# Stanford CS224W: Basics of Deep Learning

CS224W: Machine Learning with Graphs

Jure Leskovec, Stanford University

http://cs224w.stanford.edu

# Machine Learning as Optimization

- **Supervised learning:** we are given input $x$, and the goal is to predict label $y$
- **Input $x$ can be:**
  - Vectors of real numbers
  - Sequences (natural language)
  - Matrices (images)
  - Graphs (potentially with node and edge features)
- **We formulate the task as an optimization problem**

# Machine Learning as Optimization

- **Formulate the task as an optimization problem:**

$$\min_{\Theta} \boxed{\mathcal{L}(\boldsymbol{y}, f(\boldsymbol{x}))}$$

**Objective function**

- $\Theta$: a set of **parameters** we optimize
  - Could contain one or more scalars, vectors, matrices …
  - E.g. $\Theta = \{Z\}$ in the shallow encoder (the embedding lookup)

- $\mathcal{L}$: **loss function**. Example: L2 loss

$$\mathcal{L}(\boldsymbol{y}, f(\boldsymbol{x})) = \|y - f(x)\|_2$$

- Other common loss functions:
  - L1 loss, huber loss, max margin (hinge loss), cross entropy …
  - See https://pytorch.org/docs/stable/nn.html#loss-functions

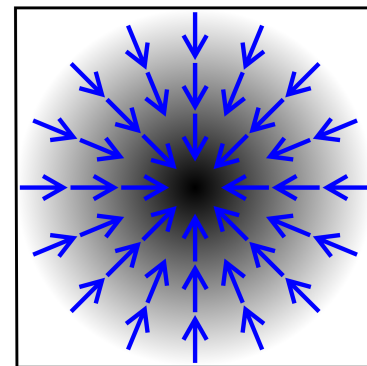# Loss Function Example

- One common loss for classification: **cross entropy (CE)**
- Label $y$ is a categorical vector (one-hot encoding)
  - e.g.  $y =$

| 0 | 0 | 1 | 0 | 0 |
|---|---|---|---|---|

  $y$ is of class "3"

- $f(x) = \text{Softmax}(g(x))$

  - Recall from lecture 3: $f(x)_i = \dfrac{e^{g(x)_i}}{\sum_{j=1}^{C} e^{g(x)_j}},$

  $g(x)_i$ denotes $i$-th coordinate of the vector output of func. $g(x)$

  where $C$ is the number of classes.
  - e.g. $f(x) =$

| 0.1 | 0.3 | 0.4 | 0.1 | 0.1 |
|-----|-----|-----|-----|-----|

- $\text{CE}\big(y, f(x)\big) = -\sum_{i=1}^{C}\big(y_i \log f(x)_i\big)$
  - $y_i, f(x)_i$ are the **actual** and **predicted** value of the $i$-th class.
  - **Intuition:** the lower the loss, the closer the prediction is to one-hot
- **Total loss over all training examples:**

  - $\mathcal{L} = \sum_{(x,y) \in \mathcal{T}} \text{CE}\big(y, f(x)\big)$
  - $\mathcal{T}$: training set containing all pairs of data and labels $(x, y)$

# Machine Learning as Optimization

- **How to optimize the objective function?**
- **Gradient vector:** Direction and rate of fastest increase

**Partial derivative**

$$\nabla_{\Theta} \mathcal{L} = (\frac{\partial \mathcal{L}}{\partial \Theta_1}, \frac{\partial \mathcal{L}}{\partial \Theta_2}, \dots)$$



- $\Theta_1, \Theta_2 \dots$ : components of $\Theta$

https://en.wikipedia.org/wiki/Gradient

- Recall **directional derivative** of a multi-variable function (e.g. $\mathcal{L}$) along a given vector represents the instantaneous rate of change of the function along the vector.
- Gradient is the directional derivative in the **direction of largest increase**

Jure Leskovec, Stanford CS224W: Machine Learning with Graphs, http://cs224w.stanford.edu

# Gradient Descent

- **Iterative algorithm:** repeatedly update weights in the (opposite) direction of gradients until convergence

$$\Theta \leftarrow \Theta - \eta \nabla_\Theta \mathcal{L}$$

- **Training:** Optimize $\Theta$ iteratively
  - **Iteration**: 1 step of gradient descent

