## Load essential libraries

```python
import numpy as np
import matplotlib.pyplot as plt
plt.style.use('dark_background')
%matplotlib inline

import tensorflow as tf
```

## Check TensorFlow version

```python
tf.__version__
```

    '2.15.0'

Answer the following questions inline using the documentation from:

- Introduction to tensors https://www.tensorflow.org/guide/tensor
- Introduction to variables https://www.tensorflow.org/guide/variable
- Introduction to gradients and automatic differentiation https://www.tensorflow.org/guide/autodiff

1. A scalar is a rank $\underbrace{0/1/2}_{\text{choose one}}$ tensor.
2. *True/false*: a scalar has no axes.
3. A matrix is a rank $\underbrace{0/1/2}_{\text{choose one}}$ tensor and has $\underbrace{1/2}_{\text{choose one}}$ axes.
4. What does the function call `tf.reshape(A, [-1]` does for a given tensor `A`?
5. *True/false*: `tf.reshape()` can be used to swap axes of a tensor such as `(patients, timestamps, features)` to `(timestamps, patients, features)`,
6. `tf.keras` uses _____ to store model parameters.
7. *True/false*: calling `assign` reuses a tensor's exisiting memory to assign the values.
8. *True/false*: creating a new tensor `b` based on the value of another tensor `a` as `b = tf.Variable(a)` will have the tensors allocated different memory.
9. *True/false*: two tensor variables can have the same name.
10. An example of a variable that would not need gradients is a _____.
11. Tensor variables are typically placed in $\underbrace{CPU/GPU}_{\text{choose one}}$.
12. *True/false*: a tensor variable is trainable by default.
13. A gradient tape tape will automatically watch a $\underbrace{\texttt{tf.Variable/tf.constant}}_{\text{choose one}}$ but not a $\underbrace{\texttt{tf.Variable/tf.constant}}_{\text{choose one}}$.
14. What attribute can be used to calculate a layer's gradient w.r.t. all its trainable variables?
15. The option `persistent = True` for a gradient tape $\underbrace{\texttt{stores/discards}}_{\text{choose one}}$ all intermediate results during the forward pass.
16. Executing the statement `print(type(x).__name__)` when `x` is a constant and when `x` is a variable results in what?
17. Which one among `tf.Tensor` and `tf.Variable` is immutable? Which one has no state but only value? Which one has a state which is actually its value?

Answers:

1. Scalar is a rank **0** tensor.
2. True.
3. A matrix is a rank **2** tensor and has 2 axes.
4. Fits the tensor in "whatever good shape" that is possible.
5. False. We use tf..transpose for this operation.
6. tf.keras use tf.Variable to store model parameters.
7. True.
8. True. Two variables will not share the same memory.
9. True.
10. Constant.

11. GPU. For better performance, tensorflow will appempt to place a variable on a faster running device compatible with its dtype. So most variables are placed on a GPU, if available. However, we can override this, and place a float tensor and a variable on a CPU, even if a GPU is available.

12. True. We can turn it off during creation by setting trainable=False.

13. A gradient tape will automatically watch a tf.Variable but not a tf.constant.

14. Module.trainable_variables.

15. stores.

16. If x is a constant, the output will be the data type of x.

    If x is a variable, the output will be the data type of the value stored in that variable at the time of execution.

17. Immutable: tf.Tensor.

    No state but only one value: tf.Tensor.

    Has a state which is actually its value: tf.Variable.

---

Consider a 1-layer neural network for a sample with 3 features: heart rate, blood pressure, and temperature and 2 possible output categories: diabetic and non-diabetic.

An individual who is diabetic has heart rate = 76 BPM, BP = 120 mm Hg, and temperature = $37.5\,^\circ\mathrm{C}$.

Here is the forward propagation:

$$\mathbf{x} \longrightarrow \mathbf{x}_B = \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix} \longrightarrow \mathbf{z}^{[1]} = \mathbf{W}^{[1]}\mathbf{x}_B \longrightarrow \underbrace{\mathbf{a}^{[1]}}_{=\hat{\mathbf{y}}} = \mathrm{softmax}\left(\mathbf{a}^{[1]}\right) \longrightarrow L = \sum_{k=0}^{1} -y_k \log(\hat{y}_k).$$

- Fill in the missing entries of the code below to calculate the gradients there-in.
- Explain why some gradient shapes do not seem to match with the usual $\mathrm{input\ shape} \times \mathrm{output\ shape}$ rule for gradient shapes using the documentation on [Gradients of non-scalar target](#) as resource.
- Try `persistent = true` and `persistent = false` in the gradient tape and observe what happens in each case.

