

CSE 446

Gradient Descent Part II

Natasha Jaques



$$\hat{w}_p = \operatorname{argmin}_w \sum_i (y_i - x_i^\top w)^2 + \lambda r(w)$$

$$\hat{w}_c = \operatorname{argmin}_w \sum_i (y_i - x_i^\top w)^2 \quad \text{subject to} \quad r(w) \leq \mu$$

Note that \hat{w}_p is a function of λ and \hat{w}_c is a function of μ . We could have written $\hat{w}_p(\lambda)$ and $\hat{w}_c(\mu)$ to make this explicit, but we'll leave it implicit to keep the notation cleaner.

Assume for convenience that the solutions to the above optimization problems are unique.

Then for any $\lambda \geq 0$, there exists μ such that $\hat{w}_p = \hat{w}_c$, and vice versa.

Proof:

Given λ , solve for \hat{w}_p , then set $\mu = r(\hat{w}_p)$. Denote $f(w) = \sum_i (y_i - x_i^\top w)^2$.

Assume for contradiction that $\hat{w}_p \neq \hat{w}_c$.

We know that \hat{w}_p is feasible (i.e., $r(\hat{w}_p) \leq \mu$) because μ was set to exactly make $r(\hat{w}_p) = \mu$.

Since \hat{w}_c is the solution of the second optimization problem, we also know that:

(1) $r(\hat{w}_c) \leq \mu = r(\hat{w}_p)$ (because \hat{w}_c is feasible by definition)

(2) $f(\hat{w}_c) < f(\hat{w}_p)$ (because \hat{w}_c minimizes $f(w)$ out of all feasible points, and \hat{w}_p is feasible).

Putting these together, we get that $f(\hat{w}_c) + \lambda r(\hat{w}_c) < f(\hat{w}_p) + \lambda r(\hat{w}_p)$.

But this is a contradiction because by definition, \hat{w}_p is supposed to minimize $f(w) + \lambda r(w)$.

Thus, $\hat{w}_p = \hat{w}_c$.

Lecture plan

- Gradient descent algorithm + examples
- Theoretical analysis ← we are here
 - When does it work?
 - How quickly does it converge?
 - How do we choose a step size?
 - Key idea: Convexity
- Not tested on proof details, but concepts are important & practical

Convexity

- Optimization problems are hard to solve in general
- The exception: convex optimization
 - Objective is a convex function
 - Constraints are convex sets



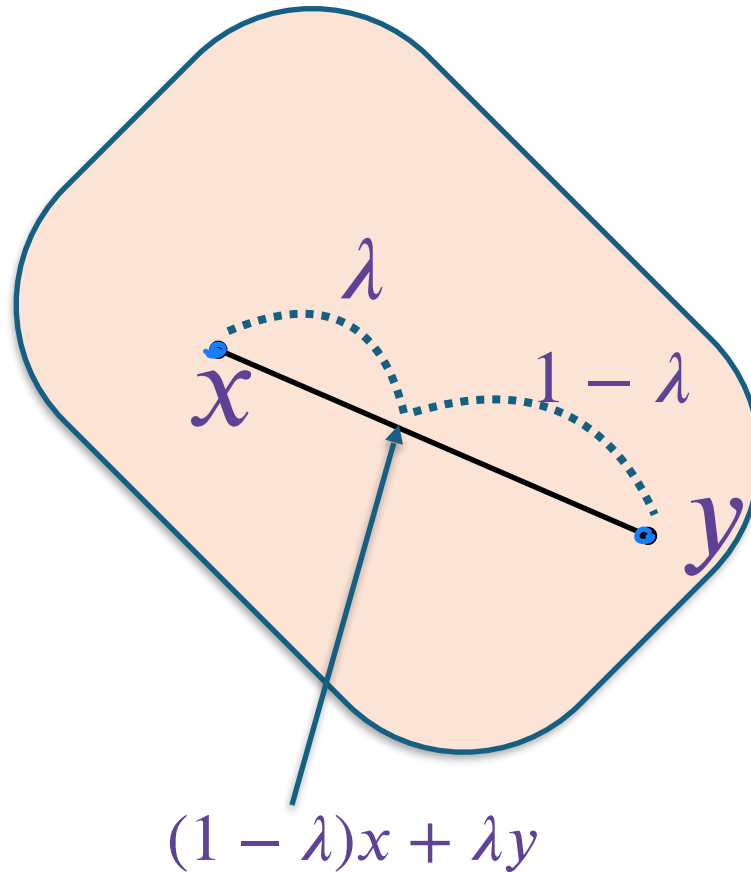
- Special class of problems that can be solved efficiently
- Surprisingly common in practice

What is a convex set?

What is a convex set?

*nothing to do
w/ regular sets*

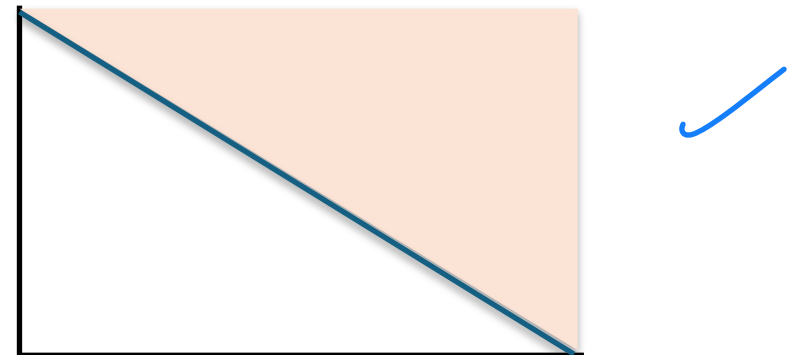
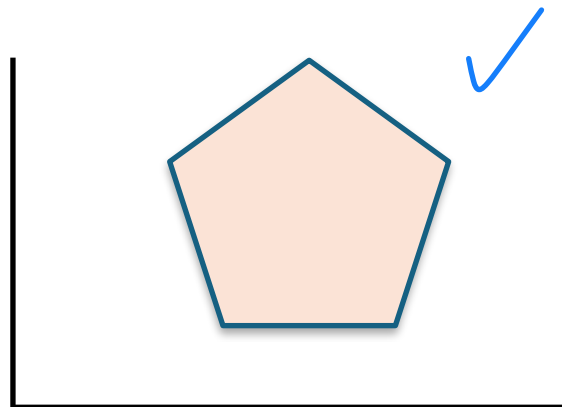
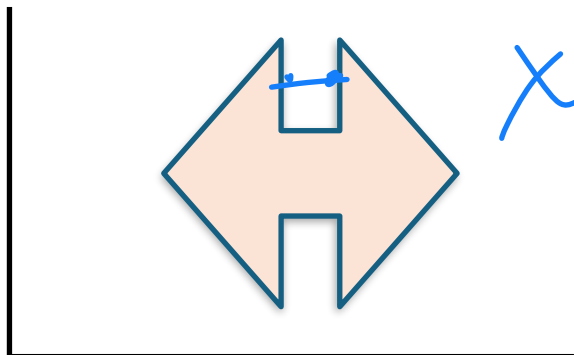
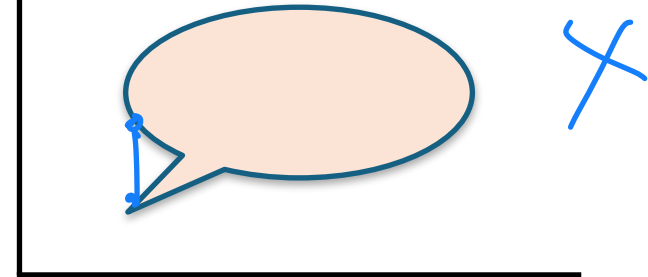
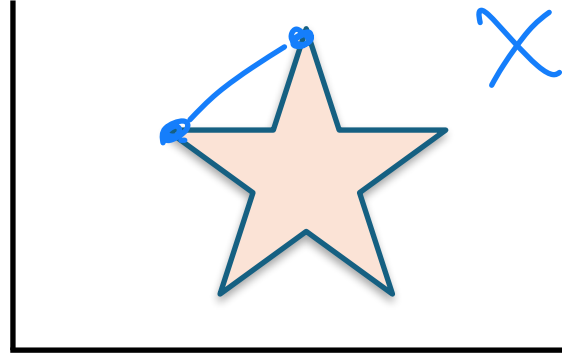
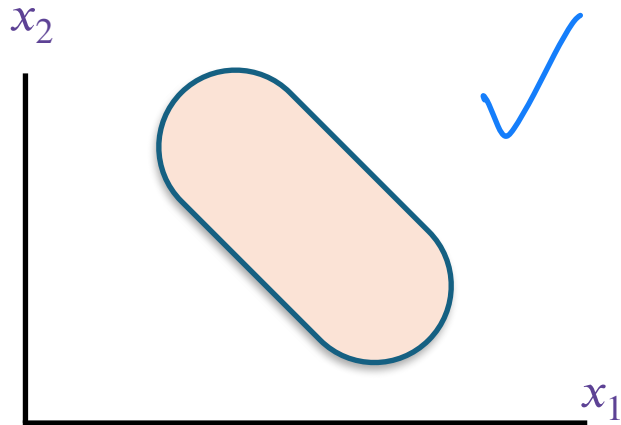
A set $K \subset \mathbb{R}^d$ is convex if $(1 - \lambda)x + \lambda y \in K$ for all $x, y \in K$ and $\lambda \in [0, 1]$



