# DATASCI 207

**Nedelina Teneva, PhD**
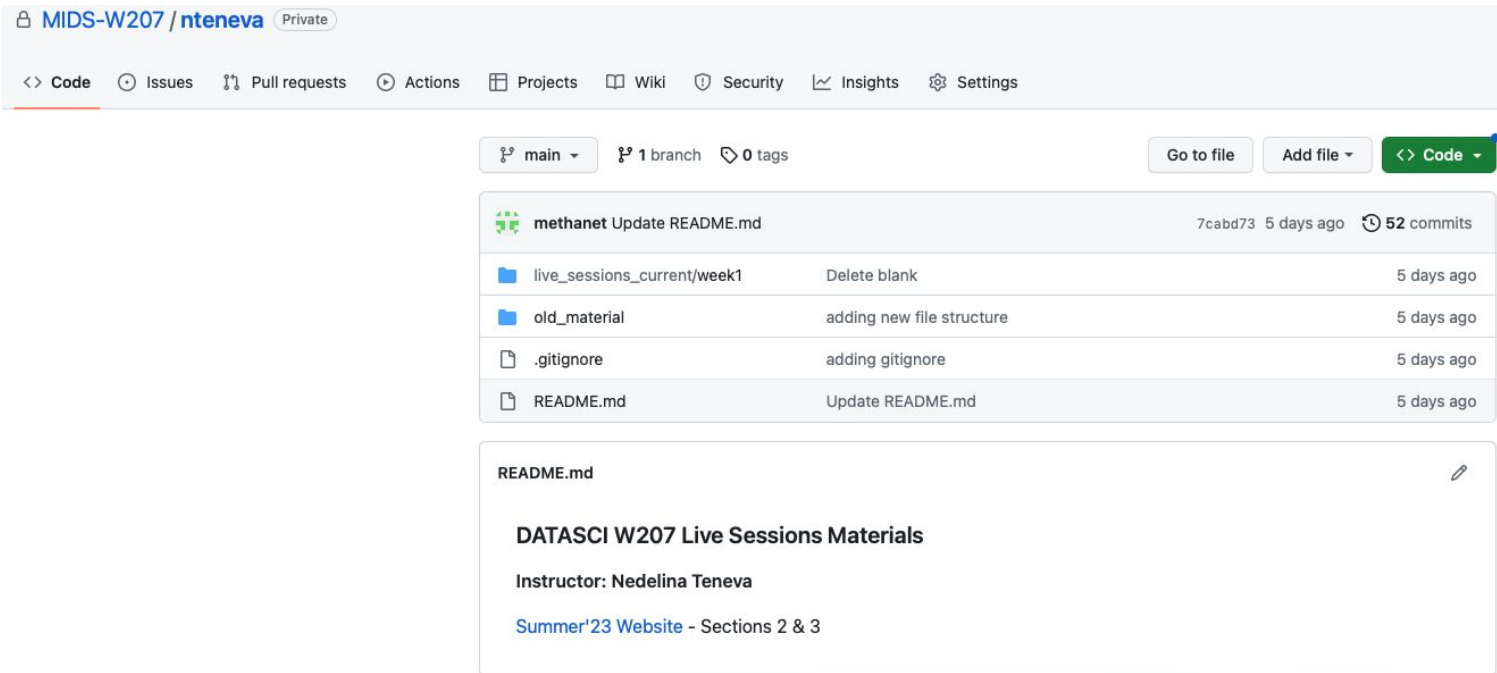
School of Information, UC Berkeley

# Bio

- Currently at Megagon Labs (R&D lab focusing on fundamental ML research)
- Previously at ML Science Manager at Amazon Alexa
- PhD from University of Chicago (focusing on optimization）
- Background in Molecular Biology

# You?

- (Under) graduate major
- Current job/occupation (if any)
- Why a Masters in Data Science?

# Announcements

- [Course Website](Course Website)
- 

# Objectives

- Intro to our first ML technique: linear regression (LR)
- Learning about how to solve LR
- Basic LR and a Tensorflow example: review after class!
  - https://github.com/MIDS-W207/nteneva/blob/main/live_sessions_current/week2/Week_2_Linear_Regression_I.ipynb

# Linear Regression

- Why do we use linear regression?
- What assumptions does LR make for the relationship between outcomes (**y**) and features (**X**)?

