

# CS6301: Optimization in Machine Learning

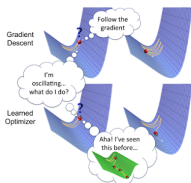
## Lecture 1: Introduction and Course Overview

Rishabh Iyer

Department of Computer Science  
University of Texas, Dallas

<https://sites.google.com/view/cs-6301-optml/home>

January 13, 2020



- Why take this course?
- Prerequisites
- Week by Week Course plan
- Course Logistics
- Basics of Linear Algebra
- Applications of Optimization in Machine Learning



# Acknowledgements

I would like to acknowledge several sources I have used to create slides

- Martin Jaggi's course at EPFL:  
[https://github.com/epfml/OptML\\_course](https://github.com/epfml/OptML_course)
- Mark Schmidt's ML and Advance ML course at UBC:  
<https://www.cs.ubc.ca/~schmidtm/Courses/LecturesOnML>
- Several Blogs: [geeksforgeeks.org](http://geeksforgeeks.org), [healthylalgorithms.com](http://healthylalgorithms.com)
- Prof. Ganesh Ramakrishnan's Convex Optimization course at IIT Bombay: <https://www.cse.iitb.ac.in/~cs709/>.
- Prof. Ramakrishnan is also the co-creator of this course which he will be teaching at IIT Bombay.



# Why take this Course?

Optimization is everywhere: Big Data and Machine Learning, Scheduling and Planning, Operations Research, control theory, data analysis, simulations, almost all technology we use, search engines, computers/laptops, smart-phones, hardware/software of all kinds, ...

- **Mathematical Modeling:**
  - defining and modeling the problem



# Why take this Course?

Optimization is everywhere: Big Data and Machine Learning, Scheduling and Planning, Operations Research, control theory, data analysis, simulations, almost all technology we use, search engines, computers/laptops, smart-phones, hardware/software of all kinds, ...

- **Mathematical Modeling:**
  - defining and modeling the problem
- **Computational Optimization:**
  - Algorithms to solve these optimization problems optimally or near optimally.



# Why take this Course?

Machine Learning and AI are embedded in practically every spear of our life: Google/Bing Search & Ads, Amazon product search/recommendation, driverless cars, Google Maps, Uber/Ola matching, Google Photos, Youtube, Facebook, Twitter, Microsoft Office, ...

- With democratization of AI, software engineers can write ML applications with a few lines of code!



# Why take this Course?

Machine Learning and AI are embedded in practically every spear of our life: Google/Bing Search & Ads, Amazon product search/recommendation, driverless cars, Google Maps, Uber/Ola matching, Google Photos, Youtube, Facebook, Twitter, Microsoft Office, ...

- With democratization of AI, software engineers can write ML applications with a few lines of code!
- Open-source libraries today offer capabilities to build products with practically zero knowledge of ML.



# Why take this Course?

Machine Learning and AI are embedded in practically every spear of our life: Google/Bing Search & Ads, Amazon product search/recommendation, driverless cars, Google Maps, Uber/Ola matching, Google Photos, Youtube, Facebook, Twitter, Microsoft Office, ...

- With democratization of AI, software engineers can write ML applications with a few lines of code!
- Open-source libraries today offer capabilities to build products with practically zero knowledge of ML.
- However to push the boundaries of research and really solve problems, you need to gain hands on experience in ML!





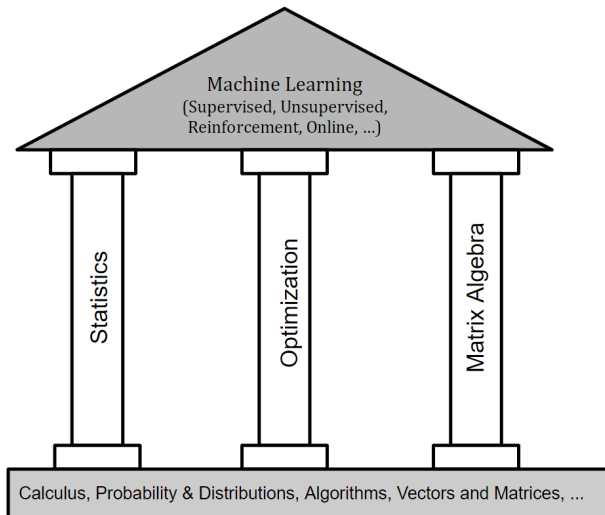
# Why take this Course?

Machine Learning and AI are embedded in practically every spear of our life: Google/Bing Search & Ads, Amazon product search/recommendation, driverless cars, Google Maps, Uber/Ola matching, Google Photos, Youtube, Facebook, Twitter, Microsoft Office, ...

- With democratization of AI, software engineers can write ML applications with a few lines of code!
- Open-source libraries today offer capabilities to build products with practically zero knowledge of ML.
- However to push the boundaries of research and really solve problems, you need to gain hands on experience in ML!
- Optimization is one of the important backbones of machine learning.



# Why take this Course?



# Why take this Course?

Optimization one of the pillars of ML!

- **Continuous Optimization:**
  - Continuous Optimization often appears as *relaxations* of risk/error minimization problem. The *Learning* problem in many parametrized models (whether supervised, semi-supervised, unsupervised, or reinforcement learning) involves **Continuous Optimization**.



# Why take this Course?

Optimization one of the pillars of ML!

- **Continuous Optimization:**

- Continuous Optimization often appears as *relaxations* of risk/error minimization problem. The *Learning* problem in many parametrized models (whether supervised, semi-supervised, unsupervised, or reinforcement learning) involves **Continuous Optimization**.

- **Discrete Optimization:**

- Discrete Optimization occurs in Inference problems in structured spaces, certain learning problems and auxiliary problems such as Feature Selection, Data subset selection, Data summarization, Architecture search etc.



# Why take this Course?

Optimization one of the pillars of ML!

- **Continuous Optimization:**

- Continuous Optimization often appears as *relaxations* of risk/error minimization problem. The *Learning* problem in many parametrized models (whether supervised, semi-supervised, unsupervised, or reinforcement learning) involves **Continuous Optimization**.

- **Discrete Optimization:**

- Discrete Optimization occurs in Inference problems in structured spaces, certain learning problems and auxilliary problems such as Feature Selection, Data subset selection, Data summarization, Architecture search etc.

- **Mixed Continuous and Discrete Optimization:**

- A growing number of problems including classical problems such as clustering, feature selection, structured sparsity occur as mixed discrete/continuous optimization problems.



# Why take this Course?

- Countless number of ML libraries available which implement all kinds of optimization algorithms (Tensorflow, pytorch, scipy, sklearn, Vowpal Wabbit, ...)



# Why take this Course?

- Countless number of ML libraries available which implement all kinds of optimization algorithms (Tensorflow, pytorch, scipy, sklearn, Vowpal Wabbit, ...)
- **This course** will give you the expertise to look inside these algorithms, understand how they work, why they work and how fast they work.



# Why take this Course?

- Countless number of ML libraries available which implement all kinds of optimization algorithms (Tensorflow, pytorch, scipy, sklearn, Vowpal Wabbit, ...)
- **This course** will give you the expertise to look inside these algorithms, understand how they work, why they work and how fast they work.
- Invariably in Research, you will come up with new optimization problems which you might need to implement custom algorithms or atleast the loss functions.





# Why take this Course?

- Countless number of ML libraries available which implement all kinds of optimization algorithms (Tensorflow, pytorch, scipy, sklearn, Vowpal Wabbit, ...)
- **This course** will give you the expertise to look inside these algorithms, understand how they work, why they work and how fast they work.
- Invariably in Research, you will come up with new optimization problems which you might need to implement custom algorithms or atleast the loss functions.
- Over 70% of the projects I worked at Microsoft involved new optimization problems (and sometimes new algorithms)



# Why take this Course?

- Countless number of ML libraries available which implement all kinds of optimization algorithms (Tensorflow, pytorch, scipy, sklearn, Vowpal Wabbit, ...)
- **This course** will give you the expertise to look inside these algorithms, understand how they work, why they work and how fast they work.
- Invariably in Research, you will come up with new optimization problems which you might need to implement custom algorithms or atleast the loss functions.
- Over 70% of the projects I worked at Microsoft involved new optimization problems (and sometimes new algorithms)
- Even if you don't implement new algorithms, you will have a better idea of which algorithm to use in which scenario.



