



Deep Learning

Lecture 2: Learning and Convex Optimization

University of Agder,
Kristiansand, Norway

Prepared by: Hadi Ghauch^{*}, Hossein S. Ghadikolaei[†]

^{*} Telecom ParisTech

[†] Royal Institute of Technology, KTH

<https://sites.google.com/view/fundl/home>

April 2019

Learning outcomes

Supervised/unsupervised learning as optimization

Definitions and properties of smooth and convex functions

Deterministic iterative algorithms for convex problems

Advantages and limitations of each

Convergence analysis

Intro to Reinforcement Learning

Outline

1. Supervised Learning
2. Optimization and Learning
3. Iterative solution approaches
4. Unsupervised and Reinforcement Learning
5. Supplements

Outline

1. Supervised Learning
2. Optimization and Learning
3. Iterative solution approaches
4. Unsupervised and Reinforcement Learning
5. Supplements

Machine Learning

- **Supervised learning**

learning from **labeled training set**, $\{(x_i, y_i)\}_{i=1}^N$

x_i is feature vector, y_i label for sample $i \in [N]$

e.g., deep learning for **regression/classification** task

- **Unsupervised learning**

learning from **unlabeled training set**, $\{x_i\}_{i=1}^N$

identify inherent structure of the data

is it low-dimensional (PCA, MF) ? does it lie on manifold (manifold learning)

- **Reinforcement learning (a.k.a online learning)**

learning by interacting with an unknown environment (modeled by a Markov decision process), e.g., Q -learning

Supervised learning: Recap

A dataset of N training samples $\{(x_i, y_i)\}_{i=1}^N$

hidden (unknown) fnc that generates training set

learned func from the training set (a.k.a **model**)

Prediction on sample i : and loss on each sample

Empirical risk (training error):

Expected risk (true risk):

necessary cond on the **shattering coeff** to ensure empirical risk and expected (true) risk are uniformly close

Outline

1. Supervised Learning
2. Optimization and Learning
3. Iterative solution approaches
4. Unsupervised and Reinforcement Learning
5. Supplements

Learning and Optimization

Linear ridge regression:

Given set of feature vectors with labels: $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$
find the hyperplane $\mathbf{w} \in \mathbb{R}^d$ that min the squared loss (board)

$$\underbrace{\frac{1}{N} \sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{x}_i)^2}_{\text{training error}} + \underbrace{\lambda \|\mathbf{w}\|_2^2}_{\text{regularizer}} := f(\mathbf{w}) \quad \text{regularized loss}$$

$f()$: **loss function** + regularization ($\ell()$ in Lec 1)

\mathbf{w} : **function/model** to be learned ($f()$ in Lec 1)

increasing λ makes training error larger while reducing overfitting

Linear LASSO regression: $f(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{x}_i)^2 + \lambda \|\mathbf{w}\|_1$

Supervised learning always solves an optimization problem

- Convex problem for linear/logistic regression, SVM

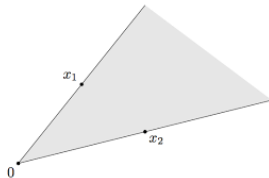
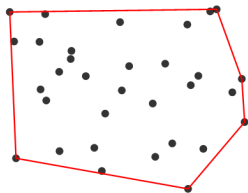
Convexity: Basic definitions

Convex combination of points $\mathcal{X} = \{\mathbf{x}_i\}_{i \in [n]}$ is $\sum_{i \in [n]} \theta_i \mathbf{x}_i$ where $\{\theta_i\}$ are *convex weights*, $\sum_i \theta_i = 1$

Conic combination of points $\mathcal{X} = \{\mathbf{x}_i\}_{i \in [n]}$ is $\sum_{i \in [n]} \theta_i \mathbf{x}_i$ where $\{\theta_i\}$ are *affine weights*, $\theta_i \geq 0$

Convex hull of \mathcal{X} : set of all convex combinations of points in \mathcal{X}

Convex cone of \mathcal{X} : set of all conic combinations of points in \mathcal{X}



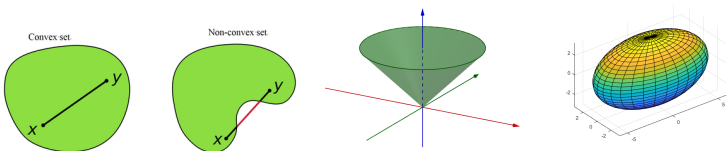
Convex Sets

Convex set \mathcal{X} : $\forall \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X}$ and $\theta \in [0, 1]$, $\theta \mathbf{x}_1 + (1 - \theta) \mathbf{x}_2 \in \mathcal{X}$.