- **Learning rate (LR) $\eta$:**
  - Hyperparameter that controls the size of gradient step
  - Can vary over the course of training (LR scheduling)
- Ideal termination condition: **0** gradient
  - In practice, we stop training if it no longer improves performance on **validation set** (part of dataset we hold out from training)

# Stochastic Gradient Descent (SGD)

- **Problem with gradient descent:**

  - Exact gradient requires computing $\nabla_\Theta \mathcal{L}(\boldsymbol{y}, f(\boldsymbol{x}))$, where $\boldsymbol{x}$ is the **entire** dataset!

    - This means summing gradient contributions over all the points in the dataset
    - Modern datasets often contain billions of data points
    - Extremely expensive for every gradient descent step

- **Solution: Stochastic gradient descent (SGD)**

  - At every step, pick a different **minibatch** $\mathcal{B}$ containing a subset of the dataset, use it as input $\boldsymbol{x}$

# Minibatch SGD

- **Concepts:**
  - **Batch size**: the number of data points in a minibatch
    - E.g. number of nodes for node classification task
  - **Iteration**: 1 step of SGD on a minibatch
  - **Epoch**: one full pass over the dataset (# iterations is equal to ratio of dataset size and batch size)

- SGD is unbiased estimator of full gradient:

  - But there is no guarantee on the rate of convergence
  - In practice often requires tuning of learning rate
- Common optimizer that improves over SGD:

  - Adam, Adagrad, Adadelta, RMSprop …

# Neural Network Function

- **Objective:** $\min_{\Theta} \mathcal{L}(y, f(x))$

- In deep learning, the function $f$ can be very complex

- To start simple, consider linear function
$$f(x) = W \cdot x, \qquad \Theta = \{W\}$$

- If $f$ returns a scalar, then $W$ is a learnable **vector**
$$\nabla_W f = \left(\frac{\partial f}{\partial w_1}, \frac{\partial f}{\partial w_2}, \frac{\partial f}{\partial w_3} \ldots\right)$$

- If $f$ returns a vector, then $W$ is the **weight matrix**
$$\nabla_W f = W^T$$

**Jacobian matrix of $f$**

# Back-propagation

- **How about a more complex function:**

$$f(x) = W_2(W_1 x), \qquad \Theta = \{W_1, W_2\}$$

- Recall **chain rule**:

$$\frac{dz}{dx} = \frac{dz}{dy} \cdot \frac{dy}{dx}$$

In other words:
$f(x) = W_2(W_1 x)$
$h(x) = W_1 x$
$g(z) = W_2 z$

- E.g. $\nabla_x f = \dfrac{\partial f}{\partial(W_1 x)} \cdot \dfrac{\partial(W_1 x)}{\partial x}$

- **Back-propagation**: Use of **chain rule** to propagate gradients of intermediate steps, and finally obtain gradient of $\mathcal{L}$ w.r.t. $\Theta$

# Back-propagation Example (1)

- **Example:** Simple 2-layer linear network
- $f(x) = g(h(x)) = W_2(W_1 x)$



- $\mathcal{L} = \sum_{(x,y) \in \mathcal{B}} \left\lVert (y, -f(x)) \right\rVert_2$ sums the L2 loss in a minibatch $\mathcal{B}$

- **Hidden layer:** intermediate representation for input $x$
  - Here we use $h(x) = W_1 x$ to denote the hidden layer
  - $f(x) = W_2 h(x)$

