

Lecture 4

- SGD Convergence Continued.
- Adam
- Neural Networks as Generalized Linear Models

Admin.

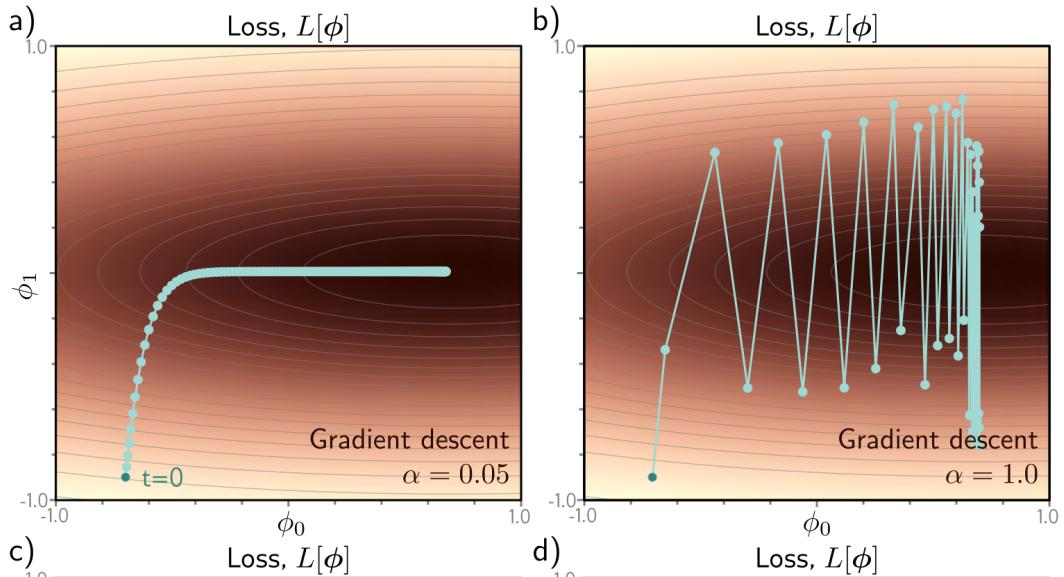
- ① Waitlist Cleared!!!
- ② 3 New staff
 - Joey, Lance, Kevin.
- ③ TA office hours start today.
- ④ 3 new discussion sections added!
- ⑤ OH + Study groups

Adam

Gradient Descent has a feature/bug → takes bigger steps in directions with large singular values and smaller steps in directions with small singular values.

- Advantage: early stopping \leftrightarrow Implicit regularization.
- Disadvantage: hard to choose a learning rate that
 - (a) makes enough progress in all directions
 - (b) doesn't blow up.

So how do you solve this? ... Brainstorm.



Prince 2025.

- Ideal world: compute SVD and normalize.

→ Too expensive :-)

→ Normalize the gradients

$$\vec{w}_{t+1} = \vec{w}_t - \eta \cdot \boxed{\nabla_{\vec{w}} L(\vec{w})}_{\vec{w}=\vec{w}_t}$$

Regular GD.

Instead:

$$\vec{m}_t \leftarrow \nabla_{\vec{w}} L(\vec{w}). \quad \text{gradient}$$

$$\vec{v}_t \leftarrow (\nabla_{\vec{w}} L(\vec{w}))^2 \quad \text{squared gradient.}$$

$$\vec{w}_{t+1} = \vec{w}_t - \eta \cdot \frac{\vec{m}_{t+1}}{\sqrt{\vec{v}_{t+1}} + \epsilon} \quad \Rightarrow \epsilon > 0 \text{ to avoid division by zero.}$$

$\sqrt{\vec{v}_{t+1}}$: Take positive square root.

$$\frac{\vec{m}_{t+1}}{\sqrt{\vec{v}_{t+1}}} = \text{sgn}(\nabla_{\vec{w}} L(\vec{w}))$$

"sign SGD"

What's the problem?

Once you are close you might keep overshooting and bouncing.

So can this be smoothed out?

→

$$\left\{ \begin{array}{l} \vec{m}_{t+1} = \beta \cdot \vec{m}_t + (1-\beta) \nabla_{\vec{w}} L(\vec{w}_t) \\ \vec{v}_{t+1} = \gamma \cdot \vec{v}_t + (1-\gamma) (\nabla_{\vec{w}} L(\vec{w}_t))^2 \end{array} \right.$$

$$\vec{w}_{t+1} = \vec{w}_t - \eta \frac{\vec{m}_{t+1}}{\sqrt{\vec{v}_{t+1}} + \epsilon}$$

But there can be a bias - Why?

Say $\beta = 0.9$. Then

$$\vec{m}_t = 0.9[0] + (0.1) \cdot \nabla_{\vec{w}} L(\vec{w}) \Big|_{\vec{w}=\vec{w}_0}$$

too small
Can lead to slow start.

