

Convex Optimization

Autumn 2022

Ye Shi



信息科学与技术学院

School of Information Science and Technology

Some slides are made by Prof Yuanming Shi

Outline

- **Data science models**

- Linear, bilinear, quadratic, low-rank, and deep models

- **Large-scale optimization**

- Constrained vs. unconstrained, convex vs. nonconvex, deterministic vs. stochastic, solvability vs. scalability

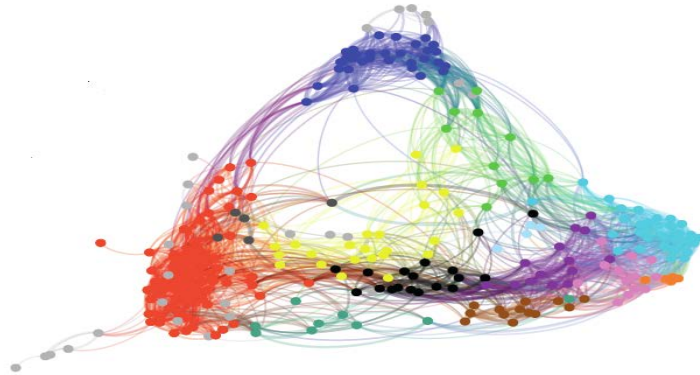
- **High-dimensional statistics**

- Convex geometry, local geometry, global geometry

- **Topics and grading**

- Theoretical foundations, first-order methods, second-order methods, stochastic methods, machine learning approaches, and applications.

*Motivations: **The Era of Big Data***



Intelligent IoT applications



Autonomous vehicles



Smart home



Smart city



Smart health



Smart agriculture



Smart drones

Financial big data

- **Financial data & AI technologies:** set up analytic models to gain valuable insights for better business decisions



Large Volume

data generating at
speed of 1TB/day in
NYSE (2013)

typically 100,000trans/s
in High-frequency Trading

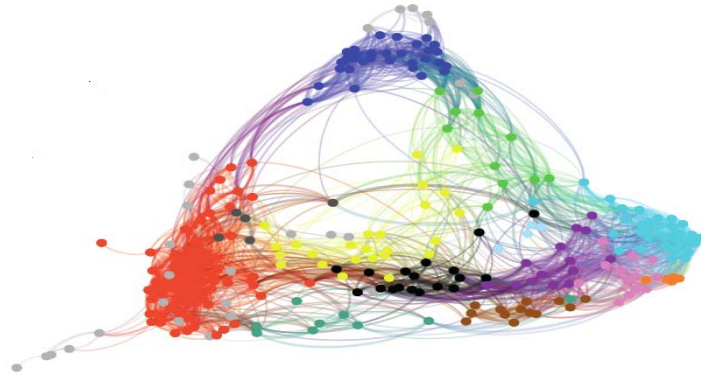
High Velocity



Wide Variety

various data sources
and types

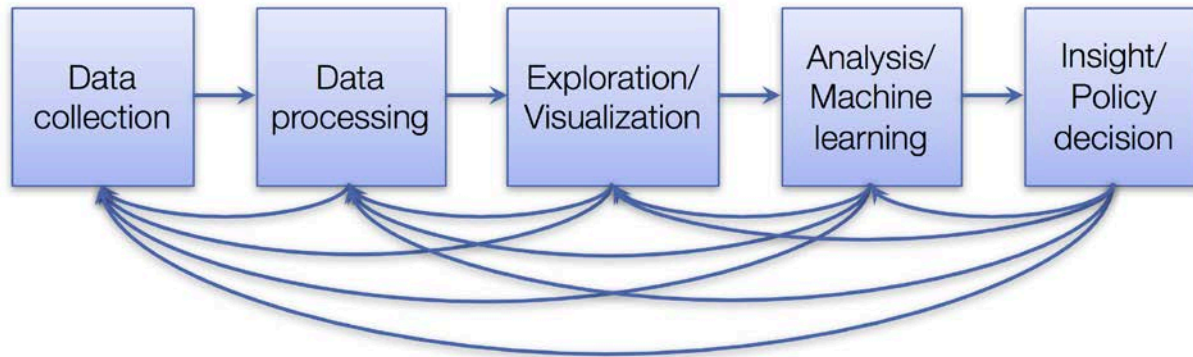
Vignettes A: Data Science Models



What is data science?

- **Some possible definitions**

- Data science is the application of **computational** and **statistical** techniques to address or gain insight into some problem in the **real world**



Challenges

- Retrieve or infer information from high-dimensional/large-scale data



limited processing ability
(computation, storage, ...)

2.5 **exabytes** of data
are generated every day (2012)

exabyte → **zettabyte** → **yottabyte...??**

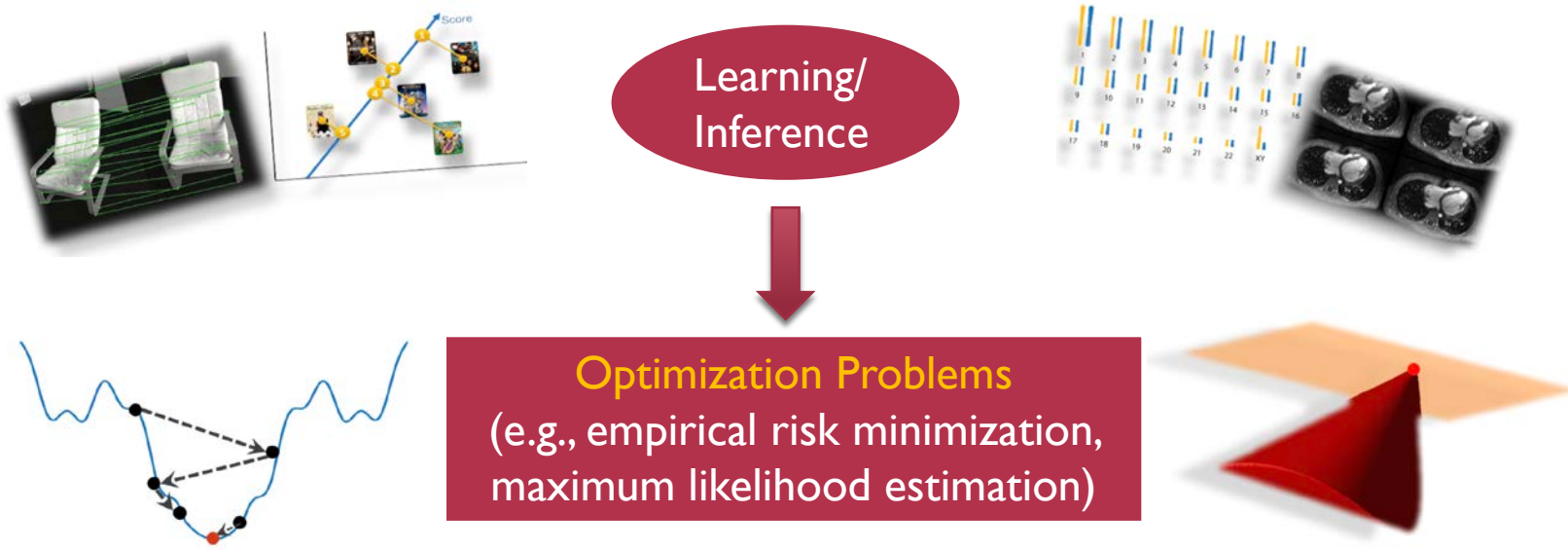
We're interested in the **information** rather
than the data

