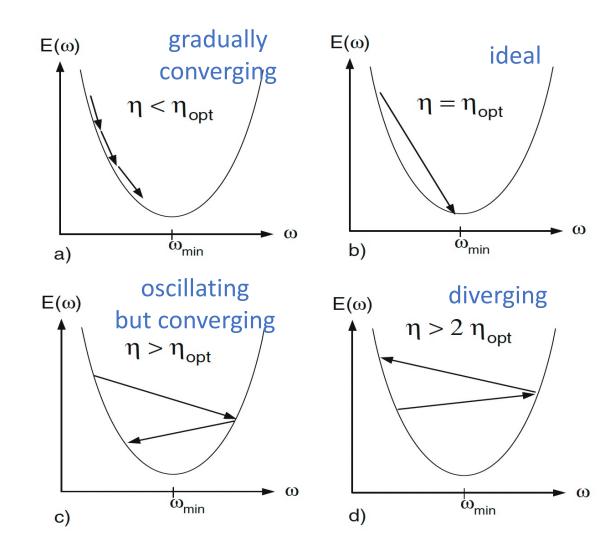
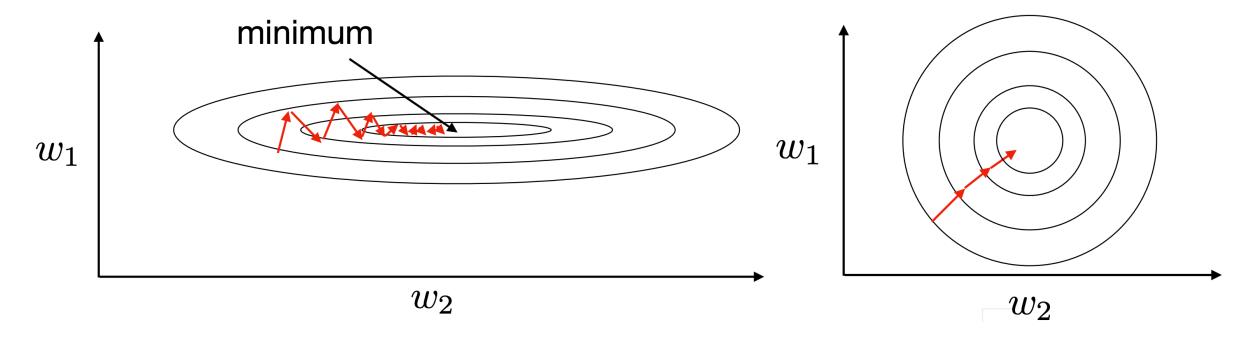
### **Learning Rate in Gradient Descent**

$$W \coloneqq W - \eta \frac{\partial \mathcal{E}}{\partial W}$$

$$\eta: \text{ learning rate}$$



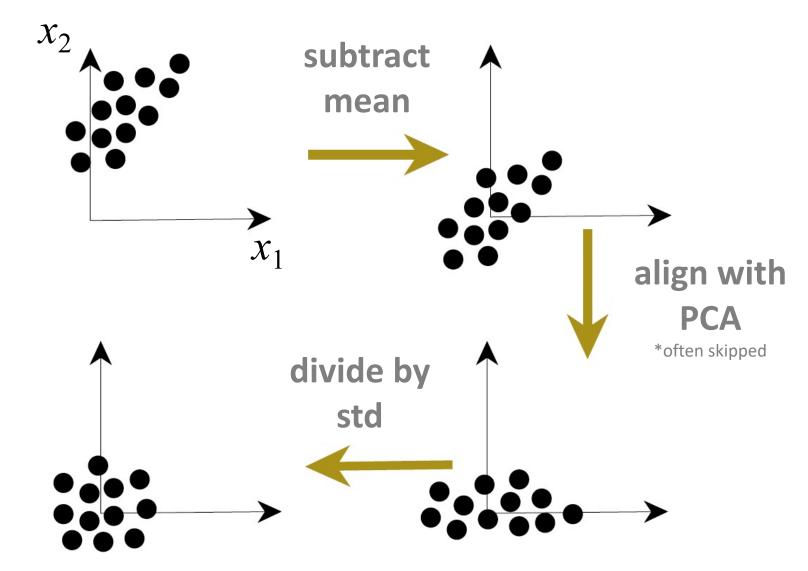
#### **Normalization**



- Same learning rate applied to all weights
- Large weights dominate updates
- Small weights oscillate (or diverge)

