75, 45, 60, 75, 35,
60, 60, 70, 10, 35, 40, 65, 50, 80, 40, 65, 55, 90,
40, 65, 90, 25, 40, 55, 55, 70, 80, 90, 50, 65, 80,
40, 80, 40, 55, 80, 50, 65, 90, 95, 95, 25, 50, 52,
35, 60, 65, 90, 80, 105, 30, 50, 30, 45, 60, 60, 35,
60, 85, 30, 55, 40, 60, 60, 95, 50, 60, 50, 50, 90,
40, 65, 80, 105, 250, 65, 105, 105, 30, 55, 45, 80, 30,
60, 40, 70, 65, 65, 65, 65, 75, 20, 95, 95, 130,
48, 55, 130, 65, 65, 65, 35, 70, 30, 60, 80, 80, 160,
90, 90, 41, 61, 91, 106, 106, 106, 100, 45, 60, 80,
39, 58, 7
70, 85, 7
90, 90, 7
65, 55, 9
100, 65, 7
69.258750
60, 90, 4
75, 75, 75, 90, 90, 80, 15, 50, 50, 50, 50, 45, 45, 45,
95, 255, 90, 115, 100, 50, 70, 100, 100, 106, 106, 100, 40,
50, 70, 70, 45, 60, 80, 80, 50, 70, 100, 100, 35, 70,
38, 78, 45, 50, 60, 50, 60, 40, 60, 80, 40, 70, 90,
40, 60, 40, 60, 28, 38, 68, 68, 40, 70, 60, 60, 60,
80, 150, 31, 61, 1, 64, 84, 104, 72, 144, 50, 30, 50,
70, 50, 50, 50, 50, 60, 70, 70, 30, 60, 60, 40,
70, 70, 60, 60, 65, 65, 50, 70, 100, 45, 70, 70, 130,
170, 60, 70, 70, 70, 60, 80, 60, 45, 50, 80, 50, 70,
45, 75, 75, 73, 73, 73, 70, 70, 50, 110, 43, 63, 40, 60,
66, 86, 45, 75, 20, 95, 70, 60, 44, 64, 64, 64, 20, 40,
99, 65, 65, 65, 95, 50, 80, 80, 70, 90, 110, 35, 55,
55, 100, 43, 45, 65, 95, 95, 40, 60, 80, 80, 80, 80,
80, 80, 80, 80, 100, 100, 100, 105, 105, 100, 50,

Data Science

Uni-Variate Statistics

Central Tendency : Mean

Natural Intuition

Average Hit Point of a Pokemon

Statistical Definition

Sum of Data / Count of Data

$$\bar{x} = \frac{x_1 + x_2 + \cdots + x_n}{n}$$

Source : Kaggle Datasets | **Pokemon with stats** by Alberto Barradas | <https://www.kaggle.com/abcsds/pokemon>



45,	60,	80,	80,	39,	58,	78,	78,	78,	44,	59,	79,	79,
45,	50,	60,	40,	45,	65,	65,	40,	63,	83,	83,	30,	55,
40,	65,	35,	60,	35,	60,	50,	75,	55,	70,	90,	46,	61,
81,	70,	95,	38,	73,	115,	140,	40,	75,	45,	60,	75,	35,
60,	60,	70,	10,	35,	40,	65,	50,	80,	40,	65,	55,	90,
40,	65,	90,	25,	40,	55,	55,	70,	80,	90,	50,	65,	80,
40,	80,	40,	55,	80,	50,	65,	90,	95,	95,	25,	50,	52,
35,	60,	65,	90,	80,	105,	30,	50,	30,	45,	60,	60,	35,
60,	85,	30,	55,	40,	60,	60,	95,	50,	60,	50,	50,	90,
40,	65,	80,	105,	250,	65,	105,	105,	30,	55,	45,	80,	30,
60,	40,	70,	65,	65,	65,	65,	75,	20,	95,	95,	130,	
48,	55,	130,	65,	65,	65,	35,	70,	30,	60,	80,	80,	160,
90,	90,	41,	61,	91,	106,	106,	106,	100,	45,	60,	80,	
39,	58,	7	--	--	--	--	--	--	--	--	--	
70,	85,	7	--	--	--	--	--	--	--	--	--	
90,	90,	7	--	--	--	--	--	--	--	--	--	
65,	55,	9	--	--	--	--	--	--	--	--	--	
100,	65,	7	--	--	--	--	--	--	--	--	--	

25.534669

Data Science

Uni-Variate Statistics

Dispersion : Standard Deviation

Natural Intuition

Average Deviation from the Mean

Statistical Definition

Sum of Deviation / Count of Data

$$\sqrt{\frac{(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_n - \bar{x})^2}{n}}$$

Source : Kaggle Datasets | **Pokemon with stats** by Alberto Barradas | <https://www.kaggle.com/abcsds/pokemon>



45, 60, 80, 80, 39, 58, 78, 78, 78, 44, 59, 79, 79,
45, 50, 60, 40, 45, 65, 65, 40, 63, 83, 83, 30, 55,
40, 65, 35, 60, 35, 60, 50, 75, 55, 70, 90, 46, 61,
81, 70, 95, 38, 73, 115, 140, 40, 75, 45, 60, 75, 35,
60, 60, 70, 10, 35, 40, 65, 50, 80, 40, 65, 55, 90,
40, 65, 90, 25, 40, 55, 55, 70, 80, 90, 50, 65, 80,
40, 80, 40, 55, 80, 50, 65, 90, 95, 95, 25, 50, 52,
35, 60, 65, 90, 80, 105, 30, 50, 30, 45, 60, 60, 35,
60, 85, 30, 55, 40, 60, 60, 95, 50, 60, 50, 50, 90,
40, 65, 80, 105, 250, 65, 105, 105, 30, 55, 45, 80, 30,
60, 40, 70, 65, 65, 65, 65, 75, 20, 95, 95, 130,
48, 55, 130, 65, 65, 65, 35, 70, 30, 60, 80, 80, 160,
90, 90, 90, 41, 61, 91, 106, 106, 106, 100, 45, 60, 80,
39, 58, 7, --, --, --, --, --, --, --, --, 55, 40,
70, 85, 7, 5, 55, 70,
90, 90, 7, 5, 30, 75,
65, 55, 9, 5, 50, 75,
100, 65, 7, 5, 80, 55,
60, 90, 4, 5, 65, 45,
75, 75, 75, 75, 80, 80, 80, 80, 80, 80, 45, 45, 45,
95, 255, 90, 115, 100, 50, 70, 100, 100, 106, 106, 100, 40,
50, 70, 70, 45, 60, 80, 80, 50, 70, 100, 100, 35, 70,
38, 78, 45, 50, 60, 50, 60, 40, 60, 80, 40, 70, 90,
40, 60, 40, 60, 28, 38, 68, 68, 40, 70, 60, 60, 60,
80, 150, 31, 61, 1, 64, 84, 104, 72, 144, 50, 30, 50,
70, 50, 50, 50, 50, 60, 70, 70, 30, 60, 60, 40,
70, 70, 60, 60, 65, 65, 50, 70, 100, 45, 70, 70, 130,
170, 60, 70, 70, 70, 60, 80, 60, 45, 50, 80, 50, 70,
45, 75, 75, 73, 73, 70, 70, 50, 110, 43, 63, 40, 60,
66, 86, 45, 75, 20, 95, 70, 60, 44, 64, 64, 20, 40,
99, 65, 65, 65, 95, 50, 80, 80, 70, 90, 110, 35, 55,
55, 100, 43, 45, 65, 95, 95, 40, 60, 80, 80, 80, 80,
80, 80, 80, 80, 100, 100, 100, 105, 105, 100, 50,

65.00000

Data Science

Uni-Variate Statistics

Central Tendency : Median

Natural Intuition

Mid-Value of Pokemon Hit Points

Statistical Definition

Marker to Divide the Data 50:50

$$P(x \leq x_M) = P(x \geq x_M) = 0.5$$



45, 60, 80, 80, 39, 58, 78, 78, 78, 44, 59, 79, 79,
45, 50, 60, 40, 45, 65, 65, 40, 63, 83, 83, 30, 55,
40, 65, 35, 60, 35, 60, 50, 75, 55, 70, 90, 46, 61,
81, 70, 95, 38, 73, 115, 140, 40, 75, 45, 60, 75, 35,
60, 60, 70, 10, 35, 40, 65, 50, 80, 40, 65, 55, 90,
40, 65, 90, 25, 40, 55, 55, 70, 80, 90, 50, 65, 80,
40, 80, 40, 55, 80, 50, 65, 90, 95, 95, 25, 50, 52,
35, 60, 65, 90, 80, 105, 30, 50, 30, 45, 60, 60, 35,
60, 85, 30, 55, 40, 60, 60, 95, 50, 60, 50, 50, 90,
40, 65, 80, 105, 250, 65, 105, 105, 30, 55, 45, 80, 30,
60, 40, 70 65 65 65 65 75 20 95, 95, 130,
48, 55, 13
90, 90, 9
39, 58, 7
70, 85, 7
90, 90, 7
65, 55, 9
100, 65, 7
60, 90, 4
75, 75, 7
95, 255, 9
50, 70, 7
38, 78, 45, 50, 60, 50, 60, 40, 60, 80, 40, 70, 90,
40, 60, 40, 60, 28, 38, 68, 68, 40, 70, 60, 60, 60,
80, 150, 31, 61, 1, 64, 84, 104, 72, 144, 50, 30, 50,
70, 50, 50, 50, 50, 60, 70, 70, 30, 60, 60, 40,
70, 70, 60, 60, 65, 65, 50, 70, 100, 45, 70, 70, 130,
170, 60, 70, 70, 70, 60, 80, 60, 45, 50, 80, 50, 70,
45, 75, 75, 73, 73, 70, 70, 50, 110, 43, 63, 40, 60,
66, 86, 45, 75, 20, 95, 70, 60, 44, 64, 64, 20, 40,
99, 65, 65, 65, 95, 50, 80, 80, 70, 90, 110, 35, 55,
55, 100, 43, 45, 65, 95, 95, 40, 60, 80, 80, 80, 80,
80, 80, 80, 80, 100, 100, 100, 100, 105, 105, 100, 50,

50.000000

80.000000

Data Science

Uni-Variate Statistics

Dispersion : Quantiles

Natural Intuition

Distribution of Pokemon Hit Points

Statistical Definition

Markers to Divide the Data 25:50:25

$$P(x \leq x_{Q1}) = 0.25, P(x \geq x_{Q3}) = 0.25$$

$$P(x_{Q1} \leq x \leq x_{Q3}) = 0.5$$



45,	60,	80,	80,	39,	58,	78,	78,	78,	44,	59,	79,	79,
45,	50,	60,	40,	45,	65,	65,	40,	63,	83,	83,	30,	55,
40,	65,	35,	60,	35,	60,	50,	75,	55,	70,	90,	46,	61,
81,	70,	95,	29	72	115	140	40	75	45,	60,	75,	35,
60,	60,	70,								65,	55,	90,
40,	65,	90,								50,	65,	80,
40,	80,	40,								25,	50,	52,
35,	60,	65,								60,	60,	35,
60,	85,	30,								50,	50,	90,
40,	65,	80,	1							45,	80,	30,
60,	40,	70,								95,	95,	130,
48,	55,	130,								80,	80,	160,
90,	90,	90,								45,	60,	80,
39,	58,	78,								40,	55,	40,
70,	85,	75,	1							65,	55,	70,
90,	90,	75,								55,	30,	75,
65,	55,	95,								70,	50,	75,
100,	65,	75,								80,	80,	55,
60,	90,	40,								65,	65,	45,
75,	75,	75,								45,	45,	45,
95,	255,	90,	1							106,	100,	40,
50,	70,	70,								100,	35,	70,
38,	78,	45,								40,	70,	90,
40,	60,	40,								60,	60,	60,
80,	150,	31,								50,	30,	50,
70,	50,	50,								60,	60,	40,
70,	70,	60,								70,	70,	130,
170,	60,	70,								80,	50,	70,
45,	75,	75,								63,	40,	60,
66,	86,	45,	75,	20,	95,	70,	60,	44,	64,	64,	20,	40,
99,	65,	65,	65,	95,	50,	80,	80,	70,	90,	110,	35,	55,
55,	100,	43,	45,	65,	95,	95,	40,	60,	80,	80,	80,	80,
80,	80,	80,	80,	100,	100,	100,	100,	105,	105,	100,	50,	