Challenges:

- ❖ High computational cost
- ❖ Only limited memory is available
- ❖ Do NOT want to compromise statistical accuracy

Optimization for data science

- Optimization has transformed algorithm design



(Convex) optimization is *almost* a tool

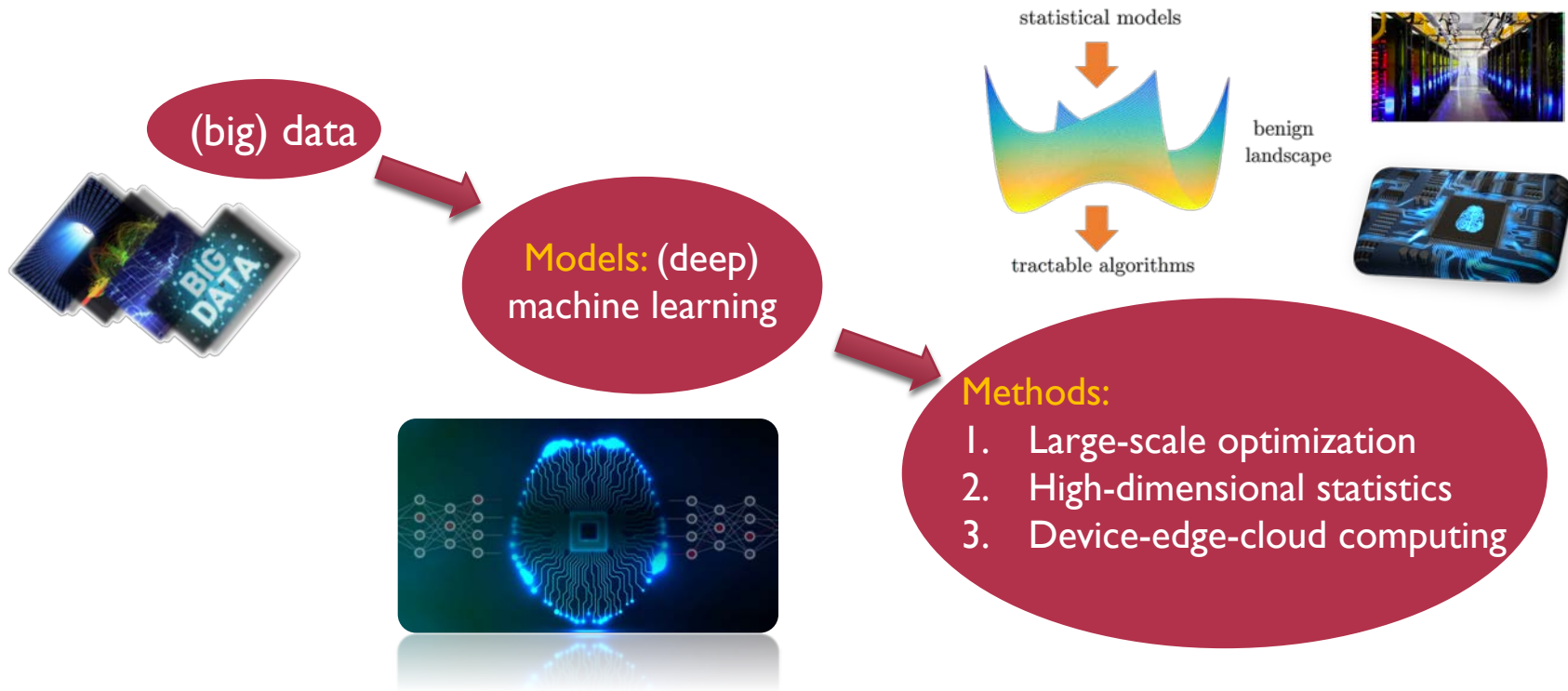
Optimization problem

- General optimization problem in standard form:

$$\begin{array}{ll}\text{minimize} & f_0(\mathbf{x}) \\ \text{subject to} & f_i(\mathbf{x}) \leq 0, i = 1, \dots, m\end{array}$$

- $\mathbf{x} = (x_1, \dots, x_n)$: optimization variables
- $f_0 : \mathbb{R}^n \rightarrow \mathbb{R}$: objective function
- $f_i : \mathbb{R}^n \rightarrow \mathbb{R}, i = 1, \dots, m$: constraint functions
- **Goal:** find optimal solution \mathbf{x}^* minimizing f_0 while satisfying constraints
- **Three basic elements:** 1) variables, 2) constraints, and 3) objective

High-dimensional data analysis



Linear model

- Let $x^\natural \in \mathbb{R}^d$ be an unknown structured sparse signal
 - Individual sparsity for compressed sensing
- Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a convex function that reflects structure, e.g., ℓ_1 -norm
- Let $A \in \mathbb{R}^{m \times d}$ be a measurement operator
- **Observe** $z = Ax^\natural$
- Find estimate \hat{x} by solving **convex program**

$$\text{minimize } f(x) \quad \text{subject to } Ax = z$$



MR scanner



MR image

- **Hope:** $\hat{x} = x^\natural$

Bilinear model

image deblurring

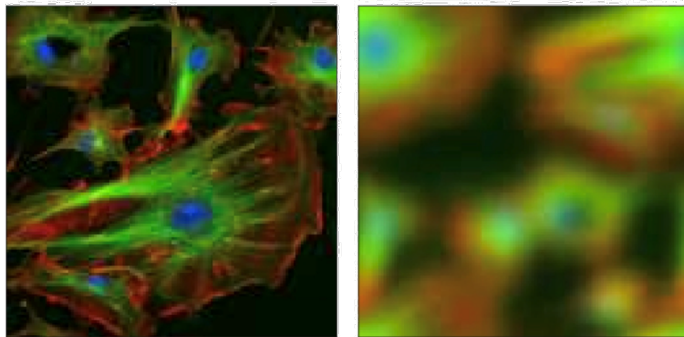
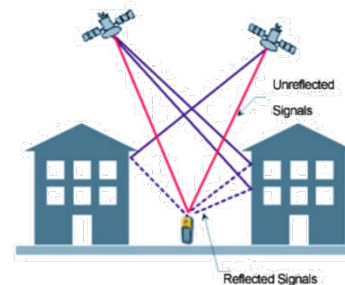


