# Week 04: Back-Prop

# &

# Optimization for

# Deep Learning

*Dr. Arghya Pal*

Lecturer

Monash University

# Classroom Setup

## Dataset
Data Calling
Data Visualization

## Network Building

## Optimization
Result Visualization

cs.toronto.edu/~kriz/cifar.html

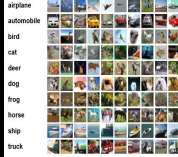< Back to Alex Krizhevsky's home page

The CIFAR-10 and CIFAR-100 datasets are labeled subsets of the 80 million tiny images dataset. CIFAR-10 and CIFAR-100 were created by Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton.

**The CIFAR-10 dataset**

The CIFAR-10 dataset consists of 60000 32x32 colour images in 10 classes, with 6000 images per class. There are 50000 training images and 10000 test images.

The dataset is divided into five training batches and one test batch, each with 10000 images. The test batch contains exactly 1000 randomly-selected images from each class. The training batches contain the remaining images in random order, but some training batches may contain more images from one class than another. Between them, the training batches contain exactly 5000 images from each class.

Here are the classes in the dataset, as well as 10 random images from each:

airplane
automobile
bird
cat
deer
dog
frog
horse
ship
truck

The classes are completely mutually exclusive. There is no overlap between automobiles and trucks. "Automobile" includes sedans, SUVs, things of that sort. "Truck" includes only big trucks. Neither includes pickup trucks.

namespace ?
allow you to use
pre-loaded datasets as
well as your own data

```python
import torch
from torch.utils.data import Dataset, DataLoader


full_train_set =
torchvision.datasets.CIFAR10 ("./data", download=True,
transform=transform)
full_test_set = torchvision.datasets.CIFAR10 ("./data",
download=True, train=False, transform=transform)
```

# Classroom Setup

## Dataset
### Data Calling
### Data Visualization

## Network Building

## Optimization
### Result Visualization

```
import torch
from torch.utils.data import Dataset
```



pip install datasets torchvision torch

# Classroom Setup

Dataset
## Data Calling
Data Visualization

Network Building

Optimization
Result Visualization

```python
import torch
from torch.utils.data import Dataset, DataLoader




full_train_set =
torchvision.datasets.CIFAR10("./data"
, download=True, transform=transform)




full_test_set =
torchvision.datasets.CIFAR10("./data",
download=True, train=False,
transform=transform)
```
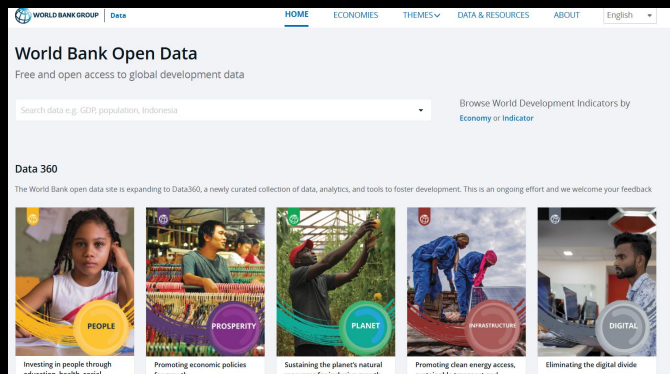
```python
import torch

from torch.utils.data import Dataset
```





```python
class CustomDataset(Dataset):
    """Face Landmarks dataset."""

    def __init__(self, csv_file, root_dir, transform=None):
        """
        Arguments:
            csv_file (string): Path to the csv file with annotations.
            root_dir (string): Directory with all the images.
            transform (callable, optional): Optional transform to be applied
                on a sample.
        """
        self.landmarks_frame = pd.read_csv(csv_file)
        self.root_dir = root_dir
        self.transform = transform

    def __len__(self):
        return len(self.landmarks_frame)

    def __getitem__(self, idx):
        if torch.is_tensor(idx):
            idx = idx.tolist()

        img_name = os.path.join(self.root_dir,
                                self.landmarks_frame.iloc[idx, 0])
        image = io.imread(img_name)
        landmarks = self.landmarks_frame.iloc[idx, 1:]
        landmarks = np.array([landmarks], dtype=float).reshape(-1, 2)
        sample = {'image': image, 'landmarks': landmarks}

        if self.transform:
            sample = self.transform(sample)

        return sample

your_dataset = CustomDataset(csv_file='data/faces/face_landmarks.csv',root_dir='data/faces/')
```