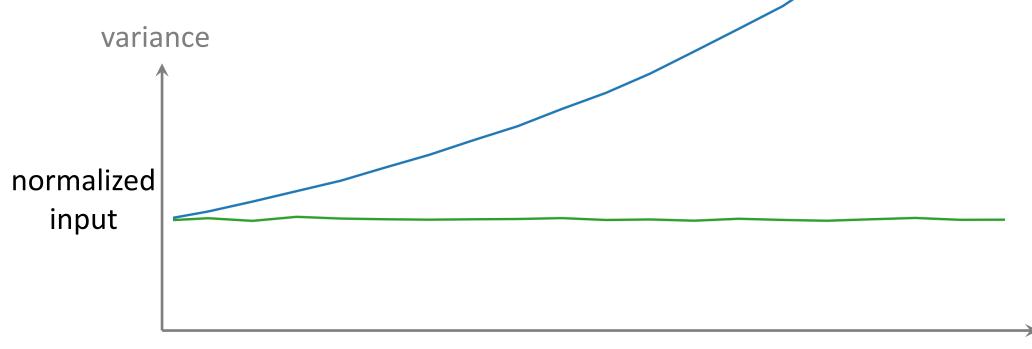
Similar pace for all weights

## **Input Normalization**



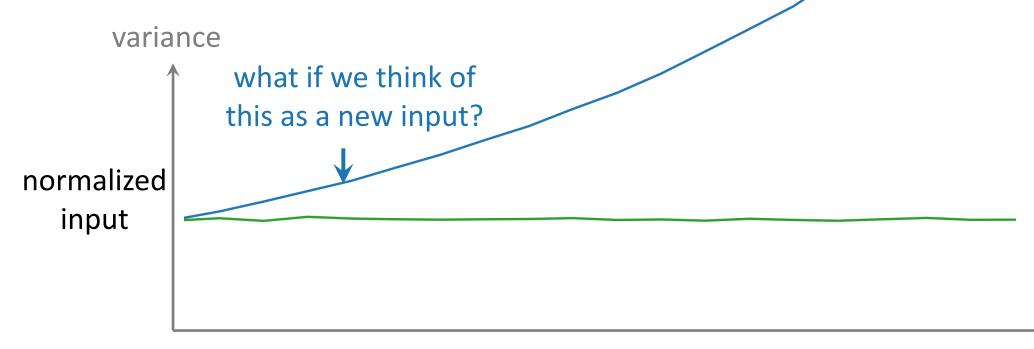
We want to maintain variance for all layers





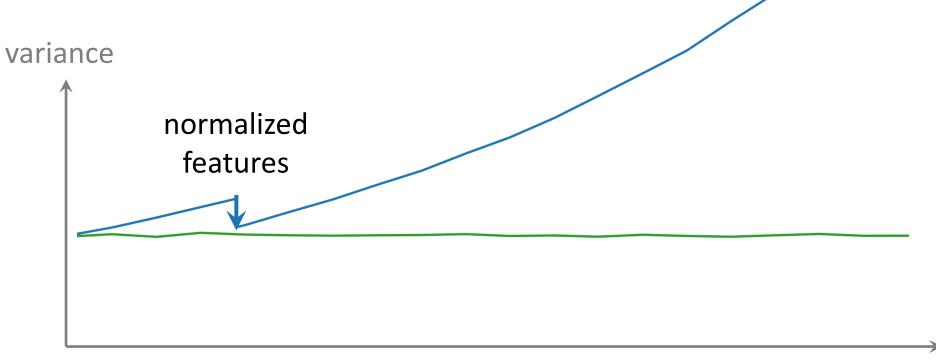
We want to maintain variance for all layers



