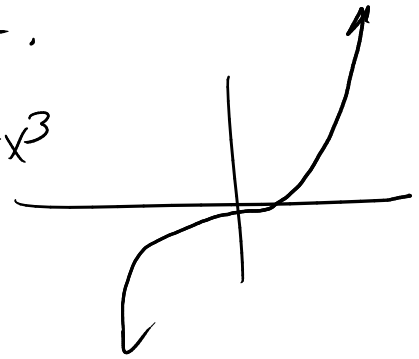


Non-convex optimization.

When objective is non-convex, generally the best we can hope for is to find a stationary point: $\|\nabla f\| = 0$, or an ϵ -approx stationary point, $\|\nabla f\| \leq \epsilon$.

Stationary points need not be local min: $y = x^3$



Recall Lemma 7 from last class:

Lemma 7. Let f be β -smooth (not nec. cvx). If $\alpha_f \equiv \alpha < \frac{1}{\beta}$, then

$$\text{G.D. satisfies } \forall T \geq 1 \quad \min_{t < T} \|\nabla f(w_t)\|^2 \leq \frac{2}{\alpha T} (f(w_0) - f(w_T))$$

So, GD on smooth objectives efficiently finds stationary points.

There are many problems where all stationary points are global optima.

Let's consider a few: Let $X^* := \{x^* : f(x^*) = \min_x f(x)\}$.

Let $\Pi(x) :=$ projection of x onto X^* .

① (Strong) convexity: for some $\mu > 0$ ($\mu > 0$: strong),
 $\forall x, y, \quad f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\mu}{2} \|x - y\|^2$

② weak strong convexity: for some $\mu > 0$, if $f^* = \min_x f(x)$,
 $\forall x, \quad f^* \geq f(x) + \langle \nabla f(x), \pi(x) - x \rangle + \frac{\mu}{2} \|x - \pi(x)\|^2$.

③ Restricted secant inequality:
 $\forall x, \quad \langle \nabla f(x), x - \pi(x) \rangle \geq \mu \|x - \pi(x)\|^2$.

④ Error bound inequality:
 $\forall x, \quad \|\nabla f(x)\| \geq \mu \|x - \pi(x)\|$.

⑤ Polyak-Łojasiewicz: (PL)
 $\forall x, \quad \frac{1}{2} \|\nabla f(x)\|^2 \geq \mu (f(x) - f^*)$

Clear that $PL \Rightarrow$ all local min are global. However:

Thm (Karimi et al '16): If f is β -smooth, then:

① \Rightarrow ② \Rightarrow ③ \Rightarrow (④ \Leftrightarrow ⑤).

So, PL inequality is quite crucial.

Thm. Suppose f is β -smooth, μ -PL, and $\exists x^*$ s.t. $f(x^*) = \min_x f(x)$.

Then GD w/ $\alpha \leq \frac{1}{\beta} \wedge \frac{1}{\mu}$ $f(x_k) - f^* \leq (1 - \alpha\mu)^k \cdot (f(x_0) - f^*)$.

Pf. By defn, $x_{k+1} = x_k - \alpha \nabla f(x_k)$.

By β -smoothness,
 $f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\beta}{2} \|x - y\|^2 \quad \forall x, y$.

$$\begin{aligned} \Rightarrow f(x_{k+1}) - f(x_k) &\leq \langle \nabla f(x_k), x_{k+1} - x_k \rangle + \frac{\beta}{2} \|x_{k+1} - x_k\|^2 \\ &= -\alpha \|\nabla f(x_k)\|^2 + \frac{\beta\alpha^2}{2} \|\nabla f(x_k)\|^2 \\ &\leq -\frac{\alpha}{2} \|\nabla f(x_k)\|^2 \quad \text{as } \alpha \leq \frac{1}{\beta}. \end{aligned}$$

$$\begin{aligned} \Rightarrow f(x_{k+1}) - f^* &\leq (1 - \alpha\mu) f(x_k) + \alpha\mu f^* - f^* \quad \text{by PL} \\ &= (1 - \alpha\mu) f(x_k) - (1 - \alpha\mu) f^* = (1 - \alpha\mu) (f(x_k) - f^*). \end{aligned}$$