Intuition: convex combination of any two points in the set is within the set.

Examples:

- Euclidean ball with radius r at \mathbf{x}_c : $\{\mathbf{x} \mid \|\mathbf{x} - \mathbf{x}_c\|_2 \leq r\}$
- Norm cone: $\{(\mathbf{x}, r) \mid \|\mathbf{x}\| \leq r\}$
- Ellipsoid centered at \mathbf{x}_c : $\{(\mathbf{x} - \mathbf{x}_c)^T \mathbf{P}^{-1}(\mathbf{x} - \mathbf{x}_c) \leq 1\}$, where \mathbf{P} is positive definite



Convex Functions

Convex function $f : \mathcal{X} \rightarrow \mathbb{R}$: if \mathcal{X} is convex and $\forall \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X}$ and $\theta \in [0, 1]$:

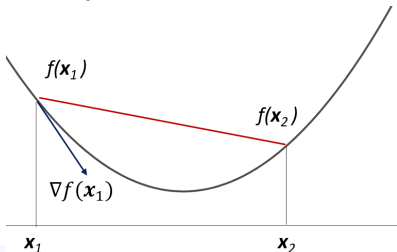
$$f(\theta \mathbf{x}_1 + (1 - \theta) \mathbf{x}_2) \leq \theta f(\mathbf{x}_1) + (1 - \theta) f(\mathbf{x}_2) \quad (1)$$

If f differentiable, (1) is equivalent to

$$f(\mathbf{x}_2) - f(\mathbf{x}_1) \geq \nabla f(\mathbf{x}_1)^T (\mathbf{x}_2 - \mathbf{x}_1), \quad \forall (\mathbf{x}_1, \mathbf{x}_2) \in \mathcal{X} \quad (2)$$

Intuition? grad is a **global linear lower bound** what about non-smooth convex ?

If f twice differentiable, (1) $\Leftrightarrow \nabla^2 f(\mathbf{x}) \succeq \mathbf{0}, \forall \mathbf{x} \in \mathcal{X}$: PSD Hessian, non-negative curvature everywhere.



Convex Functions: examples

- $\|\mathbf{x}\|_p$ for any $p \geq 1$
- $f(\mathbf{x}) = \|\mathbf{Ax} - \mathbf{b}\|_2^2$ is convex for any \mathbf{A} (since $\nabla^2 f(\mathbf{x}) := 2\mathbf{A}^T \mathbf{A} \succeq \mathbf{0}$)
- $\|\mathbf{Ax} - \mathbf{b}\|_p$ for any $p \geq 1$
- Set of PSD matrices $\mathbf{X} \succeq \mathbf{0}$ is a cone (why?). It is convex.
- $\lambda_{\max}(\mathbf{X})$ convex in \mathbf{X} , for $\mathbf{X} \succeq \mathbf{0}$
- $\text{proj}(\mathbf{x}, \mathcal{C}) := \inf_{\mathbf{y} \in \mathcal{C}} \|\mathbf{y} - \mathbf{x}\|$ for convex \mathcal{C}
- Check Boyd and Vandenberghe (2003) for more examples
Nonnegative weighted sum, composition with affine function, pointwise maximum and supremum, composition, minimization, and perspective

Convex Problems and Duality

$$\begin{aligned} \text{Primal problem:} \quad & \text{minimize } f_0(\mathbf{x}) \\ & \text{s.t. } f_i(\mathbf{x}) \leq 0, i \in [m] \\ & \quad h_i(\mathbf{x}) = 0, i \in [p] \end{aligned} \tag{3}$$

f_0 convex cost fnc, f_i convex inequality constraints, h_i affine equality constraints.