# Classroom Setup

Dataset
Data Calling
# Data
# Visualization

Network Building
Optimization
Result Visualization

```python
import math
def imshow(img):
    img = img / 2 + 0.5  # unnormalize
    plt.imshow(np.transpose(img, (1, 2, 0)))
def visualize_data(images, categories,
images_per_row = 8):
    class_names = ['airplane', 'automobile',
'bird', 'cat', 'deer',
'dog', 'frog', 'horse', 'ship', 'truck']

    n_images = len(images)
    n_rows = math.ceil(float(n_images)/images_per_row)
    fig = plt.figure(figsize=(1.5*images_per_row, 1.5*n_rows))
    fig.patch.set_facecolor('white')
    for i in range(n_images):
        plt.subplot(n_rows, images_per_row, i+1)
        plt.xticks([])
        plt.yticks([])
        imshow(images[i])
        class_index = categories[i]
        plt.xlabel(class_names[class_index])
    plt.show()
```

# Classroom Setup

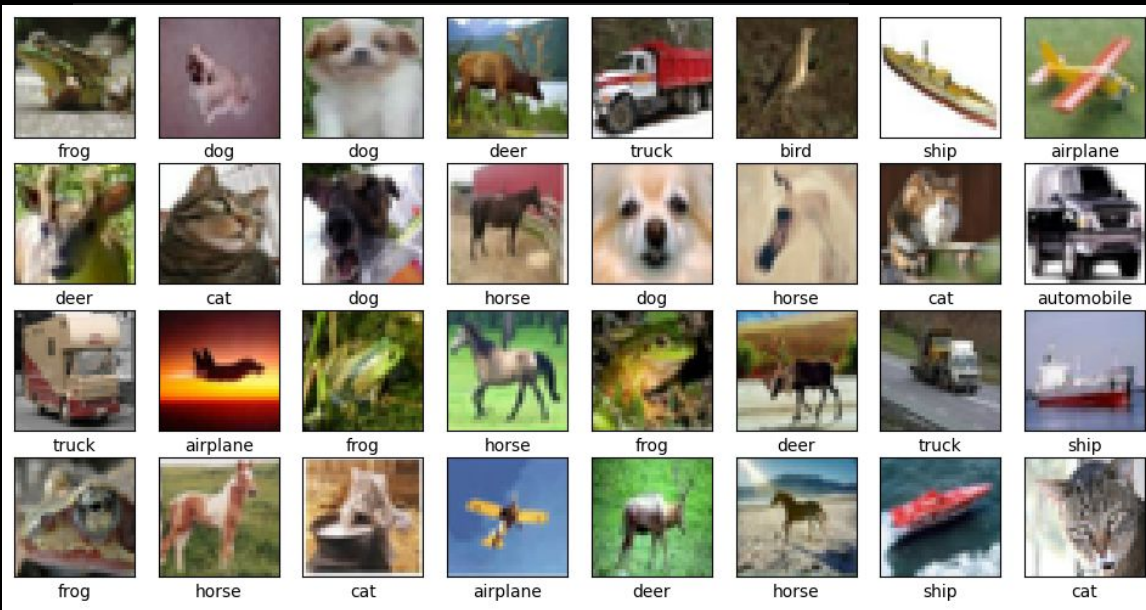Dataset
Data Calling
## Data
# Visualization

Network Building
Optimization
Result Visualization

```python
import math
def imshow(img):
    img = img / 2 + 0.5  # unnormalize
    plt.imshow(np.transpose(img, (1, 2, 0)))
```



```python
        class_index = categories[i]
        plt.xlabel(class_names[class_index])
    plt.show()
```

