

▼ Ungraded lab: Knowledge Distillation

Welcome, during this ungraded lab you are going to perform a model compression technique known as **knowledge distillation** in which a `student` model "learns" from a more complex model known as the `teacher`. In particular you will:

1. Define a `Distiller` class with the custom logic for the distillation process.
2. Train the `teacher` model which is a CNN that implements regularization via dropout.
3. Train a `student` model (a smaller version of the teacher without regularization) by using knowledge distillation.
4. Train another `student` model from scratch without distillation called `student_scratch`.
5. Compare the three students.

This notebook is based on [this](#) official Keras tutorial.

If you want a more theoretical approach to this topic be sure to check this paper [Hinton et al. \(2015\)](#).

Let's get started!

▼ Imports

```
1 # For setting random seeds
2 import os
3 os.environ['PYTHONHASHSEED']=str(42)
4
5 # Libraries
6 import random
7 import numpy as np
8 import pandas as pd
9 import seaborn as sns
10 import tensorflow as tf
11 from tensorflow import keras
12 import matplotlib.pyplot as plt
13 import tensorflow_datasets as tfds
14
15 # More random seed setup
16 tf.random.set_seed(42)
```

```
17 np.random.seed(42)
```

▼ Prepare the data

For this lab you will use the [cats vs dogs](#) which is composed of many images of cats and dogs along with their respective labels.

Begin by downloading the data:

```
1 # Define train/test splits
2 splits = ['train[:80%]', 'train[80%:90%]', 'train[90%:]']
3
4 # Download the dataset
5 (train_examples, validation_examples, test_examples), info = t
6
7 # Print useful information
8 num_examples = info.splits['train'].num_examples
9 num_classes = info.features['label'].num_classes
10
11 print(f"There are {num_examples} images for {num_classes} clas
```

Downloading and preparing dataset cats_vs_dogs/4.0.0 (download: 786.68 MiB, gene
DI Completed...: 100% 1/1 [00:09<00:00, 9.75s/ url]

DI Size...: 100% 786/786 [00:09<00:00, 84.72 MiB/s]

WARNING:absl:1738 images were corrupted and were skipped
Shuffling and writing examples to /root/tensorflow_datasets/cats_vs_dogs/4.0.0.i
100% 23261/23262 [00:02<00:00, 6899.65 examples/s]
Dataset cats_vs_dogs downloaded and prepared to /root/tensorflow_datasets/cats_v
There are 23262 images for 2 classes.

```
1 info
```

```
tfds.core.DatasetInfo(
  name='cats_vs_dogs',
  version=4.0.0,
  description='A large set of images of cats and dogs. There are 1738 corrupted
  homepage='https://www.microsoft.com/en-us/download/details.aspx?id=54765',
  features=FeaturesDict({
    'image': Image(shape=(None, None, 3), dtype=tf.uint8),
    'image/filename': Text(shape=(), dtype=tf.string),
    'label': ClassLabel(shape=(), dtype=tf.int64, num_classes=2),
  }),
  total_num_examples=23262,
```

```

splits={
    'train': 23262,
},
supervised_keys=('image', 'label'),
citation="""@Inproceedings (Conference){asirra-a-captcha-that-exploits-inter
author = {Elson, Jeremy and Douceur, John (JD) and Howell, Jon and Saul, Jar
title = {Asirra: A CAPTCHA that Exploits Interest-Aligned Manual Image Categ
booktitle = {Proceedings of 14th ACM Conference on Computer and Communicatio
year = {2007},
month = {October},
publisher = {Association for Computing Machinery, Inc.},
url = {https://www.microsoft.com/en-us/research/publication/asirra-a-captcha
edition = {Proceedings of 14th ACM Conference on Computer and Communications
}""",
redistribution_info=,
)

```

Preprocess the data for training by normalizing pixel values, reshaping them and creating batches of data:

```

1 # Some global variables
2 pixels = 224
3 IMAGE_SIZE = (pixels, pixels)
4 BATCH_SIZE = 32
5
6 # Apply resizing and pixel normalization
7 def format_image(image, label):
8     image = tf.image.resize(image, IMAGE_SIZE) / 255.0
9     return image, label
10
11 # Create batches of data
12 train_batches = train_examples.shuffle(num_examples // 4).map(
13 validation_batches = validation_examples.map(format_image).bat
14 test_batches = test_examples.map(format_image).batch(1)