---

```
x = tf.constant([76., 120., 37.5], dtype = float)
y = tf.constant([0, 1])
xB = tf.concat([x, 1.0*tf.ones(1)], axis = 0)
W = tf.Variable(0.01*(tf.random.normal((2, 4))))

with tf.GradientTape(persistent = True) as g:
  z = tf.linalg.matvec(W, xB)
  a = tf.nn.softmax(z)
  yhat = a
  L = tf.reduce_sum(-tf.one_hot(indices=y, depth=2))*tf.math.log(yhat)
  L.numpy()

print('Loss = ', L.numpy())
print('Gradient of L w.r.t. yhat')
gradL_yhat = g.gradient(L, yhat)
print(gradL_yhat)
print('-----')
print('Gradient of a w.r.t. z')
grada_z = g.gradient(a, z)
print(grada_z)
print('-----')
print('Gradient of z w.r.t. W')
gradz_W = g.gradient(z, W)
print(gradz_W)
print('-----')
print('Gradient of L w.r.t. W')
gradL_W = g.gradient(L, W)
print(gradL_W)

# Delete gradient tape and release memory
del g
```

```
    Loss =  [1.5003958 1.2783527]
    Gradient of L w.r.t. yhat
    tf.Tensor([-4.234838 -3.789839], shape=(2,), dtype=float32)
    -----
    Gradient of a w.r.t. z
    tf.Tensor([0. 0.], shape=(2,), dtype=float32)
    -----
    Gradient of z w.r.t. W
    tf.Tensor(
    [[ 76.  120.   37.5   1. ]
     [ 76.  120.   37.5   1. ]], shape=(2, 4), dtype=float32)
    -----
    Gradient of L w.r.t. W
    tf.Tensor(
    [[ -8.428984   -13.308921    -4.1590376   -0.11090767]
     [  8.428977    13.308911     4.1590347    0.11090759]], shape=(2, 4), dtype=float32)
```

Recalculate gradients pen-and-paper-way with the same weights from above using numpy. Compare the gradient results here with the ones that you had from the previous cell. Why are some gradients different? In both approaches (in this cell and in the one above), is the gradient of interest $\nabla_{\mathbf{W}^{[1]}}(L)$ the same? Note that this is the only gradient we need to update the weights matrix $\mathbf{W}^{[1]}$.

```python
x = tf.constant([76., 120., 37.5], dtype = float)
y = tf.constant([0, 1])
xB = tf.concat([x, 1.0*tf.ones(1)], axis = 0)
W = tf.Variable(0.01*(tf.random.normal((2, 4))))

with tf.GradientTape(persistent = True) as g:
  z = tf.linalg.matvec(W, xB)
  a = tf.nn.softmax(z)
  yhat = a
  L = tf.reduce_sum(-tf.one_hot(indices=y, depth=2))*tf.math.log(yhat)
  L.numpy()

print('Loss = ', L.numpy())
print('Gradient of L w.r.t. yhat')
gradL_yhat = g.gradient(L, yhat)
print(gradL_yhat)
print('-----')
print('Gradient of a w.r.t. z')
grada_z = g.gradient(a, z)
print(grada_z)
print('-----')
print('Gradient of z w.r.t. W')
gradz_W = g.gradient(z, W)
print(gradz_W)
print('-----')
print('Gradient of L w.r.t. W')
gradL_W = g.gradient(L, W)
print(gradL_W)

# Delete gradient tape and release memory
del g
```

```
Loss =  [0.06471884 6.8939757 ]
Gradient of L w.r.t. yhat
tf.Tensor([ -2.0657773 -62.811306 ], shape=(2,), dtype=float32)
-----
Gradient of a w.r.t. z
tf.Tensor([5.7706746e-08 1.8978954e-09], shape=(2,), dtype=float32)
-----
Gradient of z w.r.t. W
tf.Tensor(
[[ 76.  120.   37.5   1. ]
 [ 76.  120.   37.5   1. ]], shape=(2, 4), dtype=float32)
-----
Gradient of L w.r.t. W
tf.Tensor(
[[ 142.32019    224.7161      70.22378       1.8726342]
 [-142.32022   -224.71613    -70.22379      -1.8726344]], shape=(2, 4), dtype=float32)
```

```python
x = tf.constant([76., 120., 37.5], dtype = float)
y = tf.constant([0, 1])
xB = tf.concat([x, 1.0*tf.ones(1)], axis = 0)
W = tf.Variable(0.01*(tf.random.normal((2, 4))))

with tf.GradientTape(persistent = False) as g:
  z = tf.linalg.matvec(W, xB)
  a = tf.nn.softmax(z)
  yhat = a
  L = tf.reduce_sum(-tf.one_hot(indices=y, depth=2))*tf.math.log(yhat)
  L.numpy()

print('Loss = ', L.numpy())
print('Gradient of L w.r.t. yhat')
gradL_yhat = g.gradient(L, yhat)
print(gradL_yhat)
print('-----')
print('Gradient of a w.r.t. z')
grada_z = g.gradient(a, z)
print(grada_z)
print('-----')
print('Gradient of z w.r.t. W')
gradz_W = g.gradient(z, W)
print(gradz_W)
print('-----')
print('Gradient of L w.r.t. W')
```

```
gradL_W = g.gradient(L, W)
print(gradL_W)

# Delete gradient tape and release memory
del g
```

```
    Loss =  [1.4842352 1.2929265]
    Gradient of L w.r.t. yhat
    tf.Tensor([-4.200757  -3.8175561], shape=(2,), dtype=float32)
    -----
    Gradient of a w.r.t. z
    ---------------------------------------------------------------------------
    RuntimeError                              Traceback (most recent call last)
    <ipython-input-5-ea1c3274ae47> in <cell line: 19>()
         17 print('-----')
         18 print('Gradient of a w.r.t. z')
    ---> 19 grada_z = g.gradient(a, z)
         20 print(grada_z)
         21 print('-----')

    /usr/local/lib/python3.10/dist-packages/tensorflow/python/eager/backprop.py in
    gradient(self, target, sources, output_gradients, unconnected_gradients)
       1003      """
       1004      if self._tape is None:
    -> 1005        raise RuntimeError("A non-persistent GradientTape can only be used to
    "
       1006                           "compute one set of gradients (or jacobians)")
       1007      if self._recording:

    RuntimeError: A non-persistent GradientTape can only be used to compute one set of
    gradients (or jacobians)
```