*every point between
2 points in the
set must also be
in the set*

What is a convex set?

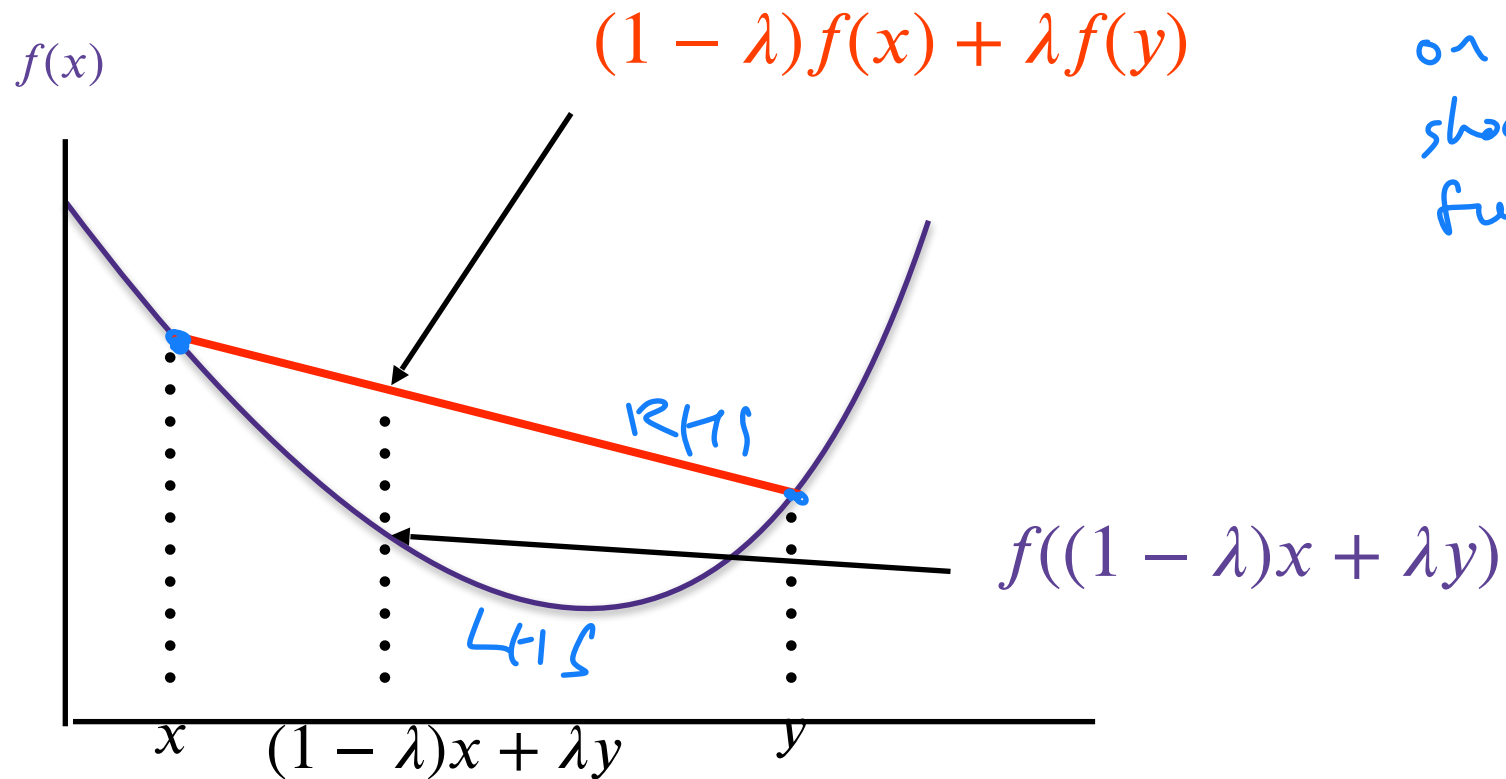
A set $K \subset \mathbb{R}^d$ is convex if $(1 - \lambda)x + \lambda y \in K$ for all $x, y \in K$ and $\lambda \in [0, 1]$



What is a convex function?

What is a convex function?

A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if $f(\underbrace{(1-\lambda)x + \lambda y}_{\text{LHS}}) \leq \underbrace{(1-\lambda)f(x) + \lambda f(y)}_{\text{RHS}}$ for all $x, y \in \mathbb{R}^d$ and $\lambda \in [0, 1]$



