# Discussion

## Comparison of Data Analysis Tools: Excel, R, and Python
### (Yujeong Kim, Ph.D)

**Shin Wha Lee**
**Asan Medical Center, University of Ulsan**

# DECLARATION OF INTERESTS
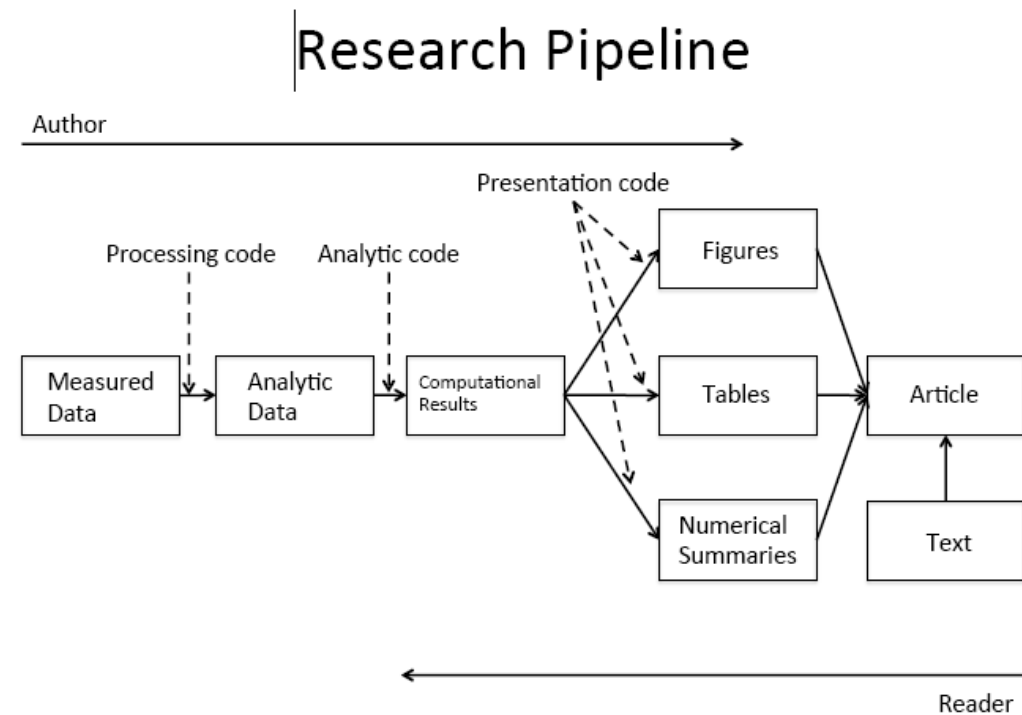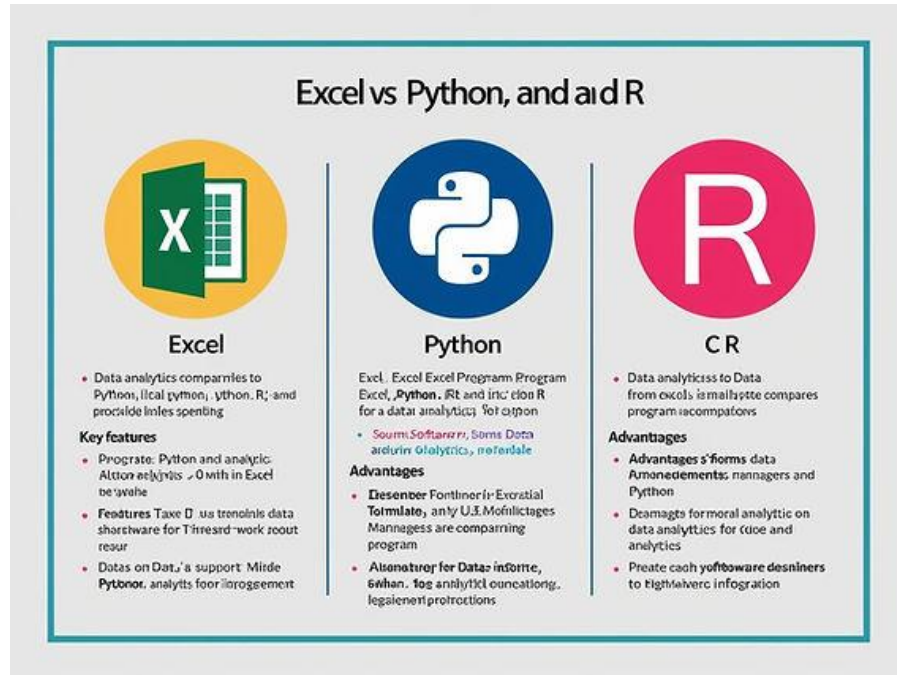
**Nothing to declare**