\Rightarrow Unrolling, we see $f(x_k) - f^* \leq (1 - \alpha\mu)^k (f(x_0) - f^*)$. \square

So, with PL on smooth objectives, can get "linear convergence" to get $f(w_t) - f^* \leq \epsilon$, only need $O(\log \frac{1}{\epsilon})$ iterations.

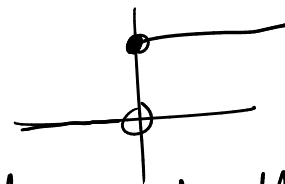
Exercise Show that there exists a β -smooth convex function on \mathbb{R} for which PL is not satisfied.

So establishing PL is clearly very powerful. Unsurprisingly, it often does not hold, and if it does, can be very hard to show.

Moreover, for many problems of interest, we do not have β -smoothness. eg, neural networks with ReLU activations,

since $\psi(t) = \max(0, t)$ has $\psi'(t) = \mathbb{1}(t \geq 0)$ (subderivative).

Clearly ψ' is not Lipschitz, so ψ is not β -smooth.



We'll next describe some alternative non-convex approaches which

- (1) allow for non-smooth objectives, i/cr
- (2) allow for situations where "all local min are good" (but not nec. global opt.)

Surrogate objectives.

In classification, we are most interested in the 0-1 loss:

$$\mathbb{1}(y \neq f(x)) = \begin{cases} 1: y \neq f(x), \\ 0: y = f(x) \end{cases} = \mathbb{1}(y \cdot f(x) < 0) \quad \text{for } y, f(x) \in \{-1, 1\}$$

[In general: $\mathbb{1}(A) = \begin{cases} 1 \text{ if } A \text{ occurs,} \\ 0 \text{ otherwise.} \end{cases}$]

Unfortunately, $\mathbb{1}(\cdot)$ is discontinuous, piecewise constant, so its derivative is zero a.e. \rightarrow gradient descent is useless.

So instead, we minimize convex surrogates for the 0-1 loss:

Call $l(t)$ a convex surrogate for $g(t)$ if:

①. l is convex

②. $g(t) \leq l(t) \quad \forall t$.

Thus, minimizing l implies minimizing g .

We will typically also require l is differentiable/smooth so we can apply ^{s.d.} _{to l.}

Common convex surrogates to 0-1 loss $\mathbb{1}(t < 0) = \begin{cases} 1: t < 0, \\ 0: t \geq 0 \end{cases}$:

• $\exp(-t)$ "exponential loss"

• $\max(0, 1-t)$ "hinge loss"

• $\log(1 + \exp(-t))$ "logistic / cross entropy"

(note ② not satisfied, but $\mathbb{1}(t < 0) \leq 2 \cdot \log(1 + \exp(-t))$)

More generally, there are a large family of surrogates to 0-1 loss.

Let $l(t)$ be convex, decreasing, and twice differentiable. Then:

Assume $l'(0) \neq 0$.

$$\textcircled{1} \quad l''(t) \geq 0 \quad (\text{so } -l' \text{ is decreasing})$$

$$\textcircled{2} \quad -l'(t) \geq 0.$$

Then: $\{t < 0\} = \{-l'(t) \geq -l'(0)\}$ by $\textcircled{1}$

If we can show $-l'(t) \rightarrow 0$, then by

$$\frac{1}{n} \sum_{i=1}^n \mathbb{1}(y_i f(x_i) < 0) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(y_i \cdot f(x_i) < 0)$$
$$= \hat{\mathbb{P}}(y \cdot f(x) < 0)$$

$$= \hat{\mathbb{P}}(-l'(y \cdot f(x)) \geq -l'(0))$$

$$\leq \frac{\hat{\mathbb{E}}(-l'(y \cdot f(x)))}{-l'(0)}$$

by Markov's $\hat{\mathbb{E}}$
 $l'(0) \neq 0$.