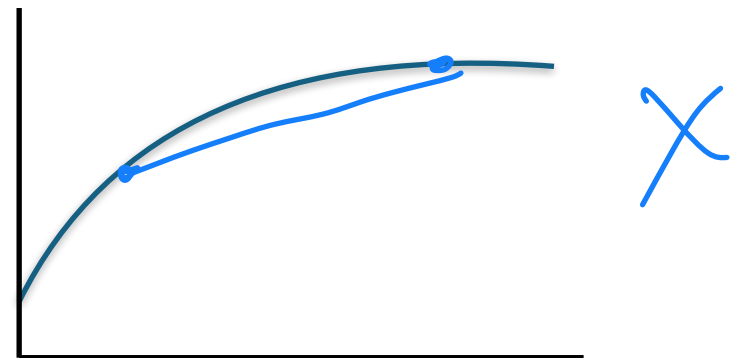
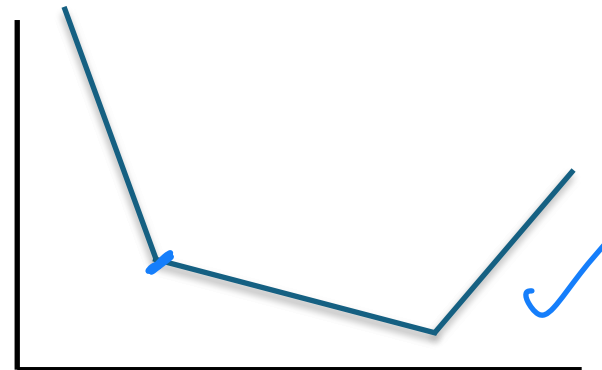
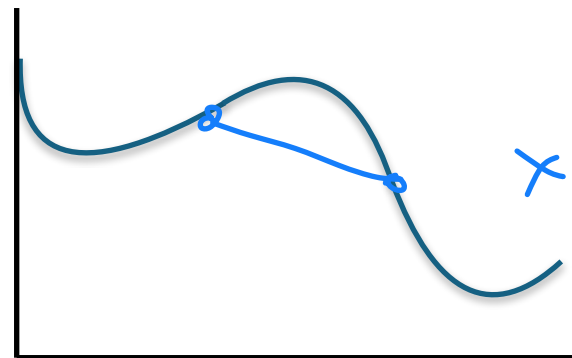
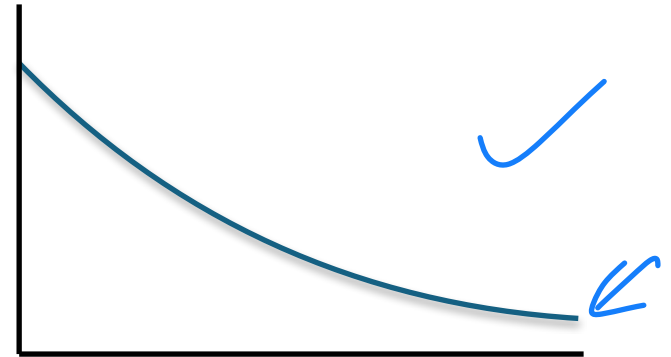
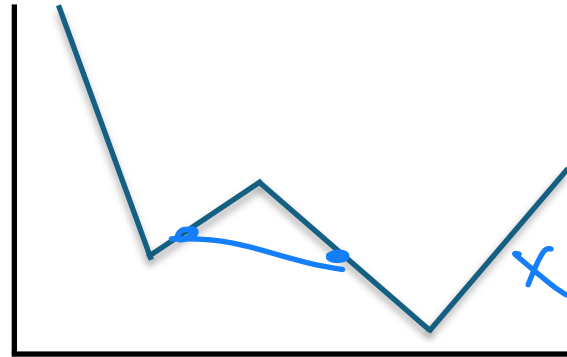
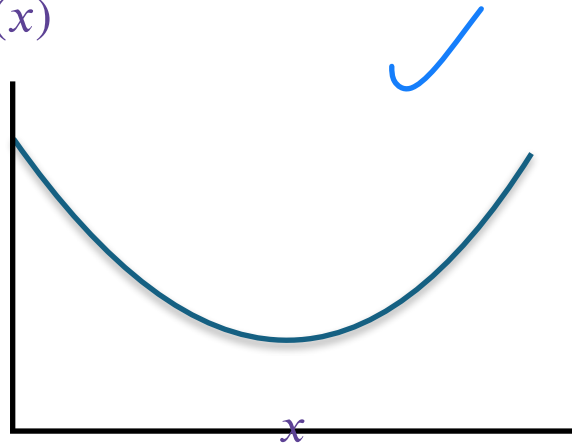
line b/w 2 points on the function should be above function itself

What is a convex function?

Smoothness
 \neq
convexity

A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if $f((1 - \lambda)x + \lambda y) \leq (1 - \lambda)f(x) + \lambda f(y)$ for all $x, y \in \mathbb{R}^d$ and $\lambda \in [0, 1]$

$f(x)$



Convex functions and convex sets



Convex functions and convex sets

A set $K \subset \mathbb{R}^d$ is convex if $(1 - \lambda)x + \lambda y \in K$ for all $x, y \in K$ and $\lambda \in [0, 1]$ ✓

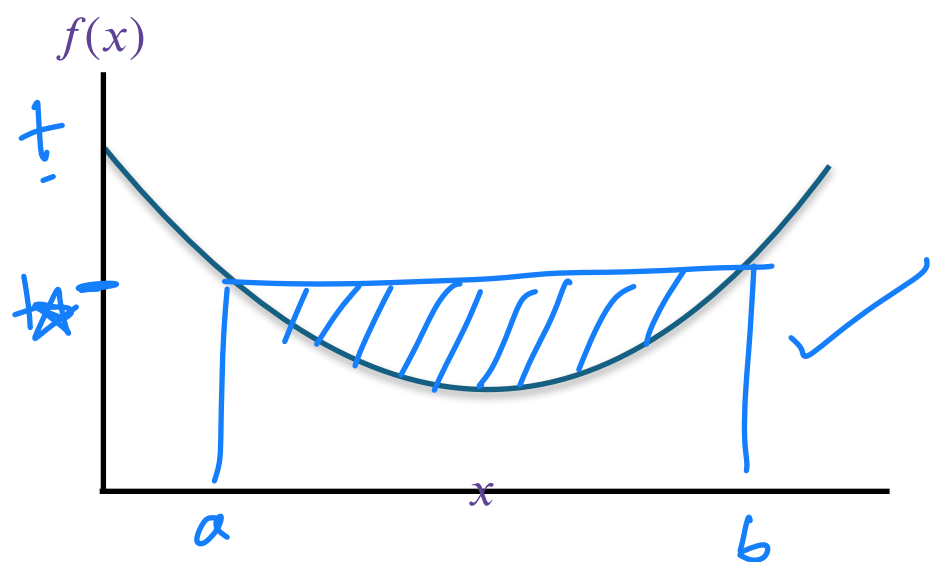
A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if $f((1 - \lambda)x + \lambda y) \leq (1 - \lambda)f(x) + \lambda f(y)$ for all $x, y \in \mathbb{R}^d$ and $\lambda \in [0, 1]$ ✓

A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if the set $\{(x, t) \in \mathbb{R}^{d+1} : f(x) \leq t\}$ is convex

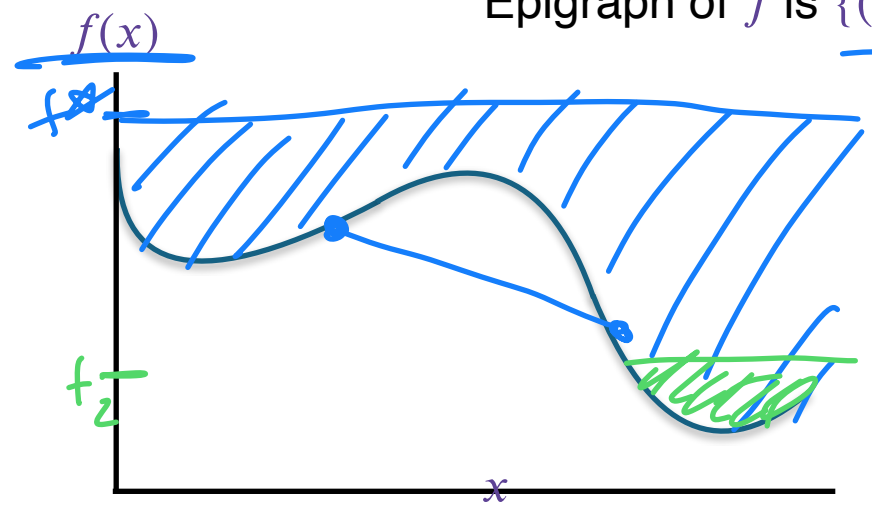
$\forall t$ ✓

epigraph of f

f is $f(x) < t$



Graph of f is $\{(x, t) : f(x) = t\}$
 Epigraph of f is $\{(x, t) : f(x) \leq t\}$