**Classroom Setup**

Dataset
Data Calling
Data Visualization

**Network
Building**
Optimization
Result Visualization
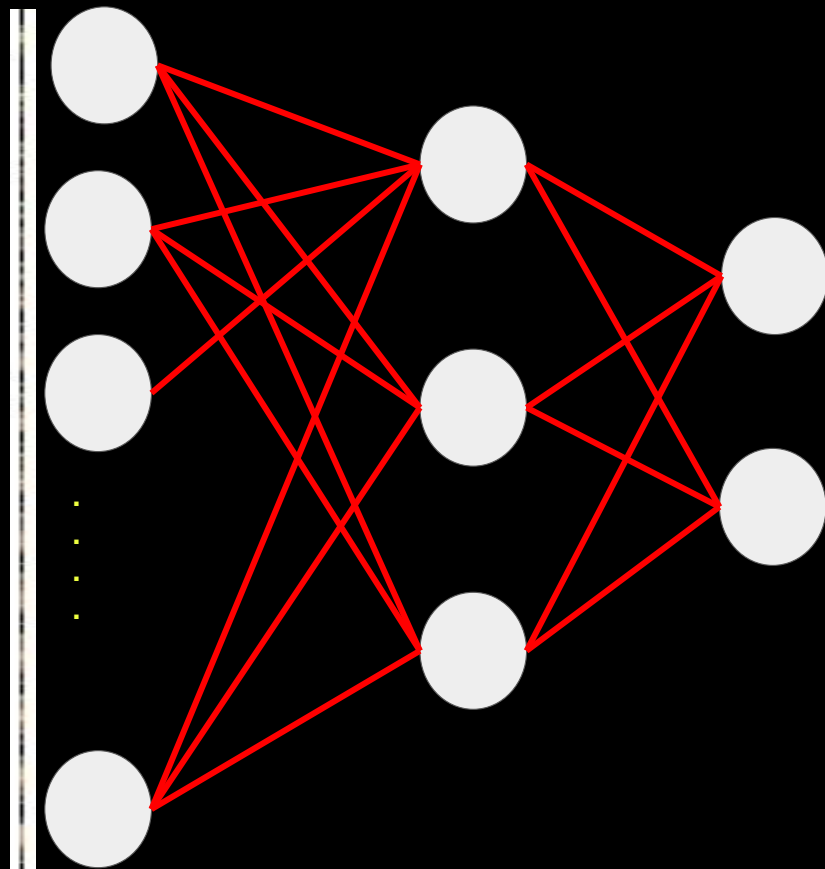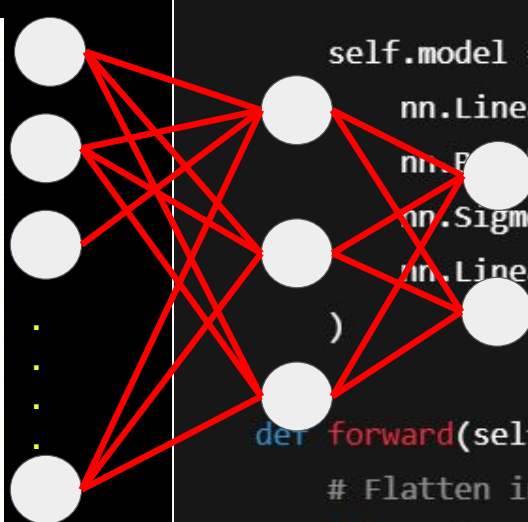
Image      Flatten

(batch, channel, H, W)     nn.Linear(32*32, 3)     nn.Linear(3, 2)

Image    Flatten

Yoshua Bengio    Geoffrey Hinton    Yann LeCun

Input Layer    Convolutional Layer    Pooling Layer    Fully Connected Layer    Output Layer

# Classroom Setup

## Dataset
## Data Calling
## Data Visualization

# Network Building

## Optimization
## Result Visualization



Input Layer · Convolutional Layer · Pooling Layer · Fully Connected Layer · Output Layer

## Complex Cell

A. Shift → Convolution

B. Size → Pooling

tolerance area

S-cell

deformed
unlearned

Classroom Setup

Dataset
Data Calling
Data Visualization
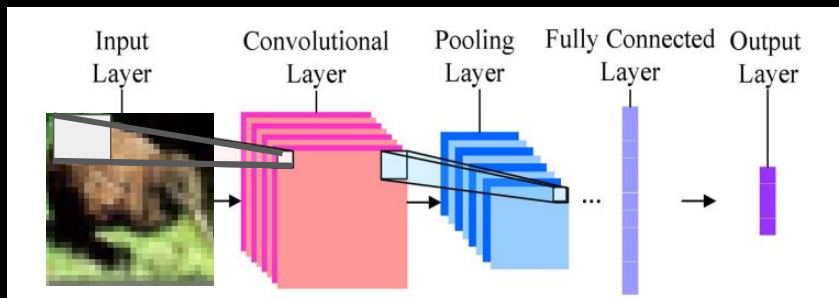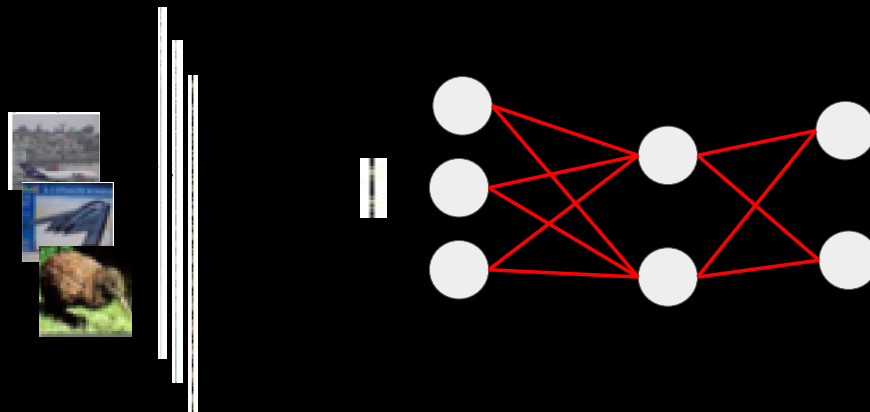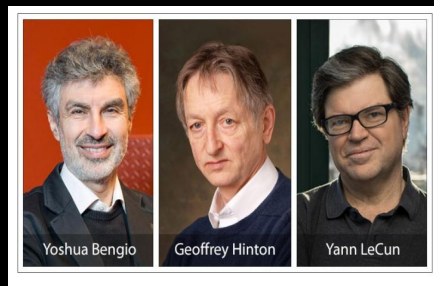
Network Building
Optimization
Result Visualization