$$= \frac{1}{-l'(0)} \cdot \frac{1}{n} \sum_{i=1}^n -l'(y_i f(x_i))$$

Thus, for any convex, decreasing, twice differentiable l s.t. $l'(0) \neq 0$,
 $-l'(t)$ is surrogate for 0-1 loss.

$$\textcircled{1} l(t) = \exp(-t) \Rightarrow -l'(t) = l(t).$$

$$\textcircled{2} l(t) = \log(1 + \exp(-t)) \Rightarrow -l'(t) = \frac{\exp(-t)}{1 + \exp(-t)} = \frac{1}{1 + \exp(t)}.$$

Here $-l'$ is not convex.

However, if we do G.D. on l (which is convex),

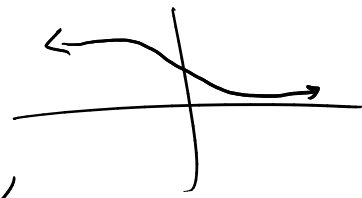
$$\text{i.e. G.D. on } \hat{L}(W) := \frac{1}{n} \sum_{i=1}^n l(y_i f(x_i; W)),$$

and we somehow found that the loss

$$\hat{G}(W) := \frac{1}{n} \sum_{i=1}^n -l'(y_i f(x_i; W)) \text{ were small, } \underline{\text{then}} \text{ we}$$

would know that the 0-1 loss is small as well.

(In the homework, we'll see why this is useful.)



Def. $f: \mathbb{R}^p$ is called (g, h, v) -proxy convex if, $\exists g, h: \mathbb{R}^p \rightarrow \mathbb{R}$ s.t.
 $\exists v$ s.t. $\forall w, \quad h(v) \geq g(w) + \langle \nabla f(w), v - w \rangle$.

Thm: If f is (g, h, v) -proxy-convex and $\|\nabla f(w)\| \leq L$, for all w ,
 Then for any $w(0), \eta > 0, T \geq 1$

$$\min_{t < T} g(w_t) \leq h(v) + \frac{\eta L^2}{2} + \frac{\|w(0) - v\|^2}{2\eta T}$$

Remark. For small step sizes η and large T , G.D. on f finds $g(w_t)$ which is close to $h(v)$.

useful if we can show there are families of related objective functions $\{g, h\}$ for which are "meaningful". i.e., $v = \arg \min f$,

$h(w) = \sqrt{f(w)}$ and $g(w) = \frac{1}{2} f(w)$ would imply

$$\min_{t < T} \frac{1}{2} f(w_t) \leq \min_w \sqrt{f(w)} + \epsilon.$$

NOT a global min, but a fun of it. If $\min_w f(w)$ is very small, then good.

In remainder, we will prove that GD converges to a global min for 2-layer leaky ReLU nets on linearly sep. data.

Two proof techniques:

- (1) Proxy convexity] lecture
- (2) Perceptron argument] homework

2-layer leaky nets: $w_j \in \mathbb{R}^d$, $W \in \mathbb{R}^{m \times d}$ (rows w_j),
 $a_j \in \{\pm 1/\sqrt{m}\}$, $m = \#$ neurons. $\varphi(t) = \max(t, \alpha t)$, $\alpha \in (0, 1]$

$$f(x; W) = \sum_{j=1}^m a_j \varphi(\langle w_j, x \rangle).$$

Only training first layer.

$$\ell(t) = \log(1 + \exp(-t)) \quad ; \quad \|x_i\| \leq 1$$

$$\exists v^* \in \mathbb{R}^d \text{ st. } \|v^*\| = 1, \forall i, \langle y_i x_i, v^* \rangle \geq \gamma, \gamma \in (0, 1].$$

GD on empirical risk $\hat{L}(W) = \frac{1}{n} \sum_{i=1}^n \ell(y_i f(x_i; W))$.