yes, all pairs x, t

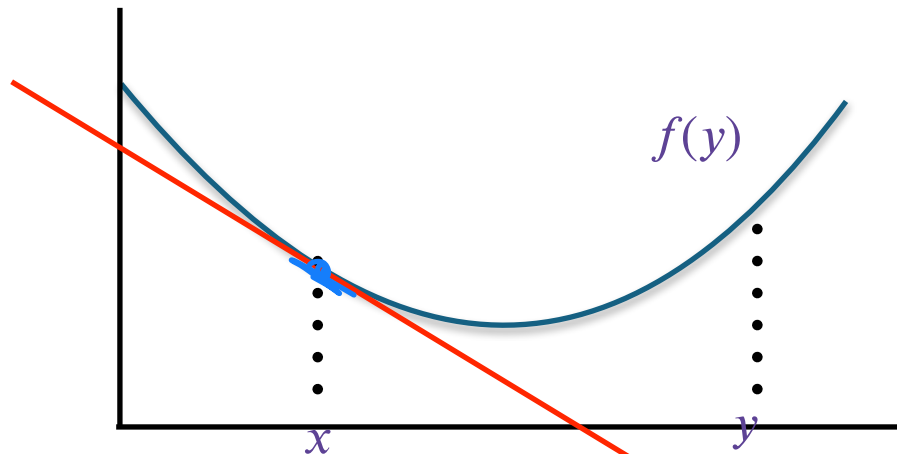
Convexity of differentiable functions

Convexity of differentiable functions

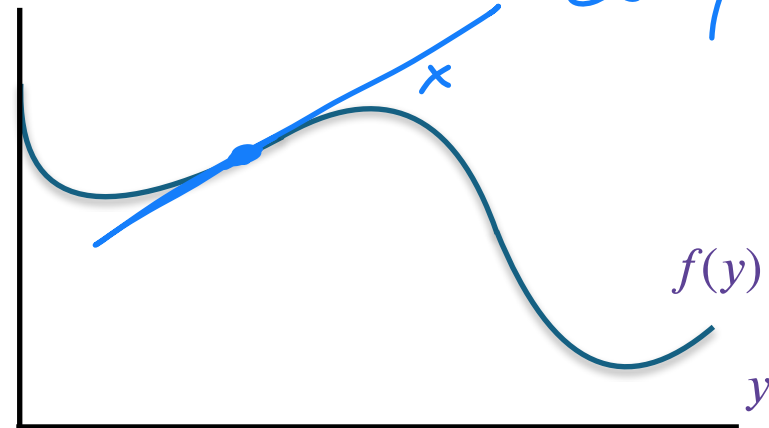
A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if the set $\{(x, t) \in \mathbb{R}^{d+1} : f(x) \leq t\}$ is convex

A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ that is differentiable everywhere is convex if $f(y) \geq f(x) + \nabla f(x)^\top (y - x)$ for all $x, y \in \text{dom}(f)$

1st order Taylor expansion \rightarrow needs to be below the curve everywhere



$$f(x) + \nabla f(x)^\top (y - x)$$



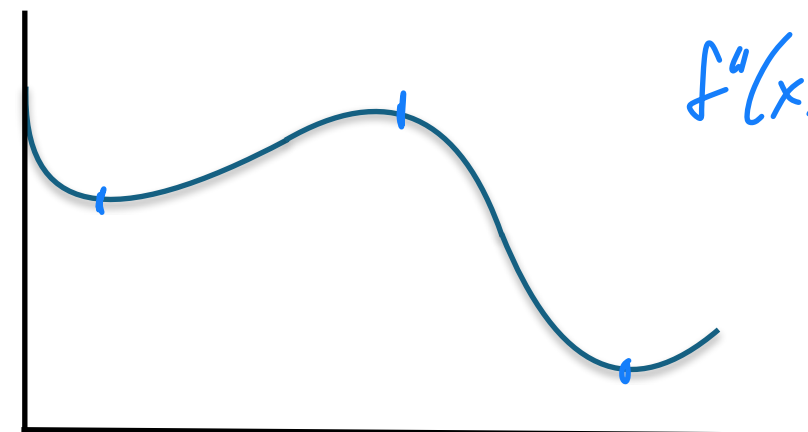
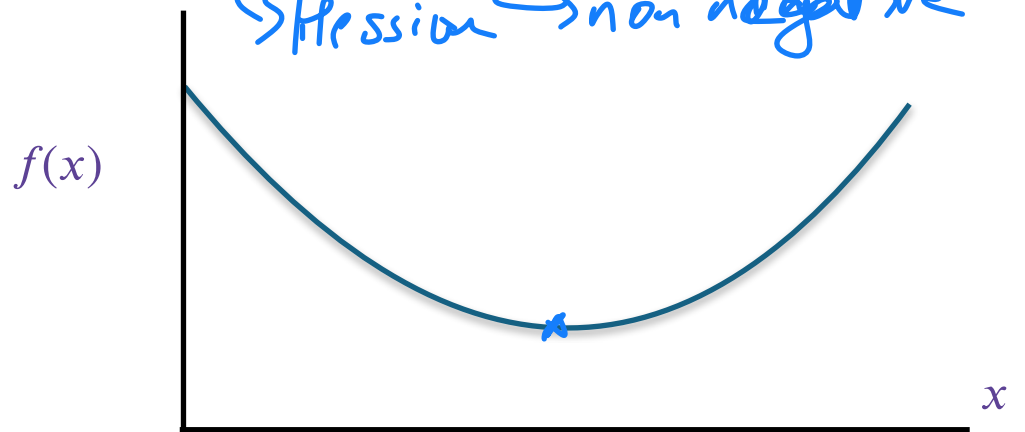
Convexity of twice-differentiable functions

Convexity of twice-differentiable functions

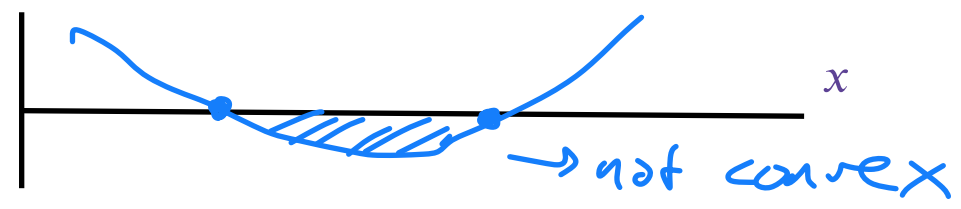
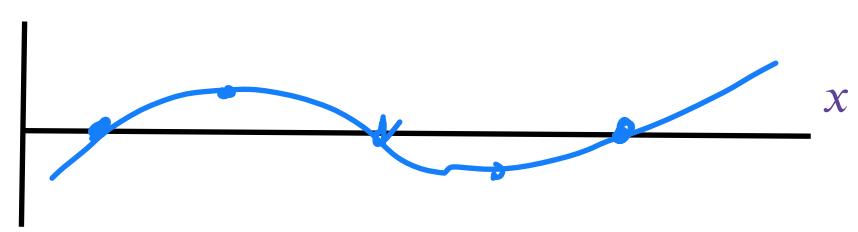
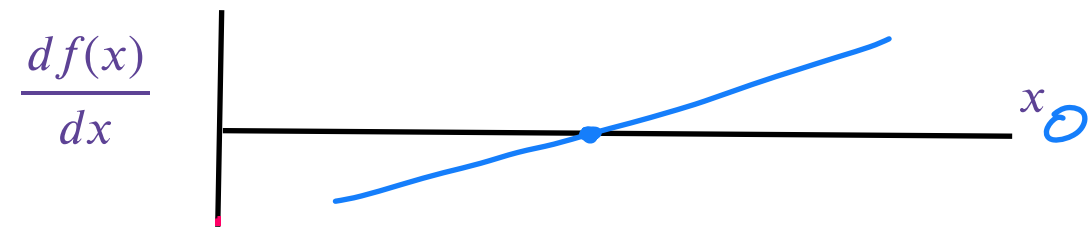
★ → should this be IA?
 YES.

