Continuing our introduction to optimization, let's see an important application:

Least Squares  Let $f(x) = \dfrac{\|Ax - b\|_2^2}{2}$

i.e we want to minimize

Then we can compute

$$\nabla f(x) = A^T(Ax - b)$$

$$\nabla^2 f(x) = A^T A$$

From this we see that $f(x)$ is strongly convex iff $A^T A$ has full rank

What does gradient descent look like, specialized to this problem?

$$X_{t+1} = X_t - \eta A^T r_t$$

where $r_t \overset{\Delta}{=} Ax_t - b$

is the <u>residual</u>

# Variants

There are some important variations on vanilla GD

## Stochastic Gradient Descent

**What do you do if the function you want to minimize is too big to fit in memory?**

e.g. $R(f(x; w)) = \sum_{i=1}^{N} \ell(f(x_i; w), y_i)$

↑
one term for each example you want to label

For $t = 1$ to $T$

    Choose $i$ at random

    Set $w_{t+1} = w_t - n \nabla_w \ell(f(w_i; w), y_i)$

Can also choose a set of $k$ random exs., called a <u>minibatch</u>, to estimate overall gradient

# Nesterov Acceleration

Can also incorporate information about past gradients to converge even faster

For $t = 1$ to $T$

$$\text{Set } y_{t+1} = x_t - \eta \nabla f(x_t)$$

$$\text{Set } x_{t+1} = y_{t+1} + \mu \underbrace{(y_{t+1} - y_t)}_{\text{momentum}}$$

This achieves a squareroot improvement:

$\text{GD} \quad \sim \dfrac{\beta}{\alpha}$ iterations to reduce distance to optimum by a constant factor

$\text{AGD} \quad \sim \sqrt{\dfrac{\beta}{\alpha}}$ iterations

# Gradient Flow

For theoretical convenience sometimes it will be helpful to work with the continuous time limit — i.e. the ODE

$$\dot{x}(t) = -\nabla f(x(t))$$

This is the natural limit of taking the stepsize to zero

$$x_{k+1} = x_k - \eta \nabla f(x_k)$$

(rearranging $\rightarrow$

$$\frac{x_{k+1} - x_k}{\eta} = -\nabla f(x_k)$$

Taking a step back, let's say out loud what is so great about <u>convex</u> optimization:

<u>many</u> <u>different</u>
<u>algorithms</u>

GD
SGD $\sim\!\!\rightarrow$ same $x^*$,
AGD provably optimal
$\vdots$

But in non-convex optimization, like when we are training a deep net, all that goes out the window

We could try to apply GD/SGD/AGD... but what goes wrong?

Issue #1: Can't find the global minimum

E.g.

Usually we find a local minimum, and in general finding a global minimum is NP-hard

Issue #2: which point we reach depends on our initialization, the randomness of our algorithm

Usually there are tons of points that achieve zero risk

And which one we reach is important for generalization

(ref: "sharp" vs. "flat" minima)

So, the journey is more important than the destination

## Neural Tangent Kernels

Next, we will prove a fundamental result about fitting an <u>overparameterized</u> deep net to data

First, fix any training set of size $N$. Then:

"For a wide enough deep net, gradient descent finds a 0-error solution"

There are some caveats:

(1) the network needs to have depth $\geq 2$, right nonlinearities

② How wide is wide enough?
Arora et al. showed that $poly(N)$ suffices

③ This all depends on having the right scaling parameters for the random initialization

Historical Notes: This line of work was kicked off by an important work of Jacot, Gabriel and Hongler

In hindsight, concurrent works of Du et al. and Allen-Zhu et al. exploited the same principles

We will follow the "lazy training" framework of Chizat and Bach, which is an abstraction of what happens as you increase the width

The main ideas are:

(1) If we take the linear approximation to a deep network, if width > #examples, we can fit perfectly

(2) As we increase the width, we need to move less in the parameter space

$\Rightarrow$ Thus the Jacobian doesn't change too much

## The Gradient Flow

For convenience, let's suppress the data points:

$$f(\omega) \overset{\Delta}{=} \begin{bmatrix} f(x_1; \omega) \\ \vdots \\ f(x_N; \omega) \end{bmatrix}, \quad y = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$$

Throughout, we want to minimize the
squared loss:

$$R(\alpha f(w)) \triangleq \frac{1}{2} \| \alpha f(w) - y \|^2$$

Here $\alpha$ is a scaling parameter

<span style="color:red">increasing
the width</span> $\longleftrightarrow$ <span style="color:red">increasing
$\alpha$</span>

More on this later.

