CIS 5200: MACHINE LEARNING CONVEX OPTIMIZATION

Content here draws from material by Sham Kakade (Harvard), Elad Hazan (Princeton)





14 February 2023

Surbhi Goel

Spring 2023



OUTLINE - TODAY

Logistics
Kernels (quick summary)
Convexity (recap)
Gradient Descent
Proof of Convergence

POWER OF KERNELS

- training points, $w = \sum \alpha_i x_i$ i=1
- Replace $x_i^T x_i \rightarrow k(x_i, x_i)$ everywhere for a valid kernel k



Show that the solution to your problem lies in the span of the

There is a general theorem called the Representer Theorem which tells us when this is true

Rewrite the algorithm and the predictor so that all training or test points are only accessed in inner-products $(x_i^{\mathsf{T}}x_i)$ with other points

Super Powerful!

CHALLENGE

* How do we choose a good feature map ϕ ? * Feature map is the same for all inputs!

Can learn the feature map itself \rightarrow deep learning!

CONVEX OPTIMIZATION

m₁n \mathcal{W}

such that

Here $\mathscr{C} \subseteq \mathbb{R}^d$ is a convex set and $F : \mathbb{R}^d \to \mathbb{R}$ is a convex function

It is a powerful sub-class of optimization problems * allows for efficient global solutions * beautiful mathematical theory that provides guarantees * very well researched with lots of available libraries



CONVEXITY - RECAP

Convex Set: For all $w, w' \in \mathcal{C}$ and $\alpha \in [0,1]$, $\alpha w + (1 - \alpha)w' \in \mathcal{C}$.







Non-convex

CONVEXITY - RECAP

Convex Function: For all $w, w' \in \mathbb{R}^d$ and $\alpha \in [0,1]$, $F(\alpha w + (1 - \alpha)w') \leq \alpha F(w) + (1 - \alpha)F(w')$



CONVEXITY - FIRST ORDER CHARACTERIZATION

F is convex and differentiable, then for all $w, w' \in \mathbb{R}^d$



 $F(w') \ge F(w) + \nabla F(w)^{\mathsf{T}}(w' - w)$

CONVEXITY - FIRST ORDER OPTIMALITY

F is convex and differentiable, any *w* that satisfies $\nabla F(w) = 0$ is a global minimum of *F*.



Convex



Non-convex

CONVEXITY - SECOND ORDER CHARACTERIZATION

F is convex and twice-differentiable, then for all $w \in \mathbb{R}^d$

Hessian is positive semi-definite (PSD)

- $\nabla^2 F(w) \geq 0$

ONVEXITY - SMOOTHNFSS *L*-smooth Function: For all $w, w' \in \mathbb{R}^d$,

 $F(w) + \nabla F(w)^{\mathsf{T}}(w' - w) + \frac{L}{2} \|w - w'\|_2^2$

 $F(w') \le F(w) + \nabla F(w)^{\mathsf{T}}(w' - w) + \frac{L}{2} \|w - w'\|_2^2$

F(w)

CONVEXITY - STRONG CONVEXITY

CONVEX FUNCTIONS - EXAMPLES

 $*\ell_2$ -norm:

* Mean of convex functions: for convex functions F_1, \ldots, F_m $F(w) = \frac{1}{m} \sum_{i=1}^{m} F_i(w)$

 $F(w) = \|x\|_{2}^{2} = x^{\mathsf{T}}x$

 $F(w; x, y) = \log(1 + \exp(-yw^{\top}x))$

GRADIENT DESCENT - MOTIVATION

Move in the opposite direction of the gradient to decrease the function

GRADIENT DESCENT - ALGORITHM

Algorithm 1: Gradient Descent (GD)

Initialize $w_1 \in \mathbb{R}^d$ while t = 1, 2, ..., T do Update $w_{t+1} = w_t - \eta_t \nabla F(w_t)$ end

Can stop when gradient becomes very small $\|\nabla F(w_t)\|_2 \leq \epsilon$

 η_t is the learning rate or step size which governs how much to move

GRADIENT DESCENT - INTERPRETATION

Consider L-smooth convex function F

We could locally minimize the quadratic $\min_{w'} F(w) + \nabla F(w)^{\mathsf{T}}(w' - w) + \frac{L}{2} ||w' - w||_2^2$

This gives
$$w' = w - \frac{1}{L}\nabla F(w)$$

Gradient step with step size $\frac{1}{r}$

GRADIENT DESCENT - STEP SIZE

Too large

Too small

Just right

GRADIENT DESCENT - CONVERGENCE

Algorithm 1: Gradient Descent (GD)

