# Lecture 4

- Feature perspective and Taylor expansion
- Adam revisited
- GD revisited.
- Initialization

feature extractor.

$$x \longrightarrow \boxed{\phantom{x}} \longrightarrow \boxed{\phantom{x}} - \cdots \longrightarrow \boxed{\text{linear}} \longrightarrow y$$

A second perspective ...

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

$\vec{\theta}$ : parameters

Focus on this as a function of $\vec{\theta}$.

$\vec{\theta} \in \mathbb{R}^m$  say.

Consider  $f: \mathbb{R}^m \to \mathbb{R}$     e.g. linear regression, binary classification

$$f(x, \theta + \Delta\theta) = f(x, \theta) + \boxed{\langle \nabla_\theta f, \Delta\theta \rangle} + \cdots \text{ higher order terms}$$
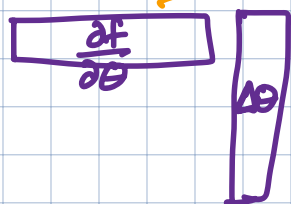
$\longrightarrow$ Taylor expansion.

"Linear approximation of the network"

Taylor approximation only holds true in a neighbourhood

often true in
practice

Lazy Training assumption:  Weights are moving by very small amounts.

$\langle \nabla_\theta f, \Delta\theta \rangle \longrightarrow$ linear term

entries are "features"!

Recall:

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

Gradient is transpose of the derivative.

So now if we revisit descent algorithms from this perspective ...

## Revisit Adam (Sign SGD).

Here $\ell$ is The loss

$$\ell(\vec{x}, \vec{\theta_i} + \Delta\vec{\theta_i}) = \ell(\vec{x}, \vec{\theta_i}) + \left.\frac{\partial \ell}{\partial \theta}\right|_{\theta=\theta_i} \vec{\Delta\theta} \qquad \left| \vec{\theta_{i+1}} = \vec{\theta_i} - \eta \cdot \text{step} \right.$$

Tension: Want large steps to be able to converge fast.

Want small steps to make sure approximation still holds.

Idea: Bound the size of your step: $\|\Delta\theta\|_{\infty} \leq \eta$

$$\hookrightarrow \max_j \|\Delta\theta[j]\| \leq \eta$$

Want:

$$\underset{\|\vec{\Delta\theta}\|_{\infty} \leq \eta}{\arg\min} \frac{1}{n}\sum_{\ell=1}^{n} \left[ \ell(\vec{x_\ell}, \cancel{\vec{\theta_i}}) + \left\langle \left.\nabla_\theta \ell\right|_{\substack{x=x_\ell \\ \theta=\theta_i}}, \vec{\Delta\theta} \right\rangle \right]$$

$$= \underset{\|\Delta\theta\|_{\infty} \leq \eta}{\arg\min} \left\langle \underbrace{\left(\frac{1}{n}\sum_{\ell=1}^{n} \left.\nabla_\theta f\right|_{\substack{x=x_\ell \\ \theta=\theta_i}}\right)}_{\nabla_\theta f_{batch}}, \vec{\Delta\theta} \right\rangle$$

$$= \operatorname*{argmin}_{\forall j \; |\Delta\theta[j]| \le \eta} \sum_j \left( \left( \frac{1}{r} \sum_{\ell=1}^{n} \nabla_\theta f \Big|_{\substack{x=x_\ell \\ \theta=\theta_i}} \right)[j], \Delta\vec{\theta}[j] \right)$$

$\underbrace{\qquad\qquad}_{\nabla_\theta f_{batch}}$

Do each optimization individually

$$\operatorname*{argmin}_{|\Delta\theta[j]| \le \eta} \quad \boxed{\nabla_\theta f_{batch}} \; \Delta\theta[j].$$

$$-\eta \le \Delta\theta[j] \le \eta$$

$$\Delta\theta_j = +\eta \qquad \text{if} \quad \nabla_\theta f_{batch}[j] < 0$$
$$= -\eta \qquad \text{if} \quad \nabla_\theta f_{batch}[j] > 0.$$

$$\Rightarrow \quad \Delta\vec{\theta}^* = -\eta \cdot \operatorname{sign}\left(\nabla_\theta f_{batch}\right).$$

Two-norm constraint: $\qquad \|\Delta\theta\|_2 \le \eta$

## Stepping back
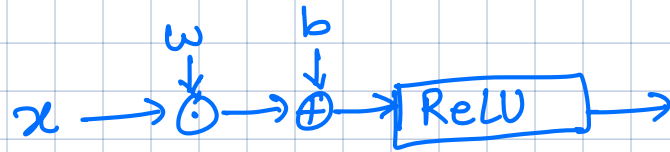
Why do we standardize data in ML?

→ Want features on similar orders of magnitude. → else large features
                                               will dominate.
    ↳ Better conditioning

→ Numerical issues and numerical precision.