- Revision of calculus

- Gradient descent and stochastic gradient descent
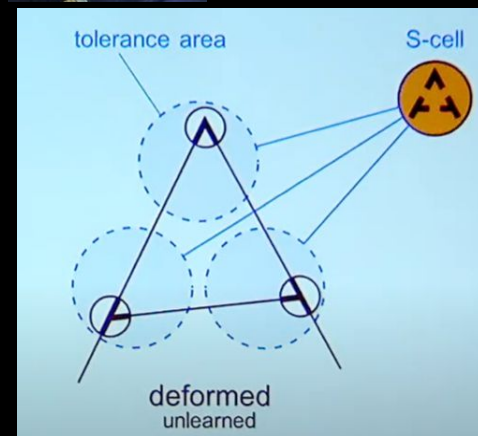
- Backpropagation in feed-forward neural networks

- Optimizers for deep learning.

| T = 0 | T = 1 | T = 2 | T = 3 | T = 4 |
|-------|-------|-------|-------|-------|
| 10 | 10 | 10 | 10 | 10 |

# Speed Vs. Velocity Vs Acceleration

Speed: 10 m/s    Average Speed: (10+10+10+10+10)/5

Velocity: 10 m/s towards your right hand side (Speed + Direction)

Acceleration:
( T = 1 )   –   ( T = 0 )        =      10  -  10  = 0
( T = 2 )   –   ( T = 1 )        =      10  -  10  = 0
( T = 3 )   –   ( T = 2 )        =      10  -  10  = 0
( T = 4 )   –   ( T = 4 )        =      10  -  10  = 0

(Observed Quantity)  /  (Controlled Quantity)

- # Revision of calculus

- Gradient descent and stochastic gradient descent

- Backpropagation in feed-forward neural networks

- Optimizers for deep learning.

| T = 0 | T = 1 | T = 2 | T = 3 | T = 4 |
|-------|-------|-------|-------|-------|
| 10 | 15 | 10 | 20 | 10 |

## Change Vs. Rate of change

Average Speed: (10+15+10+20+10) / 5       = 13 m/s
Velocity: 13 m/s towards my right hand
Acceleration
( T = 1 )  –  ( T = 0 )     =    15 - 10 =   5
( T = 2 )  –  ( T = 1 )     =    10 - 15 = - 5
( T = 3 )  –  ( T = 2 )     =    20 - 10 =   10
( T = 4 )  –  ( T = 4 )     =    10 - 20 = - 10

(Observed Quantity)  /  (Controlled Quantity)

- **Revision of calculus**

- Gradient descent and stochastic gradient descent

- Backpropagation in feed-forward neural networks

- Optimizers for deep learning.

# Numeric Values
# Numerical Methods

Speed

| T = 0 | T = 1 | T = 2 | T = 3 | T = 4 |
|:-----:|:-----:|:-----:|:-----:|:-----:|
| 10 | 15 | 10 | 20 | 10 |

Time

# Functional Values
# Functional Methods

$$f(x)=10+5 \cdot \sin(\pi x/2)+5 \cdot \sin(\pi$$



Smooth Curve Approximating [10, 15, 10, 20, 10]

$$\frac{f(t + 1) - f(t)}{(t+1) - t}$$

=

Rate of Change

$f(t)=10+5\cdot\sin(\pi t/2)+5\cdot\sin(\pi$



Smooth Curve Approximating [10, 15, 10, 20, 10]

(Observed Quantity)  /   (Controlled Quantity)

# Minimize Loss  /   ??

Network Architecture

Dataset: Batch Size

Anger

## Weight

T = 0    T = 1    T = 2

Image  Flatten    Image  Flatten    Image  Flatten

**Loss = CE(GT, Prediction)**

Ground-truth one-hot label | 0 | 1 | 0 | 0

**Image classification task**

lion    car    chair    airplane

How diverge and distant?