Fig. credit: Romberg

multipath in wireless comm



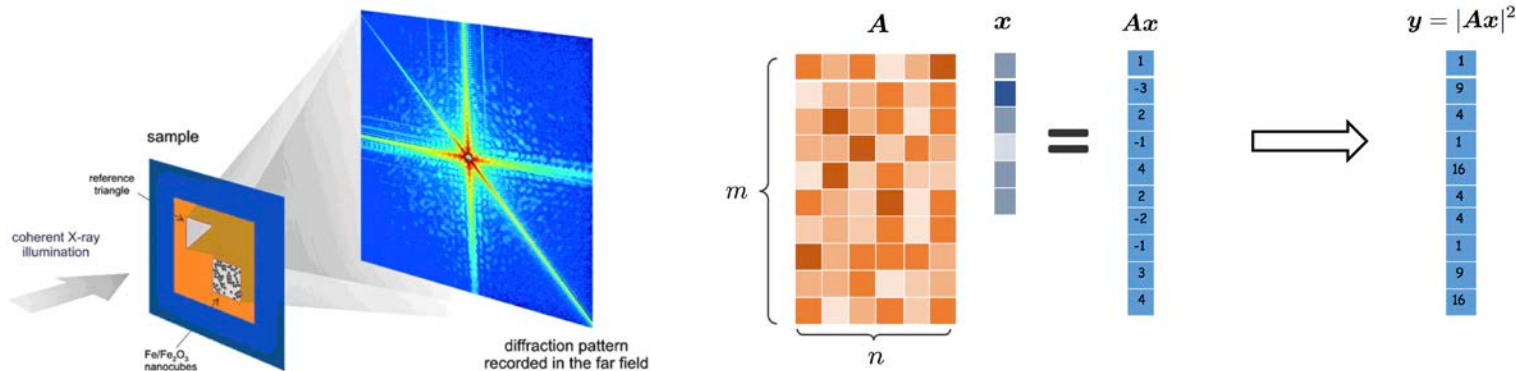
*Fig. credit:
EngineeringsALL*

- **Blind deconvolution:** reconstruct two signals from their convolution

$$\text{find } \mathbf{x}, \mathbf{h} \quad \text{subject to } z_i = \mathbf{b}_i^* \mathbf{h} \mathbf{x}^* \mathbf{a}_i, \quad 1 \leq i \leq m$$

Quadratic model

- **Phase retrieval:** recover signal from intensity (missing phase)



- Recover $z^{\natural} \in \mathbb{R}^n$ from m random quadratic measurements

$$\text{find } z \quad \text{subject to } y_r = |\langle a_r, z \rangle|^2, \quad r = 1, 2, \dots, m$$

Low-rank model



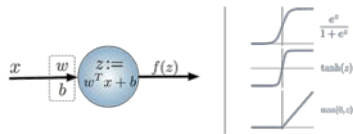
Fig. credit: Candès

- Given partial samples Ω of a low-rank matrix M^\dagger , fill in missing entries

$$\underset{M \in \mathbb{C}^{m \times n}}{\text{minimize}} \quad \text{rank}(M) \quad \text{subject to} \quad Y_{i,k} = M_{i,k}, \quad (i,k) \in \Omega$$

Deep models

- **Data:** n observations $\{\mathbf{x}_i, y_i\}_{i=1}^n \in \mathcal{X} \times \mathcal{Y}$
- **Prediction function:** $h(\mathbf{x}, \boldsymbol{\theta}) \in \mathbb{R}$ parameterized by $\boldsymbol{\theta} \in \mathbb{R}^d$
 - linear predictions: $h(\mathbf{x}, \boldsymbol{\theta}) = \boldsymbol{\theta}^T \Phi(\mathbf{x})$ using features $\Phi(\mathbf{x})$
 - neural networks: $h(\mathbf{x}, \boldsymbol{\theta}) = \boldsymbol{\theta}_m^T \sigma(\boldsymbol{\theta}_{m-1}^T \sigma(\cdots \boldsymbol{\theta}_2^T \sigma(\boldsymbol{\theta}_1^T \mathbf{x})))$
- Estimating $\boldsymbol{\theta}$ parameters is an **optimization problem** (ℓ : loss function)

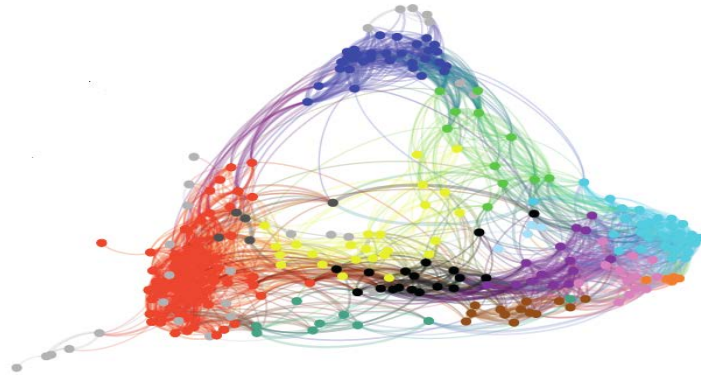


$$\text{minimize } f(\boldsymbol{\theta}) := \frac{1}{n} \sum_{i=1}^n \ell(h(\mathbf{x}_i, \boldsymbol{\theta}), y_i) \quad \text{subject to } \mathcal{R}(\boldsymbol{\theta}) \leq \tau$$

- \mathcal{R} : regularization function encoding prior information (e.g., sparse) on $\boldsymbol{\theta}$

