

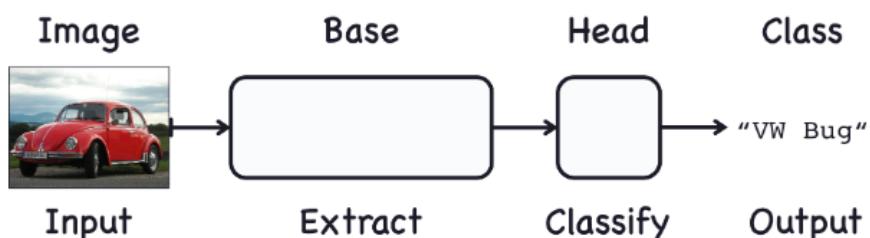


Module 1

Date	@28/08/2023
Title	Computer Vision

The Convolutional Classifier

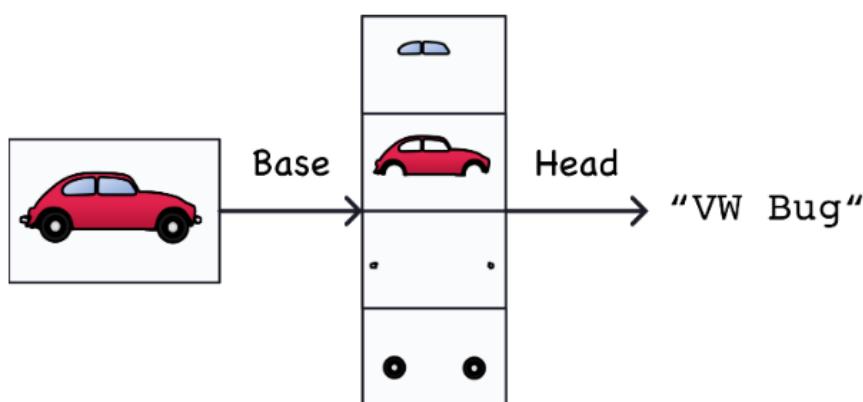
A convnet used for image classification consists of two parts: **a convolutional base** and a **dense head**.



The base is used to **extract the features** from an image. It is formed primarily of layers performing the convolution operation, but often includes other kinds of layers as well.

The head is used to **determine the class** of the image. It is formed primarily of dense layers, but might include other layers like dropout.

The whole process goes something like this:

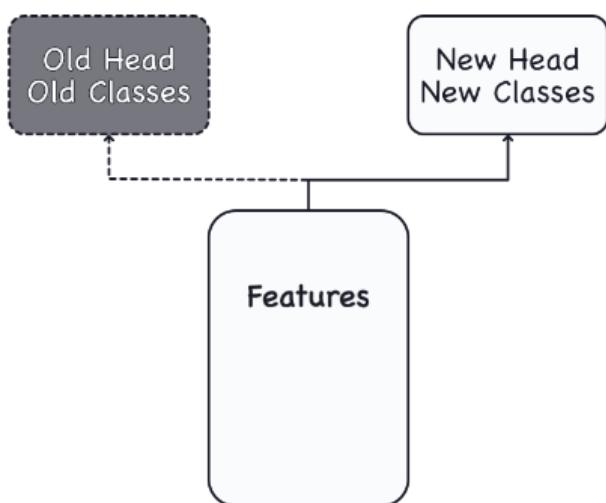


Training the Classifier

The goal of the network during training is to learn two things:

1. which features to extract from an image (base),
2. which class goes with what features (head).

These days, convnets are rarely trained from scratch. More often, we **reuse the base of a pretrained model**. To the pretrained base we then **attach an untrained head**. In other words, we reuse the part of a network that has already learned to do 1. *Extract features*, and attach to it some fresh layers to learn 2. *Classify*.



Because the head usually consists of only a few dense layers, **very accurate classifiers can be created from relatively little data**.

Reusing a pretrained model is a technique known as **transfer learning**. It is so effective, that almost every image classifier these days will make use of it.