Next steps:     **Explain error**

---

For each activation function below,

1. Sigmoid $\sigma(z)$
2. Hyperbolic tangent $\tanh(z)$
3. Rectified Linear Unit $\mathrm{ReLU}(z)$
4. Leaky rectified linear unit $\mathrm{LReLU}(z)$

- plot the activation and its gradient in the same figure for raw score values $z$ ranging between $-10$ and $10$;
- comment on whether the backward flowing gradient on the input side of the activation layer will have a smaller or bigger magnitude compared to the backward flowing gradient on the output side of the activation layer. Recall that what connects these two gradients is the local gradient of the activation layer which you may have just plotted.

---

```
z = tf.linspace(-10.0, 10.0, 129) # A tf.Tensor, not a tf.Variable

with tf.GradientTape(persistent = True) as g:
    g.watch(z)
    a_sigmoid = tf.math.sigmoid(z)
    a_tanh = tf.math.tanh(z)
    a_ReLU = tf.cast(z, tf.float64) * tf.cast((z > 0), tf.float64)
    a_LReLU = tf.cast(z, tf.float64) * tf.cast((z > 0.0), tf.float64) + 0.01 * tf.cast(z, tf.float64) * tf.cast((z <= 0.0), tf.float64)



grada_sigmoid_z = g.gradient(a_sigmoid, z)
grada_tanh_z = g.gradient(a_tanh, z)
grada_ReLU_z = g.gradient(a_ReLU, z)
grada_LReLU_z = g.gradient(a_LReLU, z)

fig, axs = plt.subplots(2, 2, figsize = (8, 8))
axs[0, 0].plot(z, a_sigmoid, label = 'activated score')
axs[0, 0].plot(z, grada_sigmoid_z, label='gradient of activated score')
axs[0, 0].legend(loc = 'upper left')
axs[0, 0].set_xlabel('z')
axs[0, 0].set_title('Sigmoid activation and gradient');

axs[0, 1].plot(z, a_tanh, label = 'activated score')
axs[0, 1].plot(z, grada_tanh_z, label='gradient of activated score')
axs[0, 1].legend(loc = 'upper left')
axs[0, 1].set_xlabel('z')
axs[0, 1].set_title('Tanh activation and gradient');

axs[1, 0].plot(z, a_ReLU, label = 'activated score')
axs[1, 0].plot(z, grada_ReLU_z, label='gradient of activated score')
axs[1, 0].legend(loc = 'upper left')
axs[1, 0].set_xlabel('z')
axs[1, 0].set_title('ReLU activation and gradient');
```

```
axs[1, 1].plot(z, a_LReLU, label = 'activated score')
axs[1, 1].plot(z, grada_LReLU_z, label='gradient of activated score')
axs[1, 1].legend(loc = 'upper left')
axs[1, 1].set_xlabel('z')
axs[1, 1].set_title('Leaky ReLU activation and gradient');
```