```

▼ Code the custom Distiller model

In order to implement the distillation process you will create a custom Keras model which you will name `Distiller`. In order to do this you need to override some of the vanilla methods of a `keras.Model` to include the custom logic for the knowledge distillation. You need to override these methods:

- `compile`: This model needs some extra parameters to be compiled such as the teacher and student losses, the alpha and the temperature.

- `train_step`: Controls how the model is trained. This will be where the actual knowledge distillation logic will be found. This method is what is called when you do `model.fit`.
- `test_step`: Controls the evaluation of the model. This method is what is called when you do `model.evaluate`.

To learn more about customizing models check out the [official docs](#).

```

1 class Distiller(keras.Model):
2
3     # Needs both the student and teacher models to create an ins
4     def __init__(self, student, teacher):
5         super(Distiller, self).__init__()
6         self.teacher = teacher
7         self.student = student
8
9
10    # Will be used when calling model.compile()
11    def compile(self, optimizer, metrics, student_loss_fn,
12                distillation_loss_fn, alpha, temperature):
13
14        # Compile using the optimizer and metrics
15        super(Distiller, self).compile(optimizer=optimizer, metr
16
17        # Add the other params to the instance
18        self.student_loss_fn = student_loss_fn
19        self.distillation_loss_fn = distillation_loss_fn
20        self.alpha = alpha
21        self.temperature = temperature
22
23
24    # Will be used when calling model.fit()
25    def train_step(self, data):
26        # Data is expected to be a tuple of (features, labels)
27        x, y = data
28
29        # Vanilla forward pass of the teacher
30        # Note that the teacher is NOT trained
31        teacher_predictions = self.teacher(x, training=False)
32
33        # Use GradientTape to save gradients
34        with tf.GradientTape() as tape:

```

```

35         # Vanilla forward pass of the student
36         student_predictions = self.student(x, training=True)
37
38         # Compute vanilla student loss
39         student_loss = self.student_loss_fn(y, student_predi
40
41         # Compute distillation loss
42         # Should be KL divergence between logits softened by
43         distillation_loss = self.distillation_loss_fn(
44             tf.nn.softmax(teacher_predictions / self.tempera
45             tf.nn.softmax(student_predictions / self.tempera
46
47         # Compute loss by weighting the two previous losses
48         loss = self.alpha * student_loss + (1 - self.alpha)
49
50         # Use tape to calculate gradients for student
51         trainable_vars = self.student.trainable_variables
52         gradients = tape.gradient(loss, trainable_vars)
53
54         # Update student weights
55         # Note that this done ONLY for the student
56         self.optimizer.apply_gradients(zip(gradients, trainable_
57
58         # Update the metrics
59         self.compiled_metrics.update_state(y, student_prediction
60
61         # Return a performance dictionary
62         # You will see this being outputted during training
63         results = {m.name: m.result() for m in self.metrics}
64         results.update({"student_loss": student_loss, "distillat
65         return results
66
67
68         # Will be used when calling model.evaluate()
69         def test_step(self, data):
70             # Data is expected to be a tuple of (features, labels)
71             x, y = data
72
73             # Use student to make predictions
74             # Notice that the training param is set to False
75             y_prediction = self.student(x, training=False)
76

```

```

77     # Calculate student's vanilla loss
78     student_loss = self.student_loss_fn(y, y_prediction)
79
80     # Update the metrics
81     self.compiled_metrics.update_state(y, y_prediction)
82
83     # Return a performance dictionary
84     # You will see this being outputted during inference
85     results = {m.name: m.result() for m in self.metrics}
86     results.update({"student_loss": student_loss})
87     return results

```

▼ Teacher and student models

For the models you will use a standard CNN architecture that implements regularization via some dropout layers (in the case of the teacher), but it could be any Keras model.

Define the `create_model` functions to create models with the desired architecture using Keras' [Sequential Model](#).