Data Science

Uni-Variate Statistics

Statistical Summary of the Data

HP

The Average Hit Point
Deviation from Average
Median and Quantiles

count 800.000000

mean 69.258750

std 25.534669

min 1.000000

25% 50.000000

50% 65.000000

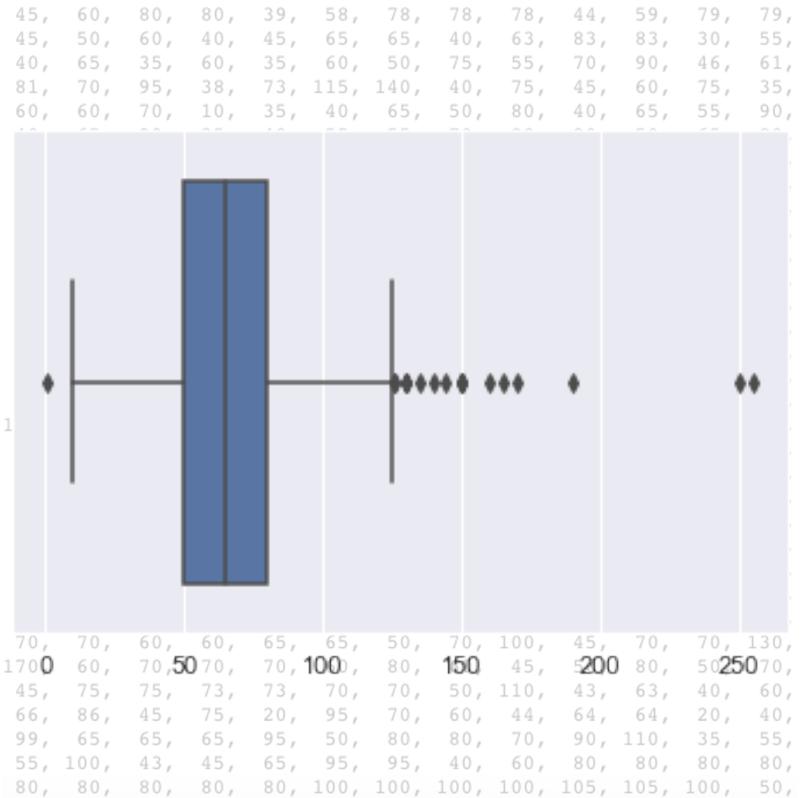
75% 80.000000

max 255.000000

Statistical Questions

- What is the Central Tendency?
- What is the Spread of the Data?





Data Science

Uni-Variate Statistics

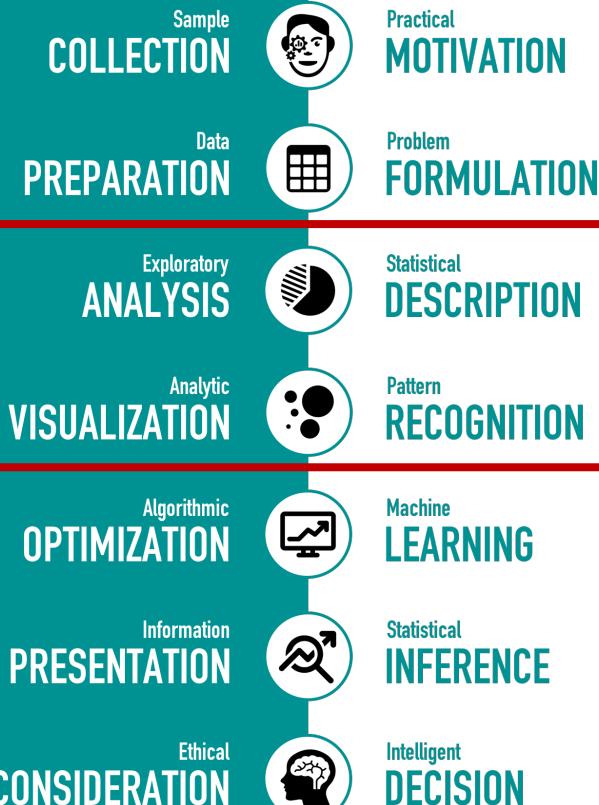
Statistical Summary of the Data

HP

The Average Hit Point
Deviation from Average
Median and Quantiles

Statistical Questions

- What is the Central Tendency?
- What is the Spread of the Data?



Data Science Pipeline Exploratory Analysis

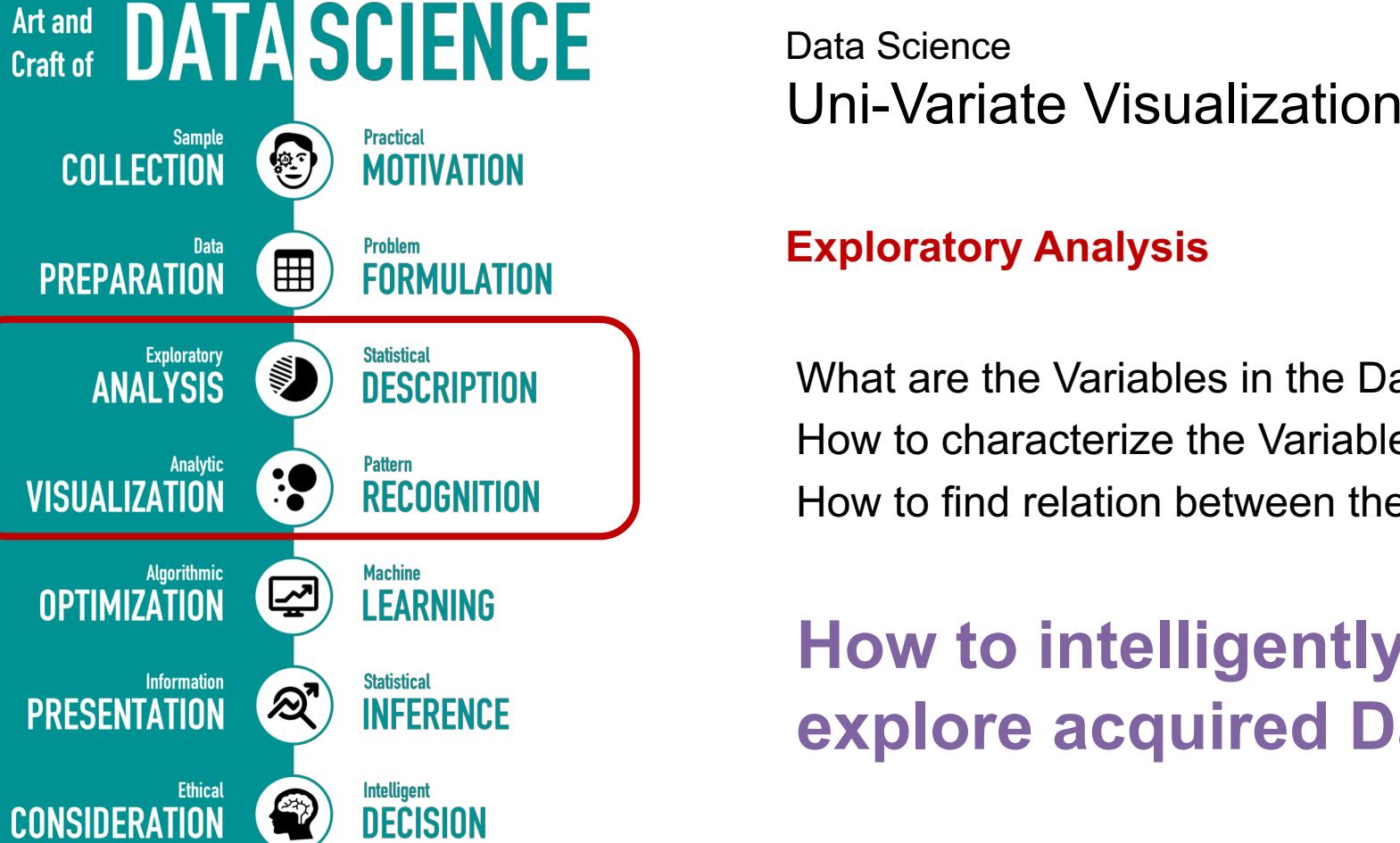
How to summarize the acquired Data?
How to visualize the acquired Data?
How to analyze the acquired Data?

How to intelligently
explore acquired Data?

Uni-Variate Visualization

Sourav SEN GUPTA
Lecturer, SCSE, NTU





Data Science Uni-Variate Visualization

Exploratory Analysis

What are the Variables in the Data?
How to characterize the Variables?
How to find relation between them?

How to intelligently explore acquired Data?



Data Science

The Pokemon Dataset

#	Name	Type 1	Type 2	Total	HP	Attack	Defense	Sp. Atk	Sp. Def	Speed	Generation	Legendary
430	Honchkrow	Dark	Flying	505	100	125	52	105	52	71	4	False
338	Solrock	Rock	Psychic	440	70	95	85	55	65	70	3	False
32	Nidoran♂	Poison	NaN	273	46	57	40	40	40	50	1	False
442	Spiritomb	Ghost	Dark	485	50	92	108	92	108	35	4	False
480	Uxie	Psychic	NaN	580	75	75	130	75	130	95	4	True
536	Palpitoad	Water	Ground	384	75	65	55	65	55	69	5	False
360	Wynaut	Psychic	NaN	260	95	23	48	23	48	23	3	False
478	Froslass	Ice	Ghost	480	70	80	70	80	70	110	4	False
76	Golem	Rock	Ground	495	80	120	130	55	65	45	1	False
177	Natu	Psychic	Flying	320	40	50	45	70	45	70	2	False

Source : Kaggle Datasets | [Pokemon with stats](#) by Alberto Barradas | <https://www.kaggle.com/abcsds/pokemon>

45, 60, 80, 80, 39, 58, 78, 78, 78, 44, 59, 79, 79,
45, 50, 60, 40, 45, 65, 65, 40, 63, 83, 30, 55,
40, 65, 35, 60, 35, 60, 50, 75, 55, 70, 90, 46, 61,
81, 70, 95, 60, 60, 70, 40, 65, 90, 50, 65, 80,
40, 65, 90, 40, 80, 40, 35, 60, 65, 30, 50, 50, 55, 90,
60, 85, 30, 40, 65, 80, 60, 40, 70, 48, 55, 130,
90, 90, 90, 39, 58, 78, 70, 85, 75, 90, 90, 90, 75,
65, 55, 95, 100, 65, 75, 60, 90, 40, 75, 255, 90,
50, 70, 70, 38, 78, 45, 40, 60, 40, 80, 150, 31, 70, 50, 50,
70, 70, 60, 170, 60, 70, 45, 75, 75, 66, 86, 45, 75, 20, 95, 70, 60, 44, 64,
99, 65, 65, 65, 95, 50, 80, 80, 70, 90, 110, 35, 55, 55, 100, 43, 45, 65, 95, 95, 40, 60, 80, 80, 80, 80,
80, 80, 80, 80, 100, 100, 100, 100, 105, 105, 100, 50,

Data Science

Uni-Variate Statistics

Statistical Summary of the Data

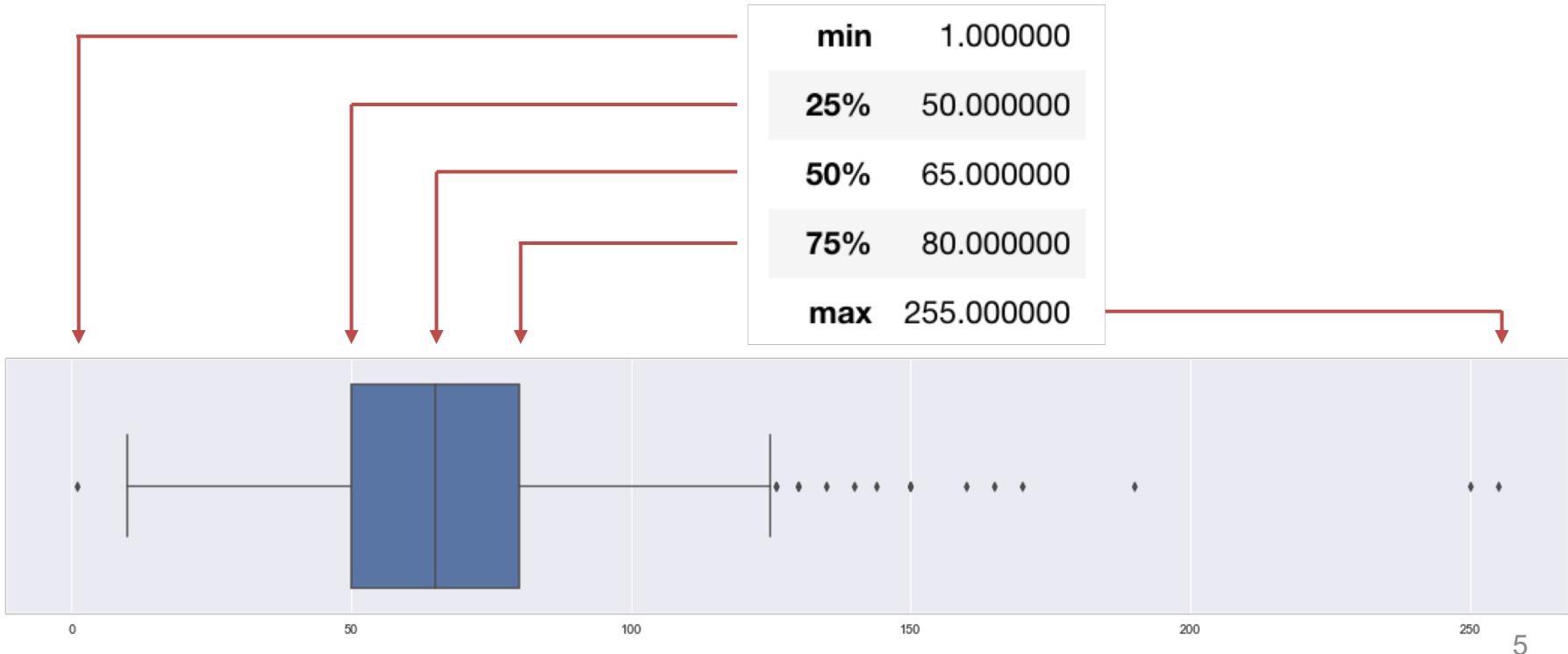
HP The Average Hit Point
 Deviation from Average
 Median and Quantiles

Statistical Questions

- What is the Central Tendency?
- What is the Spread of the Data?

