# *Neural Networks*

Lecturer: Masha Itkina

# Announcements

- Problem Set 1 is due today

- Problem Set 2 will be out later tonight; due May 4[th]

- Feedback on Project Proposals will be released within a week

# Neural Networks in the Wild

```
-------- Generated Poem 1 --------

I must have shadows on the way
If I am to walk I must have
Each step taken slowly and alone
To have it ready made

And I must think in lines of grey
To have dim thoughts to be my guide
Must look on blue and green
And never let my eye forget
That color is my friend
And purple must surround me too

The yellow of the sun is no more
Intrusive than the bluish snow
That falls on all of us.  I must have
Grey thoughts and blue thoughts walk with me
If I am to go away at all.
```

GPT-3: Brown et. Al, "Language Models are Few-Shot Learners", NeurIPS 2020.

# *Neural Networks in the Wild*



TEXT DESCRIPTION

**An astronaut**  Teddy bears   A bowl of soup

**riding a horse**  lounging in a tropical resort in space   playing basketball with cats in space

in a photorealistic style   in the style of Andy Warhol  **as a pencil drawing**

DALL·E 2

DALLE-2: Ramesh et. Al, "Hierarchical Text-Conditional Image Generation with CLIP Latents", ArXiv 2022.

# *Agenda for Today*

- Supervised learning with non-linear models

- Neural networks

# *Linear Regression Review*

$$\{x^{(i)}, y^{(i)}\}_{i=1}^{n}$$

$$h_\theta(x) = \theta^T x + b$$

$$J(\theta) = \sum_{i=1}^{n} \left(y^{(i)} - h_\theta(x^{(i)})\right)^2 = \sum_{i=1}^{n} \left(y^{(i)} - \theta^T x^{(i)} - b\right)^2$$

label    prediction

Run GD or SGD to optimize.

# Non-Linear Models: Kernels

$$\{x^{(i)}, y^{(i)}\}_{i=1}^{n} \quad \text{lin. in } \theta$$

$$\longrightarrow \text{non-linear in } x$$

$$h_\theta(x) = \theta^T \phi(x)$$

Non-linear in both $\theta$ and $x$:

$$h_\theta(x) = \theta_1^3 x_2 + \sqrt{\theta_5 x_4} + \dots$$

# Non-Linear Models

$$\{x^{(i)}, y^{(i)}\}_{i=1}^{n}$$

Assume $x^{(i)} \in \mathbb{R}^d$, $y^{(i)} \in \mathbb{R}$

$$h_\theta : \mathbb{R}^d \longrightarrow \mathbb{R}$$

Cost function for example $i$:

$$J^{(i)}(\theta) = \left(y^{(i)} - h_\theta(x^{(i)})\right)^2$$

Cost function for dataset:

$$J^{(i)}(\theta) = \frac{1}{n} \sum_{i=1}^{n} J^{(i)}(\theta)$$

$\hookrightarrow$ constant does not matter!
Minimizer will remain the same

# Non-Linear Models

We want to optimize: $\min_\theta J(\theta)$

Gradient Descent (GD): $\theta := \theta - \alpha \nabla_\theta J(\theta) = \theta - \alpha \nabla \left( \frac{1}{n} \sum_{i=1}^{n} J^{(i)}(\theta) \right)$

*update value of left side w/ right side*    $\alpha > 0$

$\hookrightarrow$ *consider all examples at once*

# *Non-Linear Models*

We want to optimize: $\min_{\theta} J(\theta)$

Stochastic Gradient Descent (SGD):

*Alternative SGD:*

*for $k = 1 : n$ epoch:*
   *shuffle dataset*
   *for $j = 1 : n$ iter:*
      $\theta = \theta - \alpha \nabla J^{(j)}(\theta)$

*no replacement*

---

**Algorithm 1** Stochastic Gradient Descent

---

1: Hyperparameter: learning rate $\alpha$, number of total iteration $n_{\text{iter}}$.
2: Initialize $\theta$ randomly.
3: **for** $i = 1$ to $n_{\text{iter}}$ **do**    → *with replacement*
4:      Sample $j$ uniformly from $\{1, \ldots, n\}$, and update $\theta$ by

$$\theta := \theta - \alpha \nabla_{\theta} J^{(j)}(\theta)$$

---

10

# Non-Linear Models

We want to optimize: $\min\limits_{\theta} J(\theta)$

Mini-batch SGD: Computing $B$ gradients $\nabla J^{(j_1)}(\theta), \ldots, \nabla J^{(j_B)}(\theta)$ simultaneously is faster than separately.