Notice that `create_small_model` returns a simplified version of the model (in terms of number of layers and absence of regularization) that `create_big_model` returns:

```

1 # Teacher model
2 def create_big_model():
3     tf.random.set_seed(42)
4     model = keras.models.Sequential([
5         keras.layers.Conv2D(32, (3, 3), activation='relu', input_s
6         keras.layers.MaxPooling2D((2, 2)),
7         keras.layers.Conv2D(64, (3, 3), activation='relu'),
8         keras.layers.MaxPooling2D((2, 2)),
9         keras.layers.Dropout(0.2),
10        keras.layers.Conv2D(64, (3, 3), activation='relu'),
11        keras.layers.MaxPooling2D((2, 2)),
12        keras.layers.Conv2D(128, (3, 3), activation='relu'),
13        keras.layers.MaxPooling2D((2, 2)),
14        keras.layers.Dropout(0.5),
15        keras.layers.Flatten(),
16        keras.layers.Dense(512, activation='relu'),
17        keras.layers.Dense(2)
18    ])
19

```

```

20 return model
21
22
23
24 # Student model
25 def create_small_model():
26     tf.random.set_seed(42)
27     model = keras.models.Sequential([
28         keras.layers.Conv2D(32, (3, 3), activation='relu', input_s
29         keras.layers.MaxPooling2D((2, 2)),
30         keras.layers.Flatten(),
31         keras.layers.Dense(2)
32     ])
33
34 return model

```

There are two important things to notice:

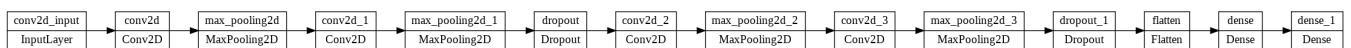
- The last layer does not have an softmax activation because the raw logits are needed for the knowledge distillation.
- Regularization via dropout layers will be applied to the teacher but NOT to the student. This is because the student should be able to learn this regularization through the distillation process.

Remember that the student model can be thought of as a simplified (or compressed) version of the teacher model.

```

1 # Create the teacher
2 teacher = create_big_model()
3
4 # Plot architecture
5 keras.utils.plot_model(teacher, rankdir="LR")

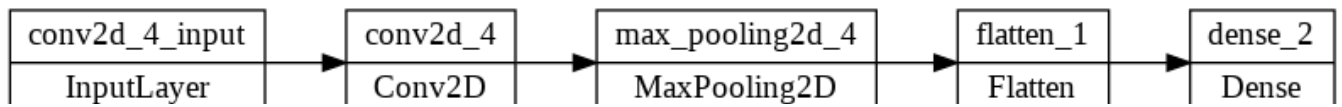
```



```

1 # Create the student
2 student = create_small_model()
3
4 # Plot architecture
5 keras.utils.plot_model(student, rankdir="LR")

```



Check the actual difference in number of trainable parameters (weights and biases) between both models:

```
1 # Calculates number of trainable params for a given model
2 def num_trainable_params(model):
3     return np.sum([np.prod(v.get_shape()) for v in model.trainab
4
5
6 student_params = num_trainable_params(student)
7 teacher_params = num_trainable_params(teacher)
8
9 print(f"Teacher model has: {teacher_params} trainable params.\
10 print(f"Student model has: {student_params} trainable params.\
11 print(f"Teacher model is roughly {teacher_params//student_para

Teacher model has: 9568898 trainable params.

Student model has: 789442 trainable params.

Teacher model is roughly 12 times bigger than the student model.
```

