

Neural Networks

Lecturer: Masha Itkina

Announcements

- Problem Set 1 is due today
- Problem Set 2 will be out later tonight; due May 4th
- 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

Non-Linear Models: Kernels

Non-Linear Models

Non-Linear Models

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

Gradient Descent (GD): $\theta := \theta - \alpha \nabla_{\theta} J(\theta)$

Non-Linear Models

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

Stochastic Gradient Descent (SGD):

Algorithm 1 Stochastic Gradient Descent

- 1: Hyperparameter: learning rate α , number of total iteration n_{iter} .
- 2: Initialize θ randomly.
- 3: **for** $i = 1$ to n_{iter} **do**
- 4: Sample j uniformly from $\{1, \dots, n\}$, and update θ by

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

Non-Linear Models

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

Mini-batch SGD:

Algorithm 2 Mini-batch Stochastic Gradient Descent

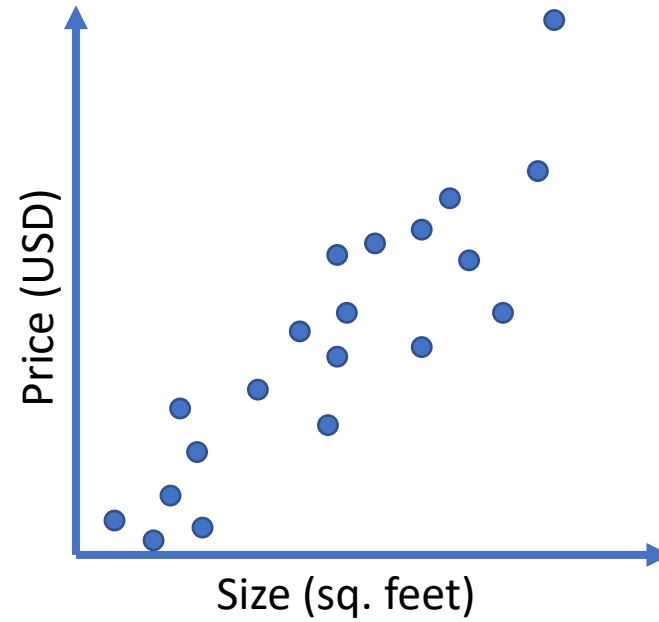
- 1: Hyperparameters: learning rate α , batch size B , # iterations n_{iter} .
- 2: Initialize θ randomly
- 3: **for** $i = 1$ to n_{iter} **do**
- 4: Sample B examples j_1, \dots, j_B (without replacement) uniformly from $\{1, \dots, n\}$, and update θ by

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

Neural Networks

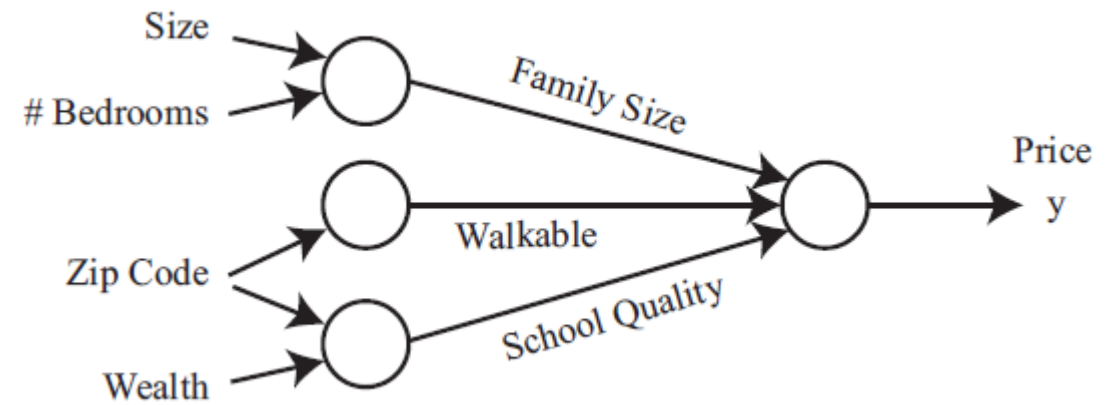
- How to define $h_{\theta}(x)$?
 - Neural network!
- How to compute $\nabla J^{(j)}(\theta)$?
 - Backpropagation (next lecture)

Housing Price Prediction



Housing Price Prediction

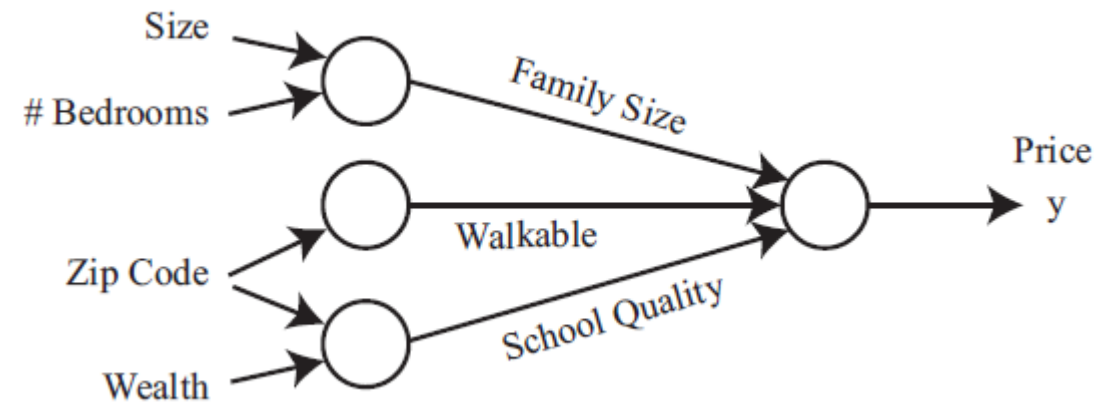
Housing Price Prediction



Two-Layer Neural Network

What if we do not have prior knowledge?

- Fully connected neural network
- Intermediate variables -> hidden units



Two-Layer Neural Network

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

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

$$a = [a_1, \dots, a_m]^\top \in \mathbb{R}^m$$

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

Vectorization

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

Vectorization

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

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

Vectorization

Multi-Layer Fully-Connected Neural Networks

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

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

...

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

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

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

Connection to Kernel Methods

Summary

- Supervised learning with non-linear models
- Neural networks
- Next time: backpropagation