GPU parallelism.

**Algorithm 2** Mini-batch Stochastic Gradient Descent

1: Hyperparameters: learning rate $\alpha$, batch size $B$, # iterations $n_{\text{iter}}$.
2: Initialize $\theta$ randomly
3: **for** $i = 1$ to $n_{\text{iter}}$ **do**    $B < n$
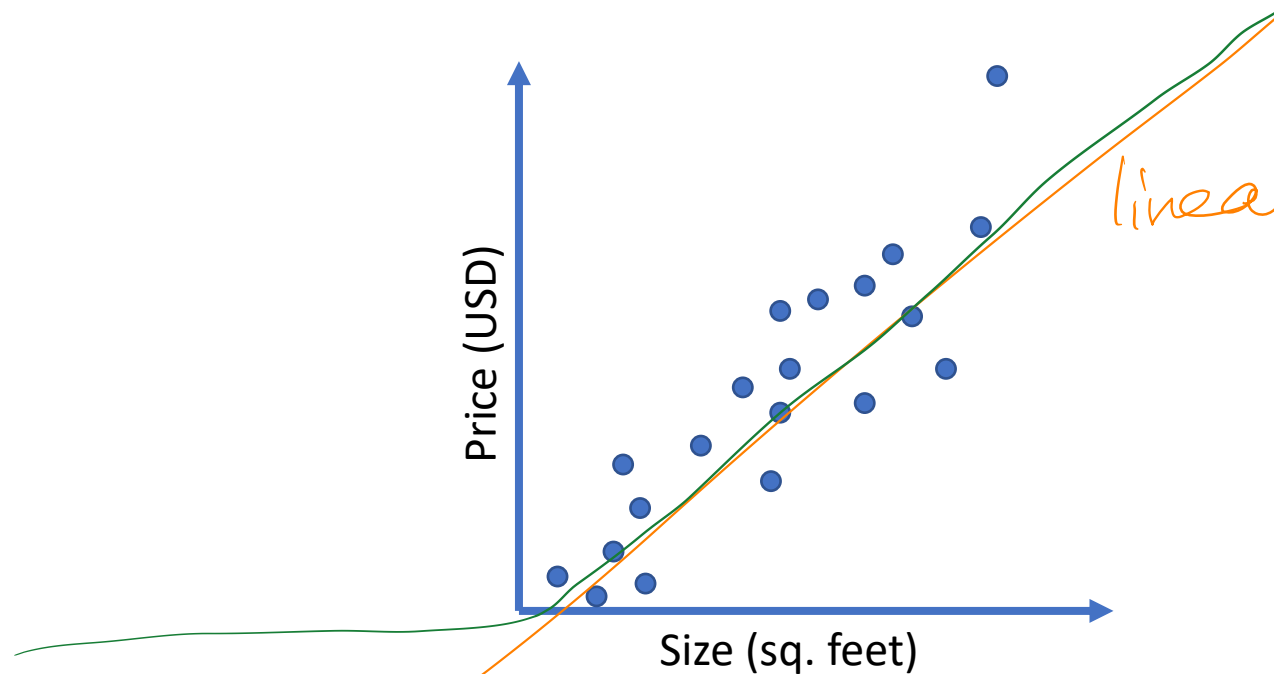4:     Sample $B$ examples $j_1, \ldots, j_B$ (without replacement) uniformly from $\{1, \ldots, n\}$, and update $\theta$ by

$$\theta := \theta - \frac{\alpha}{B} \sum_{k=1}^{B} \nabla_\theta J^{(j_k)}(\theta)$$

How large is $B$?

$\downarrow B$: better
$\uparrow B$: faster runtime

Max $B$ that you can store in GPU memory.

# Neural Networks

- How to define $h_\theta(x)$? *Expressivity, how easy to compute*
  - Neural network!

- How to compute $\nabla J^{(j)}(\theta)$?
  - Backpropagation (next lecture)

# Housing Price Prediction



linear model

① gives negative number

② prediction might have nonlinear relationship with input

$h_\theta(x) = \max \{wx + b, \theta\}$

$\theta = (w, b)$

$ReLU = \max \{0, t\}$

$h_\theta(x) = ReLU(wx + b) \leftarrow$ NN w/ one neuron (one non-linearity)

↑ activation function

13

# Housing Price Prediction

High-dimensional input : $x \in \mathbb{R}^d$, $y \in \mathbb{R}$

$h_\theta(x) = \text{ReLU}(w^T x + b)$

$x = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix} \in \mathbb{R}^d$, $w \in \mathbb{R}^d$, $b \in \mathbb{R}$

weight vector

bias

We want to stack neurons! Output of activation → input to the next.