Remark We will consider parameterizing the objective for $\hat{L}(W)$ by a matrix W , while we typically think of vectors in \mathbb{R}^d for some d .

You can cast the problem to that form by considering $\text{vec}(W) \in \mathbb{R}^{d \cdot m}$

: stack
each column
of W one after
another.

Everything works out — can check via matrix calculus.

$$\bullet \quad \langle W, V \rangle := \text{tr}(W^T V) \quad (= \text{tr}(V^T W) = \text{tr}(W V^T))$$

\uparrow $\in \mathbb{R}^{m \times d}$ \uparrow $\in \mathbb{R}^{d \times m}$

$$= \langle \text{vec}(W), \text{vec}(V) \rangle_{\mathbb{R}^{dm}}$$

$$\bullet \quad \|W\|_F^2 = \|\text{vec}(W)\|_2^2$$

Lemma. If $g: \mathbb{R}^d \rightarrow \mathbb{R}$ is L -positively homogeneous, $(g(\alpha w) = \alpha^L g(w) \text{ for } \alpha \geq 0)$
then $\langle \nabla g(w), w \rangle = L g(w)$

pf: Exercise.

Remark: $\varphi(t) = \max(t, \alpha t)$ is 1-positively homogeneous.

Our goal will be to show $\hat{L}(w)$ is proxy convex for appropriate proxies g, h , i.e. want to show, for some g, h ,

$$\langle \nabla \hat{L}(w), w - v \rangle \geq g(w) - h(v).$$

$$\nabla \hat{L}(w) = \frac{1}{n} \sum_{i=1}^n \nabla \ell(y_i f(x_i; w))$$

$$= \frac{1}{n} \sum_{i=1}^n \ell'(y_i f(x_i; w)) y_i \nabla f(x_i; w).$$

Now: $f(x_i; w) = \sum_j a_j \varphi(\langle w_j, x \rangle)$ is 1-homog. since φ is. So,

$$\langle \nabla f(x_i; w), w \rangle = f(x_i, w).$$

$$\begin{aligned} \text{Thus: } \langle \nabla \hat{L}(w), w - v \rangle &= \frac{1}{n} \sum_i^n \ell'(y_i f(x_i; w)) y_i \langle \nabla f(x_i; w), w - v \rangle \\ &= \frac{1}{n} \sum_i^n \ell'(y_i f(x_i; w)) \cdot (y_i f(x_i; w) - y_i \langle \nabla f(x_i; w), v \rangle) \end{aligned}$$

So now, which v do we want to choose?

Calculus $\Rightarrow \nabla f(x; w) = \underbrace{D_x^w a x^T}_{\in \mathbb{R}^{m \times d}}, \quad D_x^w = \text{diag}(\varphi'(kw_j, x_j))$

$$\begin{aligned} y_i \langle \nabla f(x_i; w), v \rangle &= y_i \cdot \text{tr} \left((D_{x_i}^w a x_i^T)^T v \right) \\ &= y_i \text{tr} (x_i a^T D_{x_i}^w v) \\ &= y_i \text{tr} (a^T D_{x_i}^w v x_i) \\ &= y_i a^T D_{x_i}^w v x_i \\ &= \sum_j a_j \varphi'(kw_j, x_i) \langle v_j, y_i x_i \rangle. \end{aligned}$$

know: $\langle v^*, y_i x_i \rangle \geq \gamma > 0$. So let $v_j := \rho \cdot a_j \cdot v^*$, $\rho > 0$.