Example - Train a Convnet Classifier

Throughout this course, we're going to be creating classifiers that attempt to solve the following problem: is this a picture of a *Car* or of a *Truck*? Our dataset is about 10,000 pictures of various automobiles, around half cars and half trucks.

Step 1 - Load Data

```
import os, warnings
import matplotlib.pyplot as plt
from matplotlib import gridspec

import numpy as np
import tensorflow as tf
```

```
from tensorflow.keras.preprocessing import image_dataset_from_directory

# Reproducability
def set_seed(seed=31415):
    np.random.seed(seed)
    tf.random.set_seed(seed)
    os.environ['PYTHONHASHSEED'] = str(seed)
    os.environ['TF_DETERMINISTIC_OPS'] = '1'
set_seed(31415)

# Set Matplotlib defaults
plt.rc('figure', autolayout=True)
plt.rc('axes', labelweight='bold', labelsize='large',
       titleweight='bold', titlesize=18, titlepad=10)
plt.rc('image', cmap='magma')
warnings.filterwarnings("ignore") # to clean up output cells

# Load training and validation sets
ds_train_ = image_dataset_from_directory(
    '../input/car-or-truck/train',
    labels='inferred',
    label_mode='binary',
    image_size=[128, 128],
    interpolation='nearest',
    batch_size=64,
    shuffle=True,
)
ds_valid_ = image_dataset_from_directory(
    '../input/car-or-truck/valid',
    labels='inferred',
    label_mode='binary',
    image_size=[128, 128],
    interpolation='nearest',
    batch_size=64,
    shuffle=False,
)
```

```

# Data Pipeline
def convert_to_float(image, label):
    image = tf.image.convert_image_dtype(image, dtype=tf.float32)
    return image, label

AUTOTUNE = tf.data.experimental.AUTOTUNE
ds_train = (
    ds_train_
    .map(convert_to_float)
    .cache()
    .prefetch(buffer_size=AUTOTUNE)
)
ds_valid = (
    ds_valid_
    .map(convert_to_float)
    .cache()
    .prefetch(buffer_size=AUTOTUNE)
)

```

Found 5117 files belonging to 2 classes.

Found 5051 files belonging to 2 classes.

Step 2 - Define Pretrained Base

The most commonly used dataset for pretraining is [ImageNet](#), a large dataset of many kind of natural images. Keras includes a variety models pretrained on ImageNet in its [applications module](#). The pretrained model we'll use is called **VGG16**.

```

pretrained_base = tf.keras.models.load_model(
    '../input/cv-course-models/cv-course-models/vgg16-pretrained'
pretrained_base.trainable = False

```

Step 3 - Attach Head

Next, we attach the classifier head. For this example, we'll use a layer of hidden units (the first `Dense` layer) followed by a layer to transform the outputs to a probability score for class 1, `Truck`. The `Flatten` layer transforms the two

dimensional outputs of the base into the one dimensional inputs needed by the head.

```
from tensorflow import keras
from tensorflow.keras import layers

model = keras.Sequential([
    pretrained_base,
    layers.Flatten(),
    layers.Dense(6, activation='relu'),
    layers.Dense(1, activation='sigmoid'),])
```

- `pretrained_base` : It seems like `pretrained_base` is expected to be a pre-trained neural network model. However, in the provided code, the definition of `pretrained_base` is missing. Typically, this would be a pre-trained model, such as a convolutional neural network (CNN) like VGG, ResNet, etc., that is often used as a feature extractor.
- `layers.Flatten()` : This layer is used to flatten the output of the `pretrained_base` layer. Many pre-trained models output multi-dimensional feature maps, and the `Flatten` layer converts these feature maps into a 1D array, which can be used as input for the subsequent dense layers.
- `layers.Dense(6, activation='relu')` : This is a fully connected (dense) layer with 6 neurons and ReLU activation. It takes the flattened output from the previous layer and applies a linear transformation followed by a Rectified Linear Unit (ReLU) activation function.
- `layers.Dense(1, activation='sigmoid')` : This is the final fully connected layer with 1 neuron and a sigmoid activation function. It produces a single output value between 0 and 1, which is often used for binary classification tasks.