# Linear Regression

- Why do we use linear regression?
- What assumptions does LR make for the relationship between outcome (**y**) and features (**X**)?
- **Linear Algebra notation**

    X: matrix of size (n, m) - inputs/features/covariates/independent var's…

    y: vector of size n (column by definition) - output/dependent var's/response…
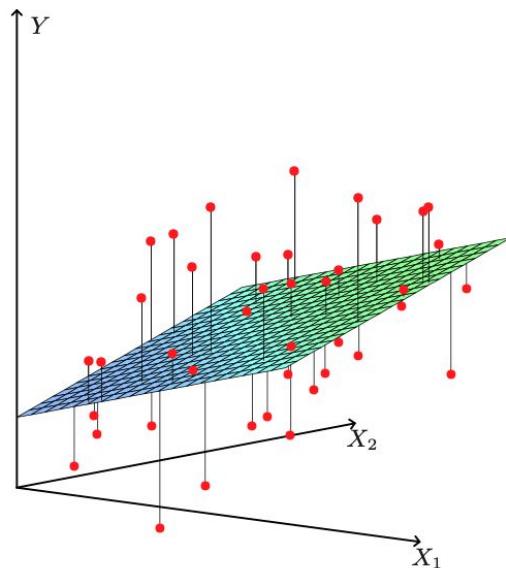
    **Model**: y = X$\beta$ + e

    $\beta$: vector of size m - parameters/weights

    e: vector of size n - error/noise

    **Goal**: estimate $\beta$  s.t. the noise/error e is minimized

# Linear Regression

Source: ESL II

# Example

Based on the data and regression line, what is:

- the actual income for the individual with 7 years of education?

$5000

- the predicted income for an individual with 7 years of education?

$35000



Relationship between income and education

Annual income / Years of Education

# How do we compute LR?

- Direct methods
  - See https://www.cs.toronto.edu/~rgrosse/courses/csc321_2017/readings/L02%20Linear%20Regression.pdf
  - 
- Iteratively

# What is gradient descent?
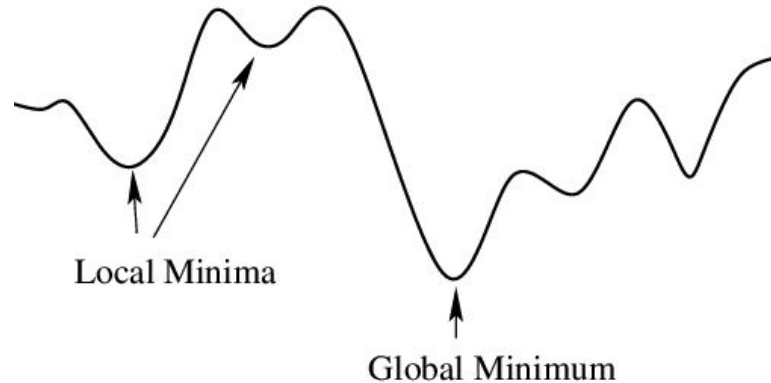
# Gradient descent: key components

- Model
- Parameters
- Cost function
- Objective: minimize the cost function

# Optimization

- What is global minim**um**?
- What is local minimum?
- How how we avoid local minim**a**?

# Optimization

- What is global minim**um**?
- What is local minimum?
- How how we avoid local minim**a**?



Local Minima

Global Minimum

# When do we stop iterating?

# When do we stop iterating?

- When the validation error stops improving (i.e., difference between step t and t +1 is below some threshold)
- Based on the loss

# Hyperparameters

- What are hyperparameter? Examples?
- How are they different from the parameters?

# Hyperparameters

- Examples
    - Batch size
    - Learning rate
    - Epochs


- How do we set their values?

# Hyperparameters

- Examples
    - Batch size
    - Learning rate
    - Epochs


- How do we set their values?
    - Hyper parameter optimization using e.g. cross validation