Lagrange dual function: $g : \mathbb{R}^m \times \mathbb{R}^p \rightarrow \mathbb{R}$

$$g(\lambda, \nu) = \inf_{\mathbf{x} \in \mathcal{X}} L(\mathbf{x}, \lambda, \nu) := f_0(\mathbf{x}) + \sum_{i \in [m]} \lambda_i f_i(\mathbf{x}) + \sum_{i \in [p]} \nu_i h_i(\mathbf{x})$$

$$\begin{aligned} \text{Lagrange dual problem:} \quad & d^* = \text{maximize } g(\lambda, \nu) \\ & \text{s.t. } \lambda_i \geq 0, i \in [m] \quad , \quad \nu_i \text{ free}, i \in [p] \end{aligned}$$

Weak duality: $d^* \leq f^*$ (holds for any (3))

Strong duality: $d^* = f^*$ (holds when (3) strongly convex)

Duality

- Check Boyd and Vandenberghe (2003) for more details

Constraint qualifications

Slater's constraint qualification

Complementary slackness

Karush-Kuhn-Tucker (KKT) conditions

Strong convexity

Differentiable function f is μ -strongly convex iff $\forall \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X}, \mu > 0$

$$\underbrace{f(\mathbf{x}_2) - f(\mathbf{x}_1) \geq \nabla f(\mathbf{x}_1)^T (\mathbf{x}_2 - \mathbf{x}_1)}_{\text{definition of convex } f} + \frac{\mu}{2} \|\mathbf{x}_2 - \mathbf{x}_1\|_2^2$$

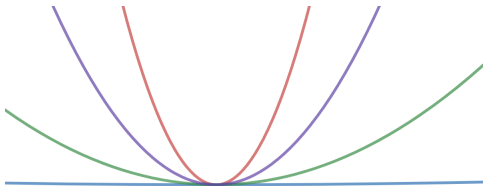
Intuition:

for convex f : gradient is global linear lowerbound

for strongly convex f : quadratic lower bound with strong convexity

Global definitions not local ($\forall \mathbf{x} \in \mathcal{X}$)

Gradient can be replaced by sub-gradient for non-smooth functions



Strong convexity: Equivalent definitions

$$f(\mathbf{x}_2) - f(\mathbf{x}_1) \geq \nabla f(\mathbf{x}_1)^T (\mathbf{x}_2 - \mathbf{x}_1) + \frac{\mu}{2} \|\mathbf{x}_2 - \mathbf{x}_1\|_2^2 \quad (4)$$

- (4) is equivalent to a **minimum positive curvature**

$$\nabla^2 f(\mathbf{x}) \succeq \mu \mathbf{I}_d, \forall \mathbf{x} \in \mathcal{X}$$

- (4) is equivalent to $(\nabla f(\mathbf{x}_2) - \nabla f(\mathbf{x}_1))^T (\mathbf{x}_2 - \mathbf{x}_1) \geq \mu \|\mathbf{x}_2 - \mathbf{x}_1\|_2^2$
- (4) implies

$$(a) \ f(\mathbf{x}) - f(\mathbf{x}^*) \leq \frac{1}{2\mu} \|\nabla f(\mathbf{x})\|_2^2, \forall \mathbf{x}$$

dist bw any points $f(\mathbf{x})$ and $f(\mathbf{x}^*)$ upperbounded by $\|\nabla f(\mathbf{x})\|$

$$(b) \ \|\mathbf{x}_2 - \mathbf{x}_1\|_2 \leq \frac{1}{\mu} \|\nabla f(\mathbf{x}_2) - \nabla f(\mathbf{x}_1)\|_2, \forall \mathbf{x}_1, \mathbf{x}_2$$

dist bw any two points upperbounded by dist of their grad

$$(c) \ (\nabla f(\mathbf{x}_2) - \nabla f(\mathbf{x}_1))^T (\mathbf{x}_2 - \mathbf{x}_1) \leq \frac{1}{\mu} \|\nabla f(\mathbf{x}_2) - \nabla f(\mathbf{x}_1)\|_2^2, \forall \mathbf{x}_1, \mathbf{x}_2$$