---

→ Want to be able to "use" the non-linearity for the
       expressive power of neural nets.

Consider ReLU.     Scalar case

$$x \longrightarrow \overset{\overset{w}{\downarrow}}{\odot} \longrightarrow \overset{\overset{b}{\downarrow}}{\oplus} \longrightarrow \boxed{\text{ReLU}} \longrightarrow$$

$\max(0, wx+b)$

↑ elbow @ $-\frac{-b}{w}$
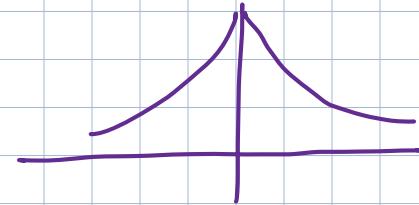
Say we initialize
$$b \sim N(0,1)$$
$$w \sim N(0,1)$$

then $-\frac{b}{w} \sim$ Cauchy distribution.

PDF: $\frac{1}{\pi}\left[\frac{1}{x^2+1}\right]$

Properties: "0-mean"

∞ variance.

But mostly around 0.

50% of probability mass in  −1  to 1

70%                                          −2 to +2

80%                                          −3 to +3

90%                                          −7 to +7

We want our values to be where the non-linearity is interesting

$\longrightarrow$ True for other non-linearities too    Tanh, sigmoid

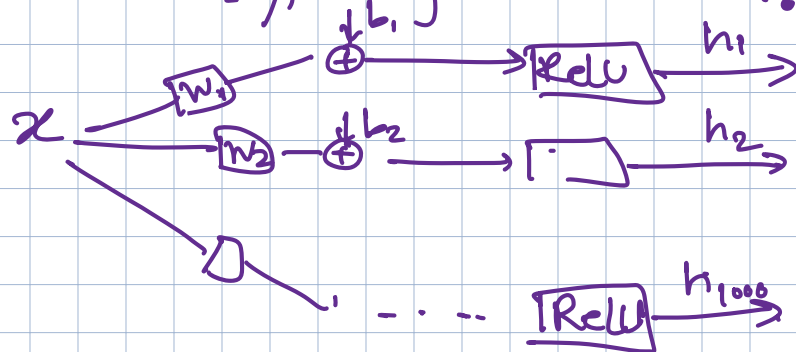What happens if we don't do this?

Suppose our data was in the range 100 to 200.

Probability of a corner landing in that range is about 0.16%

So only 1-2 units out of 1000 may have corners there.

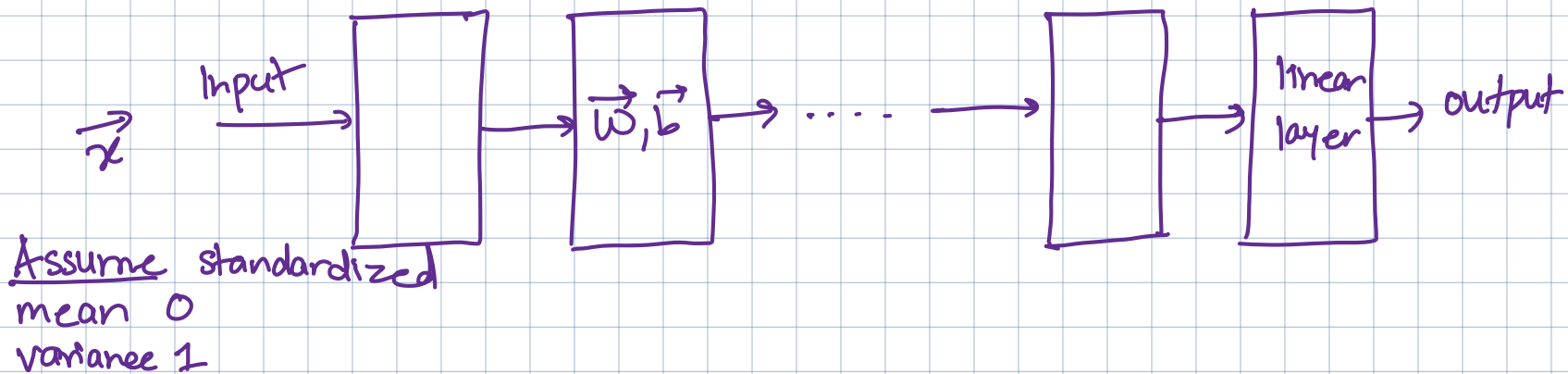Lose expressive power :'(

Mathematically, why is this bad?



