# S&DS677: Topics in High-Dimensional Statistics and Information Theory

Spring 2021

- Schedule: Tuesday 330–520pm on zoom
- Instructor: Yihong Wu yihong.wu@yale.edu
  - Office hours: by appointment
- Website:

http://www.stat.yale.edu/~yw562/teaching/SDS677/index.html
or just google S&DS677

1 Course prerequisites:

1 Course prerequisites:

Maturity with probability theory

1 Course prerequisites:

- Maturity with probability theory
- Some linear algebra

1 Course prerequisites:

- Maturity with probability theory
- Some linear algebra
- Prior knowledge on Information Theory (e.g. SDS 364) is NOT required

Course prerequisites:

- Maturity with probability theory
- Some linear algebra
- Prior knowledge on Information Theory (e.g. SDS 364) is NOT required

2 Participation (30%):

Course prerequisites:

- Maturity with probability theory
- Some linear algebra
- Prior knowledge on Information Theory (e.g. SDS 364) is NOT required

Participation (30%):

Zoom participation is highly encouraged

Course prerequisites:

- Maturity with probability theory
- Some linear algebra
- Prior knowledge on Information Theory (e.g. SDS 364) is NOT required

2 Participation (30%):

- Zoom participation is highly encouraged
- Critiques on lecture notes/maybe a few scribes towards the end

Course prerequisites:

- Maturity with probability theory
- Some linear algebra
- Prior knowledge on Information Theory (e.g. SDS 364) is NOT required
- Participation (30%):
  - Zoom participation is highly encouraged
  - Critiques on lecture notes/maybe a few scribes towards the end
- 3 Homeworks (30%): three problem sets

- Course prerequisites:
  - Maturity with probability theory
  - Some linear algebra
  - Prior knowledge on Information Theory (e.g. SDS 364) is NOT required
- Participation (30%):
  - Zoom participation is highly encouraged
  - Critiques on lecture notes/maybe a few scribes towards the end
- 3 Homeworks (30%): three problem sets
- 4 Final project (40%)

- Course prerequisites:
  - Maturity with probability theory
  - Some linear algebra
  - Prior knowledge on Information Theory (e.g. SDS 364) is NOT required
- Participation (30%):
  - Zoom participation is highly encouraged
  - Critiques on lecture notes/maybe a few scribes towards the end
- 3 Homeworks (30%): three problem sets
- 4 Final project (40%)
  - either presenting paper(s) or a standalone research project.

- Course prerequisites:
  - Maturity with probability theory
  - Some linear algebra
  - Prior knowledge on Information Theory (e.g. SDS 364) is NOT required
- Participation (30%):
  - Zoom participation is highly encouraged
  - Critiques on lecture notes/maybe a few scribes towards the end
- 3 Homeworks (30%): three problem sets
- 4 Final project (40%)
  - either presenting paper(s) or a standalone research project.
  - topics announced around week 6

- Course prerequisites:
  - Maturity with probability theory
  - Some linear algebra
  - Prior knowledge on Information Theory (e.g. SDS 364) is NOT required
- Participation (30%):
  - Zoom participation is highly encouraged
  - Critiques on lecture notes/maybe a few scribes towards the end
- 3 Homeworks (30%): three problem sets
- 4 Final project (40%)
  - either presenting paper(s) or a standalone research project.
  - topics announced around week 6
- 6 Materials: Lecture notes and additional reading materials will be posted online.

What this course is about?

Information-theoretic methods in high-dimensional statistics

What this course is about?

Information-theoretic & related methods in high-dimensional statistics

#### What this course is about?

Information-theoretic & related methods in high-dimensional statistics

• Statistical tasks: using data to make informed decisions (hypotheses testing, estimation, confidence statements)



• Statistical tasks: using data to make informed decisions (hypotheses testing, estimation, confidence statements)



• Understanding the fundamental limits:

• Statistical tasks: using data to make informed decisions (hypotheses testing, estimation, confidence statements)



• Understanding the fundamental limits:

Oharacterize statistical optimum: What is possible/impossible?

• Statistical tasks: using data to make informed decisions (hypotheses testing, estimation, confidence statements)



- Understanding the fundamental limits:
  - Of the provide the provided and the p
  - We have many samples are necessary and sufficient to achieve a prescribed goal?

• Statistical tasks: using data to make informed decisions (hypotheses testing, estimation, confidence statements)



- Understanding the fundamental limits:
  - Otheracterize statistical optimum: What is possible/impossible?
  - We have many samples are necessary and sufficient to achieve a prescribed goal?
  - Solution Can statistical limits be attained comptutationally efficiently, e.g., in poly(n, p)-time? If yes, how? If not, why?