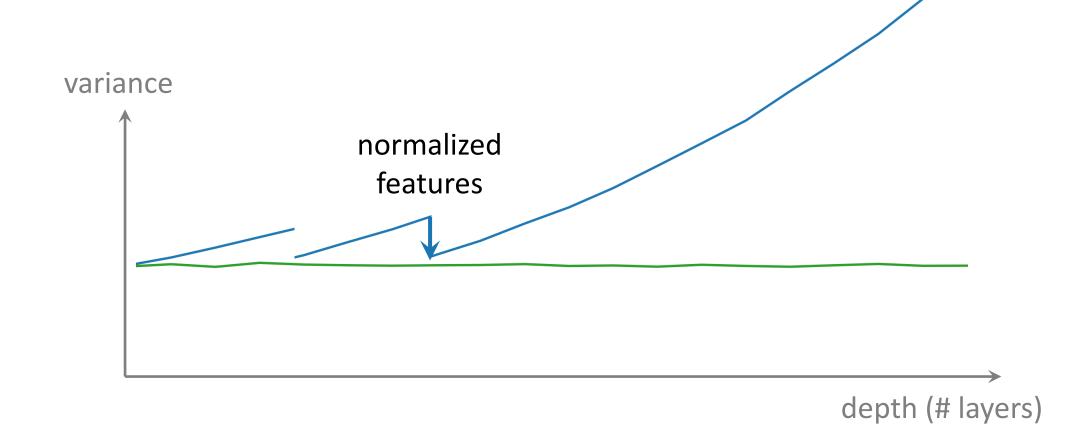


We want to maintain variance for all layers

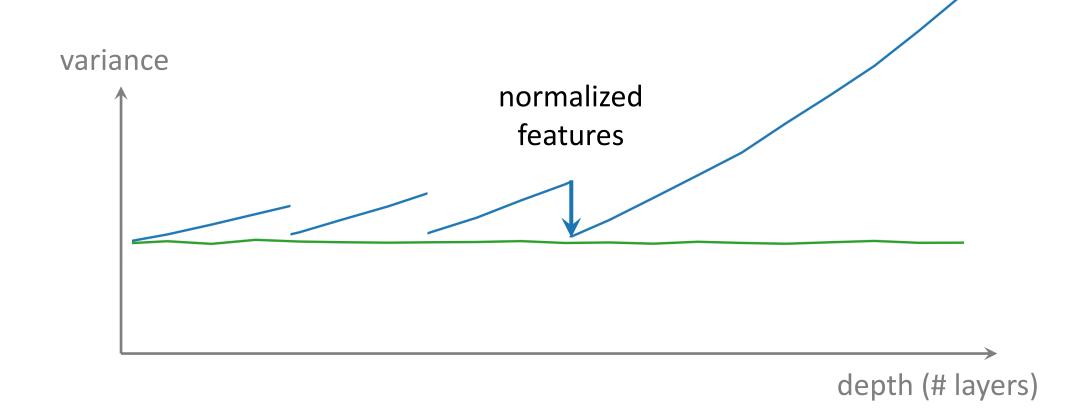




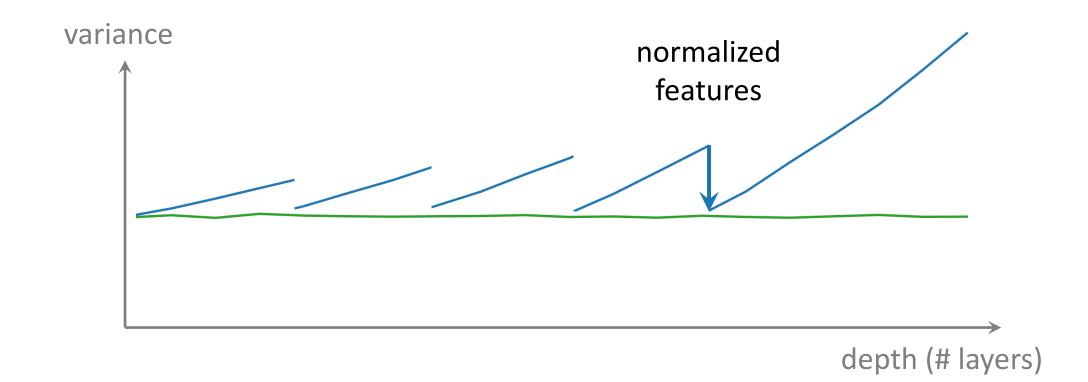
- We want to maintain variance for all layers
- normalize features in the network