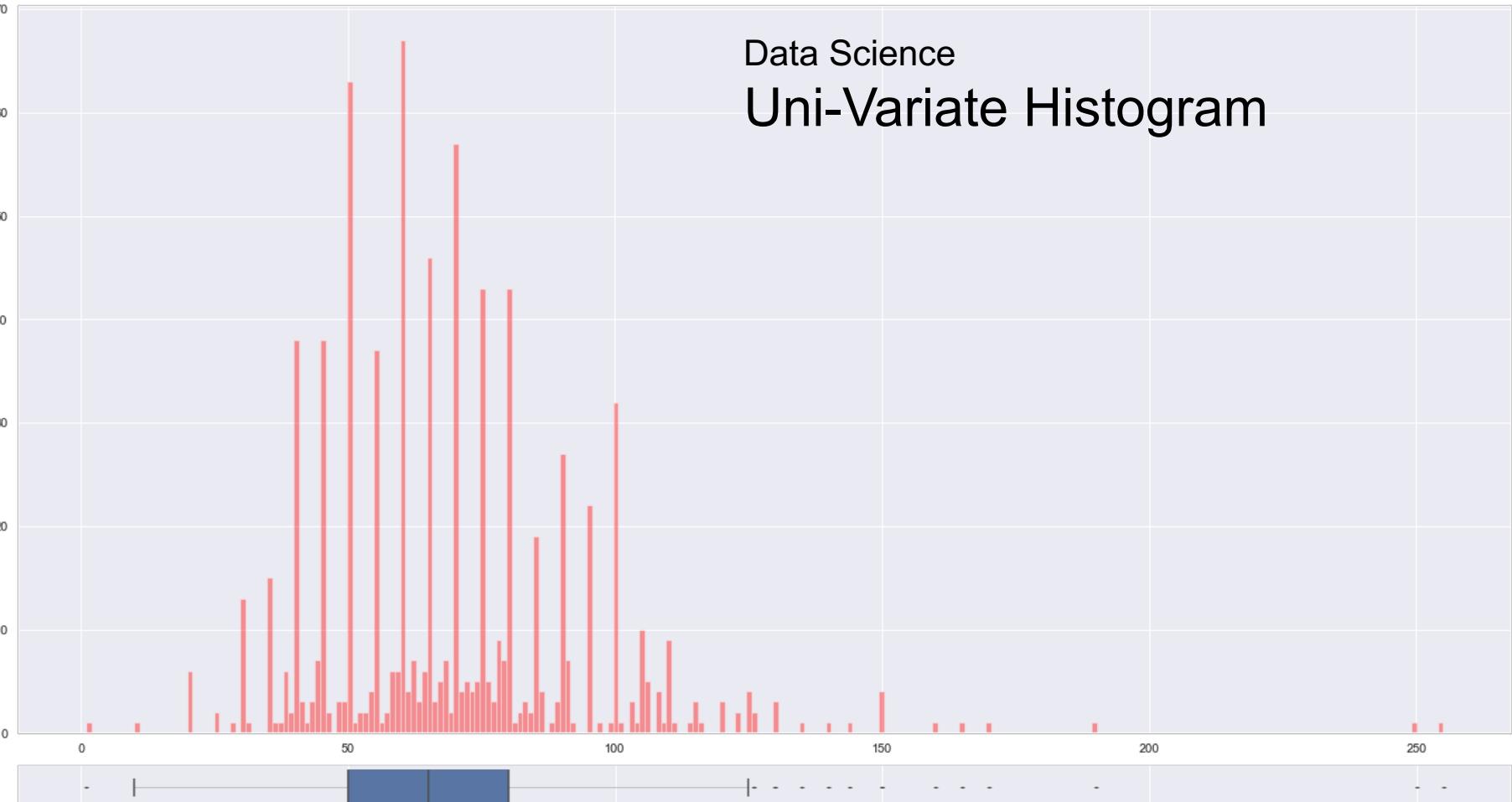
Data Science

Uni-Variate Box-Plot



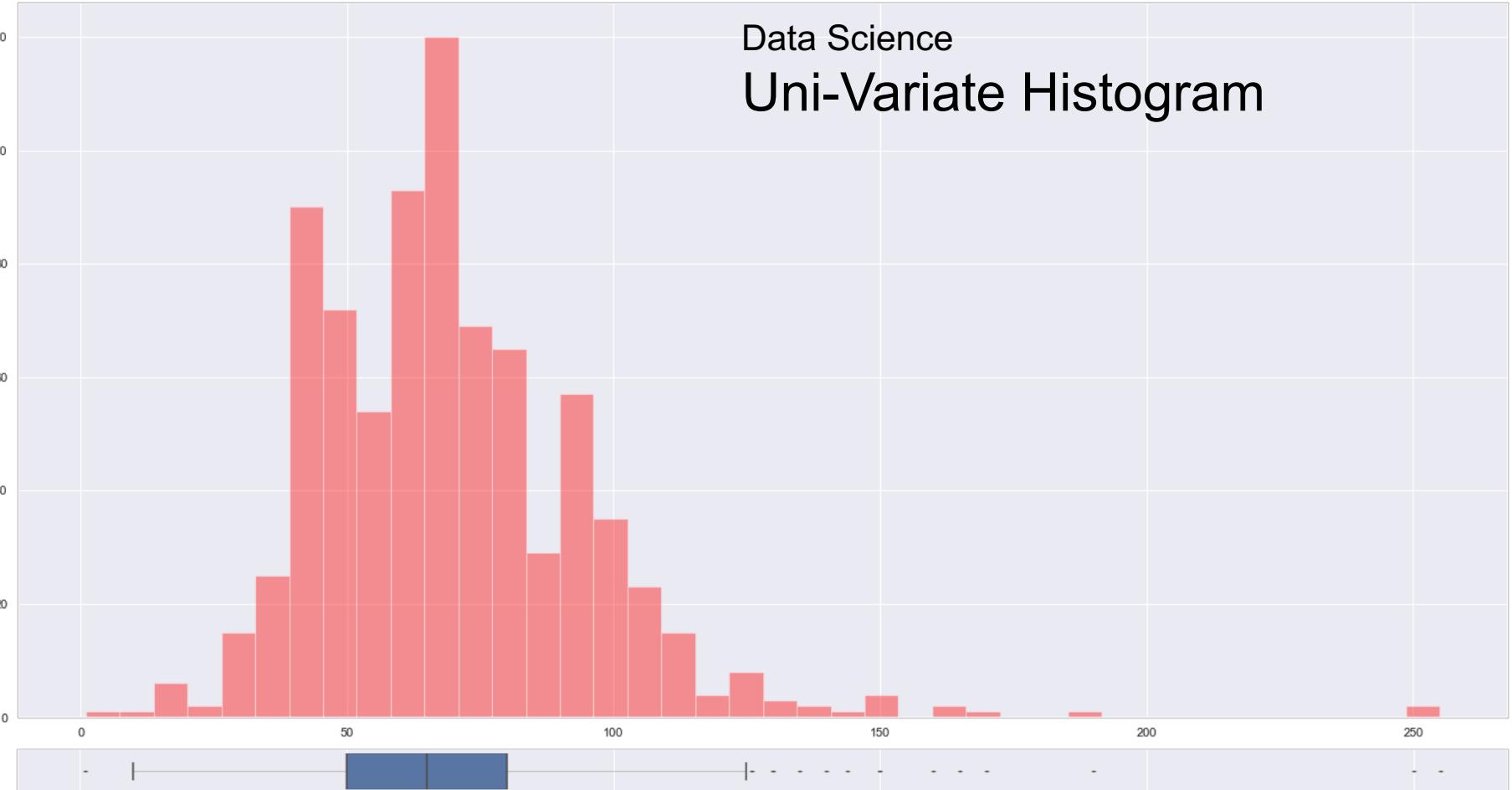
Data Science

Uni-Variate Histogram



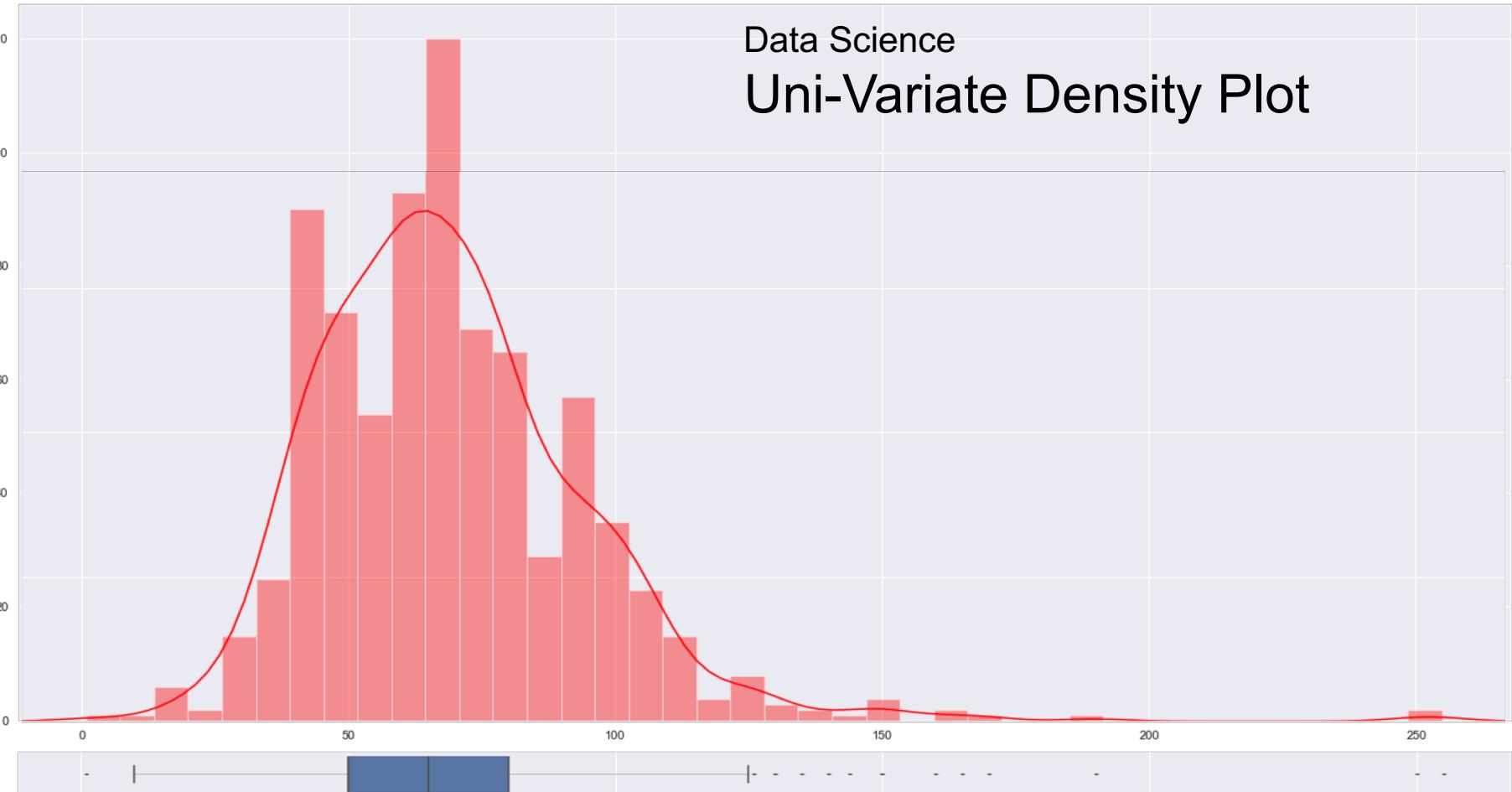
Data Science

Uni-Variate Histogram



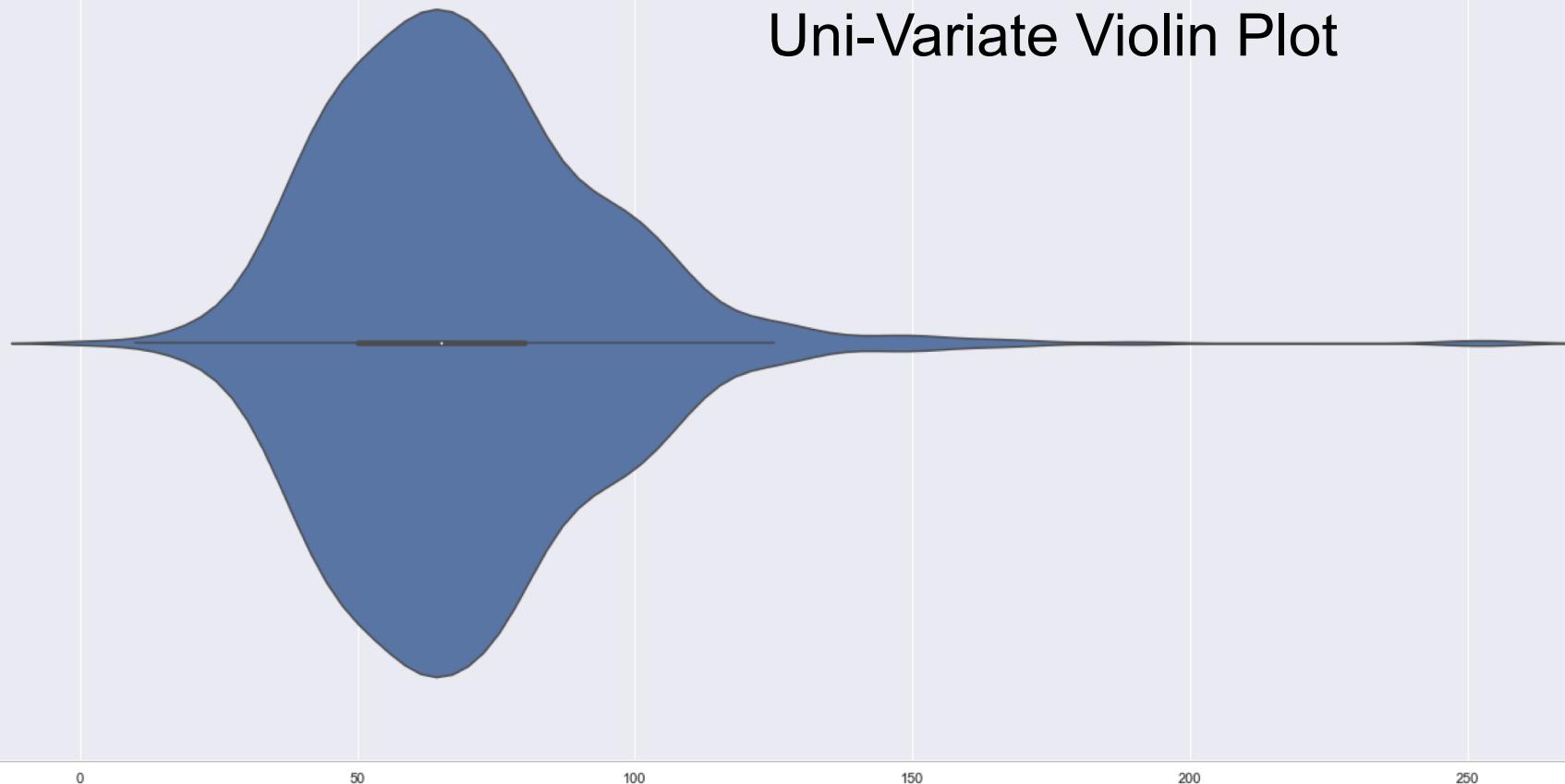
Data Science

Uni-Variate Density Plot



Data Science

Uni-Variate Violin Plot



Sample
COLLECTION



Practical
MOTIVATION

Data
PREPARATION



Problem
FORMULATION

Exploratory
ANALYSIS



Statistical
DESCRIPTION

Analytic
VISUALIZATION



Pattern
RECOGNITION

Algorithmic
OPTIMIZATION



Machine
LEARNING

Information
PRESENTATION



Statistical
INFERENCE

Ethical
CONSIDERATION



Intelligent
DECISION

Data Science Pipeline **Exploratory Analysis**

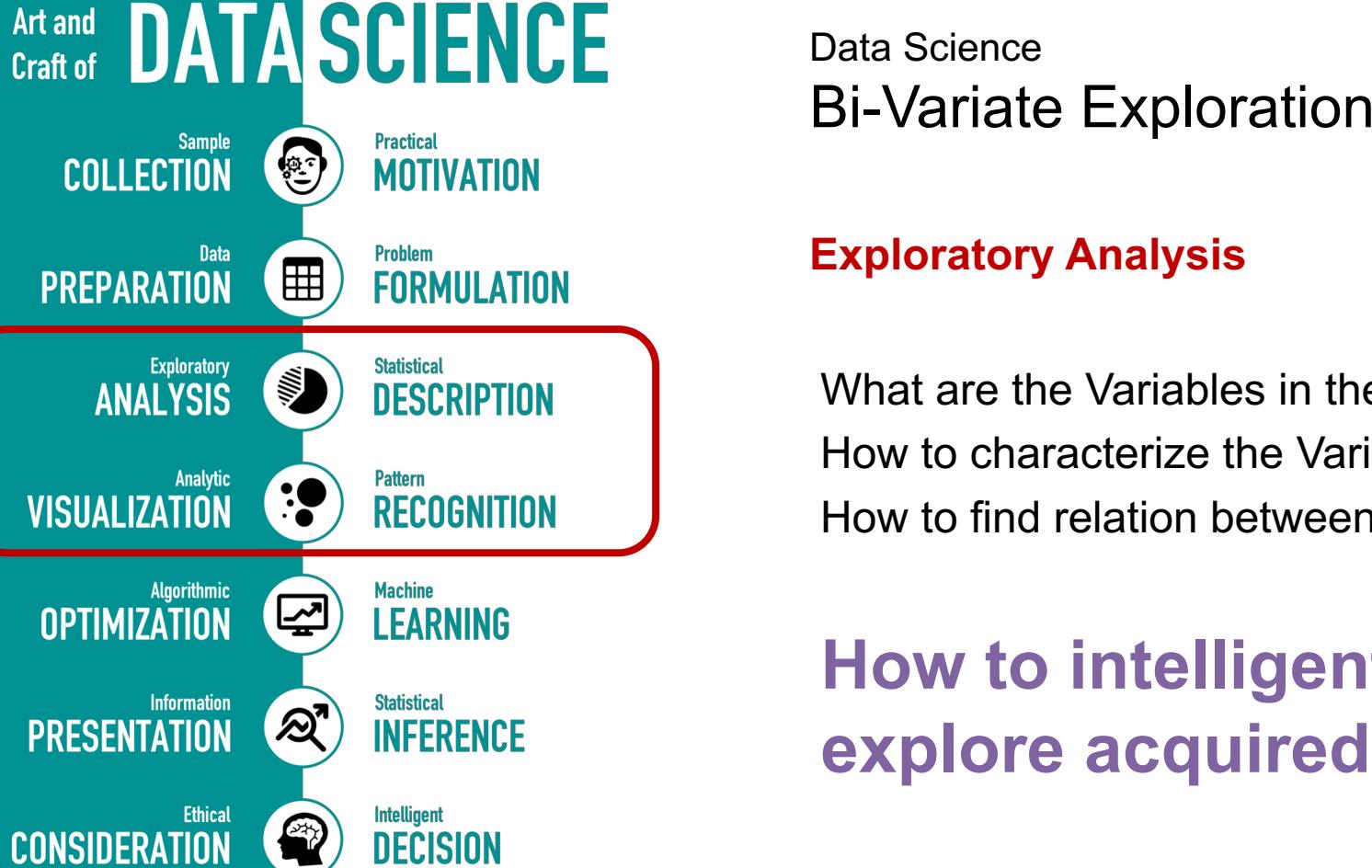
- How to summarize the acquired Data?
- How to visualize the acquired Data?
- How to analyze the acquired Data?

How to intelligently explore acquired Data?

Bi-Variate Exploration

Sourav SEN GUPTA
Lecturer, SCSE, NTU







Data Science

The Pokemon Dataset

#	Name	Type 1	Type 2	Total	HP	Attack	Defense	Sp. Atk	Sp. Def	Speed	Generation	Legendary
430	Honchkrow	Dark	Flying	505	100	125	52	105	52	71	4	False