Smoothness

A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is L -smooth iff it is differentiable and its gradient is L -Lipschitz-continuous (usually w.r.t. norm-2):

$$\forall \mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^d, \|\nabla f(\mathbf{x}_2) - \nabla f(\mathbf{x}_1)\|_2 \leq L\|\mathbf{x}_2 - \mathbf{x}_1\|_2 \quad (5)$$

If f twice differentiable and smooth:

(5) equivalent $\nabla^2 f(\mathbf{x}) \preceq L\mathbf{I}_d$: upperbound on max curvature of f
 L upperbound on largest eigenvalue of $\nabla^2 f(\mathbf{x})$

recall: for strongly convex: μ lowerbound on min eigenval of $\nabla^2 f(\mathbf{x})$

Strongly convex and L -smooth:

$$\mu\mathbf{I}_d \preceq \nabla^2 f(\mathbf{x}) \preceq L\mathbf{I}_d$$

L and μ upperbound and lowerbound on curvature of f

(5) implies

- (a) $f(\mathbf{x}_2) \leq f(\mathbf{x}_1) + \nabla f(\mathbf{x}_1)^T(\mathbf{x}_2 - \mathbf{x}_1) + \frac{L}{2}\|\mathbf{x}_2 - \mathbf{x}_1\|_2^2, \forall \mathbf{x}_1, \mathbf{x}_2$
- (b) $f(\mathbf{x}_2) \geq f(\mathbf{x}_1) + \nabla f(\mathbf{x}_1)^T(\mathbf{x}_2 - \mathbf{x}_1) + \frac{1}{2L}\|\mathbf{x}_2 - \mathbf{x}_1\|_2^2, \forall \mathbf{x}_1, \mathbf{x}_2$
- (c) $(\nabla f(\mathbf{x}_2) - \nabla f(\mathbf{x}_1))^T(\mathbf{x}_2 - \mathbf{x}_1) \geq \frac{1}{L}\|\nabla f(\mathbf{x}_2) - \nabla f(\mathbf{x}_1)\|_2^2, \forall \mathbf{x}_1, \mathbf{x}_2$

Outline

1. Supervised Learning
2. Optimization and Learning
3. Iterative solution approaches
4. Unsupervised and Reinforcement Learning
5. Supplements

Descent methods

Problem: minimize $f(\mathbf{w})$
 $\mathbf{w} \in \mathbb{R}^d$

$f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\pm\infty\}$ differentiable and convex

Intuition: iteratively move along direction where f decreases. Formally,

$$\mathbf{w}_{k+1} = \mathbf{w}_k + \alpha_k \mathbf{d}(\mathbf{w}_k),$$

- $(\alpha_k)_{k \in \mathbb{N}}$ non-negative step sizes
- $\mathbf{d}(\mathbf{w}_k)$, decent direction is such that $f(\mathbf{w}_{k+1}) < f(\mathbf{w}_k)$

How to choose $\mathbf{d}(\mathbf{w}_k)$?

$$f(\mathbf{w}_{k+1}) < f(\mathbf{w}_k) \text{ implies } -\nabla f(\mathbf{w}_k)^T \mathbf{d}(\mathbf{w}_k) > 0$$

- any direction in half-space of the negative gradient

Gradient/Steepest descent: Set $\mathbf{d}(\mathbf{w}_k) = -\nabla f(\mathbf{w}_k)$ (obvious choice)

Newton's method: Set $\mathbf{d}(\mathbf{w}_k) = -[\nabla^2 f(\mathbf{w}_k)]^{-1} \nabla f(\mathbf{w}_k)$

- adjust $\mathbf{d}(\mathbf{w}_k)$ to account for curvature: it is descent direction

Compare different iterative methods ?

- **Computational Complexity:** number of calculations/flops per iteration
- **Iteration Complexity:** number of iterations to get within ϵ -ball around optimal
- both done using order analysis: $\mathcal{O}()$

Gradient descent

Gradient descent (GD), aka, batch GD, full GD, steepest descent

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k) \quad (6)$$

for some sequence of non-negative step sizes $(\alpha_k)_{k \in \mathbb{N}}$.