Step 4 - Train

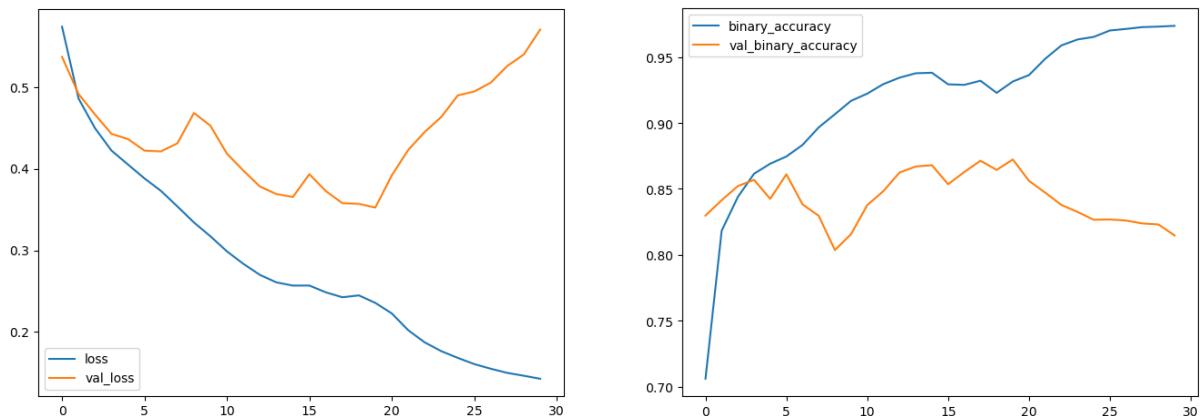
```
model.compile(
    optimizer='adam',
    loss='binary_crossentropy',
    metrics=['binary_accuracy'],
)

history = model.fit(
```

```
        ds_train,  
        validation_data=ds_valid,  
        epochs=30,  
        verbose=0,  
)
```

- `optimizer='adam'` : The Adam optimizer is being used for gradient-based optimization. Adam adapts the learning rates of each parameter individually over time to improve convergence speed.
- `loss='binary_crossentropy'` : This specifies the loss function that the model will use during training. Since the final layer of the model uses a sigmoid activation function and the task appears to be binary classification, 'binary_crossentropy' is a suitable loss function.
- `metrics=['binary_accuracy']` : During training, the accuracy of binary classification will be monitored and reported using the 'binary_accuracy' metric.
- `ds_train` : This is the training dataset that contains input data and corresponding target labels.
- `validation_data=ds_valid` : This parameter specifies the validation dataset to be used during training to monitor the model's performance on unseen data.
- `epochs=30` : The training process will run for 30 epochs, meaning the model will be exposed to the entire training dataset 30 times during training.
- `verbose=0` : This controls the verbosity of the training output. Setting it to 0 means no training progress will be printed to the console.

```
import pandas as pd  
  
history_frame = pd.DataFrame(history.history)  
history_frame.loc[:, ['loss', 'val_loss']].plot()  
history_frame.loc[:, ['binary_accuracy', 'val_binary_accuracy']]
```



Conclusion

The head, essentially, is an ordinary classifier like you learned about in the introductory course. For features, it uses those features extracted by the base. This is the basic idea behind convolutional classifiers: that we can attach a unit that performs feature engineering to the classifier itself.

This is one of the big advantages deep neural networks have over traditional machine learning models: given the right network structure, the deep neural net can learn how to engineer the features it needs to solve its problem.