Prediction probability | 0.6 | 0.2 | 0.1 | 0.1

$f(x)$

Neural Net or machine learning model $f$

Input x

| W = 0 | W = 1 | W = 2 | W = 3 | W = 4 |
|-------|-------|-------|-------|-------|
| 10 | 15 | 10 | 20 | 10 |

Loss

20

15

10

5

1    2    3    4    5    Weight

# Numeric Values
# Numerical Methods

## Loss Function

$$f(L)=10+5 \cdot sin(\pi L/2)+5 \cdot sin(\pi L)$$



Smooth Curve Approximating [10, 15, 10, 20, 10]

# Local minima-maxima and saddle point

- Given an **objective function** $J(\theta)$ with $\theta = [\theta_1, \theta_2, \ldots, \theta_P]$
  - $\theta$ is said to be a **critical point** if $\nabla J(\theta) = \mathbf{0}$ (vector $\mathbf{0}$)

- Let us denote the **set of eigenvalues** of Hessian matrix $\nabla^2 J(\theta) = H(\theta)$ by
  - $\lambda_1 \leq \lambda_2 \leq \cdots \leq \lambda_P$

- **Local minima**
  - $\nabla J(\theta) = \mathbf{0}$ and $\nabla^2 J(\theta) = H(\theta) > 0$ (positive semi-definite matrix)
  - $0 \leq \lambda_1 \leq \lambda_2 \leq \cdots \leq \lambda_P$

- **Local maxima**
  - $\nabla J(\theta) = \mathbf{0}$ and $\nabla^2 J(\theta) = H(\theta) < 0$ (negative semi-definite matrix)
  - $\lambda_1 \leq \lambda_2 \leq \cdots \leq \lambda_P \leq 0$

- **Saddle point**
  - $\nabla J(\theta) = \mathbf{0}$ and $\nabla^2 J(\theta) = H(\theta) < 0$ (indefinite matrix)
  - $\lambda_1 \leq \lambda_2 \leq \cdots < 0 < \cdots \leq \lambda_P$





local min     local max     saddle point



Loss = CE(GT, Prediction)

| | W = 0 | W = 1 | W = 2 | W = 3 | W = 4 |
|------|-------|-------|-------|-------|-------|
| Loss | 10 | 15 | 10 | 20 | 10 |

# Backprop By Hand

Goto slide 7,
Week 4

# A small detour to calculus

- Calculus = **mathematics of change** (very important for deep learning)
- Properties of derivative:
  - $f'(x) = \nabla f(x) = \lim_{h \to 0} \frac{f(x+h)-f(x)}{h}$
  - $(uv)' = u'v + uv'$
  - $\left(\frac{u}{v}\right)' = \frac{u'v - uv'}{v^2}$
  - $(e^u)' = u'e^u$
  - $(\log u)' = \frac{u'}{u}$
- Multi-variate function $f: \mathbb{R}^n \to \mathbb{R}$ with $y = f(x) = f(x_1, \ldots, x_n)$.
  - Gradient/derivative: $\frac{\partial f}{\partial x}(a) = \nabla_x f(a) = [\nabla_{x_1} f(a), \nabla_{x_2} f(a), \ldots, \nabla_{x_n} f(a)]$.
- Chain rule ∞ :
  - $\frac{\partial u}{\partial x} = \frac{\partial u}{\partial v} \times \frac{\partial v}{\partial x}$

---

# Example

- $y = f(x) = f(x_1, x_2, x_3) = (x_1^2 + x_2^2, x_2^2 + x_3^2 x_2)$
  - $f: \mathbb{R}^3 \to \mathbb{R}^2$
  - $f_1(x) = f_1(x_1, x_2, x_3) = x_1^2 + x_2^2$
  - $f_2(x) = f_2(x_1, x_2, x_3) = x_2^2 + x_3^2 x_2$
  - $\frac{\partial y}{\partial x} = \nabla f \in \mathbb{R}^{2 \times 3}$
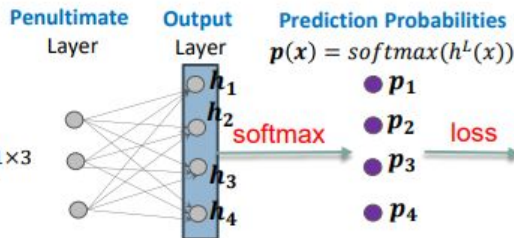
$$\frac{\partial y}{\partial x} = \nabla_x f = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \frac{\partial f_1}{\partial x_3} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \frac{\partial f_2}{\partial x_3} \end{bmatrix} = \begin{bmatrix} 2x_1 & 2x_2 & 0 \\ 0 & 2x_2 + x_3^2 & 2x_2 x_3 \end{bmatrix}$$