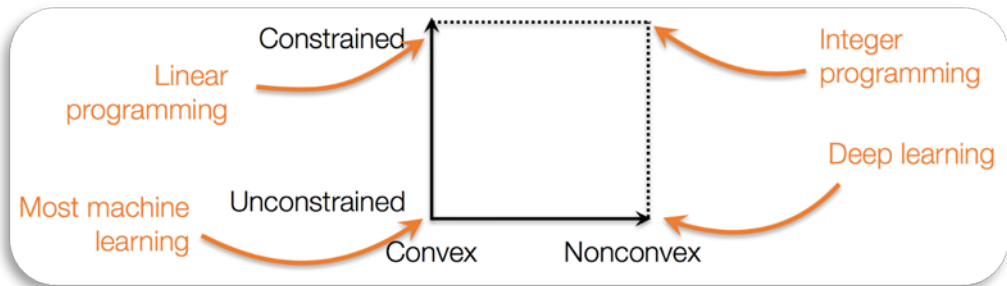
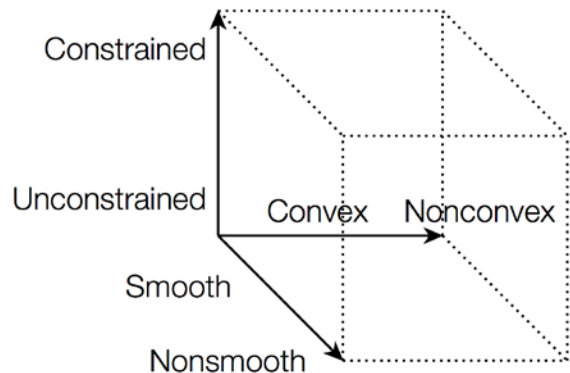
Key benefits of looking at problems in AI as optimization problems:
separate out the *definition* of the problem from the *method for solving it!*

Vignettes B: *Large-Scale Optimization*



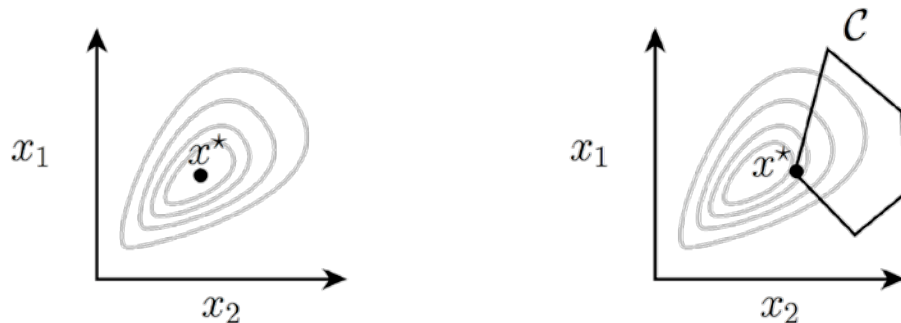
Classes of optimization problems

- Types of optimization problems: linear programming, nonlinear programming, integer programming, geometric programming, ...



- We focus on three dimensions: **unconstrained vs. constrained**, **convex vs. nonconvex**, and **smooth vs. nonsmooth**

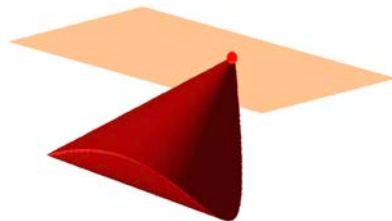
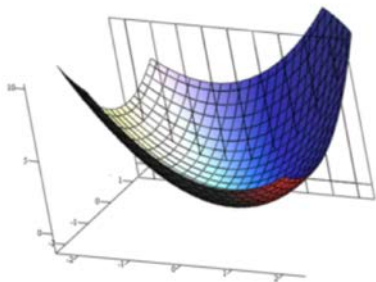
Constrained vs. unconstrained optimization



- **Unconstrained optimization:** every point $x \in \mathbb{R}^n$ is feasible, so only focus is on minimizing $f(x)$
- **Constrained optimization:** it may be difficult to even *find* a feasible point $x \in \mathcal{C}$

Typically leads to different classes of algorithms

Convex vs. nonconvex optimization



Convex optimization:

- 1) All local optima are global optima
- 2) Can be solved in polynomial-time

“... the great watershed in optimization isn’t between linearity and nonlinearity, but convexity and nonconvexity”

— R. Rockafellar ’1993



Deterministic vs. stochastic optimization

- **Stochastic optimization**

$$\text{minimize } f(\mathbf{x}) := \mathbb{E}[F(\mathbf{x}, \boldsymbol{\xi})] \quad \text{subject to } \mathbf{x} \in \mathcal{X}$$

➤ f : loss; \mathbf{x} : parameters; $\boldsymbol{\xi}$: data samples

- **Example:** supervised machine learning (finite-sum problems)

$$\text{minimize } f(\mathbf{x}) := \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \ell(b_i - \mathbf{a}_i^T \mathbf{x})$$

➤ Data observations: $(\mathbf{a}_i, b_i) \in \mathbb{R}^d \times \mathbb{R}$; loss function: $\ell : \mathbb{R}^d \rightarrow \mathbb{R}$

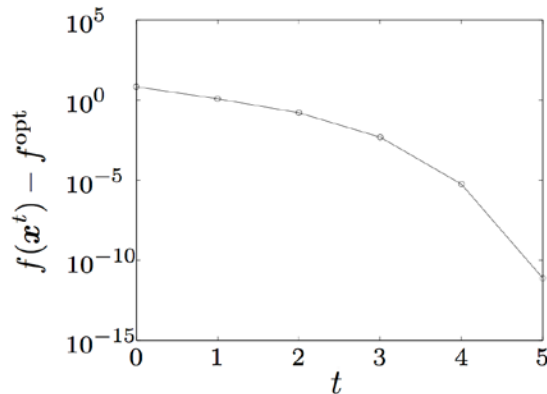
- **Stochastic gradient:** $\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \nabla f_{i(k)}(\mathbf{x}_k)$

➤ $i(k) \in \{1, 2, \dots, n\}$ uniformly at random; unbiased estimate: $\mathbb{E}[\nabla f_{i(k)}] = \nabla f$

Scaling issues: solvability vs. scalability

- Polynomial-time algorithms might be *useless* in large-scale applications
- **Example:** Newton's method

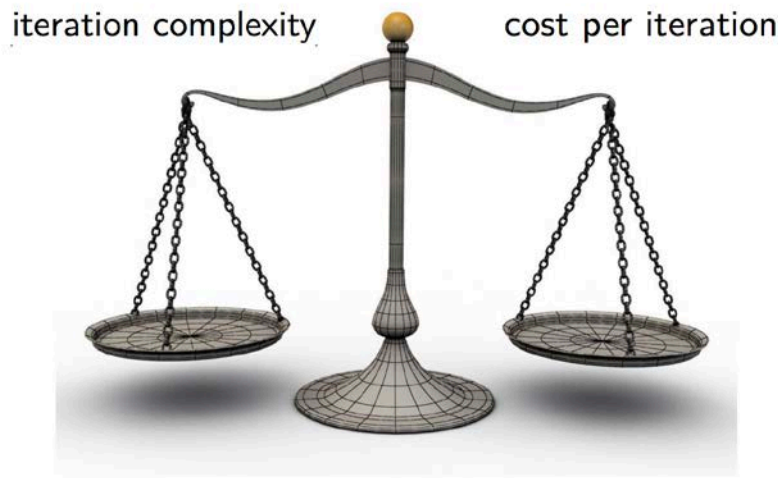
$$\begin{aligned} & \text{minimize}_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) \\ & \mathbf{x}^{t+1} = \mathbf{x}^t - (\nabla^2 f(\mathbf{x}^t))^{-1} \nabla f(\mathbf{x}^t) \end{aligned}$$