338	Solrock	Rock	Psychic	440	70	95	85	55	65	70	3	False
32	Nidoran♂	Poison	NaN	273	46	57	40	40	40	50	1	False
442	Spiritomb	Ghost	Dark	485	50	92	108	92	108	35	4	False
480	Uxie	Psychic	NaN	580	75	75	130	75	130	95	4	True
536	Palpitoad	Water	Ground	384	75	65	55	65	55	69	5	False
360	Wynaut	Psychic	NaN	260	95	23	48	23	48	23	3	False
478	Froslass	Ice	Ghost	480	70	80	70	80	70	110	4	False
76	Golem	Rock	Ground	495	80	120	130	55	65	45	1	False
177	Natu	Psychic	Flying	320	40	50	45	70	45	70	2	False

Source : Kaggle Datasets | [Pokemon with stats](#) by Alberto Barradas | <https://www.kaggle.com/abcsds/pokemon>

Data Science

Bi-Variate Statistics

Statistical Summary

HP		Attack	
count	800.000000	count	800.000000
mean	69.258750	mean	79.001250
std	25.534669	std	32.457366
min	1.000000	min	5.000000
25%	50.000000	25%	55.000000
50%	65.000000	50%	75.000000
75%	80.000000	75%	100.000000
max	255.000000	max	190.000000

HP Hit Points of a Pokemon
Attack Base Modifier for Attack

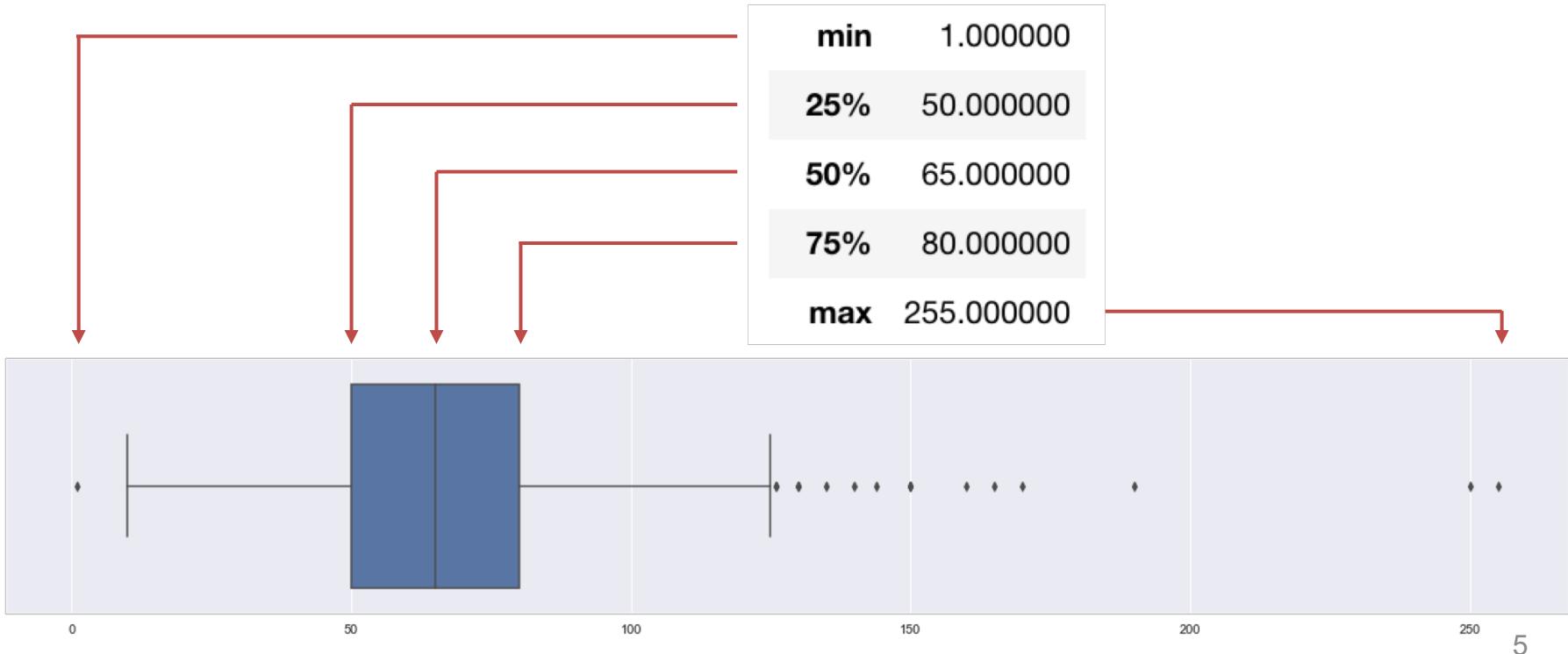
Statistical Questions

- What is the Central Tendency?
- What is the Spread of the Data?