Initialize $w_1 \in \mathbb{R}^d$ while t = 1, 2, ..., TUpdate $w_{t+1} = u$ end

Theorem: Suppose we run GD on L-smooth function F with fixed

$$f \mathbf{do}$$

 $w_t - \eta_t \nabla F(w_t)$

constant learning rate $\eta_t = 1/L$ at all time. Then at any time τ , we have $F(w_{\tau+1}) - F(w_*) \le \frac{L \|w_1 - w_*\|_2^2}{2}$ w_* is the optimal 2τ

PROOF - OVFRVIFW

iterate and the current iterate for every time t

Step I: Upper bound the difference between function value at the next $F(w_{t+1}) - F(w_t) \le -\frac{L}{2} \|w_{t+1} - w_t\|^2$

 $F(w_t)$ is non-decreasing, we are reducing the function value

PROOF - OVFRVIFW

Step 2: Upper bound the difference between function value at the next iterate and the global minimum for every time t $F(w_{t+1}) - F(w_*) \le \frac{L}{2}(|$

$$\|w_t - w_*\|_2^2 - \|w_{t+1} - w_*\|_2^2$$

 $\|w_t - w_*\|_2$ is non-decreasing, we are getting closer to the optimal

PROOF - OVERVIEW

Step 3: Use the previous to upper difference after τ iterates and the g

$F(w_{\tau+1}) - F(w_{\tau+1}) = F(w_{\tau+1}) - F(w_{\tau+1}) = F(w_{\tau+1}) - F(w_{\tau+1}) = F(w$

bound
$$F(w_{\tau+1}) - F(w_*)$$
, that is, the
global minimum.
 $V_*) \leq \frac{L \|w_1 - w_*\|_2^2}{2\tau}$

PROOF - STFP |

iterate and the current iterate for every time t

Step I: Upper bound the difference between function value at the next $F(w_{t+1}) - F(w_t) \le -\frac{L}{2} \|w_{t+1} - w_t\|^2$

 $F(w_t)$ is non-decreasing, we are reducing the function value

 $F(w') \ge F(w)$ $F(w') \le F(w) + \nabla F($

Important properties:

$$+ \nabla F(w)^{\mathsf{T}}(w' - w) (w)^{\mathsf{T}}(w' - w) + \frac{L}{2} ||w - w'||_2^2$$

PROOF - STFP 2

- **Step 2**: Upper bound the difference between function value at the next iterate and the global minimum for every time t $F(w_{t+1}) - F(w_*) \le \frac{L}{2}(\|$
 - $\|w_t w_*\|_2$ is non-decreasing, we are getting closer to the optimal

 $F(w') \ge F(w)$ $F(w') \le F(w) + \nabla F($

$$\|w_t - w_*\|_2^2 - \|w_{t+1} - w_*\|_2^2$$

Important properties:

$$+ \nabla F(w)^{\mathsf{T}}(w' - w) (w)^{\mathsf{T}}(w' - w) + \frac{L}{2} ||w - w'||_2^2$$

GRADIENT DESCENT - CONVERGENCE

Algorithm 1: Gradient Descent (GD)

Initialize $w_1 \in \mathbb{R}^d$ while t = 1, 2, ..., TUpdate $w_{t+1} = u$ end

 τ , we have

$$\|w_{\tau+1} - w_*\|_2^2 \le \left(1 - \frac{\mu}{L}\right)^{\tau} \|w_1 - w_*\|_2^2 \qquad w_* \text{ is the option}$$

$$f \mathbf{do}$$

 $v_t - \eta_t \nabla F(w_t)$

Theorem: Suppose we run GD on L-smooth μ -strongly convex function F with fixed constant learning rate $\eta_t = 1/L$ at all time. Then at any time

GRADIENT DESCENT - BENEFITS/DRAWBACKS

Algorithm 1: Gradient Descent (GD) Initialize $w_1 \in \mathbb{R}^d$ while t = 1, 2, ..., T do

Update $w_{t+1} = w_t - \eta_t \nabla F(w_t)$ end

convex problems

Requires differentiability, many problems are not convex

Easy to implement, requires only local information, very fast for strongly

GRADIENT DESCENT - RECIPE

*Write your problem as loss minimization * As long as your loss is differentiable, run gradient descent * Tune learning rate/step size to avoid divergence