A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ that is twice-differentiable everywhere is convex if $\nabla^2 f(x) \succeq 0$ for all $x \in \text{dom}(f)$

→ Hessian → non negative



$f''(x) \geq 0 \forall x \Leftrightarrow f$ is convex



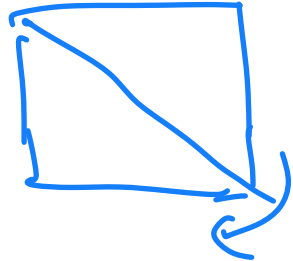
Convexity of twice-differentiable functions

A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ that is twice-differentiable everywhere is convex if $\nabla^2 f(x) \succeq 0$ for all $x \in \text{dom}(f)$

Hessian

$$\nabla^2 f(x)_{ij} = \frac{d^2 f(x)}{dx_i dx_j} = A \in \mathbb{R}^{d \times d}$$

symmetric $A_{ij} = A_{ji}$
always square
&
symmetric



A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ that is twice-differentiable everywhere is convex if $\nabla^2 f(x) \succeq 0$ for all $x \in \text{dom}(f)$

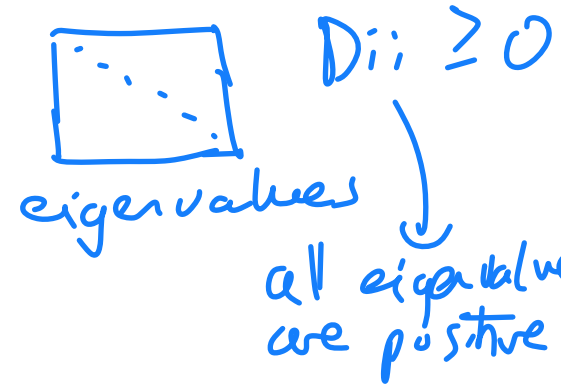
Positive semi-definite

$$\underline{A \succeq 0} \iff \forall v \in \mathbb{R}^d, \underline{v^T A v \geq 0}$$

$$\iff \underline{D_{ii} \geq 0 \text{ where } A = Q^T D Q}$$

eigen decomposition

$$A = Q^T D Q$$



Assume eigenvalues non-negative

$$D_{ii} \geq 0 \forall i$$

$$\underline{v^T A v} = v^T Q^T D Q v$$

$$\text{let } u = Qv$$

$$= u^T D u$$

$$= \sum_{i=1}^d u_i^2 D_{ii} \geq 0$$

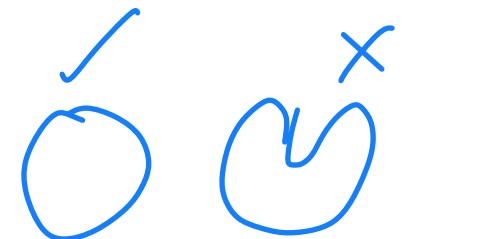
positive ≥ 0 bc squared

positive ≥ 0 by assumption

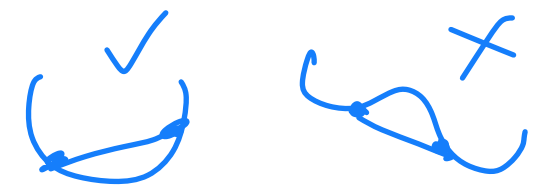
→ Check if a loss that you're given is convex

Recap: Definitions of convexity

sets
A set $K \subset \mathbb{R}^d$ is convex if $(1 - \lambda)x + \lambda y \in K$ for all $x, y \in K$ and $\lambda \in [0, 1]$



line above func
A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if $f((1 - \lambda)x + \lambda y) \leq (1 - \lambda)f(x) + \lambda f(y)$ for all $x, y \in \mathbb{R}^d$ and $\lambda \in [0, 1]$



epigraph
A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if the set $\{(x, t) \in \mathbb{R}^{d+1} : f(x) \leq t\}$ is convex



1st order Taylor expansion
A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ that is differentiable everywhere is convex if $f(y) \geq f(x) + \nabla f(x)^\top (y - x)$ for all $x, y \in \text{dom}(f)$



slope line should be below func

Hessian
A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ that is twice-differentiable everywhere is convex if $\nabla^2 f(x) \succeq 0$ for all $x \in \text{dom}(f)$



Example: Ridge regression

How to prove convex?

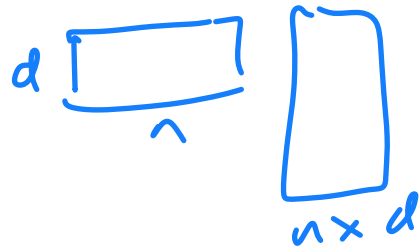
$$\operatorname{argmin}_w \underbrace{\|y - Xw\|_2^2 + \lambda \|w\|_2^2}_{f(w)}$$

$$\operatorname{argmin}_{\|w\|_2^2 \leq \rho} \|y - Xw\|_2^2$$

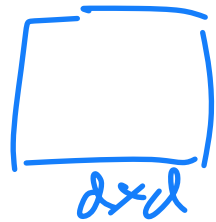
Let's use Hessian (easiest)

$$\nabla_w f(w) = 2X^T(Xw - y) + 2\lambda w$$

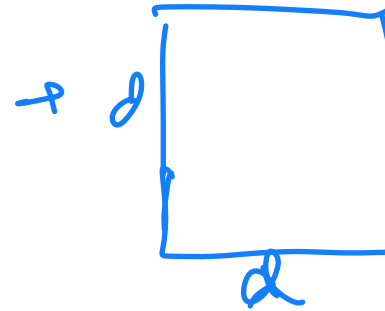
$$\nabla_w^2 f(w) = 2X^T X + 2\lambda I$$



PSD



$\lambda \geq 0$
positive



PD

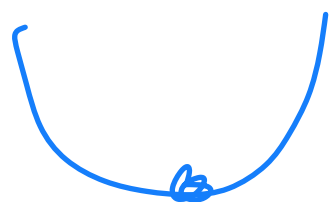
square
&
symmetric

Why $X^T X$ PSD?

$v \in \mathbb{R}^n$

$$v^T \underbrace{X^T X}_{u} v = u^T u = \underbrace{\|u\|^2} \geq 0$$

≥ 0

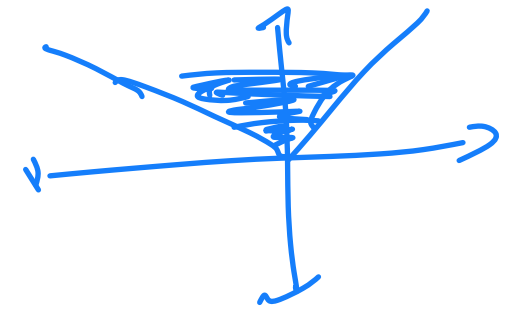


Why convex good?

local minima = global minimum

Example: Lasso

\rightarrow not twice differentiable

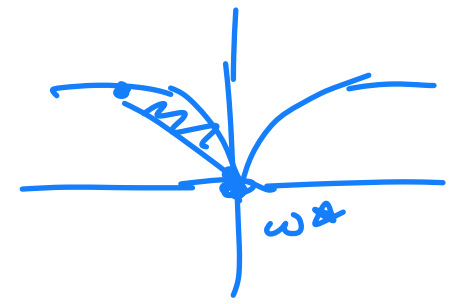


$$\text{argmin}_w \underbrace{\|y - Xw\|_2^2 + \lambda \|w\|_1}_{f(w)}$$

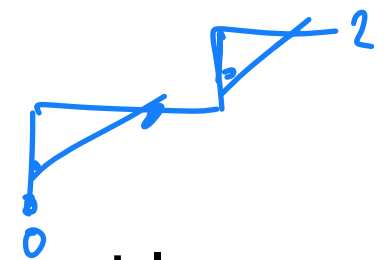
Why not directly solve $\text{argmin}_w \|y - Xw\|_2^2 + \lambda \text{card}(w)$?