HP

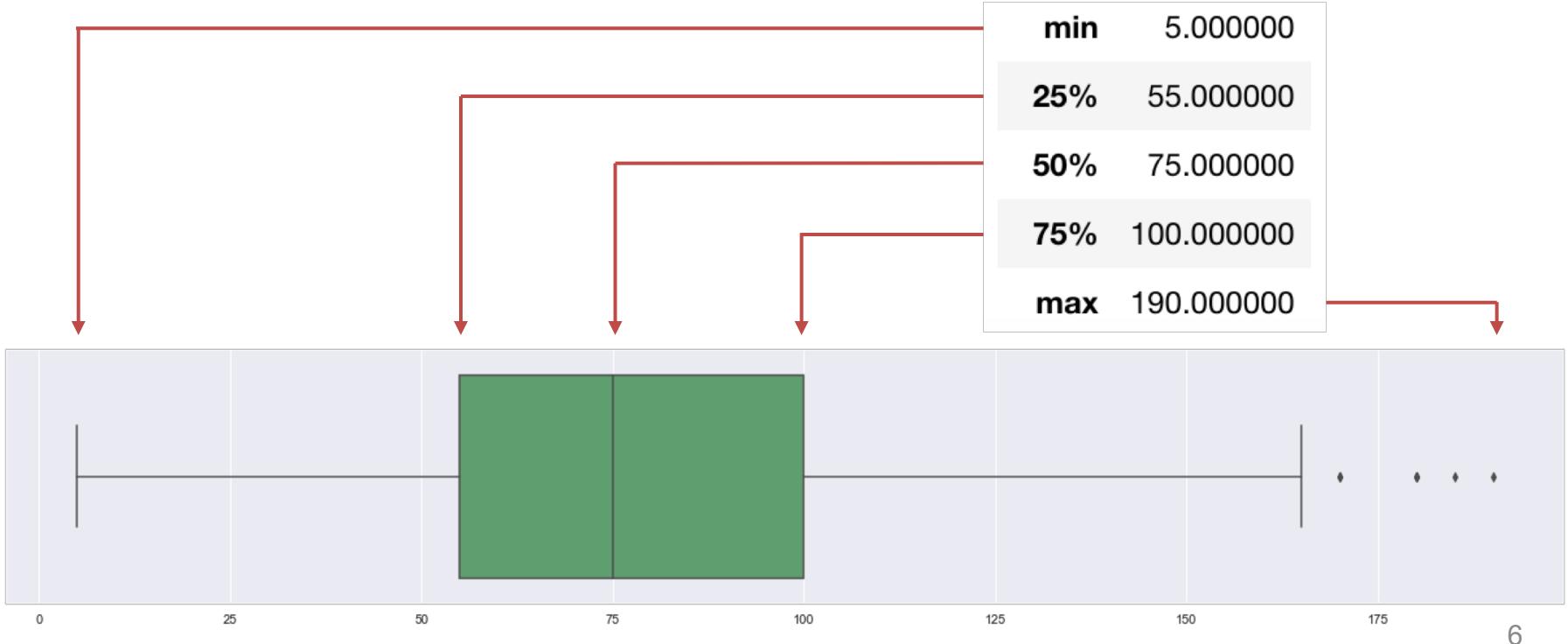
Data Science

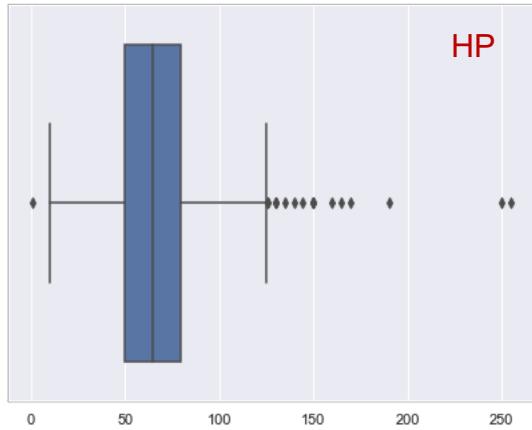
Uni-Variate Box-Plot



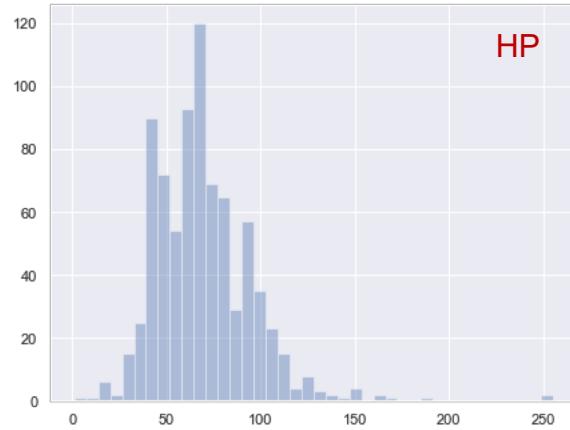
Attack

Data Science Uni-Variate Box-Plot

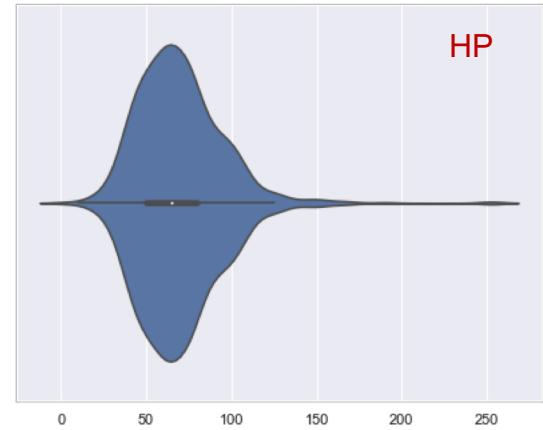




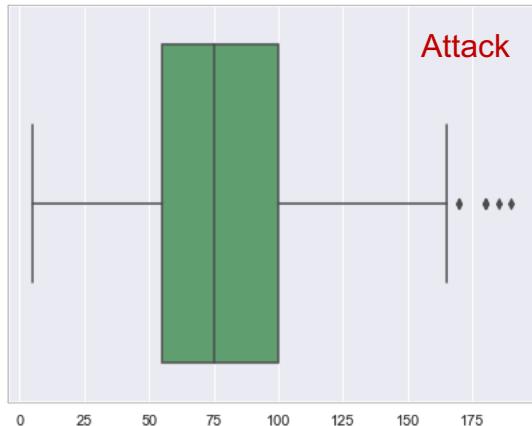
HP



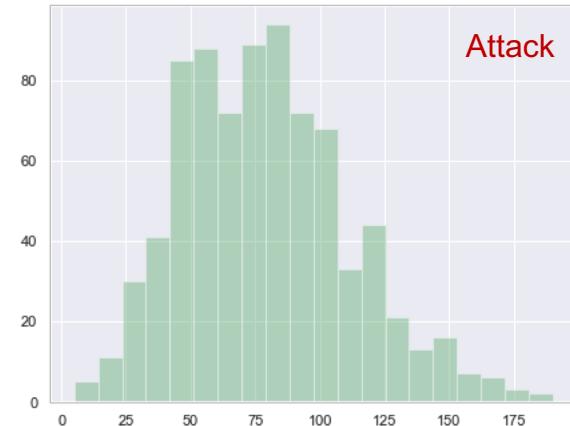
HP



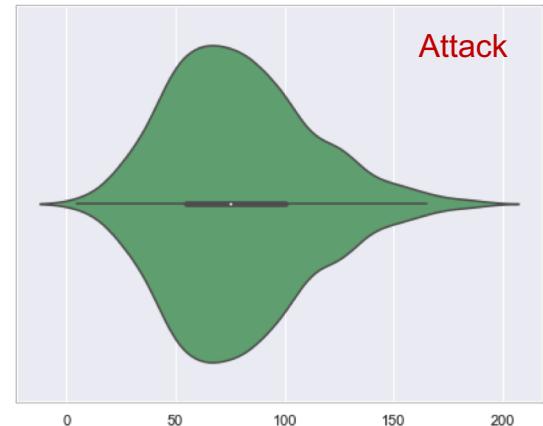
HP



Attack

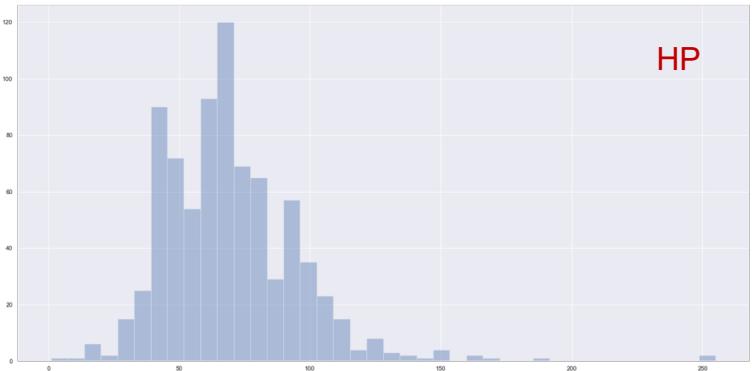


Attack



Attack



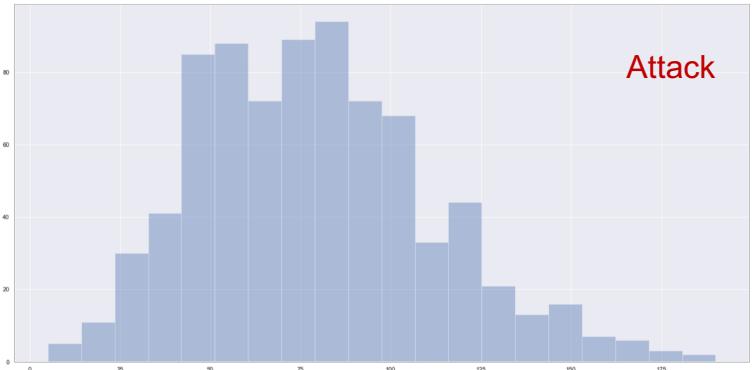


Data Science

Bi-Variate Statistics

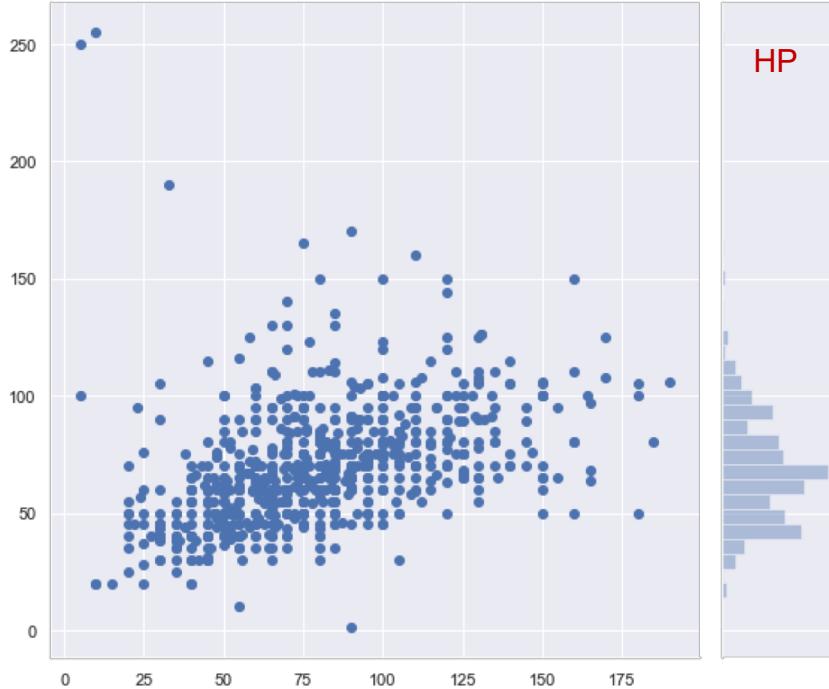
Statistical Relation

HP Hit Points of a Pokemon
Attack Base Modifier for Attack



Statistical Questions

- Is there a Mutual Dependence?
- What is the Mutual Relationship?