# About Me

- Undergrad (2007 - 2011) from IIT Bombay → Masters + Ph.D at UW Seattle (2011 - 2015) → PostDoc at UW Seattle (2015-2016) → ML Scientist & Researcher at Microsoft (2016 - 2019) → Assistant Professor at UT Dallas (2020 - )



# About Me

- Undergrad (2007 - 2011) from IIT Bombay → Masters + Ph.D at UW Seattle (2011 - 2015) → PostDoc at UW Seattle (2015-2016) → ML Scientist & Researcher at Microsoft (2016 - 2019) → Assistant Professor at UT Dallas (2020 - )
- Spent good amount of time both in Academia (writing/publishing papers, conferences, giving talks, ...) and Industry (products, impact, solving real world problems)



# About Me

- Undergrad (2007 - 2011) from IIT Bombay → Masters + Ph.D at UW Seattle (2011 - 2015) → PostDoc at UW Seattle (2015-2016) → ML Scientist & Researcher at Microsoft (2016 - 2019) → Assistant Professor at UT Dallas (2020 - )
- Spent good amount of time both in Academia (writing/publishing papers, conferences, giving talks, ...) and Industry (products, impact, solving real world problems)
- Early portion of my Ph.D (proving theorems, theoretical results) → later part of Ph.D + Industry (practical aspects of ML, solving problems, getting algorithms to work on real data)



# Philosophy of this Course

- The spirit of this course is best summarized by the quote of Thomas Cover: *Theory is only the first term of the Taylor's series of Practice*



# Philosophy of this Course

- The spirit of this course is best summarized by the quote of Thomas Cover: *Theory is only the first term of the Taylor's series of Practice*
- This course will focus on mainly the algorithmic aspects of optimization (both continuous and discrete optimization) and not so much on the modeling.



# Philosophy of this Course

- The spirit of this course is best summarized by the quote of Thomas Cover: *Theory is only the first term of the Taylor's series of Practice*
- This course will focus on mainly the algorithmic aspects of optimization (both continuous and discrete optimization) and not so much on the modeling.
- Give a flavor of the proofs and proof techniques but will try to not make this course heavily theoretical.





# Philosophy of this Course

- The spirit of this course is best summarized by the quote of Thomas Cover: *Theory is only the first term of the Taylor's series of Practice*
- This course will focus on mainly the algorithmic aspects of optimization (both continuous and discrete optimization) and not so much on the modeling.
- Give a flavor of the proofs and proof techniques but will try to not make this course heavily theoretical.
- Focus extensively on implementation aspects and as a part of assignments, we will implement many ML loss functions and algorithms.



# Why did you enroll for this Course?

Lets hear from a few of you why you took this course...



- Class is Every Mon/Wed 4 PM - 5:15 PM
- Venue: **AD 2.232**
- Office Hours: Wednesdays between 3 to 3:50 PM. Additional Office Hours (by appointment): Mondays 3 to 3:50 PM.
- TA: TBD



# Prerequisites

- Calculus Courses
- Basic Linear Algebra: Matrices, Vectors & Matrix Algebra
- Machine Learning (either an undergraduate or graduate ML course)
- Algorithms course (either in undergraduate or graduate version)
- Bonus: Linear Programming or Numerical Optimization
- It is fine if you are concurrently taking the ML course (for e.g. first year students). However you will need to put in extra effort.



# How is this course different from Machine Learning or Convex Optimization

- Many of you have taken the Machine learning course at UTD.



# How is this course different from Machine Learning or Convex Optimization

- Many of you have taken the Machine learning course at UTD.
- Similarly, many Universities offer a Convex Optimization course



# How is this course different from Machine Learning or Convex Optimization

- Many of you have taken the Machine learning course at UTD.
- Similarly, many Universities offer a Convex Optimization course
- Convex Optimization courses are mostly theoretical and focus heavily on the algorithmic/theoretical aspects



# How is this course different from Machine Learning or Convex Optimization

- Many of you have taken the Machine learning course at UTD.
- Similarly, many Universities offer a Convex Optimization course
- Convex Optimization courses are mostly theoretical and focus heavily on the algorithmic/theoretical aspects
- This course will focus:





# How is this course different from Machine Learning or Convex Optimization

- Many of you have taken the Machine learning course at UTD.
- Similarly, many Universities offer a Convex Optimization course
- Convex Optimization courses are mostly theoretical and focus heavily on the algorithmic/theoretical aspects
- This course will focus:
  - More on the *optimization algorithms* and not so much on the modeling and statistical aspects



# How is this course different from Machine Learning or Convex Optimization

- Many of you have taken the Machine learning course at UTD.
- Similarly, many Universities offer a Convex Optimization course
- Convex Optimization courses are mostly theoretical and focus heavily on the algorithmic/theoretical aspects
- This course will focus:
  - More on the *optimization algorithms* and not so much on the modeling and statistical aspects
  - Machine learning as an application domain and also emphasize the practical/implementation of algorithms,



# How is this course different from Machine Learning or Convex Optimization

- Many of you have taken the Machine learning course at UTD.
- Similarly, many Universities offer a Convex Optimization course
- Convex Optimization courses are mostly theoretical and focus heavily on the algorithmic/theoretical aspects
- This course will focus:
  - More on the *optimization algorithms* and not so much on the modeling and statistical aspects
  - Machine learning as an application domain and also emphasize the practical/implementation of algorithms,
  - A flavor of the proof techniques without going too mathematical



# How is this course different from Machine Learning or Convex Optimization

- Many of you have taken the Machine learning course at UTD.
- Similarly, many Universities offer a Convex Optimization course
- Convex Optimization courses are mostly theoretical and focus heavily on the algorithmic/theoretical aspects
- This course will focus:
  - More on the *optimization algorithms* and not so much on the modeling and statistical aspects
  - Machine learning as an application domain and also emphasize the practical/implementation of algorithms,
  - A flavor of the proof techniques without going too mathematical
- One of the first courses to provide an overview of the entire spectrum of optimization (discrete & continuous) for ML.



# First Half of this Course: Continuous optimization

- Basics of Continuous Optimization
- Convexity
- Gradient Descent
- Projected GD
- Subgradient Descent
- Accelerated Gradient Descent
- Newton & Quasi Newton
- Duality: Lagrange, Fenchel
- Coordinate Descent
- Frank Wolfe
- Optimization in Practice



# Second Half of this Course: Discrete optimization

- Linear Cost Problems
- Matroids, Spanning Trees
- s-t paths, s-t cuts
- Matchings
- Covers (Set Covers, Vertex Covers, Edge Covers)
- Optimal Transport (if time permits)
- Non-Linear Discrete Optimization
- Submodular Functions
- Submodularity and Convexity
- Submodular Minimization
- Submodular Maximization
- Optimization in Practice



# Relevant Books for this Course

- Convex Optimization: Algorithms and Complexity, by Sébastien Bubeck
- Convex Optimization, Stephen Boyd and Lieven Vandenberghe
- Introductory Lectures on Convex Optimization, Yurii Nesterov
- A Course in Combinatorial Optimization, Alexander Schrijver
- Learning with Submodular Functions: A Convex Optimization Perspective, Francis Bach
- Zhang, Lipton, Li and Smola, Dive into Deep Learning (<http://d2l.ai/>)
- Schrijver, Alexander, Combinatorial optimization: polyhedra and efficiency, Vol. 24. Springer Science & Business Media, 2003.
- Fujishige, Satoru. Submodular functions and optimization. Vol. 58. Elsevier, 2005.



- **Assignments (50%):**
  - Mix of theory and practice.
  - Simple Theory questions, e.g. proving certain loss functions are convex/non-convex, computing gradients etc.
  - Practical implementation in python
  - One assignment every two to three weeks. Around 10 assignments in total.





# Evaluation

- **Assignments (50%):**

- Mix of theory and practice.
- Simple Theory questions, e.g. proving certain loss functions are convex/non-convex, computing gradients etc.
- Practical implementation in python
- One assignment every two to three weeks. Around 10 assignments in total.

- **Project: (50%):**

- Examples: Picking a certain real world problem, modeling the problem and optimizing the loss function with different algorithms and drawing insights on its performance.
- Alternatively, also be a survey on a particular class of optimization algorithms and theoretical results
- Could also be creating a toolkit (python/c++) of the various optimization/learning algorithms.
- More on Project Later...