Intuition: Approximate f at each iteration, with the gradient

Theorem 1: Convergence of GD with constant step size

a) **Convex and L -smooth** f with $\alpha \leq 1/L$ satisfies

$$f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \frac{\|\mathbf{w}_0 - \mathbf{w}^*\|_2^2}{2k\alpha} \approx \mathcal{O}(1/k).$$

b) **μ -strongly convex and L -smooth** f with $\alpha = 2/(\mu + L)$, we have

$$\|\mathbf{w}_k - \mathbf{w}^*\|_2 \leq \left(1 - \frac{2}{1+L/\mu}\right)^k \|\mathbf{w}_0 - \mathbf{w}^*\|_2 = c^k \|\mathbf{w}_0 - \mathbf{w}^*\|_2, \text{ where } c < 1$$

dist to opt decays geometrically with k : **linear convergence**

► Proofs

GD for smooth and strongly convex functions

Iteration Complexity : Number of iterations to get within ϵ to w^*

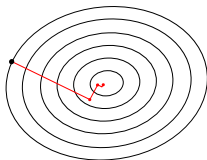
Smooth strongly-convex: $\mathcal{O}(L/\mu \log(1/\epsilon))$ iterations

Computational Complexity = $\mathcal{O}(N)$ (evaluate grad of N samples)

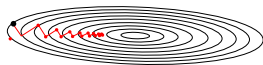
Convergence Rate $c = 1 - \frac{2}{1 + L/\mu}$ is linear

L/μ large: iteration complexity $\mathcal{O}(L/\mu \log(1/\epsilon))$ increases

Preconditioning to change the space geometry, to make sub-levels similar in all coordinates



small L/μ



large L/μ

- small L/μ : $L/\mu = 1 \Leftrightarrow$ all eigenvals of Hessian equal. **Why?**
- ideal setting for GD

Newton method

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k [\nabla^2 f(\mathbf{w}_k)]^{-1} \nabla f(\mathbf{w}_k)$$

Intuition: Approximate f at each iteration, with its Hessian

Theorem 2: Quadratic convergence of Newton's method

Assume f is a twice continuously differentiable and set $\alpha_k = 1$. If $\|\mathbf{w}_k - \mathbf{w}^*\|$ is small enough, there is a constant c ($0 < c < 1$) such that $\|\mathbf{w}_{k+1} - \mathbf{w}^*\|_2 \leq c \|\mathbf{w}_k - \mathbf{w}^*\|_2$.

Iteration complexity: $\mathcal{O}(\log \log(1/\epsilon))$ iterations for ϵ -optimality

convergence is logarithmically faster than GD

Computationally expensive. $[\nabla^2 f(\mathbf{w}_k)]^{-1}$ costs $\mathcal{O}(d^3)$. Never used in large-scale learning.

'Lighter' methods: quasi-Newton or subsampled Hessian

Accelerated First-order Methods

Heavy-ball method: (B. Polyak 64)

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k) + \underbrace{\eta_k (\mathbf{w}_k - \mathbf{w}_{k-1})}_{\text{momentum}} \quad (7)$$

Intuition: Add **momentum term** to GD to reduce zigzag

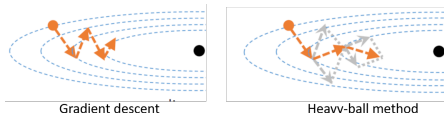
Theorem 3: Convergence of heavy-ball (strongly convex problem)

Suppose f is μ -strongly convex and L -smooth. Then,

$$\|\mathbf{w}_{k+1} - \mathbf{w}^*\|_2 + \|\mathbf{w}_k - \mathbf{w}^*\|_2 \leq \underbrace{\left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^k}_{\text{const} < 1} (\|\mathbf{w}_1 - \mathbf{w}^*\|_2 + \|\mathbf{w}_0 - \mathbf{w}^*\|_2)$$

$$\kappa := L/\mu, \alpha_k := 4/(\sqrt{L} + \sqrt{\mu})^2, \text{ and } \eta_k = \max(|1 - \sqrt{\alpha_k L}|, |1 - \sqrt{\alpha_k \mu}|)^2$$