# Housing Price Prediction

$x \in \mathbb{R}^4 \qquad x_1 \quad , \quad x_2 \quad , \quad x_3 \quad , \quad x_4$

$\uparrow \qquad \uparrow \qquad \uparrow \qquad \uparrow$

size     # bedrooms   zip code   wealth

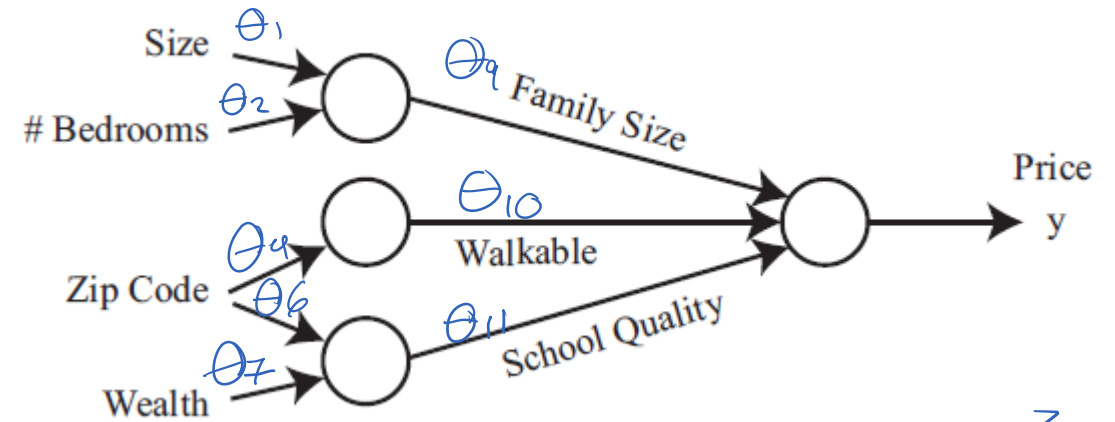intermediate variables:

max family size : $a_1$

walkable : $a_2$

school quality : $a_3$

$a_1 = ReLU(\theta_1 x_1 + \theta_2 x_2 + \theta_3)$

$a_2 = ReLU(\theta_4 x_3 + \theta_5)$

$a_3 = ReLU(\theta_6 x_3 + \theta_7 x_4 + \theta_8)$

$h_\theta(x) = ReLU(\theta_9 a_1 + \theta_{10} a_2 + \theta_{11} a_3 + \theta_{12})$

usually have linear layer at the end!

Size $\theta_1$

$\theta_2$

# Bedrooms

$\theta_9$ Family Size

$\theta_{10}$

Walkable

$\theta_4$

Zip Code

$\theta_6$

$\theta_{11}$

School Quality

$\theta_7$

Wealth

Price
y

prior knowledge

# Two-Layer Neural Network

What if we do not have prior knowledge?

- Fully connected neural network

- Intermediate variables -> hidden units

$a_1 = ReLU(\_ x_1 + \_ x_2 + \_ x_3 + \_ x_4 + \_)$

$a_2 = ReLU(\_ x_1 + \_ x_2 + \_ x_3 + \_ x_4 + \_)$

$\vdots$  ⤷ more general, depends on everything

$a_1 = ReLU(w_1^{[1]T} x + b_1^{[1]}), \quad w \in \mathbb{R}^4, x \in \mathbb{R}^4,$

$a_2 = ReLU(w_2^{[1]T} x + b_2^{[1]})$ $\quad b \in \mathbb{R}$

$a_3 = ReLU(w_3^{[1]T} x + b_3^{[1]})$

$h_\theta(x) = w^{[2]T} a + b^{[2]}, \quad w^{[2]} \in \mathbb{R}^3, a \in \mathbb{R}^3, b \in \mathbb{R}$

