- ▶ Instructor: Xiangxiong Zhang
- ▶ Course webpage:
  https://www.math.purdue.edu/∼zhan1966/teaching/574
- ▶ Selected topics from the following reference books:
  - ▶ Beck, Introduction to Nonlinear Optimization
  - ▶ Beck, First order methods in optimization
  - ▶ Ryu and Yin, Large-Scale Convex Optimization: Algorithms & Analyses via Monotone Operators
  - ▶ Nicolas Boumal, An introduction to optimization on smooth manifolds

# Differences compared to other optimization courses on campus

- ▶ There are other graduate level optimization courses offered in CS and engineering departments at Purdue.
- ▶ MA 574 is the only math graduate course on numerical optimization. Fall 2024 will be the first time that it will be taught, covering four topics:
  1. Smooth optimization methods such as gradient descent and accelerated gradient descent.
  2. Nonsmooth convex optimization such as proximal gradient and splitting methods.
  3. Randomized and stochastic methods.
  4. Riemannian optimization.
- ▶ I taught MA 598 Topics in Optimization in 2023 covering the first three topics.
- ▶ In MA 574, we focus on convergence analysis. Less than one half of MA 574 are classical ones covered in a standard optimization textbook/course, while the other content may not be covered in other optimization courses:
  - ▶ Many methods and techniques such as Nesterov's acceleration, stochastic gradient descent and Riemannian optimization became popular only after 2000, thus they were usually not covered in a book/course 20 or even 10 years ago.
  - ▶ Riemannian optimization is currently not covered in other courses on campus.

# Plan for this semester

There are many different types of optimization problems, but we mainly focus on **the convergence** of algorithms minimizing a convex function $f(x)$ with a large scale:

- ▶ Part I: some classical algorithms for minimizing a smooth function $f(x)$ such as gradient descent, accelerated gradient descent, Newton's method, quasi Newton methods, etc.

- ▶ Part II: algorithms for composite optimization of minimizing $f(x) + g(x)$ where $f(x)$ and $g(x)$ are both convex, but at least one of them is not differentiable, e.g.,

$$\min \|x\|_1 + \|Ax - b\|_2^2$$

  where $\|x\|_1 = \sum_i |x_i|$.

- ▶ Part III: stochastic type algorithms, such as stochastic gradient descent.

- ▶ Part IV: minimization over a Riemannian manifold constraint.

# Examples

▶ Part I: for $\min_x f(x)$, the gradient descent method is

$$x_{k+1} = x_k - \eta_k \nabla f(x_k)$$

When and why does gradient descent converge? How fast does it converge?
Prerequisites for Part I:
  ▶ Calculus: gradient, Hessian, Taylor Theorem...
  ▶ Linear algebra: eigenvalues, singular values and etc.

▶ Part II: we will introduce subderivatives, proximal operator, and algorithms using the subderivatives. We will use monotonicity of operators to prove convergence.
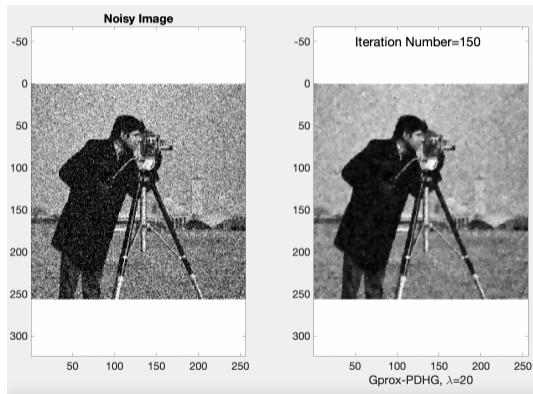
$$\min \|x\|_1 + \|Ax - b\|_2^2$$

We will need some knowledge on convex non-differentiable functions, which will be covered in the class.