## High Dimensionality of Contemporary Datasets

Fields	Data
<b>Biomedical Research</b>	microarray, ECG, fMRI,
	array sensor data,
Signal Processing	face recognition,
	hyper-spectral data,
Finance	asset returns,
:	:
:	

- Growth of data outpaced by increasing number of features
- A common feature: large d, but just comparable or smaller n

$$\theta \in \mathbb{R}^d \mapsto X_1, \dots, X_n$$

- low-dimensional structure
  - Intrinsic:  $\theta$  lies in a low-dimensional subset
  - Extrinsic:  $\theta$  has no structure but we only estimate low-dimensional functional of  $\theta$

Classical topics

## Example 1: high-dimensional linear regression

Microarray data:

- Leukaemia dataset [Golub et al. '99]: d = 7129 genes and n = 72samples
- Typically  $d \gg n$
- Interpretability (gene selection)



Ref: [Golub et al. '99, Zou-Hastie '05]

#### Example 1: high-dimensional linear regression

Statistical model

$$y = X\beta + \mathsf{noise}$$

- observation:  $y \in \mathbb{R}^n$  and  $X \in \mathbb{R}^{n \times d}$
- parameter:  $\beta \in \mathbb{R}^d$
- goal: estimate  $\beta$  or predict  $X\beta$
- assumption:  $\beta$  is sparse

#### Example 2: Covariance matrix estimation & PCA Climate Data



One observation: January average temperature in 1969 [d = 2592, n = 157]

Ref: Bickel & Levina (08)

#### Example 2: Covariance matrix estimation & PCA

Statistical model

- observation:  $X_1, \ldots, X_n \stackrel{\text{iid}}{\sim} N(0, \Sigma) \in \mathbb{R}^d$
- parameter:  $\Sigma = \mathbb{E}[XX'] \in \mathbb{R}^{d \times d}$
- goal: estimate  $\Sigma$  or its principle component (PCA)
- assumption:  $\Sigma$  is sparse/smooth(entrywise decay)/low-rank

#### Problems of combinatorial nature

Linguistics

## Estimating the number of unseen species: How many words did Shakespeare know?

Bx BRADLEY EFRON AND RONALD THISTED Department of Statistics, Stanford University, California



Ecology

#### THE RELATION BETWEEN THE NUMBER OF SPECIES AND THE NUMBER OF INDIVIDUALS IN A RANDOM SAMPLE OF AN ANIMAL POPULATION

By R. A. FISHER (Galton Laboratory), A. STEVEN CORBET (British Museum, Natural History) ND C. B. WILLIAMS (Rothamsted Experimental Station)



ACT I

SCENE I. Elsinore. A platform before the castle.

FRANCISCO at his post. Enter to him BERNARDO

BERNARDO

Who's there?

#### PRINCE FORTINBRAS

Let four captains Bear Hamlet, like a soldier, to the stage; For he was likely, had he been put on, To have proved most royally: and, for his passage, The soldiers' music and the rites of war Speak loudly for him. Take up the bodies: such a sight as this Becomes the field, but here shows much amiss. Go, bid the soldiers shoot.

A dead march. Exeunt, bearing off the dead bodies; after which a peal of ordnance is shot off

Hamlet experiment

- Starting from Act I, read a small fraction of the text
- Stop and estimate the number of distinct words in entire Hamlet

Statistical model: Distinct element problem

- observation:  $X_1, \ldots, X_n$  sampled without replacements from an urn of k colored balls
- parameter: composition of the urn (number of red, blue, etc.)
- goal: number of distinct colors
- assumption: NONE!
- Method: Estimator built from convex/LP duality



#### Example 4: Community detection in networks

Networks with community structures arise in many applications



#### Example 4: Community detection in networks

Networks with community structures arise in many applications



• Task: Discover underlying communities based on the network topology

#### Example 4: Community detection in networks

Networks with community structures arise in many applications



- Task: Discover underlying communities based on the network topology
- Applications: Friend or movie recommendation in online social networks

#### Political blogosphere

...in the 2004 U.S. election [Adamic-Glance '05]





0 n nodes are randomly partitioned into 2 equal-sized communities



- $\mathbf{0}$  n nodes are randomly partitioned into 2 equal-sized communities
- **2** For every pair of nodes in same community, add an edge w.p. p



- $\mathbf{1}$  n nodes are randomly partitioned into 2 equal-sized communities
- **2** For every pair of nodes in same community, add an edge w.p. p
- ${f 3}$  For every pair of nodes in diff. community, add an edge w.p. q



- $\mathbf{0}$  n nodes are randomly partitioned into 2 equal-sized communities
- **2** For every pair of nodes in same community, add an edge w.p. p
- ${f 3}$  For every pair of nodes in diff. community, add an edge w.p. q