▼ Train the teacher

In knowledge distillation it is assumed that the teacher has already been trained so the natural first step is to train the teacher. You will do so for a total of 8 epochs:

```
1 # Compile the teacher model
2 teacher.compile(
3     loss=tf.keras.losses.SparseCategoricalCrossentropy(from_lo
4     optimizer=keras.optimizers.Adam(),
5     metrics=[tf.keras.metrics.SparseCategoricalAccuracy()])
6 )
7
8 # Fit the model and save the training history (will take from
9 teacher_history = teacher.fit(train_batches, epochs=8, validat

Epoch 1/8
582/582 [=====] - 41s 48ms/step - loss: 0.6794 - sparse
Epoch 2/8
582/582 [=====] - 31s 46ms/step - loss: 0.5932 - sparse
Epoch 3/8
```



```

582/582 [=====] - 31s 47ms/step - loss: 0.5168 - sparse
Epoch 4/8
582/582 [=====] - 31s 46ms/step - loss: 0.4571 - sparse
Epoch 5/8
582/582 [=====] - 31s 46ms/step - loss: 0.3968 - sparse
Epoch 6/8
582/582 [=====] - 31s 46ms/step - loss: 0.3473 - sparse
Epoch 7/8
582/582 [=====] - 31s 46ms/step - loss: 0.3106 - sparse
Epoch 8/8
582/582 [=====] - 31s 46ms/step - loss: 0.2628 - sparse

```

▼ Train a student from scratch for reference

In order to assess the effectiveness of the distillation process, train a model that is equivalent to the student but without doing knowledge distillation. Notice that the training is done for only 5 epochs:

```

1 # Create student_scratch model with the same characteristics a
2 student_scratch = create_small_model()
3
4 # Compile it
5 student_scratch.compile(
6     loss=tf.keras.losses.SparseCategoricalCrossentropy(from_lo
7     optimizer=keras.optimizers.Adam(),
8     metrics=[tf.keras.metrics.SparseCategoricalAccuracy()])
9 )
10
11 # Train and evaluate student trained from scratch (will take a
12 student_scratch_history = student_scratch.fit(train_batches, e

Epoch 1/5
582/582 [=====] - 29s 42ms/step - loss: 0.7741 - sparse
Epoch 2/5
582/582 [=====] - 28s 42ms/step - loss: 0.4950 - sparse
Epoch 3/5
582/582 [=====] - 28s 42ms/step - loss: 0.3865 - sparse
Epoch 4/5
582/582 [=====] - 28s 42ms/step - loss: 0.2919 - sparse
Epoch 5/5
582/582 [=====] - 28s 42ms/step - loss: 0.1997 - sparse

```

▼ Knowledge Distillation

To perform the knowledge distillation process you will use the custom model you previously coded. To do so, begin by creating an instance of the `Distiller` class and passing in the student and teacher models. Then compile it with the appropriate parameters and train it!

The two student models are trained for only 5 epochs unlike the teacher that was trained for 8. This is done to showcase that the knowledge distillation allows for quicker training times as the student learns from an already trained model

```
1 # Create Distiller instance
2 distiller = Distiller(student=student, teacher=teacher)
3
4 # Compile Distiller model
5 distiller.compile(
6     student_loss_fn=keras.losses.SparseCategoricalCrossentropy
7     optimizer=keras.optimizers.Adam(),
8     metrics=[keras.metrics.SparseCategoricalAccuracy()],
9     distillation_loss_fn=keras.losses.KLDivergence(),
10    alpha=0.05,
11    temperature=5,
12 )
13
14 # Distill knowledge from teacher to student (will take around
15 distiller_history = distiller.fit(train_batches, epochs=5, val

Epoch 1/5
582/582 [=====] - 30s 44ms/step - sparse_categorical_ac
Epoch 2/5
582/582 [=====] - 29s 43ms/step - sparse_categorical_ac
Epoch 3/5
582/582 [=====] - 29s 44ms/step - sparse_categorical_ac
Epoch 4/5
582/582 [=====] - 29s 44ms/step - sparse_categorical_ac
Epoch 5/5
582/582 [=====] - 29s 44ms/step - sparse_categorical_ac
```

▼ Comparing the models

To compare the models you can check the `sparse_categorical_accuracy` of each one on the test set:

```
1 # Compute accuracies
2 student_scratch_acc = student_scratch.evaluate(test_batches, r
3 distiller_acc = distiller.evaluate(test_batches, return_dict=T
4 teacher_acc = teacher.evaluate(test_batches, return_dict=True)
5
```

```

6 # Print results
7 print(f"\n\nTeacher achieved a sparse_categorical_accuracy of
8 print(f"Student with knowledge distillation achieved a sparse_
9 print(f"Student without knowledge distillation achieved a spar

2326/2326 [=====] - 7s 3ms/step - loss: 0.7573 - sparse
2326/2326 [=====] - 7s 3ms/step - sparse_categorical_ac
2326/2326 [=====] - 8s 3ms/step - loss: 0.3197 - sparse

```

Teacher achieved a sparse_categorical_accuracy of 86.41%.