Data Science

Bi-Variate Joint Plot

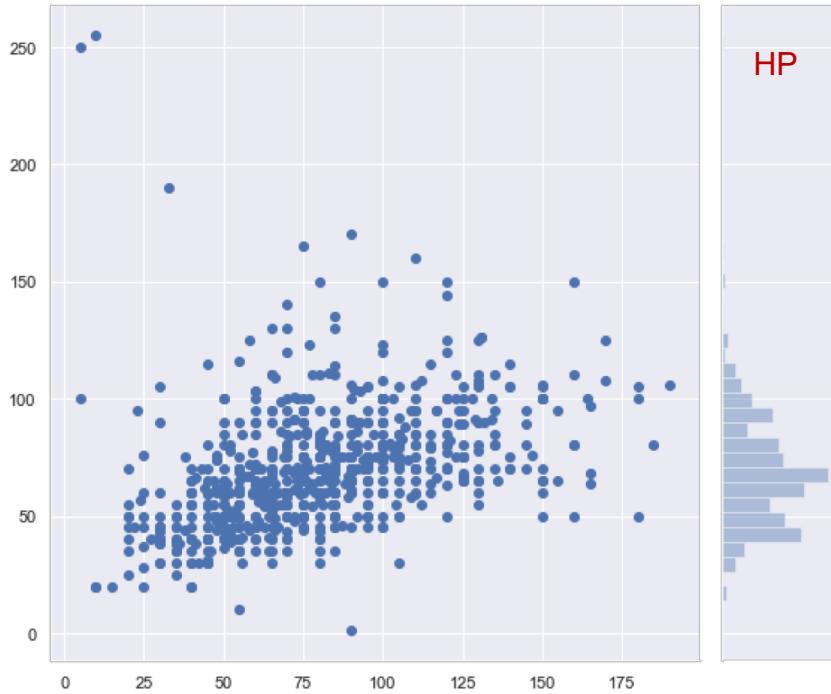
Statistical Relation

HP Plotted along Y axis

Attack Plotted along X axis

Pattern Recognition

- Is there a Mutual Dependence?
- What is the Mutual Relationship?



Data Science

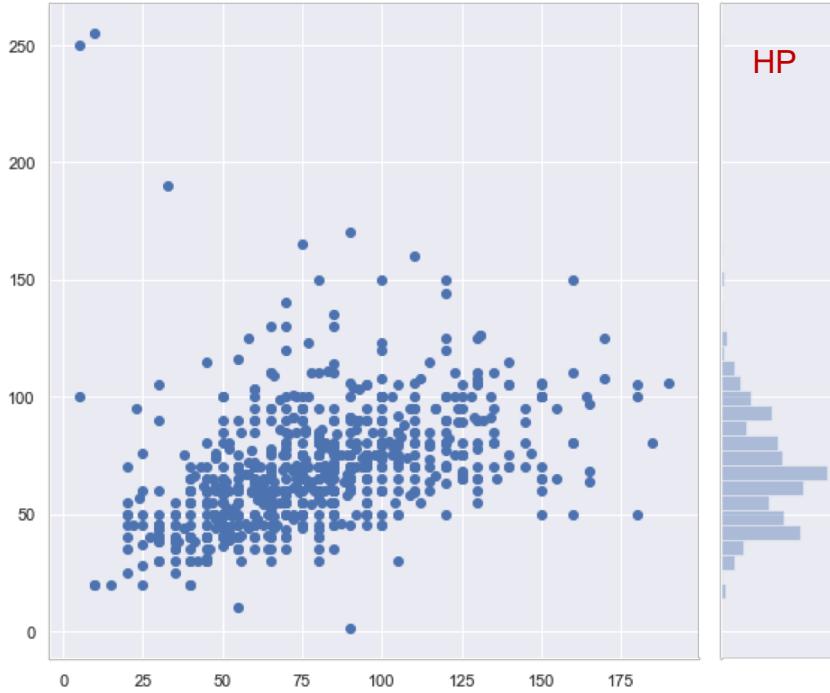
Bi-Variate Joint Plot

Statistical Relation

HP Plotted along Y axis
Attack Plotted along X axis

Pattern Recognition

- HP increases as Attack increases
- Dependence is moderately strong



Data Science Bi-Variate Relation

Correlation Coefficient

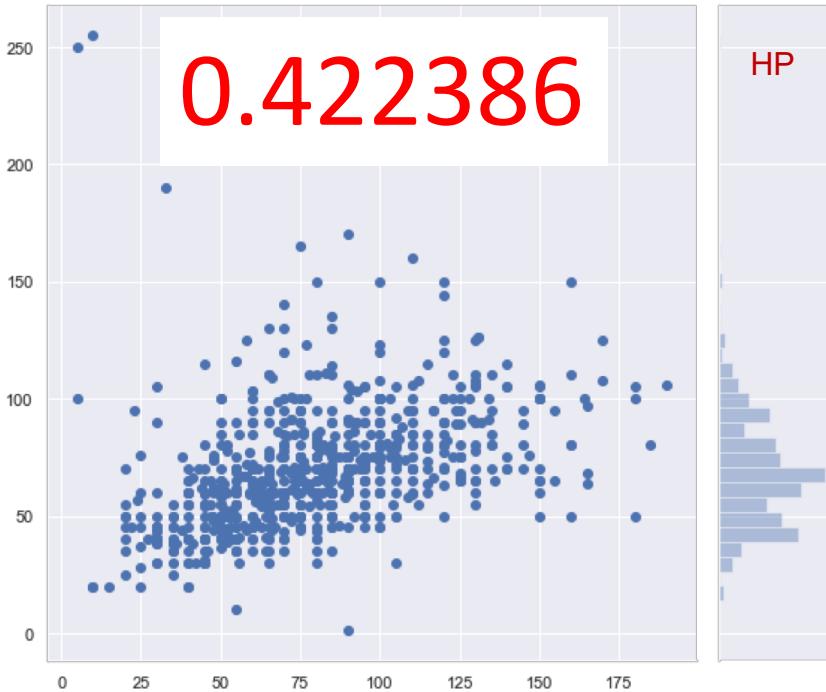
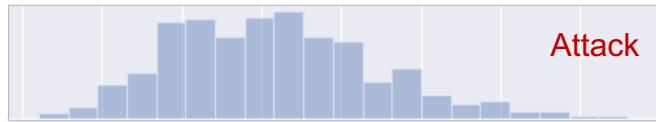
Natural Intuition

Dependence of HP and Attack

Statistical Formula

Co-Variance / St. Dev Product

$$\rho_{xy} = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum(x_i - \bar{x})^2} \sqrt{\sum(y_i - \bar{y})^2}}$$



Data Science Bi-Variate Relation

Correlation Coefficient

Natural Intuition

Dependence of HP and Attack

Statistical Intuition

No Dependence

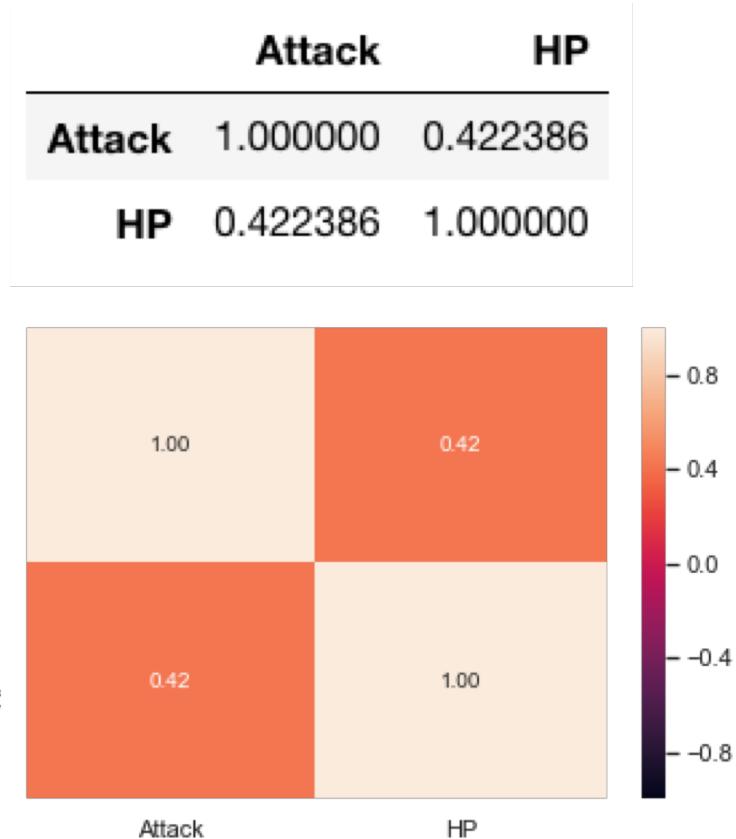
Corr = 0

Perfect Positive

Corr = + 1

Perfect Negative

Corr = - 1



Data Science Bi-Variate Relation

Correlation Matrix and Plot

Natural Intuition

Dependence of HP and Attack

Statistical Intuition

No Dependence

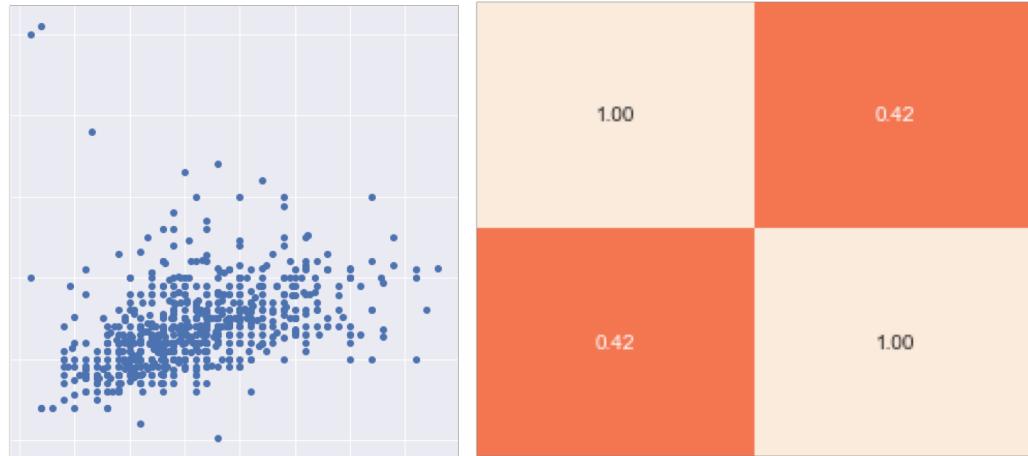
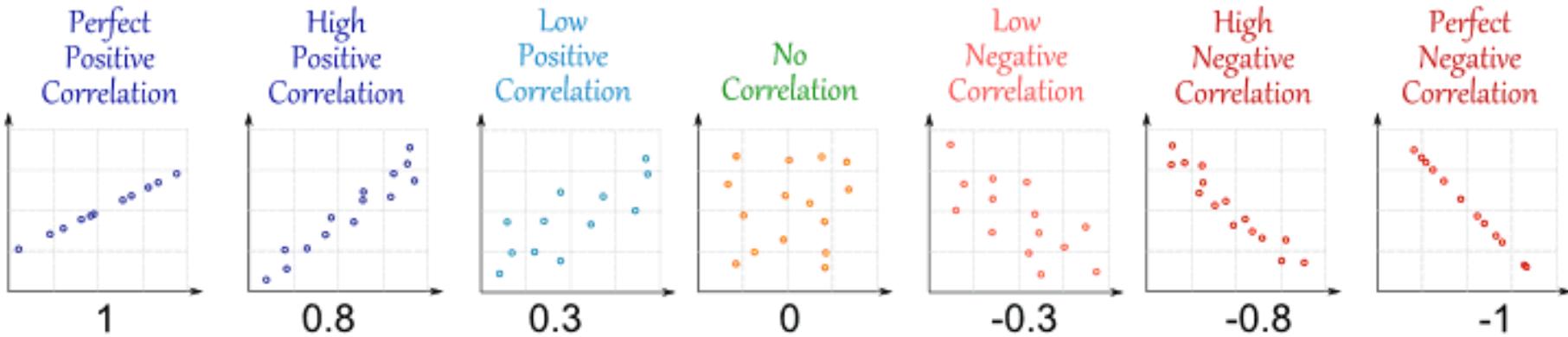
Corr = 0

Perfect Positive

Corr = + 1

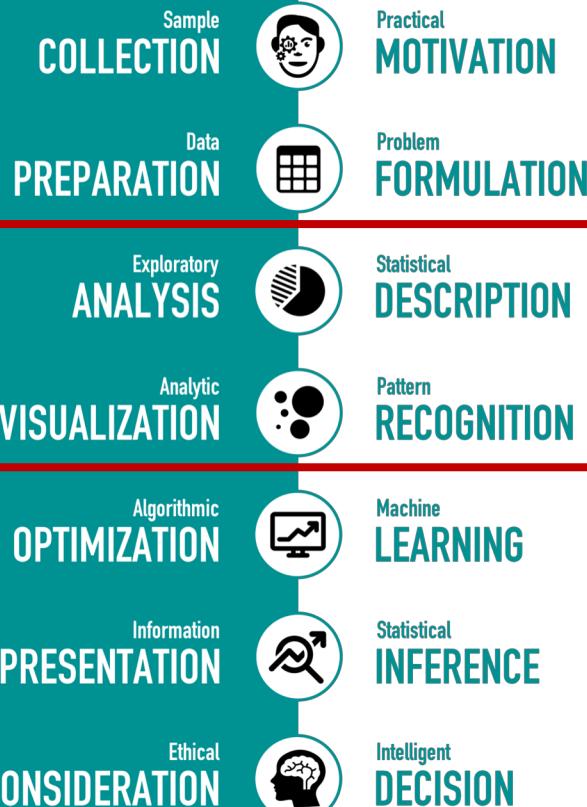
Perfect Negative

Corr = - 1



Play this to improve your
Correlation spotting skills

Guess the Correlation
<http://guessthecorrelation.com/>



Data Science Pipeline Exploratory Analysis

How to summarize the acquired Data?
How to visualize the acquired Data?
How to analyze the acquired Data?