# Attendance Policy

- Attendance is compulsory. Since this is a new course, there is not really an official textbook and the classes will be the best source of information!
- Absences due to legitimate reasons/medical etc. are fine
- Attendance policy: Two consecutive classes missed then no penalty, three consecutive classes missed then one letter downgrade, Four consecutive classes missed then an outright Fail.



# Continuous Optimization in Machine Learning

- Continuous Optimization often appears as *relaxations* of empirical risk minimization problems.
- **Supervised Learning**: Logistic Regression, Least Squares, Support Vector Machines, Deep Models
- **Unsupervised Learning**: k-Means Clustering, Principal Component Analysis
- **Contextual Bandits and Reinforcement Learning**: Soft-Max Estimators, Policy Exponential Models
- **Recommender Systems**: Matrix Completion, Non-Negative Matrix Factorization, Collaborative Filtering



# Reading Material for Linear Algebra

Good Reading Material on the basics of Linear Algebra from my Colleague Prof. Ganesh Ramakrishnan:

<https://www.cse.iitb.ac.in/~cs709/notes/LinearAlgebra.pdf>



# Notation: Vectors and Matrices

- For the most part through this course, we will use  $n$  as the number of training instances and  $m$  as the number of features.



# Notation: Vectors and Matrices

- For the most part through this course, we will use  $n$  as the number of training instances and  $m$  as the number of features.
- Denote  $x_i \in \mathbf{R}^m$  as a  $m$ -dimensional vector (feature vector).



# Notation: Vectors and Matrices

- For the most part through this course, we will use  $n$  as the number of training instances and  $m$  as the number of features.
- Denote  $x_i \in \mathbf{R}^m$  as a  $m$ -dimensional vector (feature vector).
- We shall denote the  $j$ th feature of  $x_i$  as  $x_i[j]$ .



# Notation: Vectors and Matrices

- For the most part through this course, we will use  $n$  as the number of training instances and  $m$  as the number of features.
- Denote  $x_i \in \mathbf{R}^m$  as a  $m$ -dimensional vector (feature vector).
- We shall denote the  $j$ th feature of  $x_i$  as  $x_i[j]$ .
- $y_i$  is the Label as the  $i$ th instance.





# Notation: Vectors and Matrices

- For the most part through this course, we will use  $n$  as the number of training instances and  $m$  as the number of features.
- Denote  $x_i \in \mathbf{R}^m$  as a  $m$ -dimensional vector (feature vector).
- We shall denote the  $j$ th feature of  $x_i$  as  $x_i[j]$ .
- $y_i$  is the Label as the  $i$ th instance.
- Given two vectors  $w, x \in \mathbf{R}^m$ , define the inner product  $\langle w, x \rangle = \sum_{j=1}^m x[j]w[j]$ .



# Notation: Vectors and Matrices

- For the most part through this course, we will use  $n$  as the number of training instances and  $m$  as the number of features.
- Denote  $x_i \in \mathbf{R}^m$  as a  $m$ -dimensional vector (feature vector).
- We shall denote the  $j$ th feature of  $x_i$  as  $x_i[j]$ .
- $y_i$  is the Label as the  $i$ th instance.
- Given two vectors  $w, x \in \mathbf{R}^m$ , define the inner product  $\langle w, x \rangle = \sum_{j=1}^m x[j]w[j]$ .
- The L1 Norm  $\|w\|_1 = \sum_{i=1}^m |w[i]|$



# Notation: Vectors and Matrices

- For the most part through this course, we will use  $n$  as the number of training instances and  $m$  as the number of features.
- Denote  $x_i \in \mathbf{R}^m$  as a  $m$ -dimensional vector (feature vector).
- We shall denote the  $j$ th feature of  $x_i$  as  $x_i[j]$ .
- $y_i$  is the Label as the  $i$ th instance.
- Given two vectors  $w, x \in \mathbf{R}^m$ , define the inner product  $\langle w, x \rangle = \sum_{j=1}^m x[j]w[j]$ .
- The L1 Norm  $\|w\|_1 = \sum_{i=1}^m |w[i]|$
- The Squared L2 Norm  $\|w\|_2^2 = \sum_{i=1}^m |w[i]|^2$



# Notation: Vectors and Matrices

- For the most part through this course, we will use  $n$  as the number of training instances and  $m$  as the number of features.
- Denote  $x_i \in \mathbf{R}^m$  as a  $m$ -dimensional vector (feature vector).
- We shall denote the  $j$ th feature of  $x_i$  as  $x_i[j]$ .
- $y_i$  is the Label as the  $i$ th instance.
- Given two vectors  $w, x \in \mathbf{R}^m$ , define the inner product  $\langle w, x \rangle = \sum_{j=1}^m x[j]w[j]$ .
- The L1 Norm  $\|w\|_1 = \sum_{i=1}^m |w[i]|$
- The Squared L2 Norm  $\|w\|_2^2 = \sum_{i=1}^m |w[i]|^2$
- Sometimes we shall refer to L2 as the default norm:  $\|w\| = \|w\|_2$ .



# What is a Norm

A function  $f : \mathbb{R}^n \Rightarrow \mathbb{R}_+$  is a norm if:

- $f(x) \geq 0, \forall x$
- $f(x + y) \leq f(x) + f(y)$  [Triangle Inequality]
- $f(\alpha x) = \alpha f(x)$
- $f(x) = 0 \Rightarrow x = 0$

Can you make sure the L1 and L2 norms defined above satisfy all four conditions?



# Matrix-Vector Product

- If  $A \in \mathbb{R}^{m \times n}$  and  $x \in \mathbb{R}^n$ , we can define  $y = Ax$  where  $y \in \mathbb{R}^m$  is a  $m$  dimensional vector.



# Matrix-Vector Product

- If  $A \in \mathbb{R}^{m \times n}$  and  $x \in \mathbb{R}^n$ , we can define  $y = Ax$  where  $y \in \mathbb{R}^m$  is a  $m$  dimensional vector.
- Matrix vector product is defined as below:

$$Ax = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n \\ a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n \\ \vdots \\ a_{m1}x_1 + a_{m2}x_2 + \dots + a_{mn}x_n \end{bmatrix}$$



# Matrix-Vector Product Example

For example, if

$$A = \begin{bmatrix} 1 & -1 & 2 \\ 0 & -3 & 1 \end{bmatrix}$$

and  $\mathbf{x} = (2, 1, 0)$ , then

$$\begin{aligned} A\mathbf{x} &= \begin{bmatrix} 1 & -1 & 2 \\ 0 & -3 & 1 \end{bmatrix} \begin{bmatrix} 2 \\ 1 \\ 0 \end{bmatrix} \\ &= \begin{bmatrix} 2 \cdot 1 - 1 \cdot 1 + 0 \cdot 2 \\ 2 \cdot 0 - 1 \cdot 3 + 0 \cdot 1 \end{bmatrix} \\ &= \begin{bmatrix} 1 \\ -3 \end{bmatrix}. \end{aligned}$$



# Matrix-Matrix Product

- Matrix product between Matrices  $A$  and  $B$  is defined only when the number of columns in  $A$  equals number of rows in  $B$ .



# Matrix-Matrix Product

- Matrix product between Matrices  $A$  and  $B$  is defined only when the number of columns in  $A$  equals number of rows in  $B$ .
- If  $A$  is a  $m \times n$  Matrix,  $B$  is a  $n \times p$  Matrix, the product Matrix  $C = AB$  is a  $m \times p$  Matrix.



# Matrix-Matrix Product

- Matrix product between Matrices  $A$  and  $B$  is defined only when the number of columns in  $A$  equals number of rows in  $B$ .
- If  $A$  is a  $m \times n$  Matrix,  $B$  is a  $n \times p$  Matrix, the product Matrix  $C = AB$  is a  $m \times p$  Matrix.
- We can view the Matrix-Matrix product as stacked Matrix-vector products by viewing  $B$  as a set of  $p$   $n \times 1$  vectors lined up.