Student with knowledge distillation achieved a sparse_categorical_accuracy of 73

Student without knowledge distillation achieved a sparse_categorical_accuracy of

The teacher model yields a bigger accuracy than the two student models. This is expected since it was trained for more epochs while using a bigger architecture.

Notice that the student without distillation was outperformed by the student with knowledge distillation.

Since you saved the training history of each model you can create a plot for a better comparison of the two student models.

```

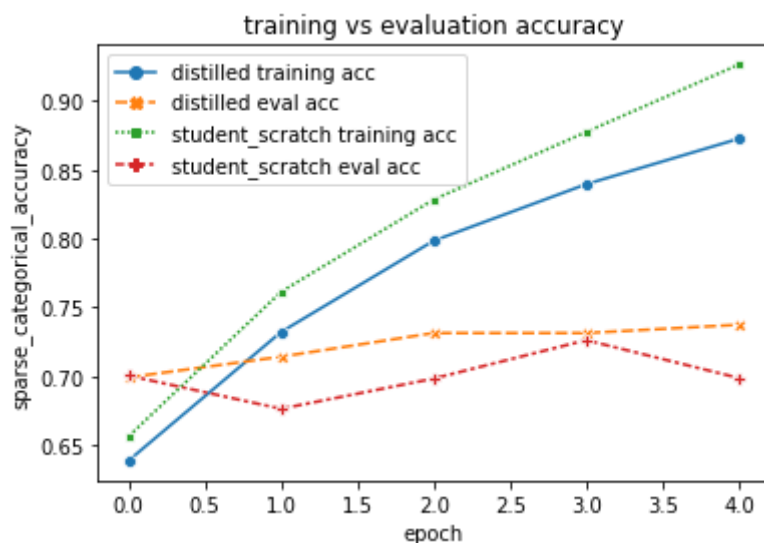
1 # Get relevant metrics from a history
2 def get_metrics(history):
3     history = history.history
4     acc = history['sparse_categorical_accuracy']
5     val_acc = history['val_sparse_categorical_accuracy']
6     return acc, val_acc
7
8
9 # Plot training and evaluation metrics given a dict of history
10 def plot_train_eval(history_dict):
11
12     metric_dict = {}
13
14     for k, v in history_dict.items():
15         acc, val_acc = get_metrics(v)
16         metric_dict[f'{k} training acc'] = acc
17         metric_dict[f'{k} eval acc'] = val_acc
18
19     acc_plot = pd.DataFrame(metric_dict)
20

```

```

20
21 acc_plot = sns.lineplot(data=acc_plot, markers=True)
22 acc_plot.set_title('training vs evaluation accuracy')
23 acc_plot.set_xlabel('epoch')
24 acc_plot.set_ylabel('sparse_categorical_accuracy')
25 plt.show()
26
27
28 # Plot for comparing the two student models
29 plot_train_eval({
30     "distilled": distiller_history,
31     "student_scratch": student_scratch_history,
32 })

```



This plot is very interesting because it shows that the distilled version outperformed the unmodified one in almost all of the epochs when using the evaluation set. Alongside this, the student without distillation yields a bigger training accuracy, which is a sign that it is overfitting more than the distilled model. **This hints that the distilled model was able to learn from the regularization that the teacher implemented!** Pretty cool, right?

Congratulations on finishing this ungraded lab! Now you should have a clearer understanding of what Knowledge Distillation is and how it can be implemented using Tensorflow and Keras.

This process is widely used for model compression and has proven to perform really well. In fact you might have heard about [DistilBert](#), which is a smaller, faster, cheaper and lighter of BERT.

Keep it up!

✓ 0s completed at 4:29 AM