- Attains ϵ accuracy within $\mathcal{O}(\log \log \frac{1}{\epsilon})$ iterations; requires $\nabla^2 f(\mathbf{x}) \in \mathbb{R}^{n \times n}$
- *A single iteration may last forever*; prohibitive storage requirement

Iteration complexity vs. per-iteration cost

computational cost = iteration complexity (#iterations) x cost per iteration



Large-scale problems call for methods with *cheap iterations*

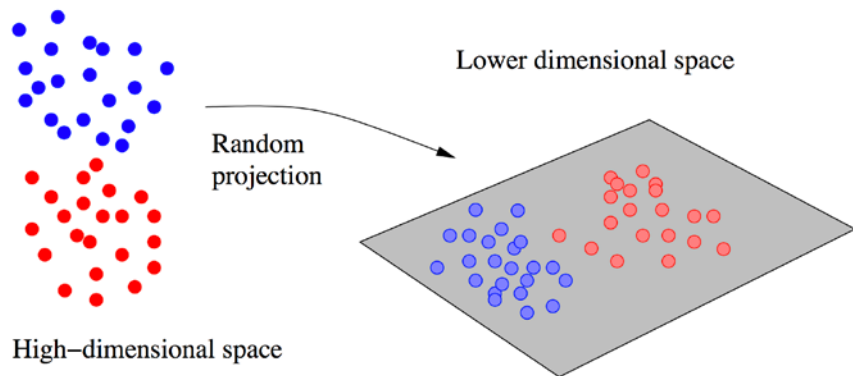
First-order methods



- **First-order methods:** methods that exploit only information on function values and (sub)gradients without using Hessian information
 - cheap iterations
 - low memory requirements

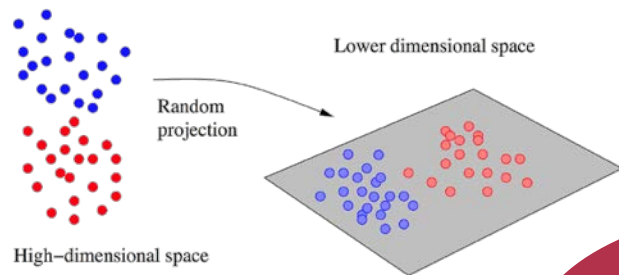
Randomized and approximation methods

- **Optimization for high-dimensional data analysis:** polynomial-time algorithms often not fast enough: further *approximations* are essential

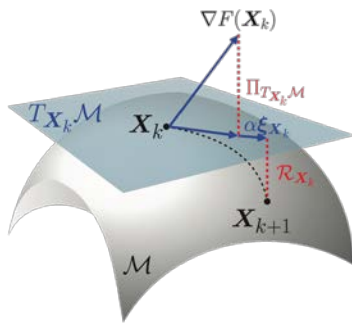


- **Randomized and stochastic methods:** project data into subspace, and solve reduced dimension problem

Advanced large-scale optimization



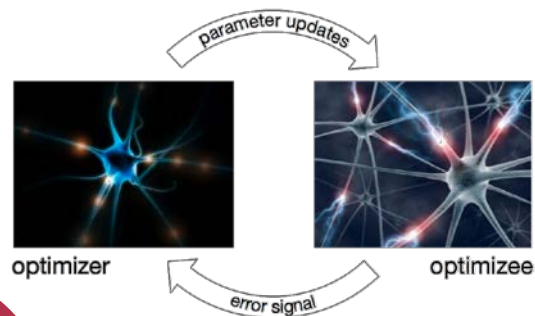
randomized methods



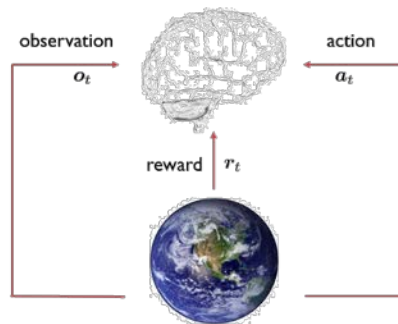
nonconvex optimization on manifold



Goal: scalable, real-time,
parallel, distributed,
automatic, etc.

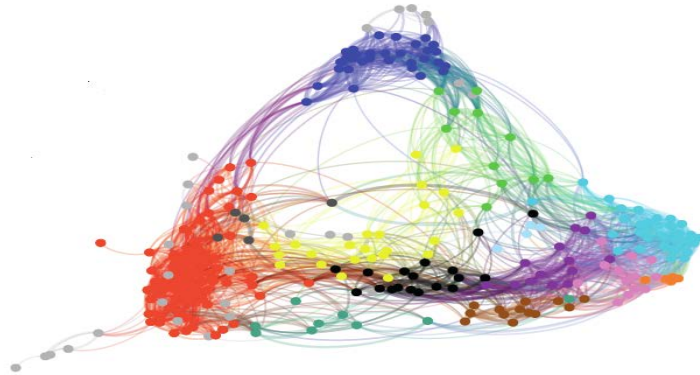


learning to optimize



deep reinforcement learning

Vignettes C: *High-Dimensional Statistics*

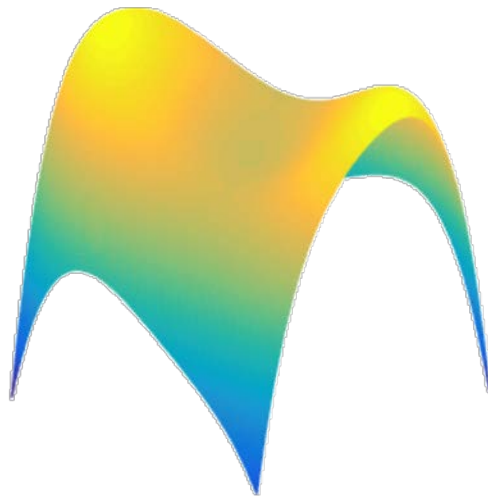


Nonconvex problems are everywhere

- Empirical risk minimization is usually nonconvex

$$\underset{\boldsymbol{x}}{\text{minimize}} \quad f(\boldsymbol{x}; \boldsymbol{\theta})$$

- low-rank matrix completion
- blind deconvolution/demixing
- dictionary learning
- phase retrieval
- mixture models
- deep learning
- ...