Since $\varphi'(t) \geq \alpha > 0 \quad \forall t$,

$$y_i \langle \nabla f(x_i; w), v \rangle = \sum_{j=1}^m a_j \varphi'(kw_j, x_i) \cdot \langle v_j, y_i x_i \rangle = \sum_{j=1}^m a_j^2 \cdot \rho \cdot \varphi'(kw_j, x_i) \langle v^*, y_i x_i \rangle \geq \frac{\rho}{m} \sum_{j=1}^m \alpha \gamma = \rho \alpha \gamma.$$

(Continuing from): Since l is decreasing, $-l'(z) < 0$. So,

$$-l'(y_i f(x_i; w)) \cdot y_i \langle \nabla f(x_i; w), V \rangle \geq -l'(y_i f(x_i; w)) \cdot \rho \alpha \gamma.$$

Thus,

$$\langle \nabla \hat{L}(w), w - V \rangle \geq \frac{1}{n} \sum_1^n l'(y_i f(x_i; w)) \cdot (y_i f(x_i; w) - \rho \alpha \gamma).$$

(since l is convex, $l'(z_1) \cdot (z_1 - z_2) \geq l(z_1) - l(z_2)$. \rightarrow

$$\rightarrow \geq \frac{1}{n} \sum_1^n [l(y_i f(x_i; w)) - l(\rho \alpha \gamma)].$$

$$= \hat{L}(w) - l(\rho \alpha \gamma).$$

Therefore, $\hat{L}(w)$ is $(\hat{L}, l(\rho \alpha \gamma), V)$ -proxy convex

with $V = V(\rho)$ having rows $v_j = \rho \alpha_j v^*$.

Since $\|\nabla f(x_i; w)\|_F = \|\text{Diag}(\psi'(kw_i, x_i) \odot x_i^+) \|_F$

$$\leq 1 \quad (\underline{\text{Exercise!}})$$

Thus the theorem guarantees that for any $T \geq 1$, $\eta > 0$, $\varepsilon > 0$,

$$\min_{t < T} \hat{L}(W_t) \leq l(\rho \alpha \gamma) + \frac{\eta}{2} + \frac{\|W(0) - V\|_F^2}{2\eta T}.$$

Since $l(t) \leq 2\exp(-t)$, taking $\rho = \alpha^{-1} \gamma^{-1} \log(\frac{6}{\varepsilon})$

guarantees $l(\rho \alpha \gamma) \leq \frac{\varepsilon}{3}$.

For $\eta \leq \frac{\varepsilon}{2}$, plus mems

$$\min_{t < T} \hat{L}(W_t) \leq \frac{2\varepsilon}{3} + \frac{\|W(0) - V\|_F^2}{2\eta T}.$$

For $T \geq 2\eta^{-1} \|W(0) - V\|_F^2 \varepsilon^{-1}$, $\min_{t < T} \hat{L}(W_t) \leq \varepsilon.$

All together:

Theorem. Suppose $\|x_i\| \leq 1$, $y_i \langle v^*, x_i \rangle \geq \gamma$ for all i for some $\|v^*\| = 1$.
 Then GD on 2-layer leaky ReLU net satisfies the following. For any $\varepsilon > 0$,
 if $\eta \leq \varepsilon/2$ then for $V \in \mathbb{R}^{m \times d}$ with rows $v_j = \alpha_j \alpha^{-1} \gamma^{-1} \log(\frac{6}{\varepsilon})$,
 if $T \geq 2\eta^{-1} \|W(0) - V\|_F^2 \varepsilon^{-1}$, then $\min_{t < T} \hat{L}(W_t) \leq \varepsilon.$



Remark I'm not aware of proofs with any of:

- ① - φ is ReLU rather than leaky ReLU
- ② - bias terms included
- ③ - training both first and second layer weights
- ④ - training > 2 layers net.

difficulty:
4 > 1 > 3 > 2.

If you can come up with a proof in any of these settings,
this would probably be sufficient for a NeurIPS paper! Requires:
(except maybe ②)

- ① Random initialization
 - ② Only assume linear separability of training data.
 - ③ No "NTK". Ideally allow for constant number of neurons.
- You can work on this instead of choosing paper for final project if you'd like.