Convolution and ReLU

```
import numpy as np
from itertools import product

def show_kernel(kernel, label=True, digits=None, text_size=28
    # Format kernel
    kernel = np.array(kernel)
    if digits is not None:
        kernel = kernel.round(digits)

    # Plot kernel
    cmap = plt.get_cmap('Blues_r')
    plt.imshow(kernel, cmap=cmap)
    rows, cols = kernel.shape
    thresh = (kernel.max() + kernel.min()) / 2
    # Optionally, add value labels
```

```

if label:
    for i, j in product(range(rows), range(cols)):
        val = kernel[i, j]
        color = cmap(0) if val > thresh else cmap(255)
        plt.text(j, i, val,
                  color=color, size=text_size,
                  horizontalalignment='center', verticalal
plt.xticks([])
plt.yticks([])

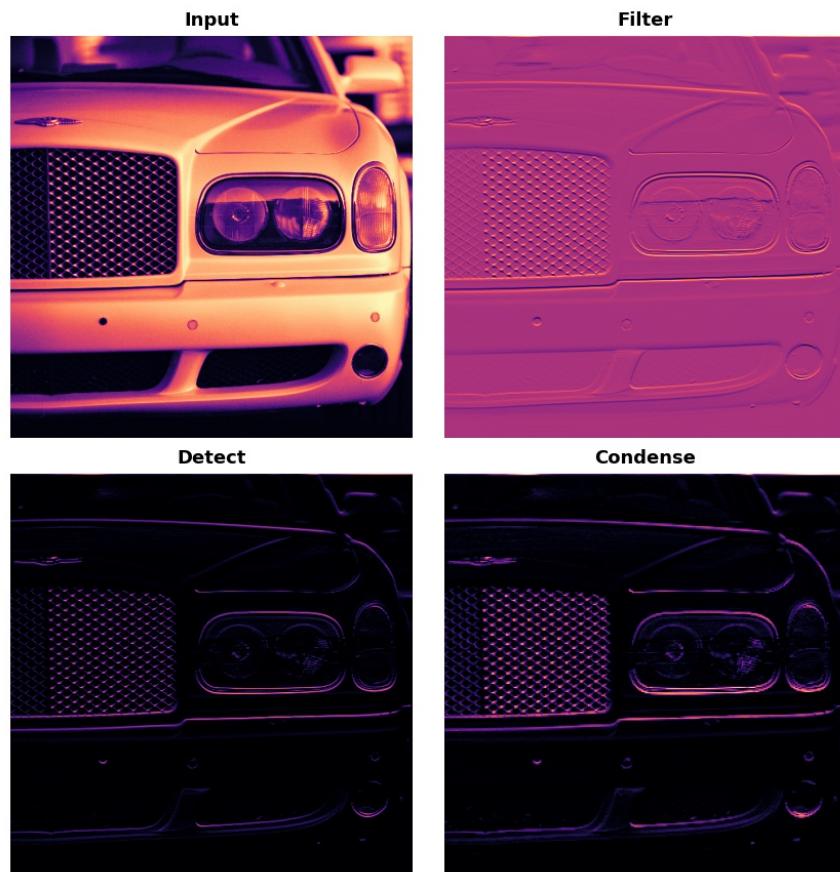
```

Feature Extraction

The **feature extraction** performed by the base consists of **three basic operations**:

1. **Filter** an image for a particular feature (convolution)
2. **Detect** that feature within the filtered image (ReLU)
3. **Condense** the image to enhance the features (maximum pooling)

The next figure illustrates this process. You can see how these three operations are able to isolate some particular characteristic of the original image (in this case, horizontal lines).