How to intelligently
explore acquired Data?

Multi-Variate Exploration

Sourav SEN GUPTA
Lecturer, SCSE, NTU



Sample
COLLECTION



Practical
MOTIVATION

Data
PREPARATION



Problem
FORMULATION

Exploratory
ANALYSIS



Statistical
DESCRIPTION

Analytic
VISUALIZATION



Pattern
RECOGNITION

Algorithmic
OPTIMIZATION



Machine
LEARNING

Information
PRESSENTATION



Statistical
INFERENCE

Ethical
CONSIDERATION



Intelligent
DECISION

Data Science Multi-Variate Exploration

Exploratory Analysis

What are the Variables in the Data?
How to characterize the Variables?
How to find relation between them?

How to intelligently explore acquired Data?



Data Science

The Pokemon Dataset

#	Name	Type 1	Type 2	Total	HP	Attack	Defense	Sp. Atk	Sp. Def	Speed	Generation	Legendary
430	Honchkrow	Dark	Flying	505	100	125	52	105	52	71	4	False
338	Solrock	Rock	Psychic	440	70	95	85	55	65	70	3	False
32	Nidoran♂	Poison	NaN	273	46	57	40	40	40	50	1	False
442	Spiritomb	Ghost	Dark	485	50	92	108	92	108	35	4	False
480	Uxie	Psychic	NaN	580	75	75	130	75	130	95	4	True
536	Palpitoad	Water	Ground	384	75	65	55	65	55	69	5	False
360	Wynaut	Psychic	NaN	260	95	23	48	23	48	23	3	False
478	Froslass	Ice	Ghost	480	70	80	70	80	70	110	4	False
76	Golem	Rock	Ground	495	80	120	130	55	65	45	1	False
177	Natu	Psychic	Flying	320	40	50	45	70	45	70	2	False

Source : Kaggle Datasets | **Pokemon with stats** by Alberto Barradas | <https://www.kaggle.com/abcsds/pokemon>

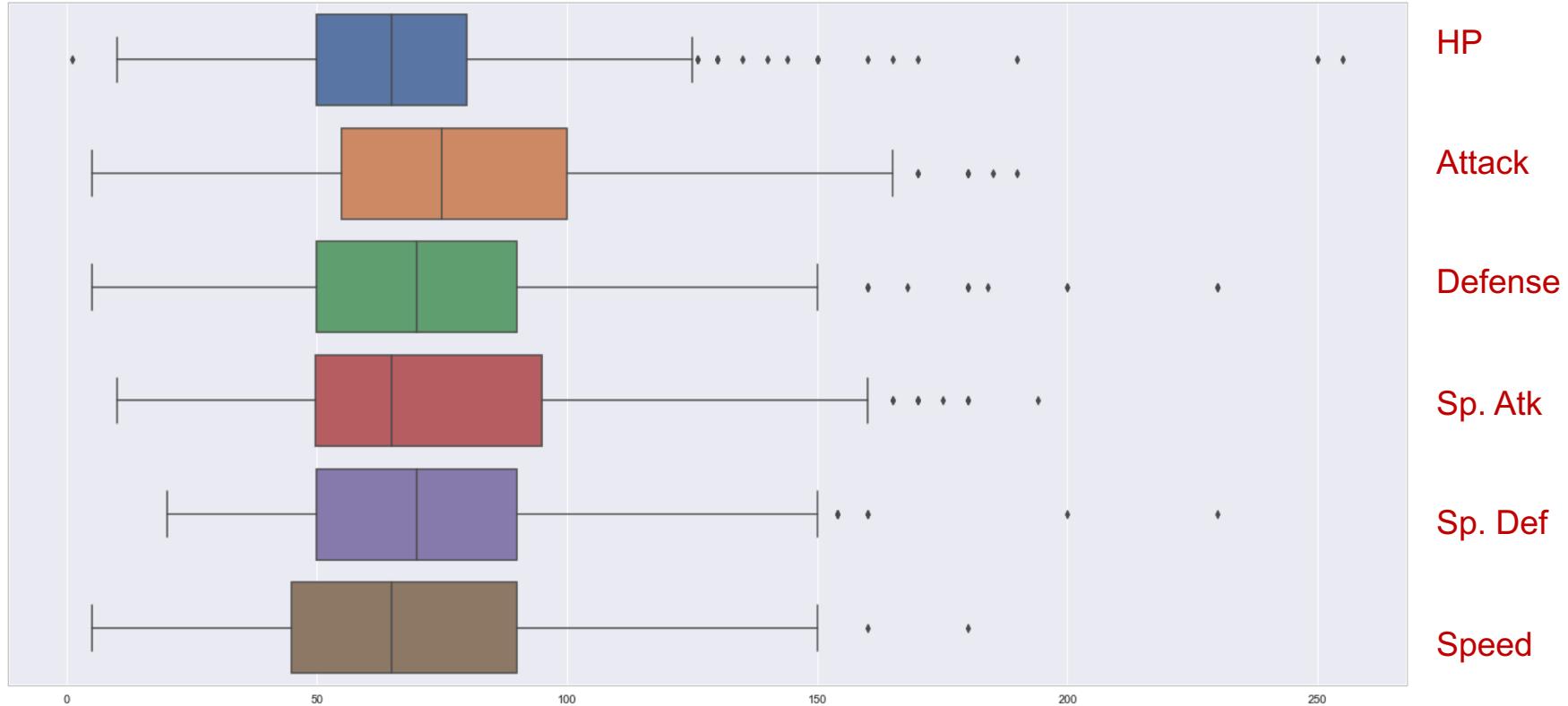
Individual Statistical Summary

Data Science Multi-Variate Statistics

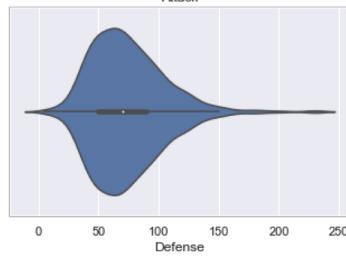
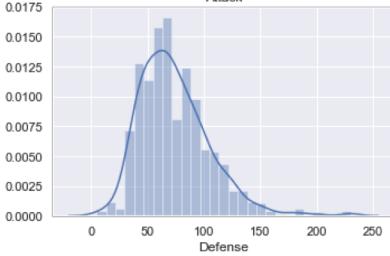
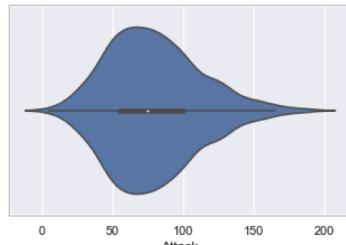
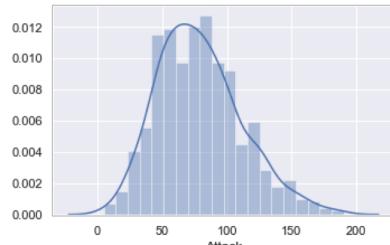
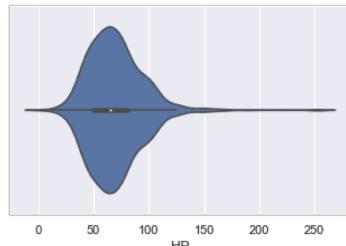
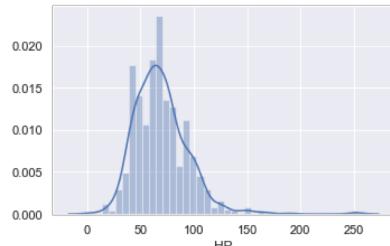
	HP	Attack	Defense	Sp. Atk	Sp. Def	Speed
count	800.000000	800.000000	800.000000	800.000000	800.000000	800.000000
mean	69.258750	79.001250	73.842500	72.820000	71.902500	68.277500
std	25.534669	32.457366	31.183501	32.722294	27.828916	29.060474
min	1.000000	5.000000	5.000000	10.000000	20.000000	5.000000
25%	50.000000	55.000000	50.000000	49.750000	50.000000	45.000000
50%	65.000000	75.000000	70.000000	65.000000	70.000000	65.000000
75%	80.000000	100.000000	90.000000	95.000000	90.000000	90.000000
max	255.000000	190.000000	230.000000	194.000000	230.000000	180.000000

Individual Box-Plots

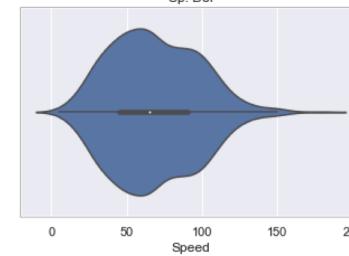
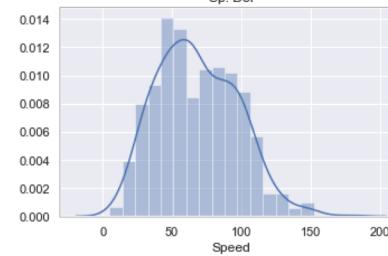
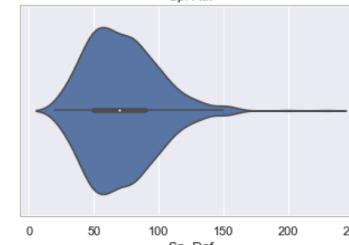
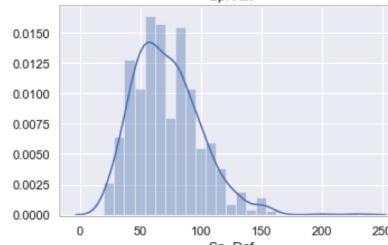
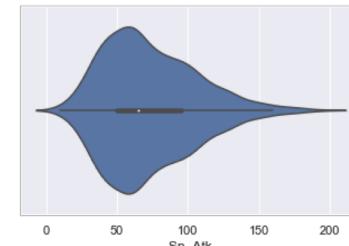
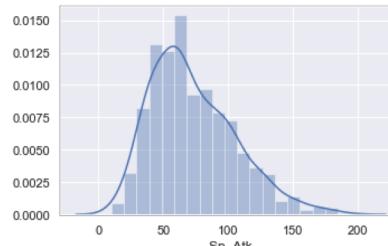
Data Science Multi-Variate Statistics



Individual Histograms and Violin Plots



Data Science Multi-Variate Statistics



Mutual Correlations

Data Science Multi-Variate Statistics

	HP	Attack	Defense	Sp. Atk	Sp. Def	Speed
HP	1.000000	0.422386	0.239622	0.362380	0.378718	0.175952
Attack	0.422386	1.000000	0.438687	0.396362	0.263990	0.381240
Defense	0.239622	0.438687	1.000000	0.223549	0.510747	0.015227
Sp. Atk	0.362380	0.396362	0.223549	1.000000	0.506121	0.473018
Sp. Def	0.378718	0.263990	0.510747	0.506121	1.000000	0.259133
Speed	0.175952	0.381240	0.015227	0.473018	0.259133	1.000000