# Q1. 도구 선택 오류는 언제 발생하는가?



- **Excel–R–Python을 역할로 구분했을 때, 임상 연구에서 가장 흔한 오류는 어느 분석 단계에서 발생합니까?**

Reference: Data analysis is an iterative process from question formulation to interpretation. *Shearer C. The CRISP-DM model. J Data Warehousing, 2000.*

대한부인종양학회
Korean Society of Gynecologic Oncology
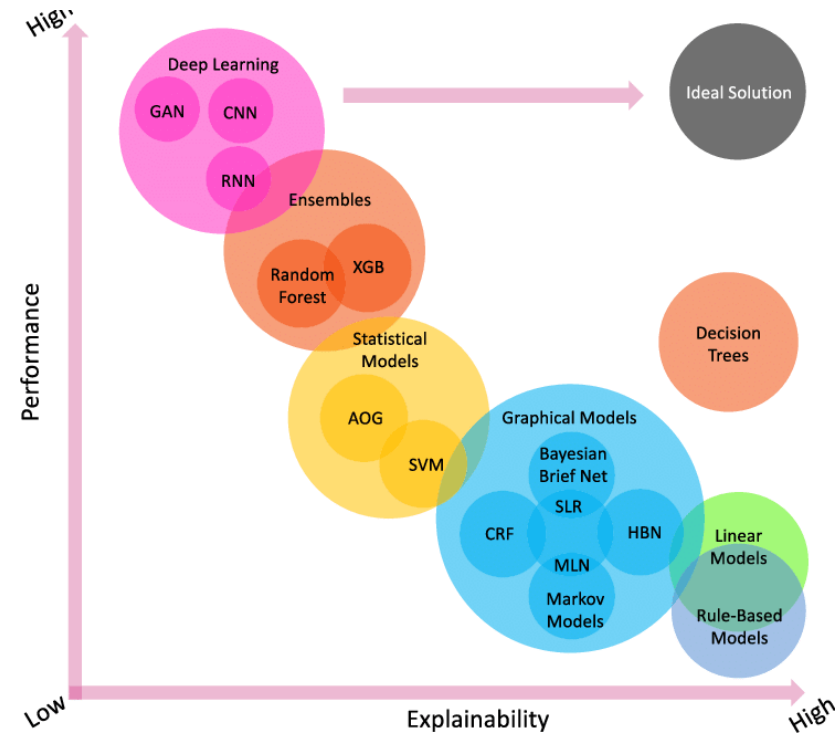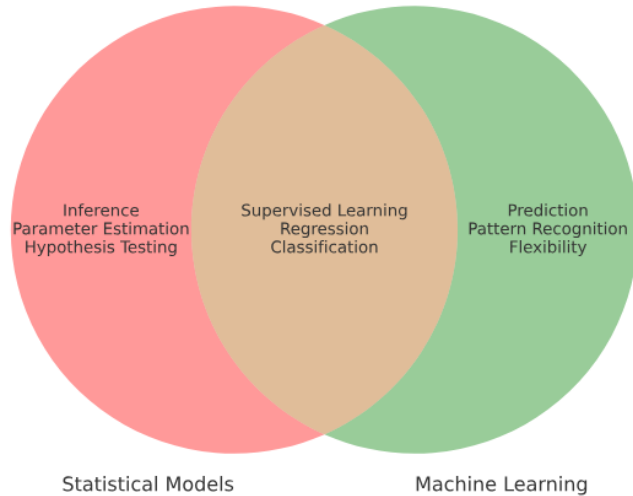
# Q2. Excel에서 R/Python으로 넘어가는 기준





- 언제부터 Excel을 넘어서 R 또는 Python으로 전환해야 할까요?

Reference: Reproducibility and transparency are improved with code-based analyses. *Peng RD. Reproducible research. Science, 2011.*.

대한부인종양학회
Korean Society of Gynecologic Oncology

# Q3. 통계 검정 vs 예측 모델



Conceptual Overlap Between Statistical Models and Machine Learning

- 통계 검정 중심 연구와 예측 모델 연구는 어떤 연구 질문에서 명확히 구분된다고 보십니까?