We will be interested in the behavior
of the gradient flow:

$$\dot{w}(t) = - \nabla_w R(\alpha f(w(t)))$$

$$= -\alpha J_t^T \nabla R(\alpha f(w(t)))$$

$$= -\alpha J_t^T (\alpha f(w(t)) - y)$$

where $J_t$ is the Jacobian at time $t$

i.e.
$$J_t \triangleq \begin{bmatrix} \nabla f(x_1; w(t))^T \\ \vdots \\ \nabla f(x_N; w(t))^T \end{bmatrix}$$

Assuming $w$ is a vector of $p$ parameters, this is an $N \times p$ matrix

Main Theorem (informal) If $\alpha$ is big enough, then

(1) $R(\alpha f(w(t))) \doteq R(\alpha f(w(0))) e^{-ct\alpha^2}$

(2) $\| w(t) - w(0) \| \leq \dfrac{\sqrt{R(\alpha f(w(0)))}}{\alpha}$

Proof Strategy: We will define an auxilliary flow, $u(t)$, using a linear approximation to $f$ and show $u(t)$ and $w(t)$ remain close

In particular, define

$$f_0(u) \stackrel{\Delta}{=} f(w(0)) + J_0(u - w(0))$$

which gives us

$$\dot{u}(t) = -\nabla_u R(\alpha f_0(u(t)))$$

$$= -\alpha J_0^T \nabla R(\alpha f_0(u(t)))$$

$$= -\alpha J_0^T (\alpha f_0(u(t)) - y) \quad (*)$$

**Lemma 1:** If $J_0 J_0^T$ is full rank, then $R(\alpha f_0(u(t)))$ goes to zero

**Proof:** Let's compute the change in predictions:

$$\frac{d}{dt} \alpha f_0(u(t)) = \frac{d}{dt} \left[ \alpha \left( f(w(0)) + J_0(u(t) - w(0)) \right) \right]$$

$$= \alpha J_0 \dot{u}(t)$$

Substituting $(*)$ we have:

$$= -\alpha^2 J_0 J_0^T (\alpha f_0(u(t)) - y)$$

Thus we conclude,

$$\frac{d}{dt} (\alpha f_0(u(t)) - y) = -\alpha^2 J_0 J_0^T (\alpha f_0(u(t)) - y)$$

Let's call $r(t) \overset{\Delta}{=} \alpha f(u(t)) - y$ the residual prediction. Then

$$\dot{r}(t) = -\alpha^2 J_0 J_0^T r(t)$$

$$\Rightarrow r(t) = e^{-\alpha^2 J_0 J_0^T t} r(0)$$

And this converges to zero by the full rank assumption ▨

Key definition: The matrix $K \overset{\Delta}{=} J_0 J_0^T$ is called the neural tangent kernel

In particular recall $J_0 J_0^T =$

$$\begin{bmatrix} \nabla f(x_1; w(0))^T \\ \vdots \\ \nabla f(x_N; w(0))^T \end{bmatrix} \begin{bmatrix} \nabla f(x_1; w(0))^T & \ldots & \nabla f(x_N; w(0))^T \end{bmatrix}$$

Thus $K(x_i, x_j) = \nabla f(x_i; w(0))^T \nabla f(x_j; w(0))$

It is known that under the right scaling, $f(\cdot; w(0))$ converges to a <u>Gaussian Process</u>

An interpretation

Actually the gradient flow for $u$ is something we have already seen, but in disguise

Recall, for least squares:

$$\min_x \frac{\|Ax - b\|^2}{2}$$

gradient descent turned into

$$X_{t+1} = X_t - \eta A^T r_t$$

$$\text{where } r_t = A X_t - b$$

Now we can rewrite this as a recurrence for $r_t$ instead

$$\Rightarrow A X_{t+1} = A X_t - \eta A A^T r_t$$

$$\Rightarrow A X_{t+1} - b = A X_t - b - \eta A A^T r_t$$

$$\Rightarrow r_{t+1} = r_t - \eta A A^T r_t$$

Thus the limit as we take the stepsize to zero is

$$\dot{r}(t) = -A A^T r(t)$$

This looks like what we had with

$$A \longleftrightarrow J_0$$

Explicitly, consider

$$\min_u \frac{\| f_0(u) - y \|^2}{2} =$$

$$\min_u \frac{\| J_0 u - \overbrace{( y + J_0 w(0) - f(w(0)))}^{b} \|^2}{2}$$

Aha! So that's what u is doing:

"u is solving a least squares
problem using the linear approx.
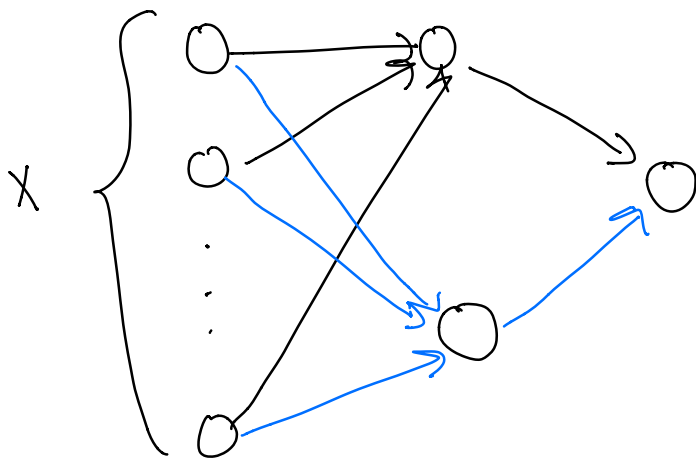to the deep net at initialization"

An Interlude

Why is ∝ an accurate abstraction
of what happens when we increase
the width?

Intuitively as $\alpha \to \infty$ we are zooming in at $w(0)$ so the linear approx. is more accurate

And as we increase the width, we need to move less to affect the output

## Thought Experiment

what if we increase the width by literally duplicating the network?

Call the new function

$$\hat{f}(x; w_1, w_2) = 2f(x; w)$$

$$\uparrow \quad \uparrow$$

duplicates of $w$

More over $\nabla_{w_1} \hat{f}(x; w_1, w_2) = \nabla_{w_2} \hat{f}(x; w_1; w_2)$

$$= \nabla_w f(x; w)$$

Thus our new gradient flow would correspond to zooming in by doubling $\alpha$

## Analyzing the Gradient Flow

Now let's return to $w(t)$, and, as we did before for $u(t)$, study the evolution in prediction space:

$$\frac{d}{dt} \alpha f(w(t)) = \alpha J_t \dot{w}(t)$$

$$= - \alpha^2 J_t J_t^T \nabla R\left(\alpha f(w(t))\right)$$

The main difference is this term now depends on $t$, so it is not exactly the gradient flow on the least squares obj.

Still it does not change much $\Rightarrow$ we can upper bound how fast risk decays

Notation

$$z(t) \triangleq \alpha f(w(t))$$

$$Q(t) \triangleq \alpha J_t J_t^T$$

Lemma 2 : Suppose

$$\dot{z}(t) = - Q(t) \nabla R\left(z(t)\right)$$

and let $\lambda = \inf_{t \in [0, \tau]} \lambda_{min}(Q) > 0$. Then

$$R(z(t)) \leq R(z(0)) \, e^{-2t\lambda}$$

for any $t \in [0, \tau]$

Detour: Our proof strategy will be to upper bound

$$\frac{d}{dt} R(z(t))$$

and show that it is at most what it would be if $Q(t) = \lambda I$, which again corresponds to least squares

**But how does this translate into a bound on $R(z(t))$?**

This is answered by Gronwall's Inequality

**Fact:** Consider an interval $[a,b]$ and suppose

$$\dot{u}(t) \leq \beta(t) u(t) \quad \forall t \in [0, \tau]$$

Then $u$ is upper bounded by the solution to

$$\dot{v}(t) = \beta(t) v(t)$$

assuming same boundary condition. In particular

$$u(t) \leq \underbrace{u(0) \, e^{\int_a^t \beta(s)\, ds}}_{v(t)}$$

Now, returning to Lemma 2:

**Proof:** Using the definition of the risk, we have

$$\frac{d}{dt} \frac{1}{2} \|z(t) - y\|^2 = \langle \dot{z}(t), z(t) - y \rangle$$

$$= \langle -\varphi(t) \nabla R(z(t)), z(t) - y \rangle$$

$$= \langle -\varphi(t)(z(t) - y), z(t) - y \rangle$$

$$\leq -2\lambda \frac{\|z(t) - y\|^2}{2}$$

Thus we have

$$\frac{d}{dt} R(z(t)) \leq -2\lambda R(z(t))$$

Now invoking Gronwall's Inequality, we have:

$$R(z(t)) \leq R(z(0)) \, e^{-\int_0^t 2\lambda \, ds}$$

$$= R(z(0)) \, e^{-2t\lambda}$$

we also have the following useful corollary:

Corollary: $\|z(t) - y\| \leq \|z(0) - y\| e^{-t\lambda}$

This follows because
$$\|z(t) - y\| = \sqrt{2R(z(t))}$$

by definition, and similarly for $\|z(0) - y\|$

Ultimately, we want to show we never venture too far from the initialization

This'll be the reason we can relate $w(t)$ and $u(t)$

Notation

$$v(t) \triangleq w(t)$$
$$g(v(t)) \triangleq \alpha f(w(t))$$
$$s(t) \triangleq \alpha \int t$$

Lemma 3: Suppose

$$\dot{v}(t) = -S(t)^T \nabla R(g(v(t)))$$

and let $Q(t) = S(t)S(t)^T$. Further suppose for all $t \in [0, \tau]$ that

$$\lambda I \preceq Q(t) \preceq \lambda_{max} I$$

Then we have

$$\| v(t) - v(0) \| \leq \frac{\sqrt{\lambda_{max}}}{\lambda} \| g(v(0)) - y \|$$

$$\leq \frac{\sqrt{2\lambda_{max} R(g(v(0)))}}{\lambda}$$

Proof: This follows from a direct computation
+ Corollary:

$$\| v(t) - v(0) \| = \left\| \int_0^t \dot{v}(t) ds \right\|$$

$$\leq \int_0^t \| \dot{v}(t) \| ds$$

$$= \int_0^t \| S(t)^T \nabla R(g(v(t))) \| ds$$

$$\leq \sqrt{\lambda_{max}} \int_0^t \| g(v(s)) - y \| \, ds$$

Now we can invoke the Corollary to get

$$\| v(t) - v(0) \| \leq \sqrt{\lambda_{max}} \int_0^t \| g(v(0)) - y \| e^{-s\lambda} \, ds$$

$$\leq \frac{\sqrt{\lambda_{max}}}{\lambda} \| g(v(0)) - y \|$$

$$= \frac{\sqrt{2 \lambda_{max} R(g(v(0)))}}{\lambda} \qquad \boxtimes$$

Putting it all together, recall that

$$Q(t) = \alpha^2 J_t J_t^T$$

So if we have bounds

$$(\Delta) \qquad \alpha^2 \sigma_{min} I \preceq Q(t) \preceq \alpha^2 \sigma_{max} I$$

we'd get

$$\text{①} \quad R(\alpha f(w(t))) \lesssim R(\alpha f(w(0))) e^{-2t\alpha^2 \sigma_{min}}$$

$$\text{②} \quad \|w(t) - w(0)\| \lesssim \frac{\sqrt{2\sigma_{max} R(\alpha f(w(0)))}}{\alpha \sigma_{min}}$$

But where do we get bounds like (Δ)?

If we make $\alpha$ large enough $\Longrightarrow$

we don't move far $\Longrightarrow$

the eigenvalues of $\phi(t)$ never change much from $\phi(0)$ $\Longrightarrow$

we can use a slightly degraded bound

(i.e. $\sigma_{min} \to \frac{\sigma_{min}}{2}$, $\sigma_{max} \to 2\sigma_{max}$)

See full details in Telgarsky's notes

Epilogue

# what random initialization puts you in the NTK regime?

Set each layer

$$X \longmapsto A_i X + b_i$$

to have $A_i^{n_i \times n_{i-1}} = \frac{1}{\sqrt{n_i}} W_i \longleftarrow$ standard Gaussian entries

Later in the course we will cover the <u>mean field view</u>, which describes the dynamics under a different scaling as a PDE

## So should you choose your architecture to be very (infinitely) wide?

No, but it does give rise to better kernel-based algorithms than we had before

e.g. ~65% accuracy on CIFAR-10

Actually you can analytically compute the kernel matrix without ever setting up the deep net

Explicit formulas in, e.g., Yang