For some $\text{argmin}_w \|y - Xw\|_2^2 + \lambda \|w\|^{1/2}$

$$+ \lambda \sum_{i=1}^d \sqrt{|w_i|}$$



$\text{card}(w) = \text{count}(w) \text{ where } w_i > 0$



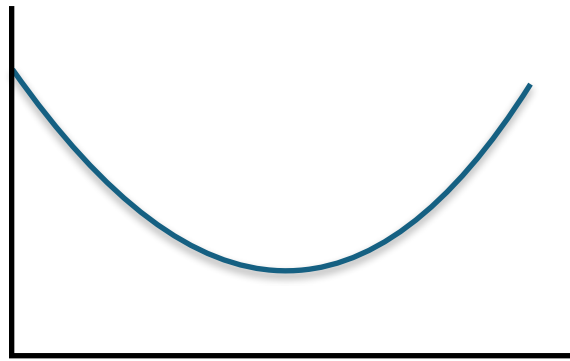
Can interpret lasso as convex relaxation of cardinality objective

Convexity and gradient descent

→ efficient to optimize (line of gradient descent)

- All local minima are global minima

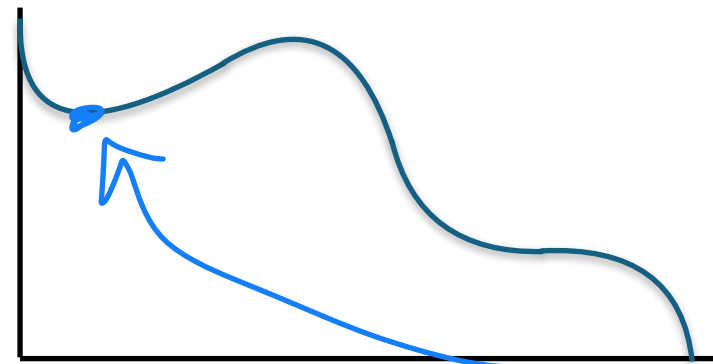
Convex function



Stationary points with $\nabla f(x) = 0$ are global minima

↙ ↘

Non-convex function



Stationary points with $\nabla f(x) = 0$ could be a local minima, a local maxima, or a saddle point

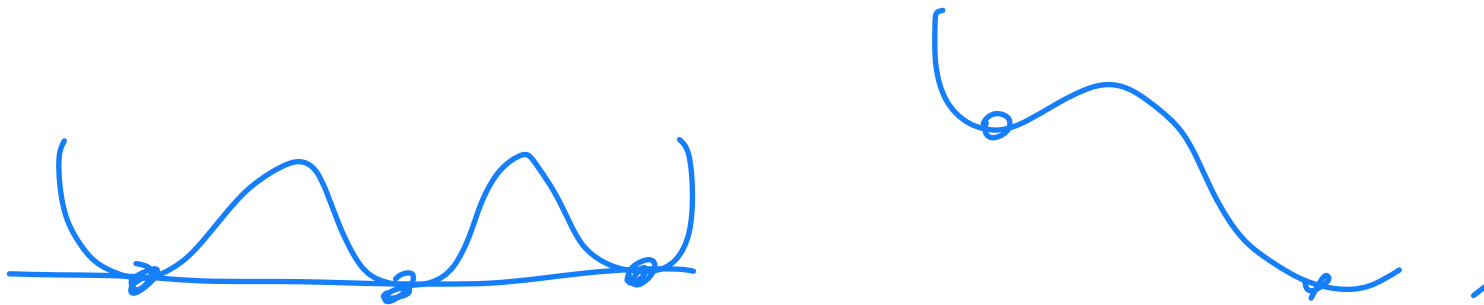
- Won't get stuck navigating the parameter constraints

Convexity and gradient descent

- Convexity \Rightarrow all local minima are global minima
 - All local minima are global minima \Rightarrow ?
- 

Convexity and gradient descent

- You can always run gradient descent whether $f(w)$ is convex or not!
- But if $f(w)$ is convex, we have guarantees on converging to the global minimum
- Linear regression, ridge regression, Lasso \rightarrow all convex!



Lecture plan

- Gradient descent algorithm + examples

- Theoretical analysis

- When does it work?

- How quickly does it converge? ← we are here

- How do we choose a step size?

- Key idea: Convexity

→ convexity

- Not tested on proof details, but concepts are important & practical

Convergence analysis steps

$$f(w_{t+1}) \leq f(w_t)$$

- 1. Study single iteration: $f(w_{t+1})$ vs. $f(w_t)$
2. Piece iterations together to study how they converge

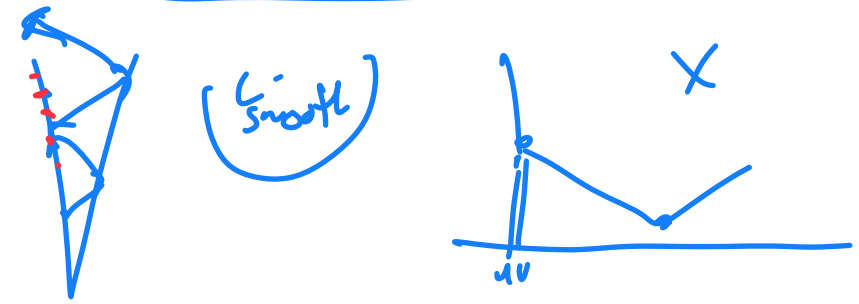
→ w^* ?
 $\|\nabla f(w)\| \approx 0$

Single-iteration progress bound 1

Don't need convexity (yet). → twice differentiable
 Assume f is C^2 , and gradient of f is Lipschitz continuous. which means

There exists L such that:

For all u, v , $\|\nabla f(u) - \nabla f(v)\| \leq L\|u - v\|$.
↙ gradient changing smoothly



For all w , $\nabla^2 f(w) \preceq LI$. ↔ $L I - \nabla^2 f(w) \succeq 0$
Hessian has bounded eigenvalues PSD

L is how fast your gradient can change

For any u, v , $v^T \nabla^2 f(u) v \leq L\|v\|^2$.

↪ use L to pick step size
small L → larger n
large L → smaller n

Single-iteration progress bound

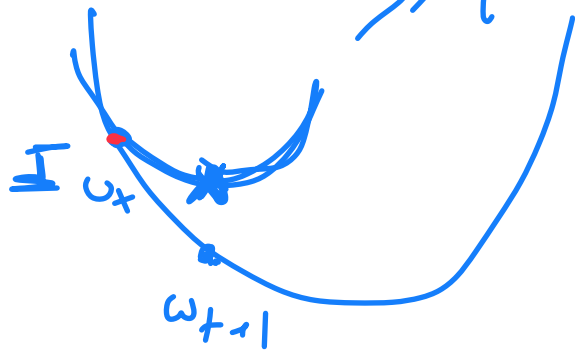
For any u, v , $v^T \nabla^2 f(u) v \leq L \|v\|^2$.