"Bias Correction"

$$\tilde{\vec{m}}_{t+1} = \frac{\vec{m}_{t+1}}{1 - \beta^{t+1}}$$

$$\tilde{\vec{v}}_{t+1}^2 = \frac{\vec{v}_{t+1}}{1 - \gamma^{t+1}}$$

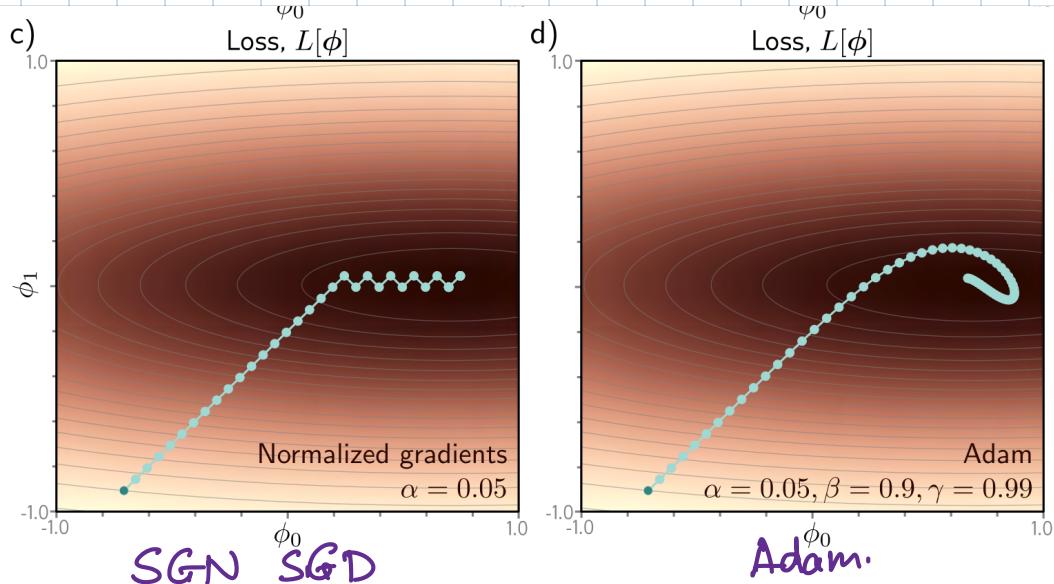
Adam:

$$\vec{w}_{t+1} = \vec{w}_t - \eta \cdot \frac{\tilde{\vec{m}}_{t+1}}{\sqrt{\tilde{\vec{v}}_{t+1}} + \epsilon}$$

Usually: Used with SGD, so compute \vec{m}, \vec{v} from mini-batches

Normalization helps with exploding / vanishing gradients

→ Even if the gradient is getting tiny, the size of the step is largely dictated by the learning rate.



Prince 2025

Memory:

SGD

SGD + Momentum

x 2

Adam

x 3

AdamW : Adam with weight decay.

Recall: Ridge regression:

$$\text{GD update: } \vec{w}_{t+1} = (1 - 2\eta\lambda)\vec{w}_t - \eta \nabla L_w(\vec{w}_t).$$

pushing weights down.

In Adam: magnitude of gradient keeps being normalized

→ So adding a regularizer to the loss doesn't work as well

AdamW = Adam with weight decay.

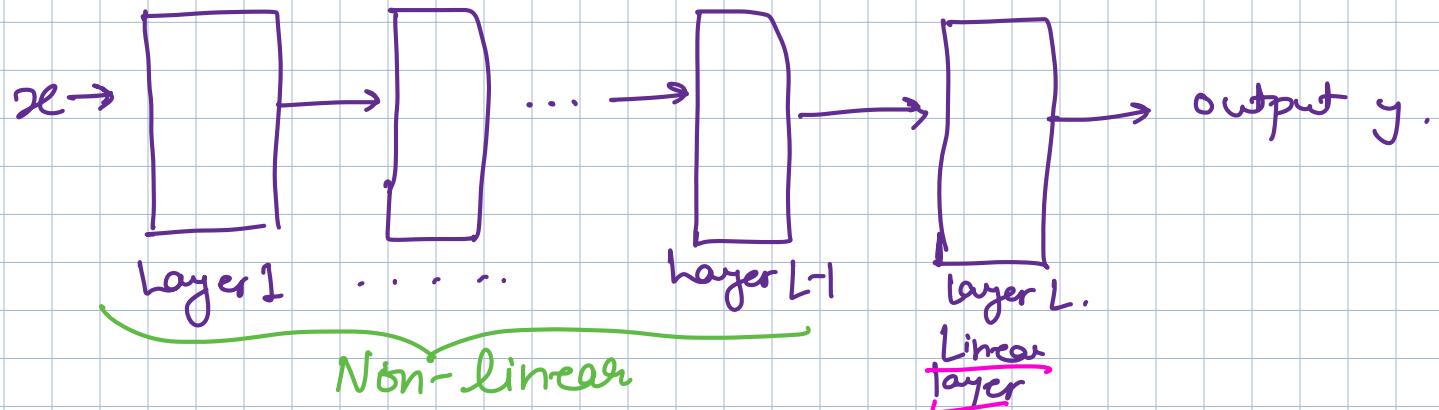
$$\text{Adam: } \overrightarrow{w}_{t+1} = (1-\lambda) \overrightarrow{w}_t - \eta \cdot \frac{\overrightarrow{m}_{t+1}}{\sqrt{\overrightarrow{\sigma}_{t+1}} + \epsilon}$$