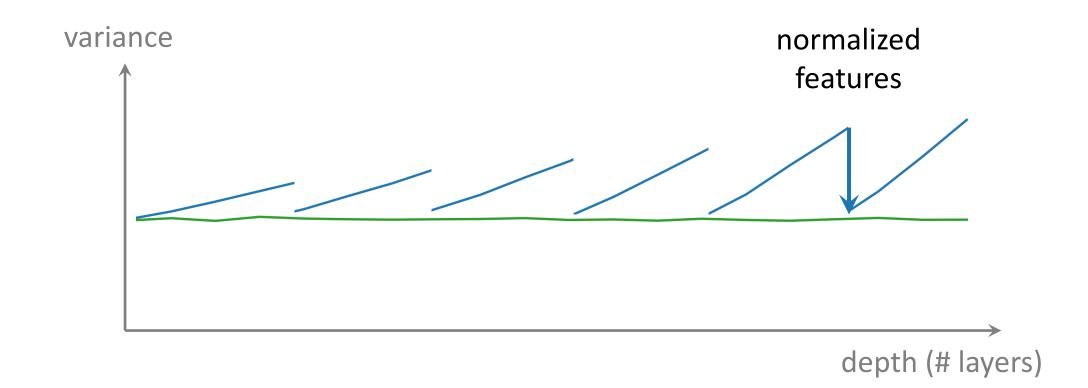
- We want to maintain variance for all layers
- normalize features in the network



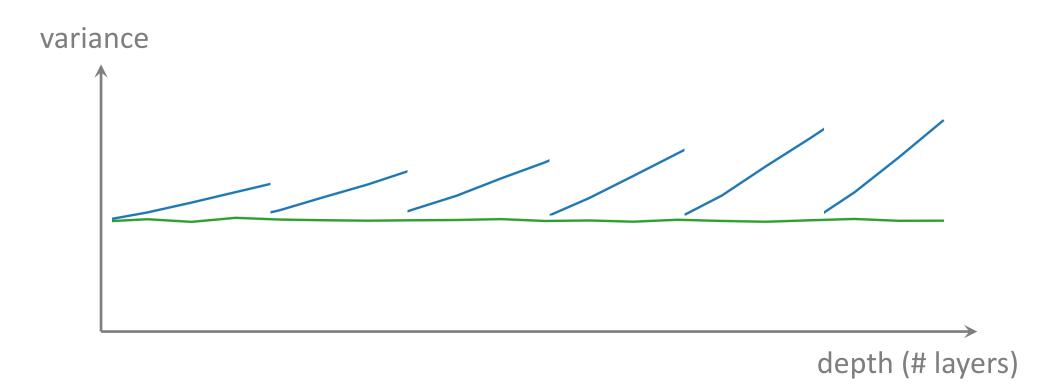
- We want to maintain variance for all layers
- normalize features in the network



- We want to maintain variance for all layers
- normalize features in the network

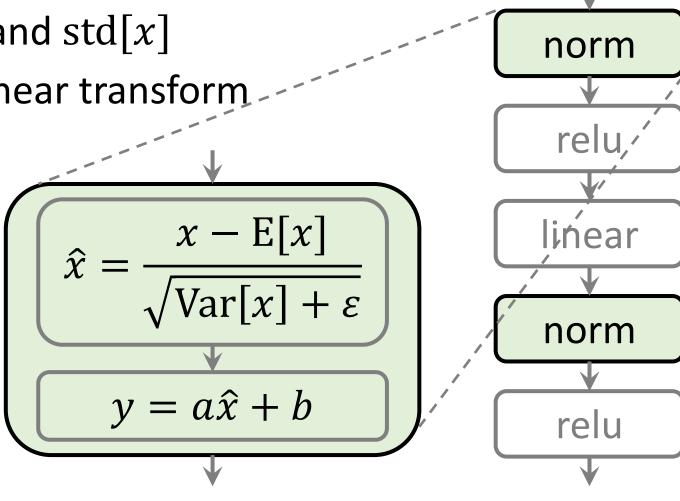


- We want to maintain variance for all layers
- normalize features in the network
- train end-to-end by BackProp



Normalization Modules: Operations

- 1. compute E[x] and Var[x]
- 2. normalize by E[x] and std[x]
- 3. compensate by a linear transform



linear

#### **Normalization Modules: Variants**

differ in support sets of E[x], Var[x]