Filter with Convolution

A convolutional layer carries out the filtering step. You might define a convolutional layer in a Keras model something like this:

```
from tensorflow import keras
from tensorflow.keras import layers

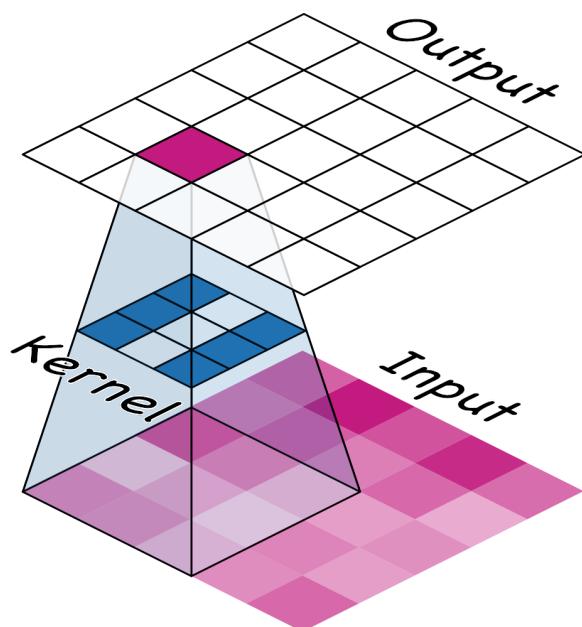
model = keras.Sequential([
    layers.Conv2D(filters=64, kernel_size=3), # activation is
])
```

Weights

The **weights** a convnet learns during training are primarily contained in its convolutional layers. These weights we call **kernels**. We can represent them as small arrays:

-1	2	-1
-1	2	-1
-1	2	-1

A kernel operates by scanning over an image and producing a *weighted sum* of pixel values. In this way, a kernel will act sort of like a polarized lens, emphasizing or deemphasizing certain patterns of information.

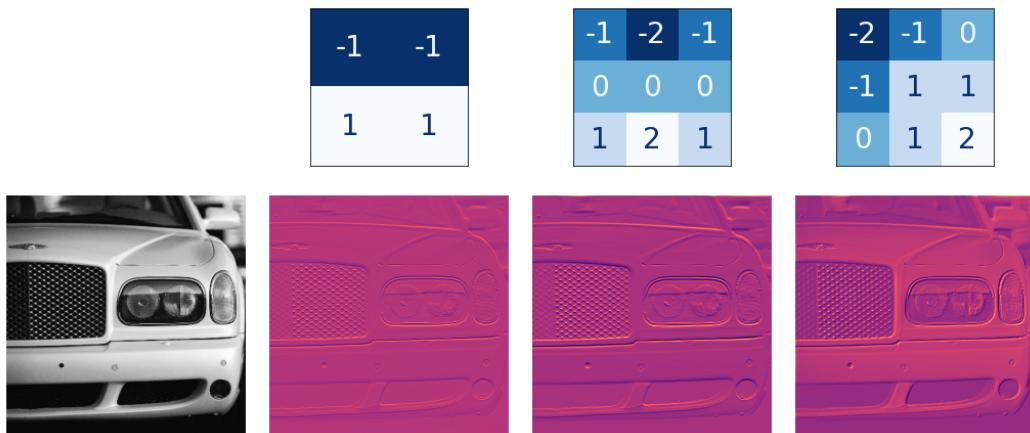


Kernels define how a convolutional layer is connected to the layer that follows. The kernel above will connect each neuron in the output to nine neurons in the input. By setting the dimensions of the kernels with `kernel_size`, you are telling the convnet how to form these connections. Most often, a kernel will have odd-numbered dimensions -- like `kernel_size=(3, 3)` or `(5, 5)` -- so that a single pixel sits at the center, but this is not a requirement.

The kernels in a convolutional layer determine what kinds of features it creates. During training, a convnet tries to learn what features it needs to solve the classification problem. This means finding the best values for its kernels.

Activations

The **activations** in the network we call **feature maps**. They are what result when we apply a filter to an image; they contain the visual features the kernel extracts. Here are a few kernels pictured with feature maps they produced.