Linearized models

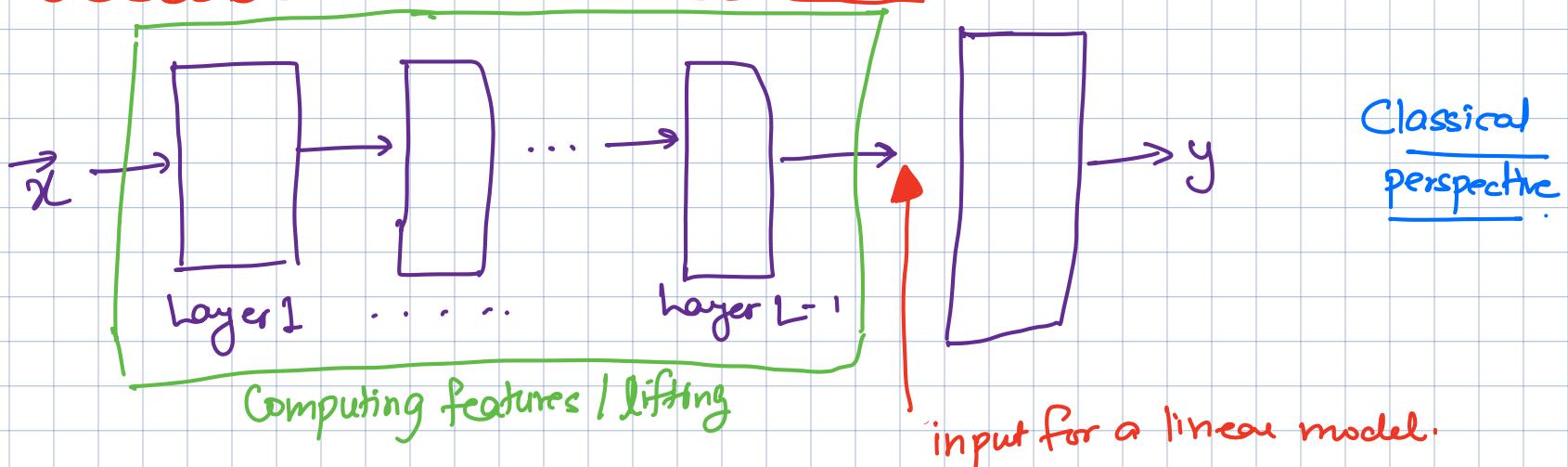
So how can we understand how these algorithms behave on neural networks?

Not nice and convex...

- Standard approach. Try to find a simpler model that can help understanding
- Linear models !!



How do I view this as a linear model?



- learn the features from data instead of committing to them in advance.
- To reuse my computed features, I can freeze all but the linear layer and retrain → Linear Probing.

Other ideas?

$$\text{Network: } f(\vec{x}, \vec{\theta})$$

\vec{x} : data

$\vec{\theta}$: parameters

$$\vec{\theta} = \begin{bmatrix} \vec{w} \\ b \end{bmatrix} \in \mathbb{R}^m$$

Linearization can come from Taylor expansion!

$\vec{\theta}$ is changing as we train the network.

$$f(x, \theta_0 + \Delta\theta) = f(x, \theta_0) + \langle \nabla_{\theta} f, \Delta\theta \rangle + \dots$$

$$f: \mathbb{R}^m \rightarrow \mathbb{R}^n$$

$$\nabla_{\theta} f \in \mathbb{R}^{m \times n}$$

$$(\nabla_{\theta} f)^T = \frac{\partial f}{\partial \theta}$$

$$\langle \nabla_{\theta} f, \Delta\theta \rangle = \frac{\partial f}{\partial \theta} \cdot \Delta\theta$$

"Linear approximation" of our network.

$\frac{\partial f}{\partial \theta}$ is multiplying the linear term. Those are features

↳ Loosy Training Assumption: Weights are moving by very small amounts.
 ↳ Why is this reasonable?

After training
practice

This perspective on linearization led to the "Neural Tangent Kernel" perspective.

→ In the limit of infinite width networks, GD will converge to zero

training loss in the over-parameterized regime.

Can show:

How NN
evolves through training.

$$\frac{df(x, \theta)}{dt} = -\frac{1}{N} \sum_{i=1}^N \underbrace{\nabla_{\theta} f(x, \theta)^T}_{\text{inner product of features @ } x \text{ with features @ } x^{(i)}} \underbrace{\nabla_{\theta} f(x^{(i)}, \theta) \nabla_f L(f, y^{(i)})}_{\text{loss @ } y.}$$

$$K(x, x^{(i)})$$

If lazy training assumption holds, parameters are not moving much and can show that $K(x, x^{(i)})$ also not changing.

→ Can show GD will converge