# Example

**Output layer**

Penultimate Layer  Output Layer  **Prediction Probabilities**

$$p(x) = softmax(h^L(x))$$

$x = [x_1\ x_2\ x_2] \in \mathbb{R}^{1\times3}$
$y = 2$

softmax → $p_1, p_2, p_3, p_4$ → loss

$$l = loss = -\log p_2$$
$$= -\log \frac{e^{h_2}}{e^{h_1}+e^{h_2}+e^{h_3}+e^{h_4}}$$

Logit $h = [h_1\ h_2\ h_3\ h_4]$   Prob $p = [p_1\ p_2\ p_3\ p_4]$

$$p_1 = \frac{e^{h_1}}{e^{h_1}+e^{h_2}+e^{h_3}+e^{h_4}}$$

$$p_2 = \frac{e^{h_2}}{e^{h_1}+e^{h_2}+e^{h_3}+e^{h_4}}$$

$$p_3 = \frac{e^{h_3}}{e^{h_1}+e^{h_2}+e^{h_3}+e^{h_4}}$$

$$p_4 = \frac{e^{h_4}}{e^{h_1}+e^{h_2}+e^{h_3}+e^{h_4}}$$

Compute $\dfrac{\partial l}{\partial h}$?

$$\overbrace{\phantom{e^{h_1} + e^{h_2} + e^{h_3} + e^{h_4}}}^{u}$$

- $l = -\log \dfrac{e^{h_2}}{e^{h_1}+e^{h_2}+e^{h_3}+e^{h_4}} = \log(e^{h_1} + e^{h_2} + e^{h_3} + e^{h_4}) - h_2$

- $\dfrac{\partial l}{\partial h_1} = \dfrac{\nabla_{h_1}u}{u} = \dfrac{e^{h_1}}{e^{h_1}+e^{h_2}+e^{h_3}+e^{h_4}} = p_1$

- $\dfrac{\partial l}{\partial h_2} = \dfrac{\nabla_{h_2}u}{u} - 1 = \dfrac{e^{h_2}}{e^{h_1}+e^{h_2}+e^{h_3}+e^{h_4}} - 1 = p_2 - 1$

- $\dfrac{\partial l}{\partial h_3} = \dfrac{\nabla_{h_3}u}{u} = \dfrac{e^{h_3}}{e^{h_1}+e^{h_2}+e^{h_3}+e^{h_4}} = p_3$

- $\dfrac{\partial l}{\partial h_4} = \dfrac{\nabla_{h_4}u}{u} = \dfrac{e^{h_4}}{e^{h_1}+e^{h_2}+e^{h_3}+e^{h_4}} = p_4$

- $\dfrac{\partial l}{\partial h} = [p_1, p_2 - 1, p_3, p_4] = [p_1, p_2, p_3, p_4] - [0,1,0,0] = \boldsymbol{p} - \mathbf{1}_2 = \boldsymbol{p} - \mathbf{1}_y$

# Example
## Intermediate layer

- $\bar{h} = xW + b$ and $h = \sigma(\bar{h})$
  - $h = \sigma(xW + b)$
  - $\sigma$ is the **activation function**