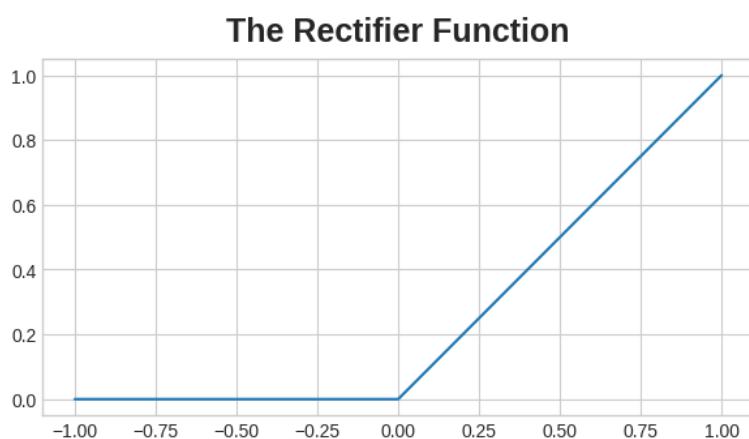
From the pattern of numbers in the kernel, you can tell the kinds of feature maps it creates. Generally, what a convolution accentuates in its inputs will match the shape of the *positive* numbers in the kernel. The left and middle kernels above will both filter for horizontal shapes.

With the `filters` parameter, you tell the convolutional layer how many feature maps you want it to create as output.

Detect with ReLU

After filtering, the feature maps pass through the activation function.

The **rectifier function** has a graph like this:



A neuron with a rectifier attached is called a *rectified linear unit*. For that reason, we might also call the rectifier function the **ReLU activation** or even the ReLU function.

The ReLU activation can be defined in its own `Activation` layer, but most often you'll just include it as the activation function of `Conv2D`.

```
model = keras.Sequential([
    layers.Conv2D(filters=64, kernel_size=3, activation='relu')
])
```

You could think about the activation function as scoring pixel values according to some measure of importance. The ReLU activation says that negative values are not important and so sets them to 0. (**"Everything unimportant is equally unimportant."**)

Here is ReLU applied the feature maps above. Notice how it succeeds at isolating the features.



Like other activation functions, the ReLU function is **nonlinear**. Essentially this means that the total effect of all the layers in the network becomes different than what you would get by just adding the effects together -- which would be the same as what you could achieve with only a single layer.

Example - Apply Convolution and ReLU

```
import tensorflow as tf
import matplotlib.pyplot as plt
plt.rc('figure', autolayout=True)
```

```
plt.rc('axes', labelweight='bold', labelsize='large',
       titleweight='bold', titlesize=18, titlepad=10)
plt.rc('image', cmap='magma')

image_path = '../input/computer-vision-resources/car_feature.'
image = tf.io.read_file(image_path)
image = tf.io.decode_jpeg(image)

plt.figure(figsize=(6, 6))
plt.imshow(tf.squeeze(image), cmap='gray')
plt.axis('off')
plt.show();
```



For the filtering step, we'll define a kernel and then apply it with the convolution. The kernel in this case is an "edge detection" kernel. You can define it with `tf.constant` just like you'd define an array in Numpy with `np.array`. This creates a *tensor* of the sort TensorFlow uses.

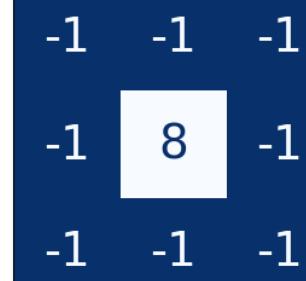
```

import tensorflow as tf

kernel =tf.constant([
    [-1, -1, -1],
    [-1, 8, -1],
    [-1, -1, -1],
])

plt.figure(figsize=(3, 3))
show_kernel(kernel)