# Matrix-Matrix Product

- Matrix product between Matrices  $A$  and  $B$  is defined only when the number of columns in  $A$  equals number of rows in  $B$ .
- If  $A$  is a  $m \times n$  Matrix,  $B$  is a  $n \times p$  Matrix, the product Matrix  $C = AB$  is a  $m \times p$  Matrix.
- We can view the Matrix-Matrix product as stacked Matrix-vector products by viewing  $B$  as a set of  $p$   $n \times 1$  vectors lined up.

$$\begin{bmatrix} b_{11} & b_{12} & \dots & b_{1p} \\ b_{21} & b_{22} & \dots & b_{2p} \\ \vdots & \vdots & \vdots & \vdots \\ b_{n1} & b_{n2} & \dots & b_{np} \end{bmatrix} = \begin{bmatrix} \begin{bmatrix} b_{11} \\ b_{21} \\ \vdots \\ b_{n1} \end{bmatrix} & \begin{bmatrix} b_{12} \\ b_{22} \\ \vdots \\ b_{n2} \end{bmatrix} & \dots & \begin{bmatrix} b_{1p} \\ b_{2p} \\ \vdots \\ b_{np} \end{bmatrix} \end{bmatrix}$$



# Matrix-Matrix Product

$$\begin{bmatrix} b_{11} & b_{12} & \dots & b_{1p} \\ b_{21} & b_{22} & \dots & b_{2p} \\ \vdots & \vdots & \vdots & \vdots \\ b_{n1} & b_{n2} & \dots & b_{np} \end{bmatrix} = \begin{bmatrix} b_{11} \\ b_{21} \\ \vdots \\ b_{n1} \end{bmatrix} \begin{bmatrix} b_{12} \\ b_{22} \\ \vdots \\ b_{n2} \end{bmatrix} \dots \begin{bmatrix} b_{1p} \\ b_{2p} \\ \vdots \\ b_{np} \end{bmatrix}$$

- We can view the Matrix-Matrix product as stacked Matrix-vector products by viewing  $B$  as a set of  $p$   $n \times 1$  vectors lined up.



# Matrix-Matrix Product

$$\begin{bmatrix} b_{11} & b_{12} & \dots & b_{1p} \\ b_{21} & b_{22} & \dots & b_{2p} \\ \vdots & \vdots & \vdots & \vdots \\ b_{n1} & b_{n2} & \dots & b_{np} \end{bmatrix} = \begin{bmatrix} b_{11} \\ b_{21} \\ \vdots \\ b_{n1} \end{bmatrix} \begin{bmatrix} b_{12} \\ b_{22} \\ \vdots \\ b_{n2} \end{bmatrix} \dots \begin{bmatrix} b_{1p} \\ b_{2p} \\ \vdots \\ b_{np} \end{bmatrix}$$

- We can view the Matrix-Matrix product as stacked Matrix-vector products by viewing  $B$  as a set of  $p$   $n \times 1$  vectors lined up.
- Let  $b_1, \dots, b_p$  be these  $p$  column vectors.



# Matrix-Matrix Product

$$\begin{bmatrix} b_{11} & b_{12} & \dots & b_{1p} \\ b_{21} & b_{22} & \dots & b_{2p} \\ \vdots & \vdots & \vdots & \vdots \\ b_{n1} & b_{n2} & \dots & b_{np} \end{bmatrix} = \begin{bmatrix} b_{11} \\ b_{21} \\ \vdots \\ b_{n1} \end{bmatrix} \begin{bmatrix} b_{12} \\ b_{22} \\ \vdots \\ b_{n2} \end{bmatrix} \dots \begin{bmatrix} b_{1p} \\ b_{2p} \\ \vdots \\ b_{np} \end{bmatrix}$$

- We can view the Matrix-Matrix product as stacked Matrix-vector products by viewing  $B$  as a set of  $p$   $n \times 1$  vectors lined up.
- Let  $b_1, \dots, b_p$  be these  $p$  column vectors.
- Note that  $c_i = Ab_i$  is a  $m \times 1$  column vector.



# Matrix-Matrix Product

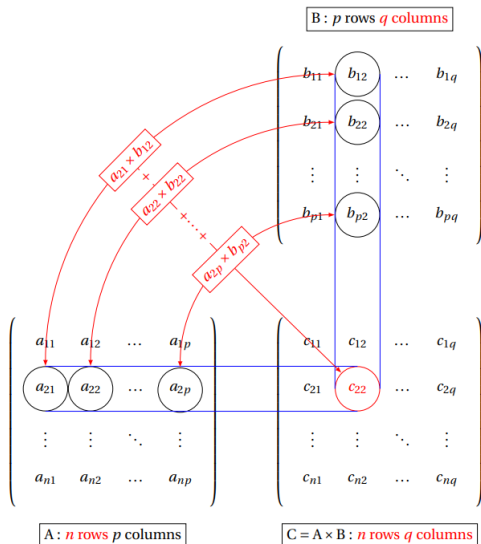
$$\begin{bmatrix} b_{11} & b_{12} & \dots & b_{1p} \\ b_{21} & b_{22} & \dots & b_{2p} \\ \vdots & \vdots & \vdots & \vdots \\ b_{n1} & b_{n2} & \dots & b_{np} \end{bmatrix} = \begin{bmatrix} b_{11} \\ b_{21} \\ \vdots \\ b_{n1} \end{bmatrix} \begin{bmatrix} b_{12} \\ b_{22} \\ \vdots \\ b_{n2} \end{bmatrix} \dots \begin{bmatrix} b_{1p} \\ b_{2p} \\ \vdots \\ b_{np} \end{bmatrix}$$

- We can view the Matrix-Matrix product as stacked Matrix-vector products by viewing  $B$  as a set of  $p$   $n \times 1$  vectors lined up.
- Let  $b_1, \dots, b_p$  be these  $p$  column vectors.
- Note that  $c_i = Ab_i$  is a  $m \times 1$  column vector.
- Hence  $AB = [Ab_1 \ Ab_2 \ \dots \ Ab_p] = [c_1 \ \dots \ c_p] = C$  is a  $m \times p$  Matrix





# Matrix-Matrix Product Illustration



# Matrix-Matrix Product Example

Compute  $BC$ , where

$$B = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \quad \text{and} \quad C = \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix}.$$

$$\begin{aligned} BC &= \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix} \\ &= \begin{bmatrix} 1 \cdot 1 + 2 \cdot 3 + 3 \cdot 5 & 1 \cdot 2 + 2 \cdot 4 + 3 \cdot 6 \\ 4 \cdot 1 + 5 \cdot 3 + 6 \cdot 5 & 4 \cdot 2 + 5 \cdot 4 + 6 \cdot 6 \end{bmatrix} \\ &= \begin{bmatrix} 22 & 28 \\ 49 & 64 \end{bmatrix} \end{aligned}$$

# Application 1: Supervised Learning

- **Data:** Given training examples  $\{(x_1, y_1), \dots, (x_n, y_n)\}$  where  $x_i \in \mathbf{R}^m$  is the feature vectors and  $y_i$  is the label.
- **Applications:** Several different models depending on the applications:
  - **Email Spam Filtering:** Features are words, phrases, regexps in the email, Label is "+1" for Spam, "0" for Not Spam.
  - **Handwritten Digit Recognition:** Features are Images of Images, Label is the Digit (say between "0" to "9").
  - **Housing price Prediction:** Features are House properties (square footage, # Bedrooms/Bathrooms, Location, ...) and Label is the Cost (continuous variable).



# Supervised Learning: Modeling

- **Data:** Given training examples  $\{(x_1, y_1), \dots, (x_n, y_n)\}$  where  $x_i \in \mathbf{R}^m$  is the feature vectors and  $y_i$  is the label.



# Supervised Learning: Modeling

- **Data:** Given training examples  $\{(x_1, y_1), \dots, (x_n, y_n)\}$  where  $x_i \in \mathbf{R}^m$  is the feature vectors and  $y_i$  is the label.
- **Model:** Denote the Model by  $F_\theta(x)$  with  $\theta$  being the parameters of the model. Model examples:  $F_\theta(x) = \theta^T x$  as a simple linear model. Deep Models are recursive functions:

$$F_{\theta_1, \theta_2, \dots, \theta_l}(x) = f_1(\theta_1^T f_2(\dots \theta_{l-1}^T f_l(\theta_l^T x)))$$



# Supervised Learning: Modeling

- **Data:** Given training examples  $\{(x_1, y_1), \dots, (x_n, y_n)\}$  where  $x_i \in \mathbf{R}^m$  is the feature vectors and  $y_i$  is the label.
- **Model:** Denote the Model by  $F_\theta(x)$  with  $\theta$  being the parameters of the model. Model examples:  $F_\theta(x) = \theta^T x$  as a simple linear model. Deep Models are recursive functions:

$$F_{\theta_1, \theta_2, \dots, \theta_l}(x) = f_1(\theta_1^T f_2(\dots \theta_{l-1}^T f_l(\theta_l^T x)))$$

- **Loss Functions:** The Loss Function  $L$  tries to measure the *distance* between  $F_\theta(x_i)$  and  $y_i$ .