If no corners:

$$h_i = \begin{cases} 0 & \forall x \in (100, 200) \\ \alpha x + \beta & \forall x \in (100, 200) \end{cases}$$

$\beta_i = 0$

$$\begin{bmatrix} \alpha_1 x_1 + \beta_1 & & \alpha_{1000} x_1 + \beta_{1000} \\ \alpha_1 x_2 + \beta_1 & & \\ & \cdots & \\ \alpha_1 x_n + \beta_1 & & \alpha_{1000} x_n + \beta_{1000} \end{bmatrix}$$

# Initialization.



$$\vec{x}$$

Input → [ ] → $\vec{W}, \vec{b}$ → . . . . → [ ] → linear layer → output

__Assume__ standardized
mean 0
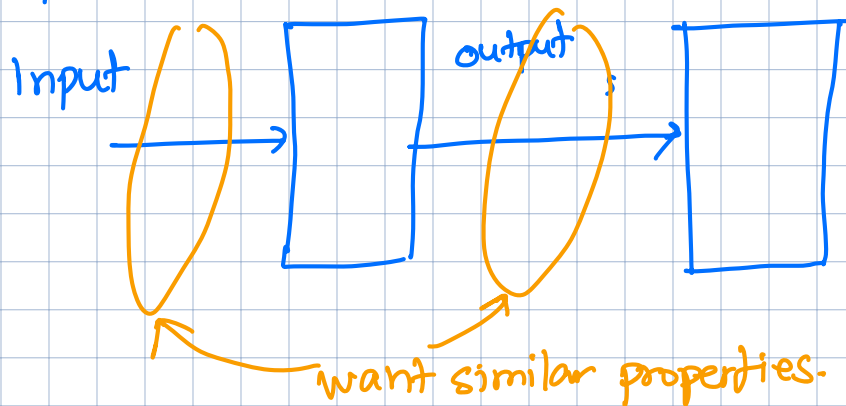variance 1

## How to initialize?

→ Set everything to 0?

But then all inputs from the previous layer get multiplied by zero
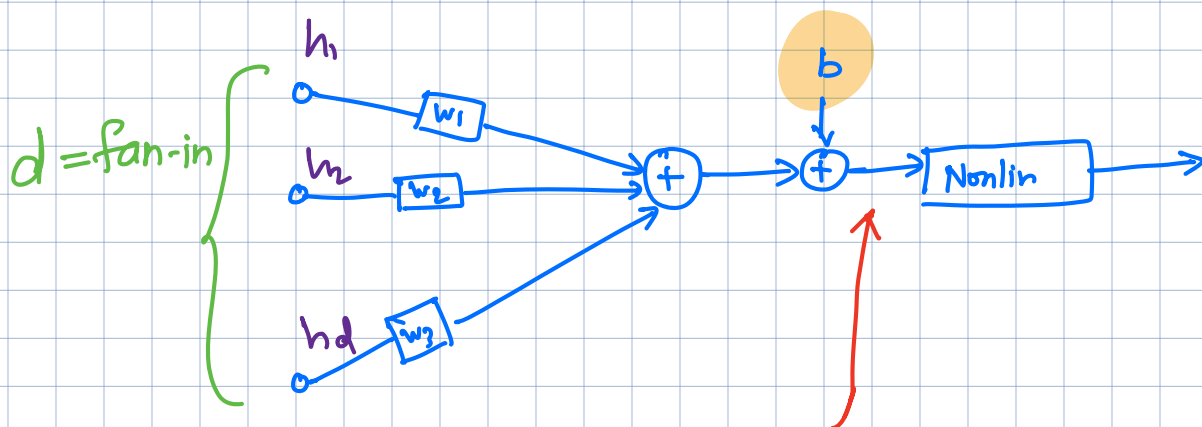
Gradients are zero, weights never more.

→ So what do we do?

Key observation: Output activations from one layer are inputs to the next layer.

Input → [ ] → output → [ ]

want similar properties.

**Goal:** Can each layer have inputs that start "standardized"?

Xavier initialization   ( weights , i.e. parameters that multiply data)



$d =$ fan-in

$h_1$

$w_1$

$h_2$

$w_2$

$h_d$

$w_d$

$b$

Nonlin

Want: Input to nonlinearity is "standardized".

Approx "0-mean"

"unit -variance"

Adding $d$ terms:     $w_1 \sim N(0, \frac{1}{d})$  ← variance.

$$\text{var}\left( \sum w_i h_i \right)$$

$$= E\left[ \left( \sum w_i h_i \right)^2 \right] - \left( E\left( \sum w_i h_i \right) \right)^2 \quad \to 0$$

$$= \sum_d E[w_i^2] \cdot E[h_i^2] \qquad \text{because independence and zero-mean.}$$

## He Initialization

ReLU non-linearities on previous layer

Observation: Output of a ReLU is zero half the time.

$$w_i \sim N(0, \frac{1}{d/2}) = N(0, \frac{2}{d})$$

## Bias initialization

(a) Initialize at $0$.

(b) Use other small number $0.01$.

(c) Treat as Xavier with initialization $d+1$.