Reference: Statistical modeling and machine learning differ in goals: inference vs prediction. *Shmueli G. Stat Sci, 2010.*

대한부인종양학회
Korean Society of Gynecologic Oncology

# Q4. Tidy data의 최소 기준



TIDY DATA is a standard way of mapping the meaning of a dataset to its structure.
—HADLEY WICKHAM

In tidy data:
• each variable forms a column
• each observation forms a row
• each cell is a single measurement

each column a variable

| id | name | color |
|----|-------|--------|
| 1 | floof | gray |
| 2 | max | black |
| 3 | cat | orange |
| 4 | donut | gray |
| 5 | merlin | black |
| 6 | panda | calico |

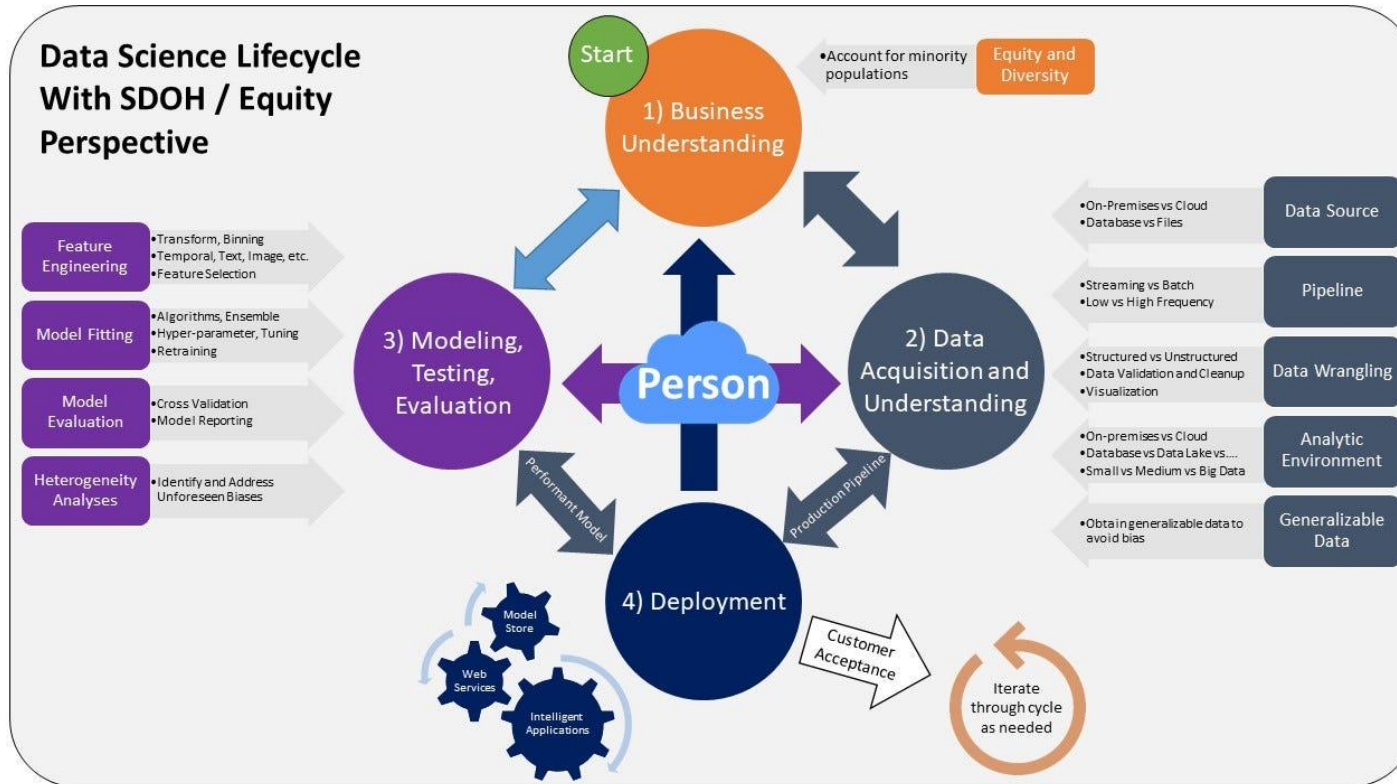each row an observation

variables | observations | values

Wickham, H. (2014). Tidy Data. Journal of Statistical Software 59 (10). DOI: 10.18637/jss.v059.i10

• 실제 임상 데이터에서 최소한 지켜야 할 Tidy data 기준은 무엇일까요?

Reference: Tidy data provides a standard way to organize datasets. *Wickham H. J Stat Softw, 2014.*

대한부인종양학회
Korean Society of Gynecologic Oncology

# Q5. Garbage In, Garbage Out



- 임상 연구에서 결과는 그럴듯하지만 위험한 대표적인 GIGO 사례에는 무엇이 있을까요?

Reference: Model performance is fundamentally limited by data quality. *Kuhn M, Johnson K. Applied Predictive Modeling, 2013.*

대한부인종양학회
Korean Society of Gynecologic Oncology