- **Iteration complexity:** $\mathcal{O}(\sqrt{\kappa} \log(1/\epsilon))$ (much better than GD)
- **Computational Complexity** = $\mathcal{O}(N)$ (evaluate grad of N samples)
- cons: need min and max eigenvalues of Hessian: L, μ



Accelerated First-order Methods

Nesterov Acceleration: (Y. Nesterov 83)

$$\begin{aligned} \mathbf{y}_k &= \mathbf{w}_k + \eta_k(\mathbf{w}_k - \mathbf{w}_{k-1}) \\ \mathbf{w}_{k+1} &= \mathbf{y}_k - \alpha_k \nabla f(\mathbf{y}_k) \end{aligned} \tag{8}$$

Intuition: mix between **gradient** and **extrapolation** steps

Theorem 4: Convergence of Nesterov accel (strongly convex problem)
Suppose f is μ -strongly convex and L -smooth. Set $\eta_k = k/(k+3)$ and $\alpha_k = \alpha = 1/L$. Then,

$$|f(\mathbf{w}_k) - f(\mathbf{w}^*)| \leq 2L \|\mathbf{w}_k - \mathbf{w}^*\|_2^2 \frac{1}{(k+1)^2} \approx \mathcal{O}(1/k^2)$$

- **convergence rate** is quadratic (like Newton): $\mathcal{O}(1/k^2)$
- **computational complexity** (like GD): $\approx \mathcal{O}(N)$
- not a descent method (iterations not necessarily monotonic)
yet converges faster than GD
- no min or max eigenvalues of Hessian needed

Choice of Method

Assume 'best-case' scenario: strong convexity.

Method	Convergence rate	Computational Complexity
GD	$\mathcal{O}(1/k)$	$\mathcal{O}(N)$
Newton	$\mathcal{O}(1/k^2)$	$\mathcal{O}(d^3)$ too high
Nesterov	$\mathcal{O}(1/k^2)$	$\mathcal{O}(N)$

Number of training samples, N , in **modern datasets**:

- **Large-scale learning** ? MovieLens 100M ($N = 10^8$), Alexnet ($N = 1.2 \times 10^6$), etc.
- even smallest computational complexity, $\mathcal{O}(N)$, is too large
 $\mathcal{O}(N)$ is due to **computing gradient over all samples**
- motivation stochastic gradient descent, workhorse for large-scale learning

Outline

1. Supervised Learning
2. Optimization and Learning
3. Iterative solution approaches
4. Unsupervised and Reinforcement Learning
5. Supplements

Unsupervised Learning

Learn from **unlabeled data**: training samples $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^N$

Goal: learn **inherent feature/structure** of the data.

- clustering, matrix factorization, dictionary learning, etc

Low-rank Matrix Factorization (MF) for recommendation systems

- Data matrix \mathbf{X} of size $M \times N$
- $[\mathbf{X}]_{u,i}$ is rating (score) that user u gave to item i
- In recommendation sys model rating as: $[\mathbf{X}]_{u,i} = \mathbf{p}_i^T \mathbf{q}_u$,
 \mathbf{p}_i vector of **item characteristics**, \mathbf{q}_u vector of **user preferences**
- factorize \mathbf{X} with low-rank matrices: $\mathbf{P} \in \mathbb{R}^{M \times d}$, $\mathbf{Q} \in \mathbb{R}^{N \times d}$, $d \ll (M, N)$ such that $\mathbf{X} = \mathbf{P}\mathbf{Q}^T$
- rows of \mathbf{Q} (resp \mathbf{P}) are user prefs (resp item characteristics)
- Unsupervised learning also solves optimization prob (sol. in Lec 4)

$$\arg \min_{\mathbf{P}, \mathbf{Q}} \|\mathbf{X} - \mathbf{P}\mathbf{Q}^T\|_F^2$$

In general MF represents training set \mathbf{X} into low-dim subspaces,

- compression of large data matrix, and prediction

Sequential Decision Making

Markov process (\mathcal{S}, p) , \mathcal{S} is state space, p transition probability

Markov decision process $(\mathcal{S}, \mathcal{A}, p, r)$, \mathcal{A} is action space

reward, r , depends on state and action, $r : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}$

maximize reward fnc: closed-form (high complexity), iteratively (low complexity)

Policy: $\pi : \mathcal{S} \mapsto \mathcal{A}$

Optimal value function:

Let $\gamma \in [0, 1]$. Define the **state-value function (discounted reward)** starting from state s following policy π is $v_\pi(s) = \mathbb{E} \left[\sum_{i=1}^T \gamma^{i-1} r_i \mid s_0 = s \right]$.