# Supervised Learning: Optimization Problem

- "Loss plus Regularizer" Framework:

$$\min_{\theta} G(\theta) = \sum_{i=1}^n L(F_{\theta}(x_i), y_i) + \lambda \Omega(\theta)$$



# Supervised Learning: Optimization Problem

- "Loss plus Regularizer" Framework:

$$\min_{\theta} G(\theta) = \sum_{i=1}^n L(F_{\theta}(x_i), y_i) + \lambda \Omega(\theta)$$

- $L$ : Loss function,  $\Omega$ : Regularizer. Example  $F_{\theta}(x) = \theta^T x$





# Supervised Learning: Optimization Problem

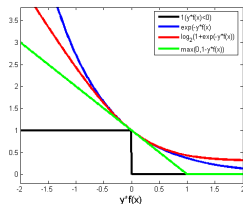
- "Loss plus Regularizer" Framework:

$$\min_{\theta} G(\theta) = \sum_{i=1}^n L(F_{\theta}(x_i), y_i) + \lambda \Omega(\theta)$$

- $L$ : Loss function,  $\Omega$ : Regularizer. Example  $F_{\theta}(x) = \theta^T x$

- Examples of  $L$ :

- Logistic Loss:  $\log(1 + \exp(-y_i F_{\theta}(x_i)))$
- Hinge Loss:  $\max\{0, 1 - y_i F_{\theta}(x_i)\}$
- Softmax Loss:  
 $-F_{\theta_{y_i}}(x_i) + \log(\sum_{c=1}^k \exp(F_{\theta_c}(x_i)))$
- Absolute Error:  $|F_{\theta}(x_i) - y_i|$
- Least Squares:  $(F_{\theta}(x_i) - y_i)^2$



# Supervised Learning: Optimization Problem

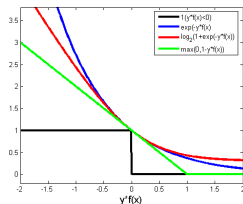
- "Loss plus Regularizer" Framework:

$$\min_{\theta} G(\theta) = \sum_{i=1}^n L(F_{\theta}(x_i), y_i) + \lambda \Omega(\theta)$$

- $L$ : Loss function,  $\Omega$ : Regularizer. Example  $F_{\theta}(x) = \theta^T x$

- Examples of  $L$ :

- Logistic Loss:  $\log(1 + \exp(-y_i F_{\theta}(x_i)))$
- Hinge Loss:  $\max\{0, 1 - y_i F_{\theta}(x_i)\}$
- Softmax Loss:  
 $-F_{\theta_{y_i}}(x_i) + \log(\sum_{c=1}^k \exp(F_{\theta_c}(x_i)))$
- Absolute Error:  $|F_{\theta}(x_i) - y_i|$
- Least Squares:  $(F_{\theta}(x_i) - y_i)^2$



- Examples of  $\Omega$ :

- L1 Regularizer:  $\|\theta\|_1 = \sum_{i=1}^m |\theta[i]|$
- L2 Regularizer:  $\|\theta\|_2^2 = \sum_{i=1}^m \theta[i]^2$



# Some Concrete Supervised Learning Instances

- L1 Regularized Logistic Regression:

$$\min_{\theta} \sum_{i=1}^n \log(1 + \exp(-y_i F_{\theta}(x_i))) + \lambda \|\theta\|_1$$



# Some Concrete Supervised Learning Instances

- L1 Regularized Logistic Regression:  
$$\min_{\theta} \sum_{i=1}^n \log(1 + \exp(-y_i F_{\theta}(x_i))) + \lambda \|\theta\|_1$$
- L2 Regularized Logistic Regression:  
$$\min_{\theta} \sum_{i=1}^n \log(1 + \exp(-y_i F_{\theta}(x_i))) + \lambda \|\theta\|_2^2$$



# Some Concrete Supervised Learning Instances

- L1 Regularized Logistic Regression:  
$$\min_{\theta} \sum_{i=1}^n \log(1 + \exp(-y_i F_{\theta}(x_i))) + \lambda \|\theta\|_1$$
- L2 Regularized Logistic Regression:  
$$\min_{\theta} \sum_{i=1}^n \log(1 + \exp(-y_i F_{\theta}(x_i))) + \lambda \|\theta\|_2^2$$
- L2 Regularized SVMs:  $\min_{\theta} \sum_{i=1}^n \max\{0, 1 - y_i F_{\theta}(x_i)\} + \lambda \|\theta\|_2^2$



# Some Concrete Supervised Learning Instances

- L1 Regularized Logistic Regression:

$$\min_{\theta} \sum_{i=1}^n \log(1 + \exp(-y_i F_{\theta}(x_i))) + \lambda \|\theta\|_1$$

- L2 Regularized Logistic Regression:

$$\min_{\theta} \sum_{i=1}^n \log(1 + \exp(-y_i F_{\theta}(x_i))) + \lambda \|\theta\|_2^2$$

- L2 Regularized SVMs:  $\min_{\theta} \sum_{i=1}^n \max\{0, 1 - y_i F_{\theta}(x_i)\} + \lambda \|\theta\|_2^2$

- L2 Regularized Multi-class Logistic Regression:

$$\min_{\theta_1, \dots, \theta_k} \sum_{i=1}^n \{-F_{\theta_{y_i}}(x_i) + \log(\sum_{c=1}^k \exp(F_{\theta_c}(x_i)))\} + \sum_{i=1}^c \lambda \|\theta_i\|_2^2$$



# Some Concrete Supervised Learning Instances

- L1 Regularized Logistic Regression:  
$$\min_{\theta} \sum_{i=1}^n \log(1 + \exp(-y_i F_{\theta}(x_i))) + \lambda \|\theta\|_1$$
- L2 Regularized Logistic Regression:  
$$\min_{\theta} \sum_{i=1}^n \log(1 + \exp(-y_i F_{\theta}(x_i))) + \lambda \|\theta\|_2^2$$
- L2 Regularized SVMs:  $\min_{\theta} \sum_{i=1}^n \max\{0, 1 - y_i F_{\theta}(x_i)\} + \lambda \|\theta\|_2^2$
- L2 Regularized Multi-class Logistic Regression:  
$$\min_{\theta_1, \dots, \theta_k} \sum_{i=1}^n \{-F_{\theta_{y_i}}(x_i) + \log(\sum_{c=1}^k \exp(F_{\theta_c}(x_i)))\} + \sum_{i=1}^c \lambda \|\theta_i\|_2^2$$
- L1 Regularized Least Squares (Lasso):  
$$\min_{\theta} \sum_{i=1}^n (F_{\theta}(x_i) - y_i)^2 + \lambda \|\theta\|_1$$



# Some Concrete Supervised Learning Instances

- L1 Regularized Logistic Regression:  
$$\min_{\theta} \sum_{i=1}^n \log(1 + \exp(-y_i F_{\theta}(x_i))) + \lambda \|\theta\|_1$$
- L2 Regularized Logistic Regression:  
$$\min_{\theta} \sum_{i=1}^n \log(1 + \exp(-y_i F_{\theta}(x_i))) + \lambda \|\theta\|_2^2$$
- L2 Regularized SVMs:  $\min_{\theta} \sum_{i=1}^n \max\{0, 1 - y_i F_{\theta}(x_i)\} + \lambda \|\theta\|_2^2$
- L2 Regularized Multi-class Logistic Regression:  
$$\min_{\theta_1, \dots, \theta_k} \sum_{i=1}^n \{-F_{\theta_{y_i}}(x_i) + \log(\sum_{c=1}^k \exp(F_{\theta_c}(x_i)))\} + \sum_{i=1}^c \lambda \|\theta_i\|_2^2$$
- L1 Regularized Least Squares (Lasso):  
$$\min_{\theta} \sum_{i=1}^n (F_{\theta}(x_i) - y_i)^2 + \lambda \|\theta\|_1$$
- L2 Regularized Least Squares (Ridge):  
$$\min_{\theta} \sum_{i=1}^n (F_{\theta}(x_i) - y_i)^2 + \lambda \|\theta\|_2^2$$





# Application 2: Clustering

- This is an instance of unsupervised learning.



