

Lecture 1

Implications of Big Data for ML

Lecture I – Big Data in Machine Learning

Marina Meilă

Department of Statistics
University of Washington

STAT 548/CSE 547
Spring, 2025

Big Data

How Big?

- Lobster 100's TB 7000\$ $n \times D$
~10 GB (free)
- cell painting ~100's TB
- real time 50 TB
- social media
- brain phys data + connectome 14 TB
- 10000 in financial
- Youtube ~100's hours
- health care claims ~10⁹
- 10⁵ particle trajectories - in brain
- $n \Rightarrow$ ◦ 10⁵ London marathon data

How used?

- pattern recognition → find pattern
- discover dependencies classification
- relationship detection
- find outliers ~~for~~ completion
- find neighbours testing
- similar points to regression
some x
(retrieval)

High dimensional

- images 4K x 4K = D
1 pixel = 3 channels^u
x 256 levels
- medical, bio, astro
- econ
- financial
- biological sequences
10 - 10⁶
- music, audio
- video
- atmospheric
- sports data

Time Scales?

Repeated?

Big data and Machine Learning I

Big Data has implications for ML at many levels

- Storage

- may not fit in local memory
- expensive/slow to move around ←
- I/O expensive/slow → *output large*

- Access

- serial/by block, not random

- Indexing

- Preprocessing steps that allow faster access during

← *pointers to data*

- Computing

- Parallelization when possible $p \leq P$
- Automation of resource management (Hadoop, Spark)

- Algorithms

- predominantly sub-quadratic, i.e. $\mathcal{O}(n)$, $\tilde{\mathcal{O}}(n)$
- sub-linear, i.e. $o(n)$ when possible – sampling, Stochastic Gradient Descent (SGD)