Nonconvex optimization may be super scary

- **Challenges:** saddle points, local optima, bumps,...

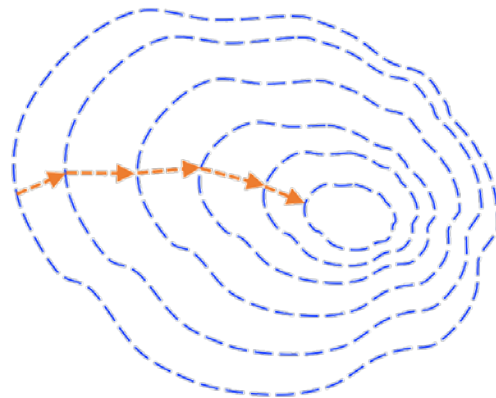
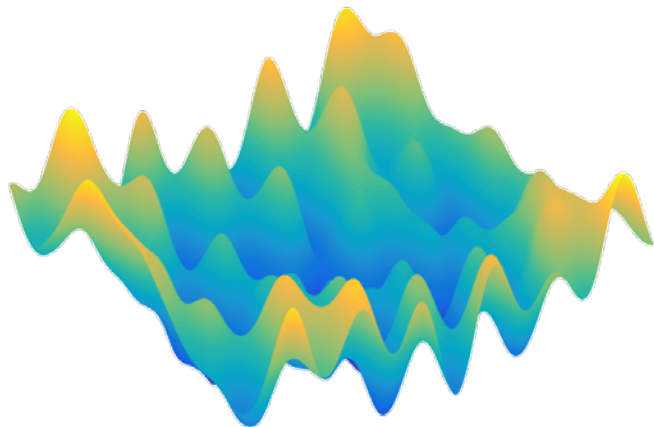


Fig. credit: Chen

- **Fact:** they are usually solved on a daily basis via simple algorithms like (stochastic) gradient descent

Statistical models come to rescue

- **Blessings:** when data are generated by certain statistical models, problems are often much nicer than worst-case instances

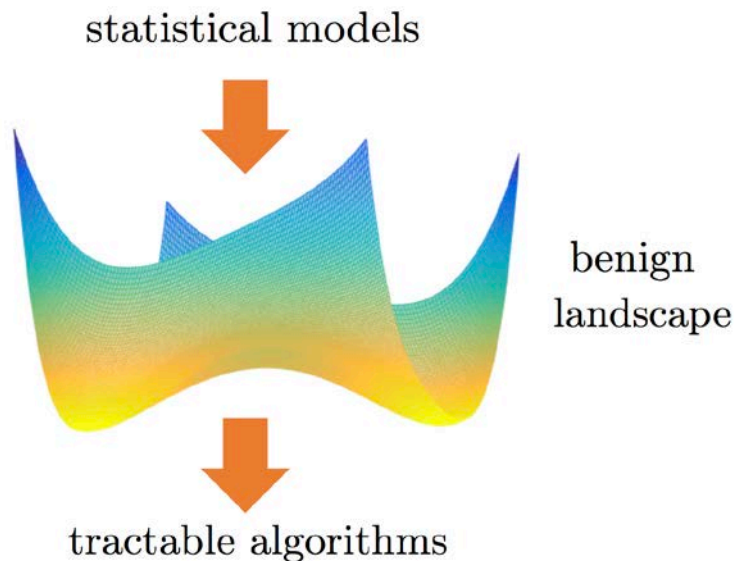


Fig. credit: Chen

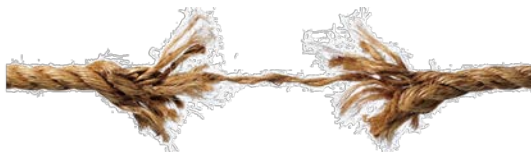
Global geometry

- **Proposal:** separation of landscape analysis and generic algorithm design

landscape analysis
(statistics)

all local minima are
global minima

- dictionary learning (Sun et al. '15)
- phase retrieval (Sun et al. '16)
- matrix completion (Ge et al. '16)
- synchronization (Bandeira et al. '16)
- inverting deep neural nets (Hand et al. '17)
- ...



generic algorithms
(optimization)

all the saddle points
can be escaped

- gradient descent (Lee et al. '16)
- trust region method (Sun et al. '16)
- perturbed GD (Jin et al. '17)
- cubic regularization (Agarwal et al. '17)
- Natasha (Allen-Zhu '17)
- ...

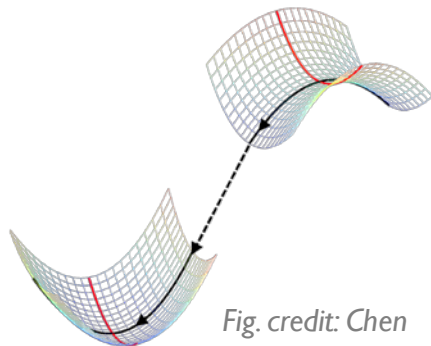


Fig. credit: Chen

Local geometry

- **Initialize** within local basin sufficiently close to ground-truth (i.e., strongly convex, no saddle points/ local minima)
- **Iterative refinement** via some iterative optimization algorithms

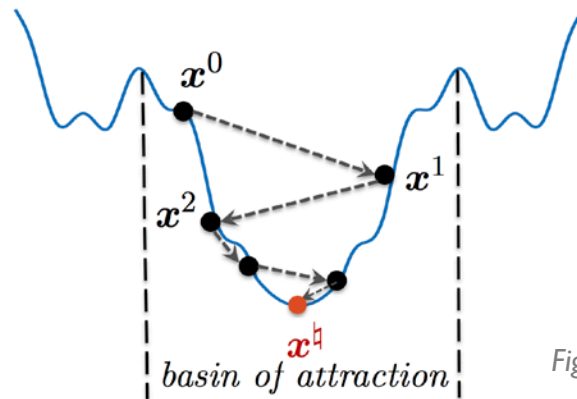
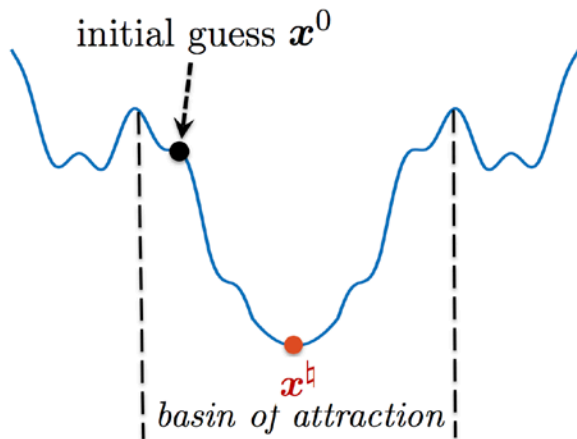
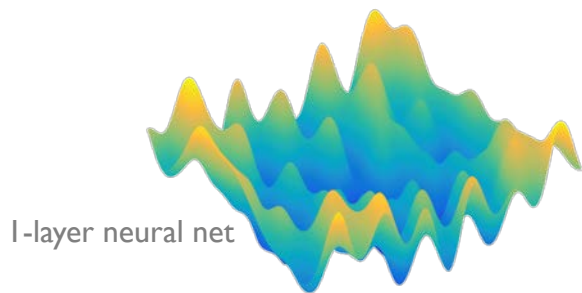