▶ Part II: here is another example of nonsmooth convex optimization for denoising a given noisy image $A$ via TV (total variation) norm minimization

$$\min_x \|x\|_{TV} + \lambda \|x - A\|_2^2,$$

where $\|x\|_{TV} = \sum_{i,j} \sqrt{|x_{i,j} - x_{i+1,j}|^2 + |x_{i,j} - x_{i,j+1}|^2}$



The algorithm PDHG will be covered in part II, and the paper on this method would be a good choice for the final presentation.

Large scale means: if the dimension of $x$ is $n$ then only $\mathcal{O}(n)$ storage is acceptable. What is $n^2$?

# Examples

▶ Part III: for minimizing $f(x) := \sum\limits_{i=1}^{N} f_i(x)$, the full gradient is $\nabla f(x) = \sum\limits_{i=1}^{N} \nabla f_i(x)$, we can use the stochastic gradient like

$$\nabla_S f(x) := \sum_{i \in S} \nabla f_i(x)$$

where $S$ is a random small subset of $\{1, 2, \cdots, N\}$. The stochastic gradient descent can be defined as:
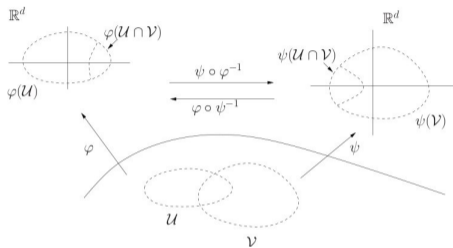
$$x_{k+1} = x_k - \eta_k \nabla_{S_k} f(x_k).$$

In order to analyze the convergence, we need some probability knowledge, which will be introduced.
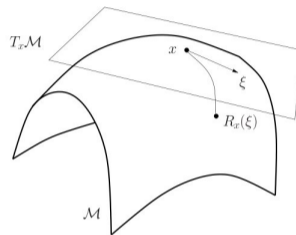
An example where $N$ is too large: recommendation systems for customers rating products (movies, merchandise, etc).

## Examples

Part IV: consider minimizing $f(x)$ with $x \in \mathcal{M} \subset \mathbb{R}^N$ where $\mathcal{M}$ is a Riemannian manifold. If you have not heard of manifolds, just think of $\mathcal{M}$ being a surface, e.g., a unit sphere.
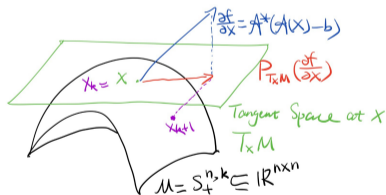


(a) Charts

(b) Tangent Space and Retraction

► Manifold over $\mathbb{R}$: $M$ is a set and it is locally diffeomorphic to $\mathbb{R}^d$.

► Tangent Space: a tangent vector is tangent to a curve on $M$.

► For $f(X)$ defined on $M$, the Riemannian gradient grad $f(X)$ is a tangent vector.

# An example of Riemannian Gradient

▶ Consider $\min_{X \in \mathcal{M}} f(X) = \frac{1}{2}\|\mathcal{A}(X) - b\|^2$ where $\mathcal{A}$ is a linear operator and $\mathcal{M}$ is an embedded manifold in $\mathbb{R}^N$.

▶ Gradient: $\nabla f(X) = \mathcal{A}^*(\mathcal{A}(X) - b)$.

▶ Riemannian gradient is the projection of $\frac{\partial f(X)}{\partial X} = \mathcal{A}^*(\mathcal{A}(X) - b)$ onto $T_X \mathcal{M}$

# Focuses and learning outcomes of this course

- ▶ We focus on analysis of classical algorithms, i.e., why and how fast they converge. Applications will be barely mentioned, though questions about applications are always welcome.
- ▶ A final presentation/report (depending on our schedule) is required by reading a paper and/or implementing some classical/novel algorithms. Examples of possible choices of papers:
  - ▶ Convergence of nonlinear conjugate gradient method.
  - ▶ Convergence analysis of Adam.
  - ▶ Stochastic gradient Langevin dynamics.
- ▶ Learning outcome: by the end of the semester, I expect you to ??