Take Taylor expansion:

$$\frac{f(\omega_{t+1}) \leq f(\omega_t)}{\omega_t, f(\omega_t), \nabla f(\omega_t), f(\omega_{t+1})}$$

Recall: 2nd order Taylor expansion around point x_0
 $f(x) \leq f(x_0) + f'(x_0)(x - x_0) + \frac{1}{2} f''(x_0)(x - x_0)^2$

→ quadratic upper bound



$f(\omega_{t+1}) \leq$ Taylor expansion

Single-iteration progress bound

For any u, v , $v^T \nabla^2 f(u) v \leq L \|v\|^2$. ①

$u \in [\omega_{t+1}, \omega_t]$
 $u \in \alpha \omega_t + (1-\alpha)\omega_{t+1}$

Take Taylor expansion:

$$\begin{aligned} \underline{f(\omega_{t+1})} &= f(\omega_t) + (\omega_{t+1} - \omega_t)^T \nabla f(\omega_t) + \frac{1}{2} (\omega_{t+1} - \omega_t)^T \nabla^2 f(u) (\omega_{t+1} - \omega_t) \\ &\leq f(\omega_t) + (\omega_{t+1} - \omega_t)^T \nabla f(\omega_t) + \frac{L}{2} \|\omega_{t+1} - \omega_t\|^2 \end{aligned}$$

by assumption ①

Please fill out course feedback!
<https://uw.iasystem.org/survey/306248>

Let $\Delta = \underline{w_{t+1} - w_t}$ \rightarrow the step we take

Pick Δ to minimize

$$\nabla f(w_t)^T \Delta + \frac{L}{2} \|\Delta\|_2^2$$

\rightarrow take deriv, set to zero:

$$\nabla f(w_t) + L \Delta = 0$$

$$\Delta = -\frac{1}{L} \nabla f(w_t)$$

recall: gradient descent

$$w_{t+1} = w_t - \eta \nabla f(w_t)$$

$$w_{t+1} = w_t - \frac{1}{L} \nabla f(w_t)$$

$$\eta = \frac{1}{L} \leftarrow$$

Single-iteration progress bound

$$w_{t+1} = w_t - \frac{1}{L} \nabla f(w_t) \rightarrow \eta = \frac{1}{L} \text{ C.D.}$$

recall: $w_{t+1} - w_t = -\eta \nabla f(w_t)$
①

$$\begin{aligned} f(w_{t+1}) &\leq f(w_t) + \nabla f(w_t)^\top (w_{t+1} - w_t) + \frac{L}{2} \|w_{t+1} - w_t\|_2^2 \\ &\leq f(w_t) - \eta \|\nabla f(w_t)\|^2 + \frac{1}{2} \eta^2 \|\nabla f(w_t)\|^2 \end{aligned}$$

// by ①

$$f(w_{t+1}) \leq f(w_t) - \frac{1}{2L} \|\nabla f(w_t)\|^2$$

→ one step progress bound
↳ amount of guaranteed progress every step

more progress when?
→ large gradients
→ smooth function (small L)

Single-iteration progress bound

When $\eta = \frac{1}{L}$,

$$f(w_{t+1}) \leq f(w_t) - \frac{1}{2L} \left\| \nabla f(w_t) \right\|_2^2$$

for any η in $0 < \eta < \frac{2}{L}$
still make progress

✓ valid

Same argument shows any $\eta < \frac{2}{L}$ will decrease f .

Convergence analysis steps

1. Study single iteration: $f(w_{t+1})$ vs. $f(w_t)$
2. Piece iterations together to study how they converge

Convergence rate of gradient descent

Converge to
a
minima
or
saddle point

$$\|\nabla f(w_t)\| \rightarrow 0$$

assume
convexity



$$f(w_t) - f(w^*) \rightarrow 0$$

the minimum
you converge
to is the
global
minimum

Convergence rate of gradient descent

For some ϵ , how many iterations before $\|\nabla f(w_t)\|^2 \leq \epsilon$?

Assumptions:

- Gradient is Lipschitz continuous (as before) ✓
- Step size is small enough (assume $1/L$) ✓
- f is bounded below by $f(w^*)$

$f(w^*) \geq c$ ✓
(otherwise, ϵ could have ∞ steps) MSE

Proof sketch:

- Each iteration decreases f by at least $\frac{1}{2L} \|\nabla f(w_t)\|^2$ ✓
- Can't decrease below $f(w^*)$ ✓
- So $\|\nabla f(w_t)\|^2$ must be decaying fast enough

Convergence rate of gradient descent

$$f(w_{t+1}) \leq f(w_t) - \frac{1}{2L} \|\nabla f(w_t)\|_2^2$$

(*) previously shown

$$\sum_{t=0}^T \|\nabla f(w_t)\|_2^2 \leq 2L \sum_{t=0}^T [f(w_t) - f(w_{t+1})]$$

// telescoping sum

$$= 2L [f(w_0) - \cancel{f(w_1)} + \cancel{f(w_1)} - \cancel{f(w_2)} + \dots + \cancel{f(w_{T-1})} + f(w_T)]$$

$$= 2L [f(w_0) - \underline{f(w_T)}]$$

$$\leq \underline{2L [f(w_0) - f(w^*)]}$$

// if you have convexity

← initial gap (depends on w_0)

analyze LHS on next slide



LMS?

$$\sum_{t=0}^T \|\nabla f(\omega_t)\|_2^2 \leq 2L [f(\omega_0) - f(\omega^*)]$$

$$T \cdot \min_{t=0}^T \|\nabla f(\omega_t)\|_2^2 \leq 2L [f(\omega_0) - f(\omega^*)]$$

$$\rightarrow \underline{\underline{\min \|\nabla f(\omega_t)\|_2^2}} \leq \frac{2L}{T} [f(\omega_0) - f(\omega^*)] \leq \varepsilon$$

$$\sum_{t=0}^T \|\nabla f(\omega_t)\|_2^2 \geq T \cdot \min \|\nabla f(\omega_t)\|_2^2$$

Convergence rate of gradient descent

$\rightarrow L$ -smooth, smooth or $\mu > 0$, faster convergence

$$T \geq \frac{2L(f(w_0) - f(w^*))}{\epsilon}$$

\rightarrow initial gap

\leftarrow closeness to $f(w^*)$

Gradient descent requires

$$T = O(1/\epsilon) \text{ iterations}$$

to achieve $\|\nabla f(w_t)\|^2 \leq \epsilon$

$\exists C$ for $\epsilon > \bar{\epsilon}$ s.t.

$$T_\epsilon \leq \frac{C}{\epsilon}$$

Convergence rate of gradient descent

$$\|\nabla f(w_t)\| \leq \epsilon$$

$T_\epsilon = O(1/\epsilon)$
if L -smooth



$$f(w_t) - \underline{f(w^*)} \leq \epsilon$$

if L -smooth & convex
 $T_\epsilon = O(1/\epsilon)$

Strongly convex

$$\nabla^2 f \succeq \mu I$$

$$\hookrightarrow T_\epsilon = O(\log \frac{1}{\epsilon})$$

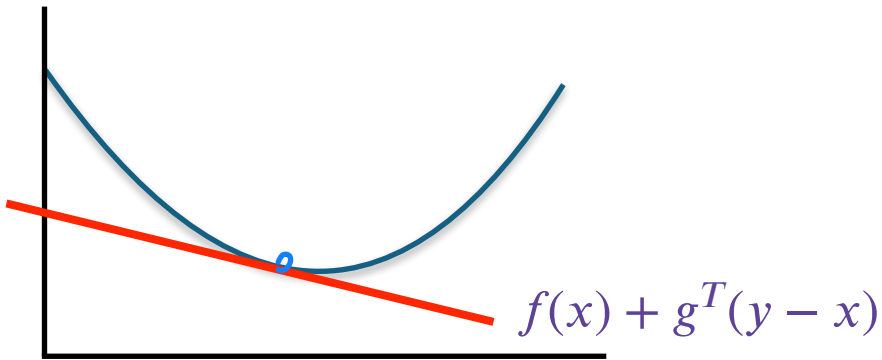
Lasso revisited

Subgradients

→ set of $g \in \mathbb{R}^d$ s.t.