## Application 2: Clustering

- This is an instance of unsupervised learning.
- **Data:** Given unsupervised data  $\{x_1, x_2, \dots, x_n\}$  where  $x_i \in \mathbf{R}^m$  is the feature vectors.



# Application 2: Clustering

- This is an instance of unsupervised learning.
- **Data:** Given unsupervised data  $\{x_1, x_2, \dots, x_n\}$  where  $x_i \in \mathbf{R}^m$  is the feature vectors.
- **Goal:** Find clusters (sets)  $C_1, C_2, \dots, C_k$  with each cluster consisting of *similar* instances. Denote  $V = \{1, \dots, n\}$ . Then  $\cup_{i=1}^k C_i = V$ .



# Application 2: Clustering

- This is an instance of unsupervised learning.
- **Data:** Given unsupervised data  $\{x_1, x_2, \dots, x_n\}$  where  $x_i \in \mathbf{R}^m$  is the feature vectors.
- **Goal:** Find clusters (sets)  $C_1, C_2, \dots, C_k$  with each cluster consisting of *similar* instances. Denote  $V = \{1, \dots, n\}$ . Then  $\cup_{i=1}^k C_i = V$ .
- **Optimization Problem:** The k-means optimization problem is:

$$\min_{C_1, C_2, \dots, C_k} \sum_{i=1}^k \sum_{x \in C_i} |x - 1/|C_i| \sum_{x_j \in C_i} x_j|_2^2$$



# Application 2: Clustering

- This is an instance of unsupervised learning.
- **Data:** Given unsupervised data  $\{x_1, x_2, \dots, x_n\}$  where  $x_i \in \mathbf{R}^m$  is the feature vectors.
- **Goal:** Find clusters (sets)  $C_1, C_2, \dots, C_k$  with each cluster consisting of *similar* instances. Denote  $V = \{1, \dots, n\}$ . Then  $\cup_{i=1}^k C_i = V$ .
- **Optimization Problem:** The k-means optimization problem is:

$$\min_{C_1, C_2, \dots, C_k} \sum_{i=1}^k \sum_{x \in C_i} |x - \mu_i|^2$$

- This problem can actually be viewed as a joint discrete and continuous problem.



# Application 3: Principal Component Analysis

- This is an instance of unsupervised learning.



# Application 3: Principal Component Analysis

- This is an instance of unsupervised learning.
- **Data:** Given unsupervised data  $\{x_1, x_2, \dots, x_n\}$  where  $x_i \in \mathbf{R}^m$  is the feature vectors.



# Application 3: Principal Component Analysis

- This is an instance of unsupervised learning.
- **Data:** Given unsupervised data  $\{x_1, x_2, \dots, x_n\}$  where  $x_i \in \mathbf{R}^m$  is the feature vectors.
- **Goal:** Find compressors  $V \in \mathbf{R}^{k \times m}$  (and correspondingly decompresses  $U \in \mathbf{R}^{m \times k}$ ) such that  $x_i$  is *close* to  $UVx_i$ .





# Application 3: Principal Component Analysis

- This is an instance of unsupervised learning.
- **Data:** Given unsupervised data  $\{x_1, x_2, \dots, x_n\}$  where  $x_i \in \mathbf{R}^m$  is the feature vectors.
- **Goal:** Find compressors  $V \in \mathbf{R}^{k \times m}$  (and correspondingly decompresses  $U \in \mathbf{R}^{m \times k}$ ) such that  $x_i$  is *close* to  $UVx_i$ .
- It turns out the compressor must satisfy  $V = U^T$  such that  $U^T U = I$ .



# Application 3: Principal Component Analysis

- This is an instance of unsupervised learning.
- **Data:** Given unsupervised data  $\{x_1, x_2, \dots, x_n\}$  where  $x_i \in \mathbf{R}^m$  is the feature vectors.
- **Goal:** Find compressors  $V \in \mathbf{R}^{k \times m}$  (and correspondingly decompresses  $U \in \mathbf{R}^{m \times k}$ ) such that  $x_i$  is *close* to  $UVx_i$ .
- It turns out the compressor must satisfy  $V = U^T$  such that  $U^T U = I$ .
- **Optimization Problem:** The PCA optimization problem is:

$$\min_{U: U^T U = I} \sum_{i=1}^n \|x_i - UU^T x_i\|_2^2$$



# Application 4: Matrix Completion

- **Data:** Given observations  $y_1, \dots, y_n$ , such that each  $y_j = A_j(X)$  where  $A_j$  could be a single element or a combination of elements in  $X \in \mathbb{R}^{m \times n}$ . Consider for example  $X$  being product recommendation matrix.



## Application 4: Matrix Completion

- **Data:** Given observations  $y_1, \dots, y_n$ , such that each  $y_j = A_j(X)$  where  $A_j$  could be a single element or a combination of elements in  $X \in \mathbb{R}^{m \times n}$ . Consider for example  $X$  being product recommendation matrix.
- **Goal:** Find the *simplest* matrix  $X$  s.t  $A_j(X) \approx y_j, \forall j \in 1, \dots, n$



## Application 4: Matrix Completion

- **Data:** Given observations  $y_1, \dots, y_n$ , such that each  $y_j = A_j(X)$  where  $A_j$  could be a single element or a combination of elements in  $X \in \mathbb{R}^{m \times n}$ . Consider for example  $X$  being product recommendation matrix.
- **Goal:** Find the *simplest* matrix  $X$  s.t  $A_j(X) \approx y_j, \forall j \in 1, \dots, n$
- **Optimization Problem:** Matrix completion optimization problem is:

$$\min_X \sum_{j=1}^n \|y_j - A_j(X)\|_2^2 + \|X\|_*$$

(The nuclear norm tries to ensure the Matrix  $X$  is low-rank)



## Application 4: Matrix Completion

- **Data:** Given observations  $y_1, \dots, y_n$ , such that each  $y_j = A_j(X)$  where  $A_j$  could be a single element or a combination of elements in  $X \in \mathbb{R}^{m \times n}$ . Consider for example  $X$  being product recommendation matrix.
- **Goal:** Find the *simplest* matrix  $X$  s.t  $A_j(X) \approx y_j, \forall j \in 1, \dots, n$
- **Optimization Problem:** Matrix completion optimization problem is:

$$\min_X \sum_{j=1}^n \|y_j - A_j(X)\|_2^2 + \|X\|_*$$

(The nuclear norm tries to ensure the Matrix  $X$  is low-rank)

- Another way is to explicitly model this is by assuming  $X = LR$  where  $L \in \mathbb{R}^{m \times k}$  and  $R \in \mathbb{R}^{k \times n}$  (and hence  $X$  is rank  $r$ ), and optimize for  $L$  and  $R$ .



# Application 5: Low Rank and Non Negative Matrix Factorization

- **Goal:** Find low rank matrices  $L, R$  with  $L \in \mathbb{R}^{m \times k}$  and  $R \in \mathbb{R}^{k \times n}$  s.t  $A_j(LR) \approx y_j, \forall j \in 1, \dots, n$



# Application 5: Low Rank and Non Negative Matrix Factorization

- **Goal:** Find low rank matrices  $L, R$  with  $L \in \mathbb{R}^{m \times k}$  and  $R \in \mathbb{R}^{k \times n}$  s.t  $A_j(LR) \approx y_j, \forall j \in 1, \dots, n$
- **Optimization Problem:** Matrix Factorization optimization problem is:

$$\min_{L, R} \sum_{i=1}^n \|y_i - A_i(LR)\|_2^2$$

No need of matrix regularization.





# Application 5: Low Rank and Non Negative Matrix Factorization

- **Goal:** Find low rank matrices  $L, R$  with  $L \in \mathbb{R}^{m \times k}$  and  $R \in \mathbb{R}^{k \times n}$  s.t  $A_j(LR) \approx y_j, \forall j \in 1, \dots, n$
- **Optimization Problem:** Matrix Factorization optimization problem is:

$$\min_{L, R} \sum_{i=1}^n \|y_i - A_i(LR)\|_2^2$$

No need of matrix regularization.

- We can also add non-negativity constraints and this becomes non-negative Matrix Factorization:

$$\min_{L, R: L \geq 0, R \geq 0} \sum_{i=1}^n \|y_i - A_i(LR)\|_2^2$$



# Application 5: Low Rank and Non Negative Matrix Factorization

- **Goal:** Find low rank matrices  $L, R$  with  $L \in \mathbb{R}^{m \times k}$  and  $R \in \mathbb{R}^{k \times n}$  s.t  $A_j(LR) \approx y_j, \forall j \in 1, \dots, n$
- **Optimization Problem:** Matrix Factorization optimization problem is:

$$\min_{L, R} \sum_{i=1}^n \|y_i - A_i(LR)\|_2^2$$

No need of matrix regularization.