- $\dfrac{\partial h}{\partial x} = \dfrac{\partial h}{\partial \bar{h}} \times \dfrac{\partial \bar{h}}{\partial x} = diag\left(\sigma'(\bar{h})\right) W^T \in \mathbb{R}^{4\times 3}$
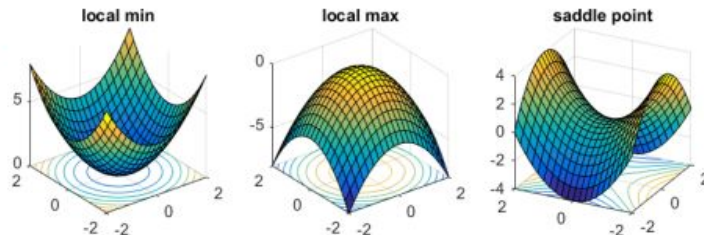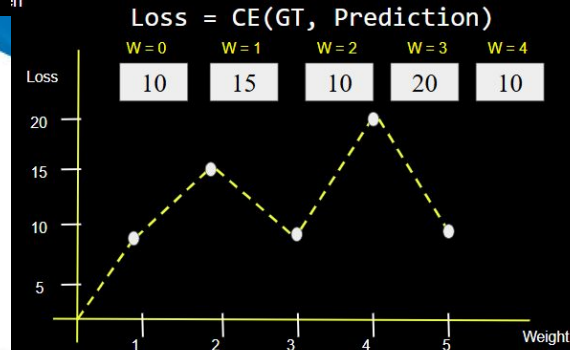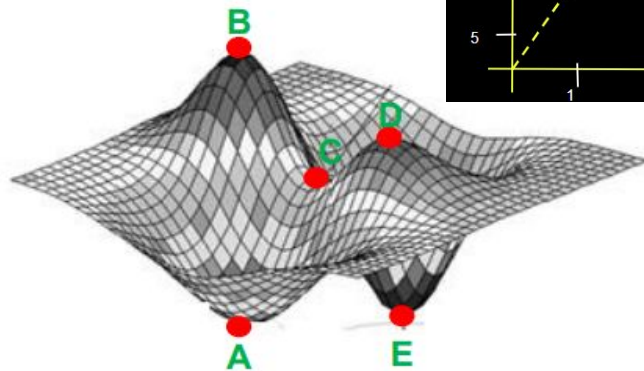
- $\dfrac{\partial h}{\partial \bar{h}} = \begin{bmatrix} \frac{\partial h_1}{\partial \bar{h}_1} & \frac{\partial h_1}{\partial \bar{h}_2} & \frac{\partial h_1}{\partial \bar{h}_3} & \frac{\partial h_1}{\partial \bar{h}_4} \\ \frac{\partial h_2}{\partial \bar{h}_1} & \frac{\partial h_2}{\partial \bar{h}_2} & \frac{\partial h_2}{\partial \bar{h}_3} & \frac{\partial h_2}{\partial \bar{h}_4} \\ \frac{\partial h_3}{\partial \bar{h}_1} & \frac{\partial h_3}{\partial \bar{h}_2} & \frac{\partial h_3}{\partial \bar{h}_3} & \frac{\partial h_3}{\partial \bar{h}_4} \\ \frac{\partial h_4}{\partial \bar{h}_1} & \frac{\partial h_4}{\partial \bar{h}_2} & \frac{\partial h_4}{\partial \bar{h}_3} & \frac{\partial h_4}{\partial \bar{h}_4} \end{bmatrix} = \begin{bmatrix} \sigma'(\bar{h}_1) & 0 & 0 & 0 \\ 0 & \sigma'(\bar{h}_2) & 0 & 0 \\ 0 & 0 & \sigma'(\bar{h}_3) & 0 \\ 0 & 0 & 0 & \sigma'(\bar{h}_4) \end{bmatrix} = diag(\sigma'(\bar{h}))$

- $\dfrac{\partial \bar{h}}{\partial x} = W^T$

$h_1\ h_2 h_3 h_4$

$h = \sigma(\bar{h})$

$\sigma\ \sigma\ \sigma\ \sigma$

$\bar{h}_1\ \bar{h}_2 \bar{h}_3 \bar{h}_4$

$\bar{h} = xW + b$

$W \in \mathbb{R}^{3\times 4}$
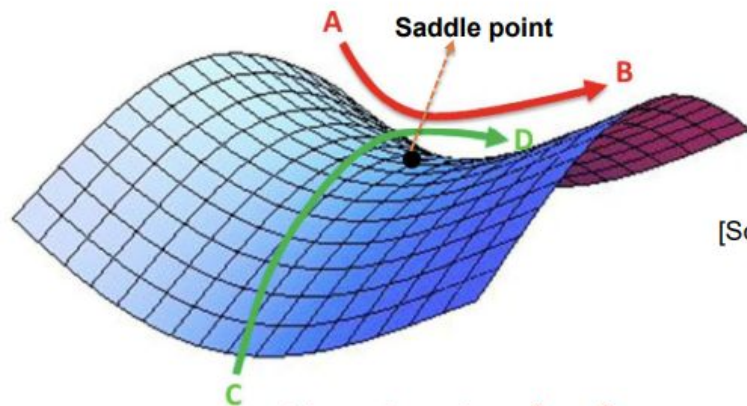$b \in \mathbb{R}^{1\times 4}$

$x \in \mathbb{R}^{1\times 3}$

# Local minima-maxima and saddle point

- Given an **objective function** $J(\theta)$ with $\theta = [\theta_1, \theta_2, \ldots, \theta_P]$
  - $\theta$ is said to be a **critical point** if $\nabla J(\theta) = \mathbf{0}$ (vector $\mathbf{0}$)

- Let us denote the **set of eigenvalues** of Hessian matrix $\nabla^2 J(\theta) = H(\theta)$ by
  - $\lambda_1 \leq \lambda_2 \leq \cdots \leq \lambda_P$

