# CASOS

# A Cautionary Note about Saving Unicode Text from Excel

Neal Altman

na@cmu.edu

**Carnegie Mellon**

Center for Computational Analysis of
Social and Organizational Systems
http://www.casos.cs.cmu.edu/

---

**Carnegie Mellon**
isr institute for SOFTWARE RESEARCH
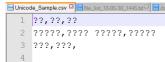
# Microsoft Excel's CSV Save

- Non-ASCII characters may be lost when exporting in CSV (comma separated value) format:

| | A | B | C |
|---|---|---|---|
| 1 | 漢語 | 汉语 | 中文 |
| 2 | | زبان فارسی فارسی | پارسی |
| 3 | 한국어 | 韓國語 | |

Unicode_Sample.xlsx

Save as: CSV

Unicode_Sample.csv

```
1  ??,??,??
2  ?????,???? ?????,?????
3  ???,???,
4
```

- Why?  The content was saved as "Western European (Windows)" (aka Windows-1252) a limited extension to the standard ASCII character set.

- Good news! The CSV UTF-8 option was announced for MS Office build 16.0.7466.2023 and after.
  (But only if you pay for an upgrade.)

CASOS

June 2020

CASOS