# Non-linearity

- Note that in $f(\boldsymbol{x}) = W_2(W_1\boldsymbol{x})$, $W_2W_1$ is another matrix (vector, if we do binary classification)
- Hence $f(\boldsymbol{x})$ is still linear w.r.t. $\boldsymbol{x}$ no matter how many weight matrices we compose
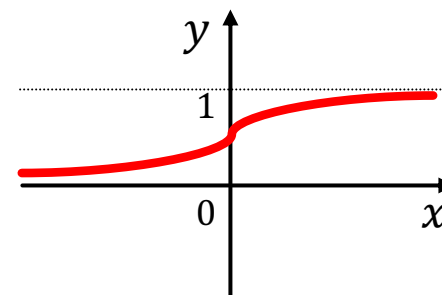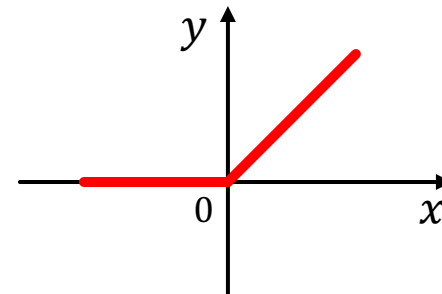- **Introduce non-linearity:**
  - **Rectified linear unit (ReLU)**
    $ReLU(x) = \max(x, 0)$
  - **Sigmoid**
    $\sigma(x) = \dfrac{1}{1 + e^{-x}}$
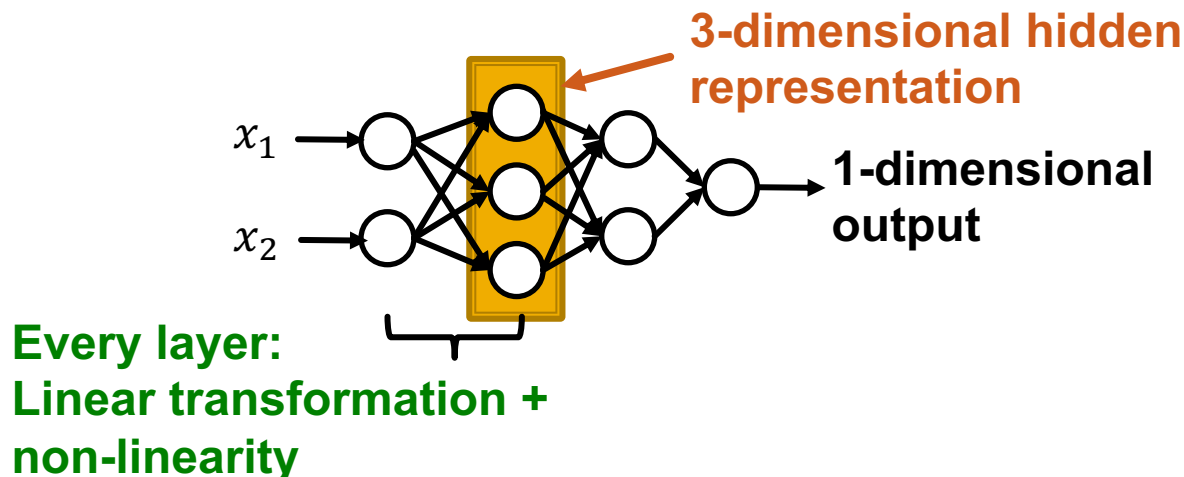
# Multi-layer Perceptron (MLP)

- **Each layer of MLP combines linear transformation and non-linearity:**

$$\boldsymbol{x}^{(l+1)} = \sigma(W_l \boldsymbol{x}^{(l)} + b^l)$$

  - where $W_l$ is weight matrix that transforms hidden representation at layer $l$ to layer $l + 1$

  - $b^l$ is bias at layer $l$, and is added to the linear transformation of $\boldsymbol{x}$

  - $\sigma$ is non-linearity function (e.g., sigmod)

- Suppose $\boldsymbol{x}$ is 2-dimensional, with entries $x_1$ and $x_2$



**3-dimensional hidden representation**

$x_1$

$x_2$

**1-dimensional output**

**Every layer: Linear transformation + non-linearity**

# Summary

- **Objective function:**

$$\min_{\Theta} \mathcal{L}(\boldsymbol{y}, f(\boldsymbol{x}))$$

- $f$ can be a simple linear layer, an MLP, or other neural networks (e.g., a GNN later)
- Sample a minibatch of input $\boldsymbol{x}$
- **Forward propagation**: compute $\mathcal{L}$ given $\boldsymbol{x}$
- **Back-propagation:** obtain gradient $\nabla_{\Theta}\mathcal{L}$ using a chain rule
- Use **stochastic gradient descent (SGD)** to optimize for $\Theta$ over many iterations