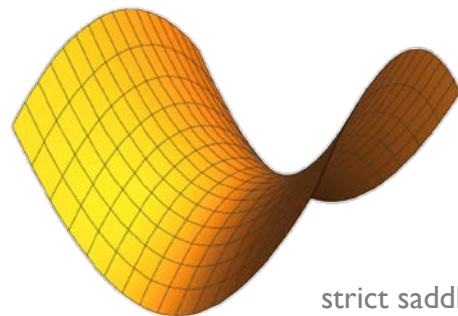
Fig. credit: Chen

Optimization meets statistics



big data & deep learning

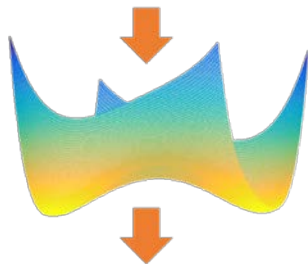
statistical models



nonconvex optimization may be super scary

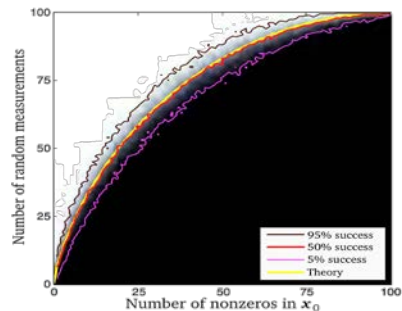
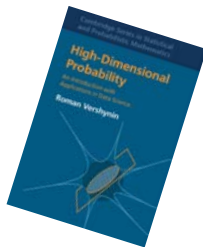
benign geometry: no spurious local optima

statistical models



tractable algorithms

high-dimensional
probability & statistics



framework: high-dimensional data analysis

Goals: data sizes, expressivity,
information propagations, etc.

Case study: bilinear model

■ Demixing from bilinear measurements

$$\begin{aligned} &\text{find } \{x_i\}, \{h_i\} \\ &\text{subject to } z_j = \sum_{i=1}^s b_j^* h_i x_i^* a_{ij} \end{aligned}$$

■ Applications

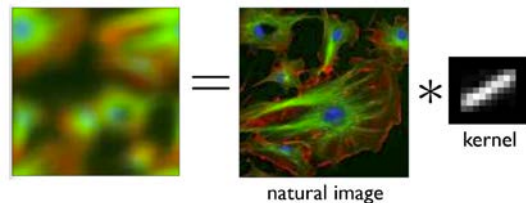
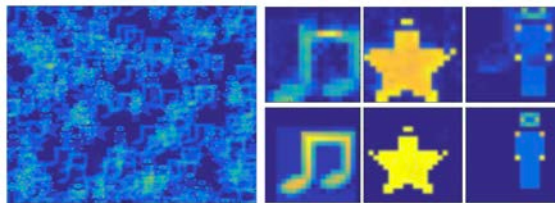
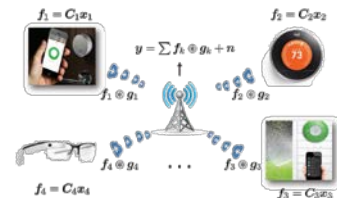


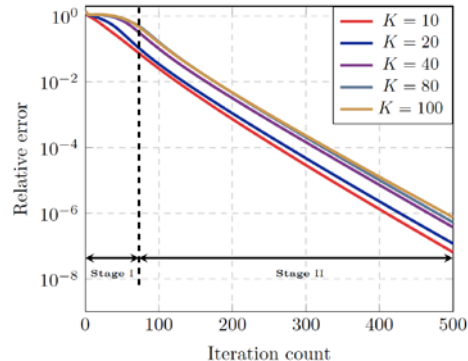
image deblurring



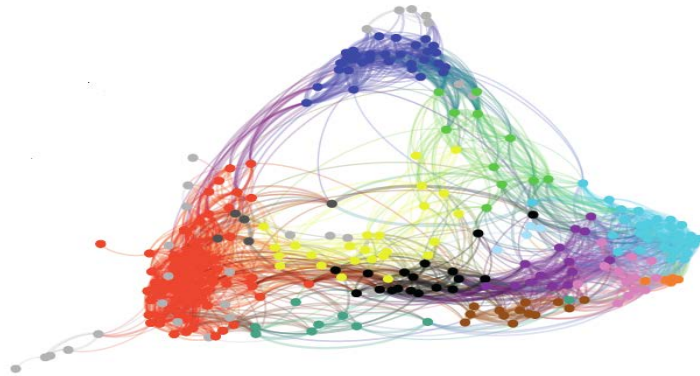
convolutional dictionary learning



low-latency communication

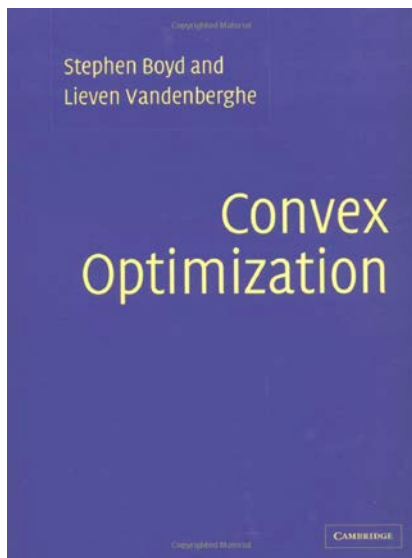


Topics and Grading



Theoretical foundations

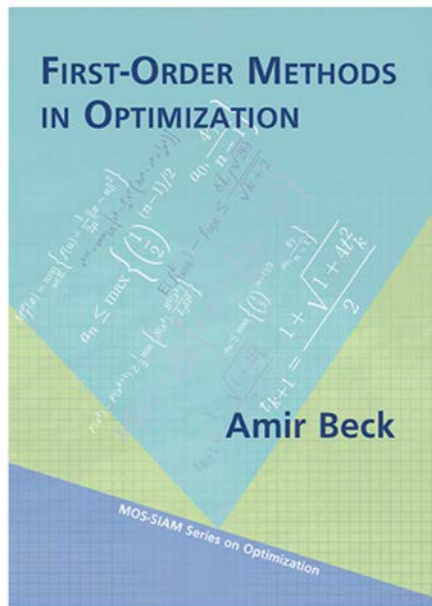
- **Main topics:** convex sets, convex functions, convex problems, Lagrange duality and KKT conditions, disciplined convex programming



Convex Optimization, by S. Boyd and L. Vandenberghe, Cambridge University Press, 2003.