For any given policy and state, we can compute $v_\pi(s)$ in closed-form.

The optimal value function is $v^* = \max_\pi v_\pi(s)$. Computing v^* needs inverse of large matrix (scales with $|\mathcal{S}|$): impossible in modern application (Go)

More efficient way using **Bellman's optimality eqt**

$$v^*(s) = \max_a \sum_{s', r} p(s', r \mid s, a) [r + \gamma v^*(s')]$$

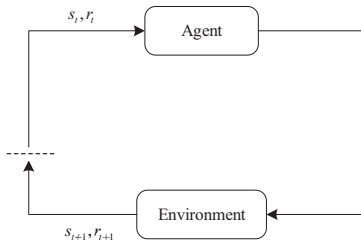
equivalent to optimal value func. fixed point eqt (both sides depend of s)

Solve iteratively until converge (fixed point exist)

So far no learning, just sequential decision making

Reinforcement Learning (RL)

In SDM: p , r , are known. Only need the policy that maximizes value function



- In reinforcement learning (RL). r and p unknown
- Agent interacts with the environment:
In state, s_t , agent takes action (a_t), gets a reward, r_t (fnc of s_t and a_t)
agent observes next state, s_{t+1}
- **Objective:** find the optimal policy (with unknown r and p)

Agent interacts with dataset. Compare with SL where learning done with a fixed dataset.

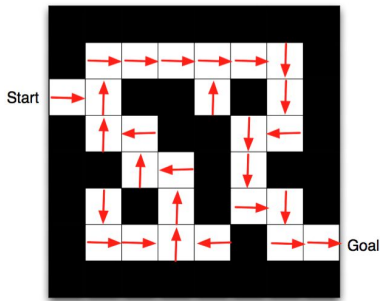
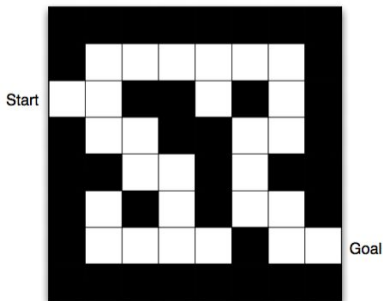
RL more adaptive/efficient: online learning

Reinforcement Learning (example)

Actions: taking one step (top/bottom/left/right) at a time

Objective: reach destination with a minimum number of steps

Assume reward known: -1 per step



Reinforcement Learning

Sample mean and stochastic approximation

estimate sample mean of random process, $\{r_n\}$, using stochastic approx

$$q_{n+1} = \frac{r_1 + r_2 + \cdots + r_{n+1}}{n+1} = q_n + \frac{1}{n+1} [r_{n+1} - q_n]$$

NewEstimate \leftarrow OldEstimate + StepSize [Target – OldEstimate]

Reinforcement Learning

$$\text{NewEstimate} \leftarrow \text{OldEstimate} + \text{StepSize} [\text{Target} - \text{OldEstimate}]$$

TD: temporal-difference

$$v(s_t) \leftarrow v(s_t) + \alpha_{n(s)} [r_{t+1} + \gamma v(s_{t+1}) - v(s_t)]$$

SARSA: state-action-reward-state-action

From state s , take action a , observe immediate reward r and next state s' . Then choose action a' using policy derived from Q , update

$$q(s, a) \leftarrow q(s, a) + \alpha_{n(s,a)} [r + \gamma q(s', a') - q(s, a)]$$

Q-learning

$$q(s, a) \leftarrow q(s, a) + \alpha_{n(s,a)} \left[r + \gamma \max_j q(s', j) - q(s, a) \right]$$