#### Stochastic block model - adjacency matrix view



#### Stochastic block model - adjacency matrix view



#### Example 4: Community detection

Statistical model: Stochastic block model SBM(n, p, q)

- observation: a single graph  ${\cal G}$
- parameter: partition of two communities (subsets of [n])
- goal: locate the community (under various criteria)
- assumption: low-rankness of  $\mathbb{E}[adjancency matrix]$

#### Example 5: spiked Wigner model

Noisy observation of rank-one matrix:

$$Y = \lambda x x^{\top} + Z,$$

where

- signal: x uniform on the hypercube  $\{\pm \frac{1}{\sqrt{n}}\}^n$
- noise: Z iid  $N(0, \frac{1}{n})$
- goal: recover x better than chance
  - Find unit vector  $\hat{x} = \hat{x}(Y)$ , s.t.  $\mathbb{E}|\langle \hat{x}, x \rangle| = \Omega(1)$

#### Example 5: spiked Wigner model

Noisy observation of rank-one matrix:

$$Y = \lambda x x^\top + Z,$$

where

- signal: x uniform on the hypercube  $\{\pm \frac{1}{\sqrt{n}}\}^n$
- noise: Z iid  $N(0, \frac{1}{n})$
- goal: recover x better than chance

Find unit vector  $\hat{x} = \hat{x}(Y)$ , s.t.  $\mathbb{E}|\langle \hat{x}, x \rangle| = \Omega(1)$ 

• Random matrix theory: PCA works iff  $\lambda > 1$  [Baik-Ben Arous-Peche '04]

#### Example 5: spiked Wigner model

Noisy observation of rank-one matrix:

$$Y = \lambda x x^\top + Z,$$

where

- signal: x uniform on the hypercube  $\{\pm \frac{1}{\sqrt{n}}\}^n$
- noise: Z iid  $N(0, \frac{1}{n})$
- goal: recover x better than chance

Find unit vector  $\hat{x} = \hat{x}(Y)$ , s.t.  $\mathbb{E}|\langle \hat{x}, x \rangle| = \Omega(1)$ 

- Random matrix theory: PCA works iff  $\lambda > 1$  [Baik-Ben Arous-Peche '04]
- We will show  $\lambda > 1$  is needed by any algo (information-percolation method)

#### What is information theory

Information theory: theory of fundamental limits

- Information measures: How to measure randomness, dependency, dissimilarity (entropy, mutual information, divergence...)
- Coding theorems: Operational problems (data compression, data transmission, etc)

information measures <u>coding theorems</u> fundamental limits operational meaning

#### What is information theory

Information theory: theory of fundamental limits

- Information measures: How to measure randomness, dependency, dissimilarity (entropy, mutual information, divergence...)
- Coding theorems: Operational problems (data compression, data transmission, etc)

information measures <u>coding theorems</u> fundamental limits operational meaning

#### Information-theoretic methods

- Negative results (converse, impossibility results, lower bound):
  - Conceptually: quantify "information" and "dissimilarity"
    - two distributions too "close"  $\Rightarrow$  impossible to distinguish
    - $I(\text{observation}; \text{parameter}) \text{ too "small"} \Rightarrow \text{impossible to estimate}$
    - dimension/entropy too "high"  $\Rightarrow$  need large sample size

#### Information-theoretic methods

- Negative results (converse, impossibility results, lower bound):
  - Conceptually: quantify "information" and "dissimilarity"
    - two distributions too "close"  $\Rightarrow$  impossible to distinguish
    - $I(\text{observation}; \text{parameter}) \text{ too "small"} \Rightarrow \text{impossible to estimate}$
    - dimension/entropy too "high"  $\Rightarrow$  need large sample size
  - More advanced techniques:
    - area theorem
    - strong data processing inequality and information-percolation method (Broadcasting on trees, spiked Wigner model...)
    - (truncated) second moment method

#### Information-theoretic methods

- Negative results (converse, impossibility results, lower bound):
  - Conceptually: quantify "information" and "dissimilarity"
    - two distributions too "close"  $\Rightarrow$  impossible to distinguish
    - $I(\text{observation}; \text{parameter}) \text{ too "small"} \Rightarrow \text{impossible to estimate}$
    - dimension/entropy too "high"  $\Rightarrow$  need large sample size
  - More advanced techniques:
    - area theorem
    - strong data processing inequality and information-percolation method (Broadcasting on trees, spiked Wigner model...)
    - (truncated) second moment method
- Positive results (achievability, constructive results, upper bound):
  - maximal likelihood estimate
  - entropy method (estimators based on pairwise comparison)
  - duality method
  - aggregation
  - efficient procedures/algorithms