Mutual Correlations

No Dependence

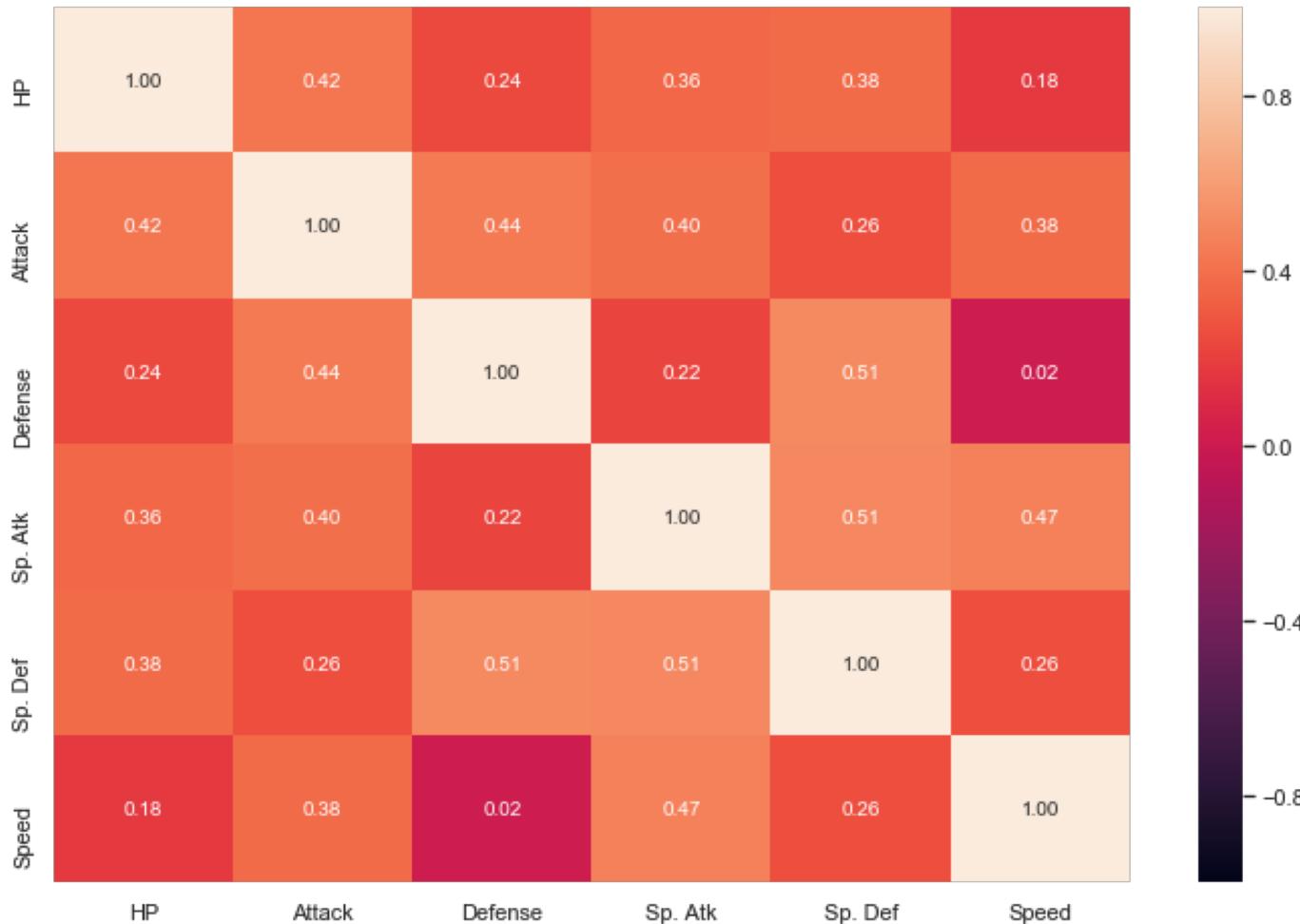
Corr = 0

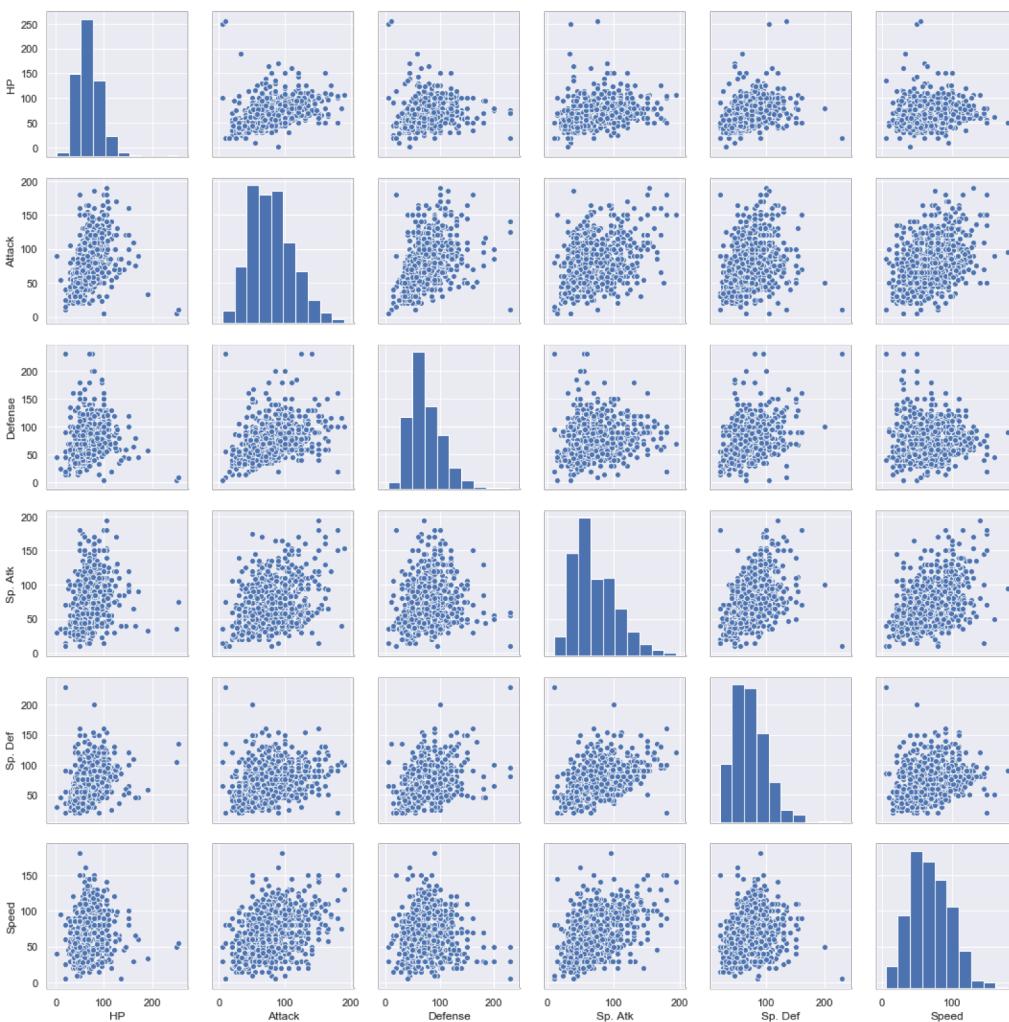
Perfect Positive

Corr = + 1

Perfect Negative

Corr = - 1





Data Science Multi-Variate Dist.

Pair-Plot of Multi-Variate Data

Histograms in the Diagonals
Scatter-Plot in Off-Diagonals

- Pattern recognition
- Distributions of Variables
 - Inter-Variable Dependence

Sample
COLLECTION



Practical
MOTIVATION

Data
PREPARATION



Problem
FORMULATION

Exploratory
ANALYSIS



Statistical
DESCRIPTION

Analytic
VISUALIZATION



Pattern
RECOGNITION

Algorithmic
OPTIMIZATION



Machine
LEARNING

Information
PRESENTATION



Statistical
INFERENCE

Ethical
CONSIDERATION



Intelligent
DECISION

Data Science Pipeline **Exploratory Analysis**

- How to summarize the acquired Data?
- How to visualize the acquired Data?
- How to analyze the acquired Data?

How to intelligently explore acquired Data?

The Normal Distribution

Sourav SEN GUPTA
Lecturer, SCSE, NTU





Data Science

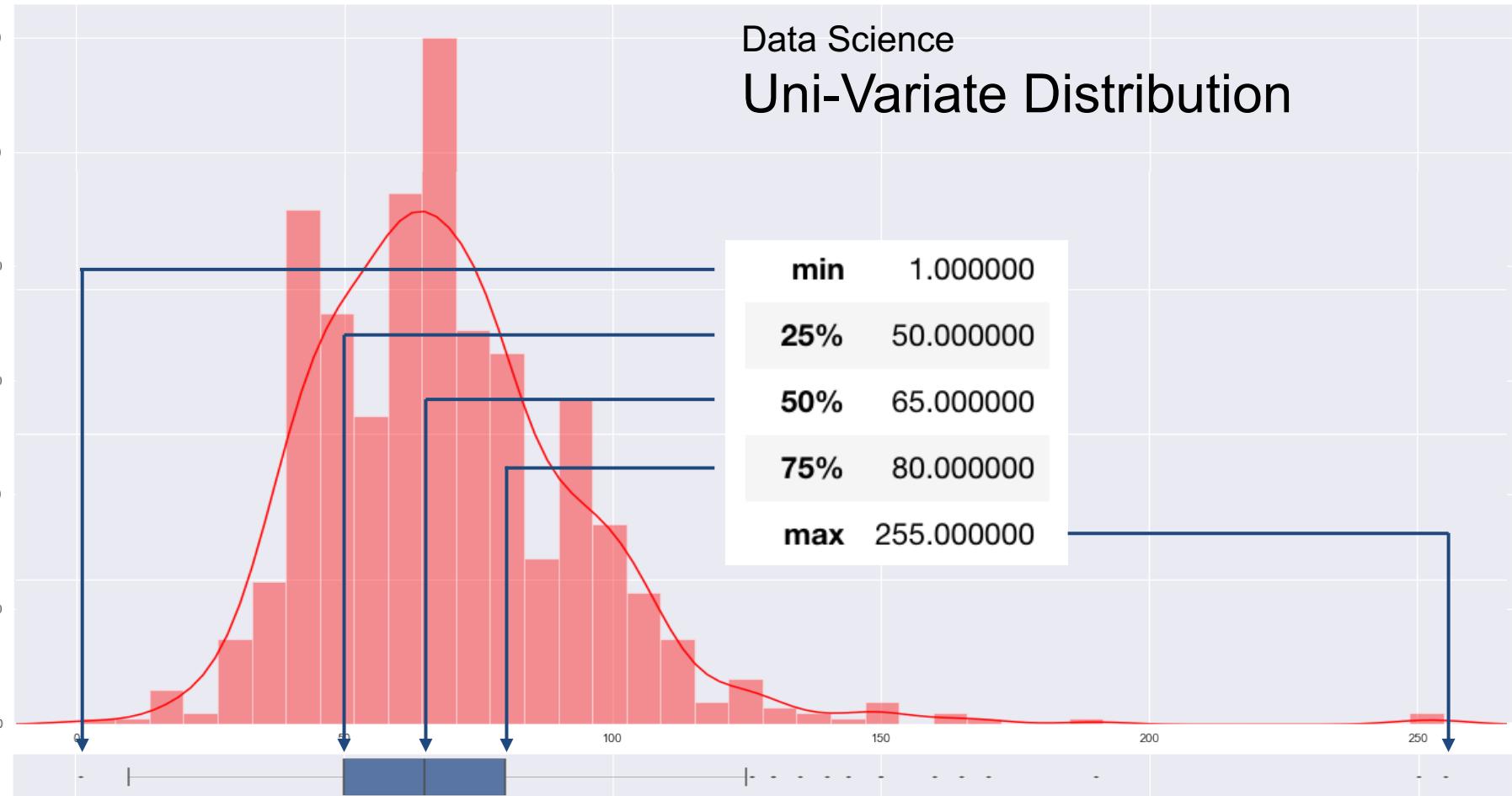
The Pokemon Dataset

#	Name	Type 1	Type 2	Total	HP	Attack	Defense	Sp. Atk	Sp. Def	Speed	Generation	Legendary
430	Honchkrow	Dark	Flying	505	100	125	52	105	52	71	4	False
338	Solrock	Rock	Psychic	440	70	95	85	55	65	70	3	False
32	Nidoran♂	Poison	NaN	273	46	57	40	40	40	50	1	False
442	Spiritomb	Ghost	Dark	485	50	92	108	92	108	35	4	False
480	Uxie	Psychic	NaN	580	75	75	130	75	130	95	4	True
536	Palpitoad	Water	Ground	384	75	65	55	65	55	69	5	False
360	Wynaut	Psychic	NaN	260	95	23	48	23	48	23	3	False
478	Froslass	Ice	Ghost	480	70	80	70	80	70	110	4	False
76	Golem	Rock	Ground	495	80	120	130	55	65	45	1	False
177	Natu	Psychic	Flying	320	40	50	45	70	45	70	2	False

Source : Kaggle Datasets | **Pokemon with stats** by Alberto Barradas | <https://www.kaggle.com/abcsds/pokemon>

Data Science

Uni-Variate Distribution



Data Science

Uni-Variate Distribution

mean 69.258750

std 25.534669