- **Local minima**
  - $\nabla J(\theta) = \mathbf{0}$ and $\nabla^2 J(\theta) = H(\theta) > 0$ (positive semi-definite matrix)
  - $0 \leq \lambda_1 \leq \lambda_2 \leq \cdots \leq \lambda_P$

- **Local maxima**
  - $\nabla J(\theta) = \mathbf{0}$ and $\nabla^2 J(\theta) = H(\theta) < 0$ (negative semi-definite matrix)
  - $\lambda_1 \leq \lambda_2 \leq \cdots \leq \lambda_P \leq 0$

- **Saddle point**
  - $\nabla J(\theta) = \mathbf{0}$ and $\nabla^2 J(\theta) = H(\theta) < 0$ (indefinite matrix)
  - $\lambda_1 \leq \lambda_2 \leq \cdots < 0 < \cdots \leq \lambda_P$





local min          local max          saddle point



Loss = CE(GT, Prediction)

| | W = 0 | W = 1 | W = 2 | W = 3 | W = 4 |
|------|-------|-------|-------|-------|-------|
| Loss | 10 | 15 | 10 | 20 | 10 |

# More on saddle point



[Source: Internet]

$$f(\theta) = f(\theta_1, \theta_2) = \theta_1^2 - \theta_2^2$$

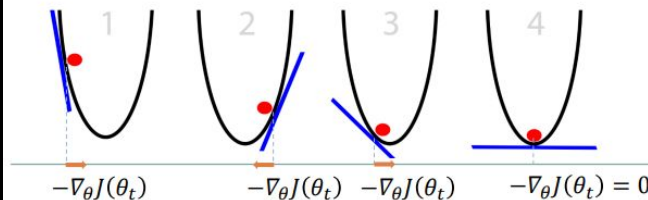**Gradient** $g = \begin{bmatrix} \frac{\partial f}{\partial \theta_1} \\ \frac{\partial f}{\partial \theta_2} \end{bmatrix} = \begin{bmatrix} 2\theta_1 \\ -2\theta_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \rightarrow$ a **critical point** $\begin{bmatrix} 0 \\ 0 \end{bmatrix}$.

**Hessian matrix** is $H = \begin{bmatrix} \frac{\partial^2 f}{\partial \theta_1^2} & \frac{\partial^2 f}{\partial \theta_1 \partial \theta_2} \\ \frac{\partial^2 f}{\partial \theta_2 \partial \theta_1} & \frac{\partial^2 f}{\partial \theta_2^2} \end{bmatrix} = \begin{bmatrix} 2 & 0 \\ 0 & -2 \end{bmatrix}$
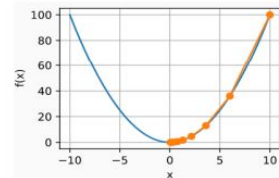
**Two eigenvalues** $\lambda_1 = -2 < 0 < 2 = \lambda_2 \rightarrow \begin{bmatrix} 0 \\ 0 \end{bmatrix}$ is a **saddle point**.

# Gradient descend



$-\nabla_\theta J(\theta_t)$     $-\nabla_\theta J(\theta_t)$   $-\nabla_\theta J(\theta_t)$    $-\nabla_\theta J(\theta_t) = 0$

❑ We need to solve
- $\min_\theta J(\theta)$

❑ Follow to **the opposite side** of the current gradient
- $\theta_{t+1} = \theta_t - \eta\nabla_\theta J(\theta_t)$ where $\boldsymbol{\eta > 0}$ is the **learning rate**.

❑ Guarantee to converge to a **global minima** if $J(.)$ is **convex**.

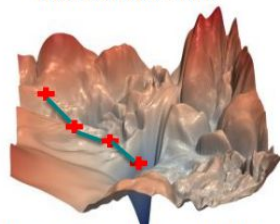❑ Get stuck in a **local minima** or **saddle points** if $J(.)$ is non-convex.



**Convex case**



(Source: www.cs.ubc.ca)

**Non-convex case**



**DL case: easy to get stuck in saddle points**

21

# Gradient descend

## Algorithm

❑ **Input**: objective function $J(\boldsymbol{\theta})$

❑ **Output**: optimal solution $\boldsymbol{\theta}^*$

1. Initialize parameters $\theta_0$ randomly $\sim N(0, \sigma^2)$.

2. for t=1 to T

3.        Compute gradients $\nabla_\theta J(\theta_t) = \frac{\partial J}{\partial \theta}(\theta_t)$

4.        Update $\theta_{t+1} = \theta_t - \eta_t \nabla_\theta J(\theta_t)$

5. Return $\theta^* = \theta_{T+1}$