$[1]$ $\quad$ $[2]$ $\quad$ two layers

⇩

one hidden layer

# Two-Layer Neural Network

$$\forall j \in [1, ..., m], \quad z_j = w_j^{[1]\top} x + b_j^{[1]} \text{ where } w_j^{[1]} \in \mathbb{R}^d, b_j^{[1]} \in \mathbb{R}, \; x \in \mathbb{R}^d$$

$$a_j = \text{ReLU}(z_j),$$

$$a = [a_1, \ldots, a_m]^\top \in \mathbb{R}^m \quad m \text{ hidden units}$$

$$h_\theta(x) = w^{[2]\top} a + b^{[2]} \text{ where } w^{[2]} \in \mathbb{R}^m, b^{[2]} \in \mathbb{R},$$

# *Vectorization*

$$W^{[1]} = \begin{bmatrix} — & w_1^{[1]^\top} & — \\ — & w_2^{[1]^\top} & — \\ & \vdots & \\ — & w_m^{[1]^\top} & — \end{bmatrix} \in \mathbb{R}^{m \times d}$$

# *Vectorization*

$$\underbrace{\begin{bmatrix} z_1 \\ \vdots \\ z_m \end{bmatrix}}_{z \in \mathbb{R}^{m \times 1}} = \underbrace{\begin{bmatrix} - w_1^{[1]^\top} - \\ - w_2^{[1]^\top} - \\ \vdots \\ - w_m^{[1]^\top} - \end{bmatrix}}_{W^{[1]} \in \mathbb{R}^{m \times d}} \underbrace{\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}}_{x \in \mathbb{R}^{d \times 1}} + \underbrace{\begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ \vdots \\ b_m^{[1]} \end{bmatrix}}_{b^{[1]} \in \mathbb{R}^{m \times 1}}$$

$$\boxed{z = W^{[1]} x + b^{[1]}}$$

# *Vectorization*

Pre-activation: $z = W^{[1]} x + b^{[1]} \in \mathbb{R}^m$

$$a = \begin{bmatrix} a_1 \\ \vdots \\ a_m \end{bmatrix} = \begin{bmatrix} \text{ReLU}(z_1) \\ \vdots \\ \text{ReLU}(z_m) \end{bmatrix} \triangleq \text{ReLU}(z)$$

overloaded notation for multi-dim vector inputs

$W^{[2]} = \begin{bmatrix} w^{[2]T} \end{bmatrix} \in \mathbb{R}^{1 \times m}, \ b^{[2]} \in \mathbb{R}$

$h_\theta(x) = W^{[2]} a + b^{[2]}$

Vectorization helps us parallelize on GPU!

# *Multi-Layer Fully-Connected Neural Networks*

weight matrix

$a^{[1]} = \text{ReLU}(W^{[1]}x + b^{[1]})$ → bias  $\dim(a^{[k]}) = m_k$

hidden units

$a^{[2]} = \text{ReLU}(W^{[2]}a^{[1]} + b^{[2]})$

$\dots$

$a^{[r-1]} = \text{ReLU}(W^{[r-1]}a^{[r-2]} + b^{[r-1]})$

$h_\theta(x) = W^{[r]}a^{[r-1]} + b^{[r]}$

$W^{[1]} \in \mathbb{R}^{m_1 \times d}$

$W^{[2]} \in \mathbb{R}^{m_2 \times m_1}$

$W^{[k]} \in \mathbb{R}^{m_k \times m_{k-1}}$

$b^{[k]} \in \mathbb{R}^{m_k}$

# Why do we need an activation function (e.g., ReLU)?

$$a^{[1]} = W^{[1]}x + b^{[1]}$$

$$h_\theta(x) = W^{[2]}a + b^{[2]} = W^{[2]}(W^{[1]}x + b^{[1]}) + b^{[2]}$$

$$= W^{[2]}W^{[1]}x + W^{[2]}b^{[1]} + b^{[2]}$$

$$\underbrace{W^{[2]}W^{[1]}}_{\tilde{W}} \qquad \underbrace{W^{[2]}b^{[1]} + b^{[2]}}_{}$$

linear in these parameters

# Connection to Kernel Methods

Kernel method: $h_\theta(x) = \theta^T \phi(x)$    linear in parameters, not $x$

$$a^{[r-1]} = \phi_\beta(x), \quad \beta = (W^{[1]}, \ldots, W^{[r-1]}, b^{[1]}, \ldots, b^{[r-1]})$$

$\hookrightarrow$ fix beta

$$h_\theta(x) = W^{[r]} \phi_\beta(x) + b^{[r]}$$

Neural network: $\phi_\beta(x)$ is learned - best that works
for data.

$a^{[r-1]} \rightarrow$ features/representations

# Summary

- Supervised learning with non-linear models

- Neural networks

- Next time: backpropagation