Some references

- S. Bubeck, "Convex optimization: Algorithms and complexity," Foundations and Trends in Machine Learning, 2015.
- Y. Nesterov, "Introductory lectures on convex optimization: A basic course, " Springer Science & Business Media, 2004.
- F. Bach, "Large-scale machine learning and convex optimization," Machine Learning Summer School, 2018.
- Y. Chen, " Large-Scale Optimization for Data Science," ELE538B course, spring 2018.
- S. Trivedi and R. Kondor, "Deep Learning", CMSC 35246 course, spring 2017.
- L. Bottou, F. Curtis, and J. Norcedal, "Optimization methods for large-scale machine learning," SIAM Rev., 2018.
- N. Parikh and S. Boyd, "Proximal algorithms", gradient methods," Foundations and Trends in Optimization, 2013.
- A. Beck and M. Teboulle, "A fast iterative shrinkage-thresholding algorithm for linear inverse problems," SIAM Journal on Imaging Sciences, 2009.

Outline

1. Supervised Learning
2. Optimization and Learning
3. Iterative solution approaches
4. Unsupervised and Reinforcement Learning
5. Supplements

Proof sketch for convex and L -smooth f

By convexity of f , $f(\mathbf{x}_i) \leq f(\mathbf{x}^*) + \langle \nabla f(\mathbf{x}_i), \mathbf{x}_i - \mathbf{x}^* \rangle$ (*)

Use smoothness property $f(\mathbf{x}_{i+1}) \leq f(\mathbf{x}_i) + \nabla f(\mathbf{x}_i)^T (\mathbf{x}_{i+1} - \mathbf{x}_i) + \frac{L}{2} \|\mathbf{x}_{i+1} - \mathbf{x}_i\|_2^2$
and $\alpha \leq 1/L$ to conclude $f(\mathbf{x}_{i+1}) \leq f(\mathbf{x}_i) - \frac{\alpha}{2} \|\nabla f(\mathbf{x}_i)\|_2^2$ (**)

Substitute (*) into (**) and note that from (6): $\nabla f(\mathbf{x}_i) = \frac{\mathbf{x}_i - \mathbf{x}_{i+1}}{\alpha}$

Observe $f(\mathbf{x}_{i+1}) \leq f^* + \frac{1}{2\alpha} (\|\mathbf{x}_i - \mathbf{x}^*\|_2^2 - \|\mathbf{x}_{i+1} - \mathbf{x}^*\|_2^2)$

Summing over iterations: $\frac{1}{k} \sum_{i \in [k]} (f(\mathbf{x}_i) \leq f^*) \leq \frac{1}{2\alpha k} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2$

Conclude from the non-increasing property of GD iterates:

$$f(\mathbf{x}_k) - f^* \leq \frac{1}{k} \sum_{i \in [k]} f(\mathbf{x}_i) - f^* \leq \frac{1}{2\alpha k} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2.$$

Proof sketch for strongly-convex and L -smooth f

From smoothness and vanishing gradient of the optimal point, conclude

$$f(\mathbf{x}_i) - f(\mathbf{x}^*) \leq \frac{L}{2} \|\mathbf{x}_i - \mathbf{x}^*\|_2^2 \quad (*)$$

Use the coercivity of the gradient

$$(\nabla f(\mathbf{x}) - \nabla f(\mathbf{y}))^T (\mathbf{x} - \mathbf{y}) \geq \frac{\mu L}{\mu + L} \|\mathbf{x} - \mathbf{y}\|_2^2 + \frac{1}{\mu + L} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2^2$$

Use $\alpha < 2/(L + \mu)$ to obtain $\|\mathbf{x}_i - \mathbf{x}^*\|^2 \leq \left(1 - 2t\alpha \frac{\mu L}{\mu + L}\right) \|\mathbf{x}_i - \mathbf{x}^*\|_2^2$

Iterate over i and use $(*)$ to obtain

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \frac{L}{2} \prod_{i \in [k]} \left(1 - 2t\alpha \frac{\mu L}{\mu + L}\right) \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2$$

Use $e^{-x} \geq 1 - x$ to conclude $f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \frac{L}{2} e^{-ck} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2$ for $c = \frac{2\alpha\mu L}{\mu + L}$