- We can also add non-negativity constraints and this becomes non-negative Matrix Factorization:

$$\min_{L, R: L \geq 0, R \geq 0} \sum_{i=1}^n \|y_i - A_i(LR)\|_2^2$$

- Sometimes  $Y$  is fully observed and we want a non-negative low rank factorization of  $Y \approx LR$ . The optimization problem is:

$$\min_{L, R: L \geq 0, R \geq 0} \|Y - LR\|_2^2.$$



# Application 6: Contextual Bandits and Learning from Logged Data

- **Scenario:** Learn from logged contextual bandit data. Example: We need to show  $k$  ads to users with each ad comprising of features (title, ad text, query etc.), and given an online policy which (with certain randomization) picks ads to show to users and the system logs feedback (whether the user clicks on the ad or not).



# Application 6: Contextual Bandits and Learning from Logged Data

- **Scenario:** Learn from logged contextual bandit data. Example: We need to show  $k$  ads to users with each ad comprising of features (title, ad text, query etc.), and given an online policy which (with certain randomization) picks ads to show to users and the system logs feedback (whether the user clicks on the ad or not).
- **Data:** We are given bandit logged data in the form of  $\{(x_1, a_1, r_1, p_1), \dots, (x_n, a_n, r_n, p_n)\}$ . Notation::



# Application 6: Contextual Bandits and Learning from Logged Data

- **Scenario:** Learn from logged contextual bandit data. Example: We need to show  $k$  ads to users with each ad comprising of features (title, ad text, query etc.), and given an online policy which (with certain randomization) picks ads to show to users and the system logs feedback (whether the user clicks on the ad or not).
- **Data:** We are given bandit logged data in the form of  $\{(x_1, a_1, r_1, p_1), \dots, (x_n, a_n, r_n, p_n)\}$ . Notation:
  - Assume we can take fixed number of actions  $1 : k$



# Application 6: Contextual Bandits and Learning from Logged Data

- **Scenario:** Learn from logged contextual bandit data. Example: We need to show  $k$  ads to users with each ad comprising of features (title, ad text, query etc.), and given an online policy which (with certain randomization) picks ads to show to users and the system logs feedback (whether the user clicks on the ad or not).
- **Data:** We are given bandit logged data in the form of  $\{(x_1, a_1, r_1, p_1), \dots, (x_n, a_n, r_n, p_n)\}$ . Notation::
  - Assume we can take fixed number of actions  $1 : k$
  - Denote  $x_i = \{x_i^1, \dots, x_i^k\}$  as the feature vectors of the  $k$  actions



# Application 6: Contextual Bandits and Learning from Logged Data

- **Scenario:** Learn from logged contextual bandit data. Example: We need to show  $k$  ads to users with each ad comprising of features (title, ad text, query etc.), and given an online policy which (with certain randomization) picks ads to show to users and the system logs feedback (whether the user clicks on the ad or not).
- **Data:** We are given bandit logged data in the form of  $\{(x_1, a_1, r_1, p_1), \dots, (x_n, a_n, r_n, p_n)\}$ . Notation::
  - Assume we can take fixed number of actions  $1 : k$
  - Denote  $x_i = \{x_i^1, \dots, x_i^k\}$  as the feature vectors of the  $k$  actions
  - Denote  $a_i$  as the chosen action by the current online policy



# Application 6: Contextual Bandits and Learning from Logged Data

- **Scenario:** Learn from logged contextual bandit data. Example: We need to show  $k$  ads to users with each ad comprising of features (title, ad text, query etc.), and given an online policy which (with certain randomization) picks ads to show to users and the system logs feedback (whether the user clicks on the ad or not).
- **Data:** We are given bandit logged data in the form of  $\{(x_1, a_1, r_1, p_1), \dots, (x_n, a_n, r_n, p_n)\}$ . Notation::
  - Assume we can take fixed number of actions  $1 : k$
  - Denote  $x_i = \{x_i^1, \dots, x_i^k\}$  as the feature vectors of the  $k$  actions
  - Denote  $a_i$  as the chosen action by the current online policy
  - Denote  $p_i$  as the probability of the logged action (by the current online policy)





# Application 6: Contextual Bandits and Learning from Logged Data

- **Scenario:** Learn from logged contextual bandit data. Example: We need to show  $k$  ads to users with each ad comprising of features (title, ad text, query etc.), and given an online policy which (with certain randomization) picks ads to show to users and the system logs feedback (whether the user clicks on the ad or not).
- **Data:** We are given bandit logged data in the form of  $\{(x_1, a_1, r_1, p_1), \dots, (x_n, a_n, r_n, p_n)\}$ . Notation::
  - Assume we can take fixed number of actions  $1 : k$
  - Denote  $x_i = \{x_i^1, \dots, x_i^k\}$  as the feature vectors of the  $k$  actions
  - Denote  $a_i$  as the chosen action by the current online policy
  - Denote  $p_i$  as the probability of the logged action (by the current online policy)
  - Denote  $r_i$  as the Reward obtained by choosing action  $a_i$



# Application 6: Contextual Bandits and Learning from Logged Data

- **Scenario:** Learn from logged contextual bandit data. Example: We need to show  $k$  ads to users with each ad comprising of features (title, ad text, query etc.), and given an online policy which (with certain randomization) picks ads to show to users and the system logs feedback (whether the user clicks on the ad or not).
- **Data:** We are given bandit logged data in the form of  $\{(x_1, a_1, r_1, p_1), \dots, (x_n, a_n, r_n, p_n)\}$ . Notation::
  - Assume we can take fixed number of actions  $1 : k$
  - Denote  $x_i = \{x_i^1, \dots, x_i^k\}$  as the feature vectors of the  $k$  actions
  - Denote  $a_i$  as the chosen action by the current online policy
  - Denote  $p_i$  as the probability of the logged action (by the current online policy)
  - Denote  $r_i$  as the Reward obtained by choosing action  $a_i$
  - Define our policy as  $\pi_\theta(x) = \operatorname{argmax}_{i=1:k} F_\theta(x^i)$ . Again the simplest example of  $F_\theta(x) = \theta^T x$ .



# Application 6: Contextual Bandits and Learning from Logged Data

- **Optimization Problem:** The Inverse Propensity Estimate of the Reward (which is an unbiased estimate of the Reward function is):

$$\max_{\theta} \text{IPS}(\theta) = \max_{\theta} \sum_{i=1}^n r_i / p_i I(\pi_{\theta}(x_i) == a_i)$$



# Application 6: Contextual Bandits and Learning from Logged Data

- **Optimization Problem:** The Inverse Propensity Estimate of the Reward (which is an unbiased estimate of the Reward function is):

$$\max_{\theta} \text{IPS}(\theta) = \max_{\theta} \sum_{i=1}^n r_i / p_i I(\pi_{\theta}(x_i) == a_i)$$

- **SoftMax Relaxation:** The IPS objective above is not continuous and non differentiable. We can define a softmax relaxation as:

$$\max_{\theta} \text{SM}(\theta) = \max_{\theta} \sum_{i=1}^n r_i / p_i \frac{\exp(F_{\theta}(x_i^{a_i}))}{\sum_{j=1}^k \exp(F_{\theta}(x_i^j))}$$



# Discrete Optimization in Machine Learning

- MAP inference in Probabilistic Models: Ising Models, DPPs
- Feature Subset Selection
- Data Partitioning
- Data Subset Selection
- Data Summarization: Text, Images, Video Summarization
- Social networks, Influence Maximization



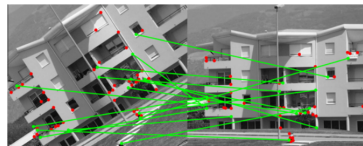
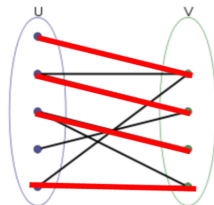
# Discrete Optimization in Machine Learning

- MAP inference in Probabilistic Models: Ising Models, DPPs
- Feature Subset Selection
- Data Partitioning
- Data Subset Selection
- Data Summarization: Text, Images, Video Summarization
- Social networks, Influence Maximization
- Natural Language Processing: words, phrases, n-grams, syntax trees, semantic structures
- Computer Vision: Image Segmentation, Image Correspondence
- Genomics and Computational Biology: cell types or assay selection, selecting peptides and proteins

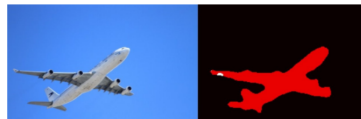
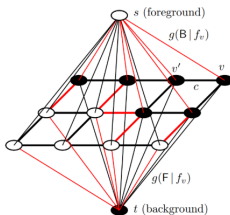


# Application 1: Image Segmentation and Correspondence

Bipartite Matchings



Minimum Cuts



## Application 2: Feature Selection

- **Data:** Given random variables  $X_1, X_2, \dots, X_n$  as features of a given ML task. Denote  $I(A, B)$  as the Mutual Information between variables  $A$  and  $B$ .