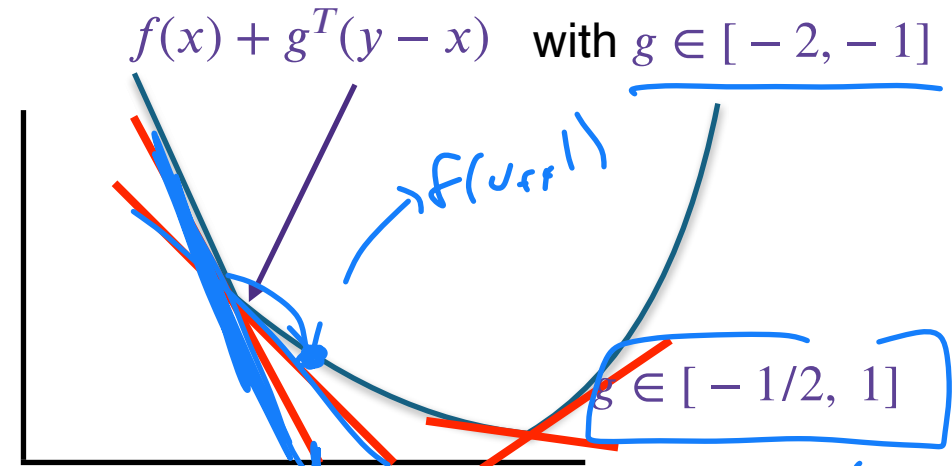
A vector $g \in \mathbb{R}^d$ is a **subgradient** at x if it satisfies $f(y) \geq f(x) + g^T(y - x)$ for all $y \in \mathbb{R}^d$

Smooth convex function



Gradient is unique sub-gradient
Minimum at points where gradient is 0

Non-smooth convex function



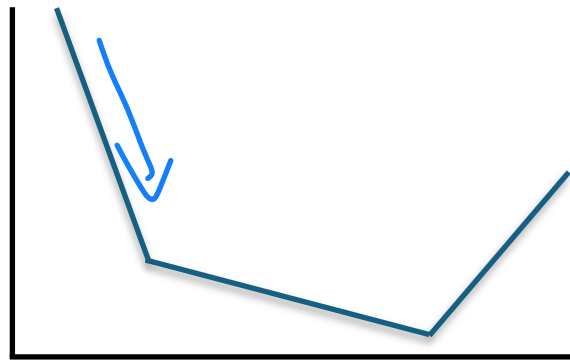
Minimum achieved at points where sub-gradient set includes 0 vector

→ 0 is in this set this point is the min

Subgradient descent for non-smooth functions

For each t ,

Find any subgradient g_t , then set $w_{t+1} \leftarrow w_t - \eta_t g_t$



Works on non-smooth convex functions

Slower compared to smooth convex functions

Gradients don't get smaller near global minima

- Instead of last iterate w_t , keep track of best one
- Step size needs to decrease with t

$$w_t \sim f(w)$$

Stochastic gradient descent (SGD)

$$\hat{w} = \operatorname{argmin}_w \frac{1}{n} \sum_{i=1}^n \ell_i(w)$$

MSE

Gradient descent: $w_{t+1} = w_t - \eta \nabla_w \left(\frac{1}{n} \sum_{i=1}^n \ell_i(w) \right) \Big|_{w=w_t}$

all examples

Stochastic gradient descent: $w_{t+1} = w_t - \eta \nabla_w \ell_{I_t}(w) \Big|_{w=w_t}$

I_t drawn uniformly at random from {1, ..., n}
1 example at a time

→ n times faster per iteration! → *n practice parallelization*

→ And can even be better minimizer.
→ *function value decreases faster in terms of gradient computation*

Minibatch stochastic gradient descent

→ workhorse of modern DL

"SGD"

1 sample

$3L$

s^1 batch

GD

η → whole dataset

$B \sim \text{i.i.d. data}$

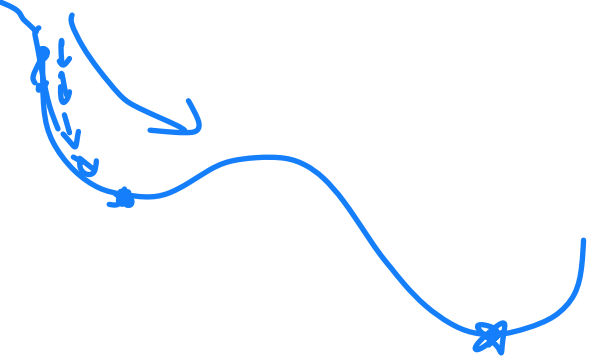
- Instead of one iterate, average B stochastic gradients together
- Advantages:
 - Smaller variance (by $1/B$) (vs 1-sample SGD)
 - Parallelization: Each gradient in the minibatch can be computed in parallel
- This is very widely used!



Summary

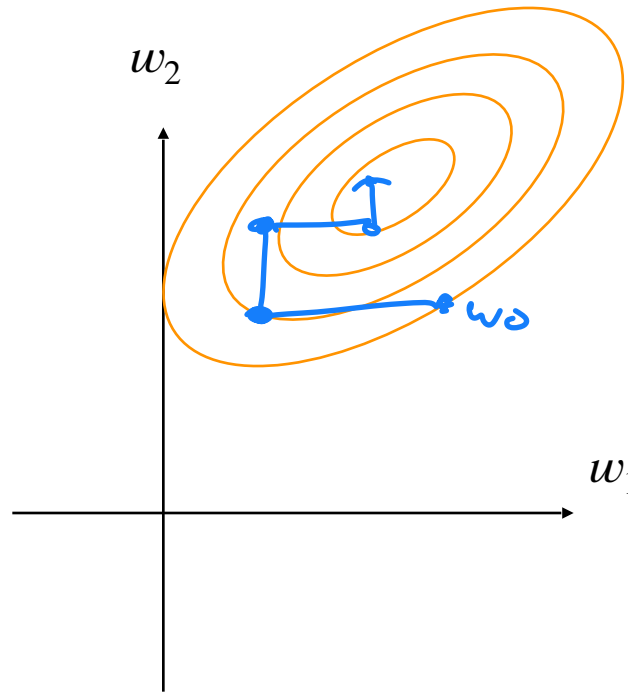
- Closed form -> iterative methods
- (Minibatch stochastic) gradient descent as a general-purpose optimizer
- Key theoretical tool: Convexity → local minima → global minimum
- • Many many variants. Highly active research area!
 - Schedulers
 - • Adaptive step sizes
 - Momentum
 - Higher-order methods Hessian
 - Non-convex analysis
 - ...

ADAM → 300k citations



Bonus: Coordinate descent

- ↳ optimize 1 weight coordinate at a time
- ↳ solve 1D optimization each time



eventually reach
 w^* if
convex &
differentiable

Coordinate descent for lasso

$$\hat{w} = \underset{w}{\operatorname{argmin}} \frac{1}{2} \sum_{i=1}^n (y_i - x_i^\top w)^2 + \lambda \|w\|_1$$

$$\begin{aligned} & \frac{d}{dw_k} f(w) \\ &= \sum_{i=1}^n (x_i^\top w - y_i) x_{ik} + \lambda \operatorname{sign}(w_k) \\ &= \sum_{i=1}^n \left(\sum_{j \neq k} x_{ij} w_j + x_{ik} w_k - y_i \right) x_{ik} + \lambda \operatorname{sign}(w_k) \\ &= \sum_{i=1}^n \left(\sum_{j \neq k} x_{ij} w_j - y_i \right) x_{ik} + w_k \sum_{i=1}^n x_{ik} + \lambda \operatorname{sign}(w_k) \ni 0 \\ & \dots \end{aligned}$$

Further reading

- Example gradient descent code on class website
- Boyd and Vandenberghe, Convex Optimization <https://stanford.edu/~boyd/cvxbook/>
- Mark Schmidt's CPSC 540 notes: <https://www.cs.ubc.ca/~schmidtm/Courses/540-W18/L4.pdf>
- 3Blue1Brown