First-order methods

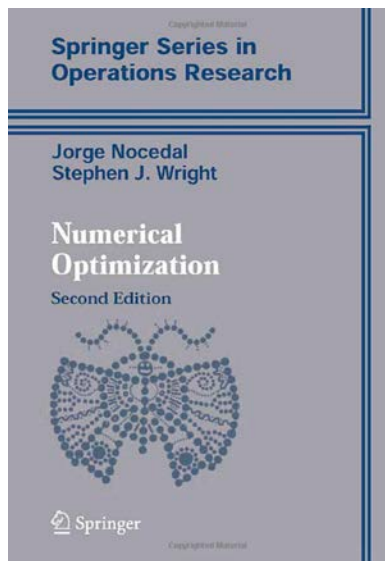
- **Main topics:** gradient methods, subgradient methods, proximal methods



First-order Methods in Optimization, by A. Beck,
MOS-SIAM Series on Optimization, 2017.

Second-order methods

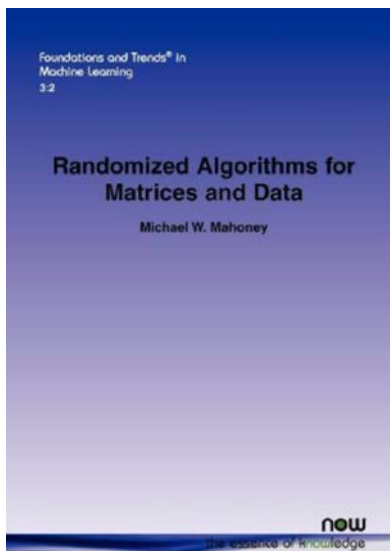
- **Main topics:** Newton method, interior-point methods, quasi-Newton methods



Numerical Optimization, by J. Nocedal and S. Wright, Springer-Verlag, 2006.

Stochastic and randomized methods

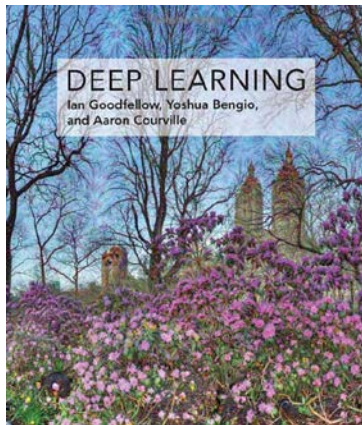
- **Main topics:** stochastic gradient methods, stochastic Newton methods, randomized sketching methods, randomized linear algebra



Lecture Notes on Randomized Linear Algebra
, by Mahoney, Michael, 2016.

Machine learning for optimization

- **Main topics:** sparse optimization (deep neural networks), mixed integer nonlinear programming (imitation learning), nonconvex optimization (statistical learning), multiple objective optimization (active learning)



Deep Learning, by I. Goodfellow, Y. Bengio and A. Courville, MIT Press, 2016.

Applications

- **Machine learning:** Optimization Layers, Learning to Optimize, Federated learning, etc.
- **Smart grids:** Sparse and low rank optimization (Optimal power flow), mixed integer linear/nonlinear programming (Electric vehicle charging/discharging, Demand response, PMU placement)
- **Finance Engineering:** portfolio optimization, factor models, time series modeling, robust portfolio optimization, risk-parity portfolio, index tracking, pairs trading

Prerequisites

- **Warning:** there will be quite a few THEOREMS and PROOFS ...
- Basic linear algebra
- Basic probability
- A programming language (e.g. Matlab, Python, ...)

Somewhat surprisingly, most proofs rely only on basic linear algebra and elementary recursive formula

Grading

- **Homework:** 3 homework assignments

- Quiz: Random in class

- **Course project:**

- either individually or in groups of two/three
- list of topics; final report & slides

- **Final Exam:**

$$\text{Grade} = 0.2H + 0.1Q + 0.3P + 0.4E$$

H: homework; Q: quiz; P: course project; E: final exam

Grading

- **Regrade Requests**

- If you feel you deserved a better grade on an assignment, you may submit a regrade request by email within **3 days** of the grade release. Your request should briefly summarize why you feel the original grade was unfair. Your TA will re-evaluate your assignment as soon as possible, and then issue a decision.

- **Late Policy**

- All students have 4 free late days for the quarter.
- You may use up to 2 late days per assignment with no penalty.
- You may use late days for the assignments, project.

Once you have exhausted your free late days, we will deduct a late penalty of 25% per additional late day. For example: you submit A1 one day late, submit A2 three days late, and submit A3 two days late. You receive no penalty for A1, and exhaust one of your free late days. For A2 the first two late days exhaust two of your free late days; the third day late incurs a 25% penalty. For A3 the first late day exhausts your final free late day; the second late day incurs a 25% penalty.

Course information

- **Instructor:** Ye Shi (<http://faculty.sist.shanghaitech.edu.cn/faculty/shiye>)
 - Email: shiye@shanghaitech.edu.cn
 - Office location: Room 1A-404A, SIST Building
 - Office hours: Wednesday 15:00-16:00 (or by appointments)
- **TAs:**
 - Haixiang Sun (sunhx@shanghaitech.edu.cn)
 - Pengchao Tian (tianpch@shanghaitech.edu.cn)
 - Hongxia Li (lihx2@shanghaitech.edu.cn)
 - Shutong Ding (dingsht@shanghaitech.edu.cn)
 - Office location: Room 1A-405, SIST Building
 - Office hours: Thursday 19:00-20:00 (or by appointments)

Course information

- Use **WeChat** as the main mode of electronic communication; please post (and answer) questions there!



SI251凸优化—2022 Fall



该二维码7天内(9月11日前)有效, 重新进入将更新

- Course website: **BlackBoard**

[https://elearning.shanghaitech.edu.cn:8443/webapps/blackboard/execute/modulepage/view?course_id= 3088 _1&cmp_tab_id= 7476 _1&editMode=true&mode=cpview](https://elearning.shanghaitech.edu.cn:8443/webapps/blackboard/execute/modulepage/view?course_id=3088_1&cmp_tab_id=7476_1&editMode=true&mode=cpview)

Thanks