## Application 2: Feature Selection

- **Data:** Given random variables  $X_1, X_2, \dots, X_n$  as features of a given ML task. Denote  $I(A, B)$  as the Mutual Information between variables  $A$  and  $B$ .
- **Goal:** Select a subset of features  $A \subseteq \{1, \dots, n\}$  such that the subset of features are *as good* as the original set.



## Application 2: Feature Selection

- **Data:** Given random variables  $X_1, X_2, \dots, X_n$  as features of a given ML task. Denote  $I(A, B)$  as the Mutual Information between variables  $A$  and  $B$ .
- **Goal:** Select a subset of features  $A \subseteq \{1, \dots, n\}$  such that the subset of features are *as good* as the original set.
- **Optimization Problem:** Maximize the Mutual Information between the set of features and the label  $Y$ :

$$\max_{A: |A| \leq k} I(X_A; Y)$$



## Application 2: Feature Selection

- **Data:** Given random variables  $X_1, X_2, \dots, X_n$  as features of a given ML task. Denote  $I(A, B)$  as the Mutual Information between variables  $A$  and  $B$ .
- **Goal:** Select a subset of features  $A \subseteq \{1, \dots, n\}$  such that the subset of features are *as good* as the original set.
- **Optimization Problem:** Maximize the Mutual Information between the set of features and the label  $Y$ :

$$\max_{A: |A| \leq k} I(X_A; Y)$$

- We will see in the second part of this course that this is related to the concept of *submodularity*.



## Application 3: Training Data Subset Selection

- **Data:** Given a training dataset  $\{(x_1, y_1), \dots, (x_n, y_n)\}$  and a budget  $k$ .



# Application 3: Training Data Subset Selection

- **Data:** Given a training dataset  $\{(x_1, y_1), \dots, (x_n, y_n)\}$  and a budget  $k$ .
- **Goal:** Select a subset of data-points  $A \subseteq \{1, \dots, n\}$  such that the model trained on the subset of data is as good as the entire dataset.



# Application 3: Training Data Subset Selection

- **Data:** Given a training dataset  $\{(x_1, y_1), \dots, (x_n, y_n)\}$  and a budget  $k$ .
- **Goal:** Select a subset of data-points  $A \subseteq \{1, \dots, n\}$  such that the model trained on the subset of data is as good as the entire dataset.
- This setting makes even more sense when the labels are missing on some or all of the given data-points.



# Application 3: Training Data Subset Selection

- **Data:** Given a training dataset  $\{(x_1, y_1), \dots, (x_n, y_n)\}$  and a budget  $k$ .
- **Goal:** Select a subset of data-points  $A \subseteq \{1, \dots, n\}$  such that the model trained on the subset of data is as good as the entire dataset.
- This setting makes even more sense when the labels are missing on some or all of the given data-points.
- **Optimization Problem:** Let  $D_{KL}(\cdot)$  denote the KL divergence between two distributions. The optimization problem can be cast as:

$$\max_{A: |A| \leq k} D_{KL}(p(X_A) \| p(X_V))$$



# Application 3: Training Data Subset Selection

- **Data:** Given a training dataset  $\{(x_1, y_1), \dots, (x_n, y_n)\}$  and a budget  $k$ .
- **Goal:** Select a subset of data-points  $A \subseteq \{1, \dots, n\}$  such that the model trained on the subset of data is as good as the entire dataset.
- This setting makes even more sense when the labels are missing on some or all of the given data-points.
- **Optimization Problem:** Let  $D_{KL}(\cdot)$  denote the KL divergence between two distributions. The optimization problem can be cast as:

$$\max_{A: |A| \leq k} D_{KL}(p(X_A) \| p(X_V))$$

- We will see in the second part of this course that this is also related to the concept of *submodularity*.





# Application 4: k-Medoids Clustering

- **Data:** Given a set of datapoints  $\{x_1, x_2, \dots, x_n\}$ , a similarity function  $s_{ij}, i, j \in \{1, \dots, n\}$  and a budget  $k$ .



# Application 4: k-Medoids Clustering

- **Data:** Given a set of datapoints  $\{(x_1, x_2, \dots, x_n)\}$ , a similarity function  $s_{ij}, i, j \in \{1, \dots, n\}$  and a budget  $k$ .
- **Goal:** Select a subset of data-points  $A \subseteq \{1, \dots, n\}$  which can act as k-medoids (similar to k-means except that the means are *a part* of the original set of points).



# Application 4: k-Medoids Clustering

- **Data:** Given a set of datapoints  $\{(x_1, x_2, \dots, x_n)\}$ , a similarity function  $s_{ij}, i, j \in \{1, \dots, n\}$  and a budget  $k$ .
- **Goal:** Select a subset of data-points  $A \subseteq \{1, \dots, n\}$  which can act as k-medoids (similar to k-means except that the means are *a part* of the original set of points).
- **Optimization Problem:** The optimization problem can be cast as:

$$\max_{A: |A| \leq k} \sum_{i=1}^n \max_{j \in A} s_{ij}$$



# Application 4: k-Medoids Clustering

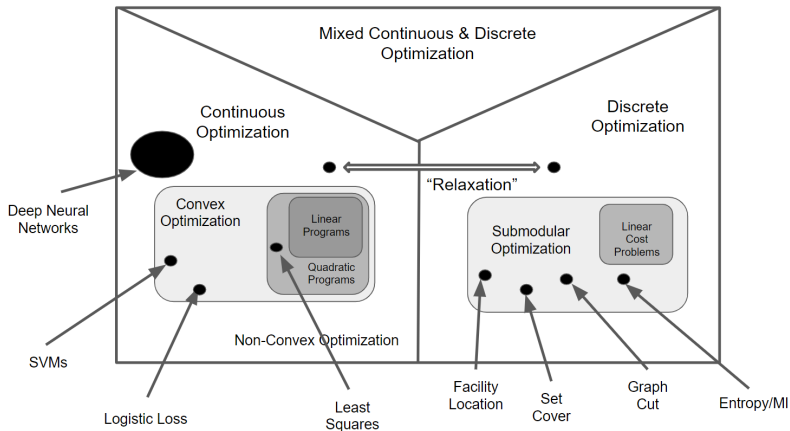
- **Data:** Given a set of datapoints  $\{(x_1, x_2, \dots, x_n)\}$ , a similarity function  $s_{ij}, i, j \in \{1, \dots, n\}$  and a budget  $k$ .
- **Goal:** Select a subset of data-points  $A \subseteq \{1, \dots, n\}$  which can act as k-medoids (similar to k-means except that the means are *a part* of the original set of points).
- **Optimization Problem:** The optimization problem can be cast as:

$$\max_{A: |A| \leq k} \sum_{i=1}^n \max_{j \in A} s_{ij}$$

- We will see in the second part of this course that this problem is called *Facility Location* also related to the concept of *submodularity*.



# Big Picture: Types of Optimization Problems



# Course Project Ideas

- Let's spend a few minutes discussing some ideas for course project(s)
- Why not we all jointly create a *OptML* python toolkit which implements several (discrete and continuous) loss functions, optimization algorithms along with wrappers to machine learning models (e.g. classification, recommender systems, regression etc.)
- Why another toolkit when there are already so many out there?
- A lot of the base for this toolkit will already be covered in this course
- Each group can take on a particular component of the toolkit: with components as a) linear classification/regression, b) non-linear classification/regression, c) recommendation and matrix factorization, d) contextual bandits, e) submodular minimization, f) submodular maximization g) graph algorithms and so on...
- We can discuss ideas on this as the class progresses.