```



TensorFlow includes many common operations performed by neural networks in its `tf.nn` module. The two that we'll use are `conv2d` and `relu`. These are simply function versions of Keras layers.

```

# Reformat for batch compatibility.
image = tf.image.convert_image_dtype(image, dtype=tf.float32)
image = tf.expand_dims(image, axis=0)
kernel = tf.reshape(kernel, [*kernel.shape, 1, 1])
kernel = tf.cast(kernel, dtype=tf.float32)

image_filter = tf.nn.conv2d(
    input=image,
    filters=kernel,
    strides=1,
    padding='SAME',
)

plt.figure(figsize=(6, 6))
plt.imshow(tf.squeeze(image_filter))
plt.axis('off')
plt.show();

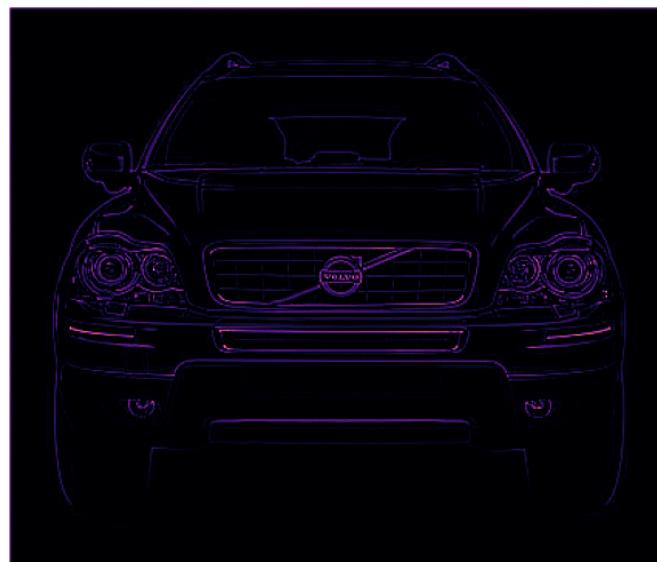
```



Next is the detection step with the ReLU function. This function is much simpler than the convolution, as it doesn't have any parameters to set.

```
image_detect = tf.nn.relu(image_filter)

plt.figure(figsize=(6, 6))
plt.imshow(tf.squeeze(image_detect))
plt.axis('off')
plt.show();
```



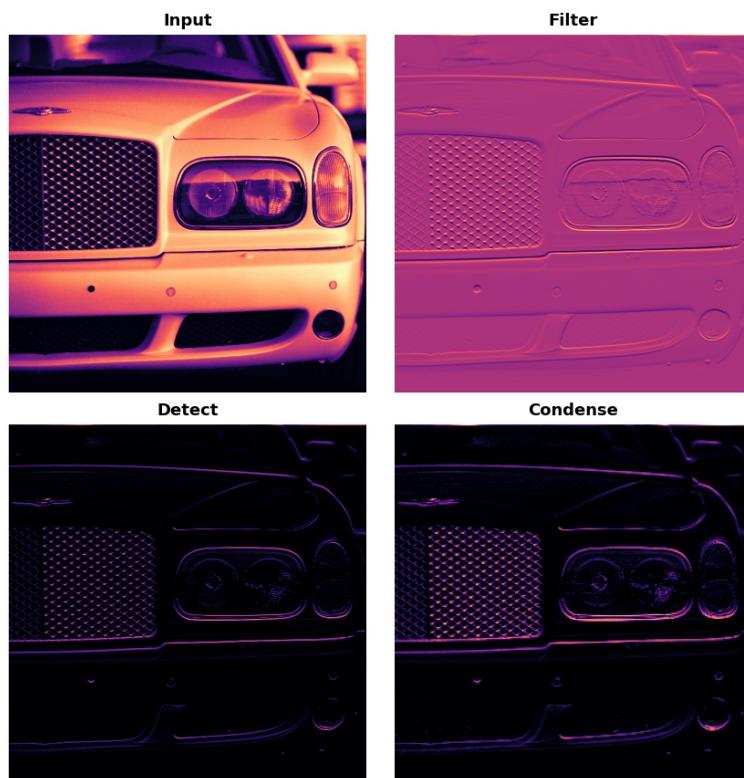
Maximum Pooling

Condense with Maximum Pooling

```
from tensorflow import keras
from tensorflow.keras import layers

model = keras.Sequential([
    layers.Conv2D(filters=64, kernel_size=3),
    layers.MaxPool2D(pool_size=2),
])
```

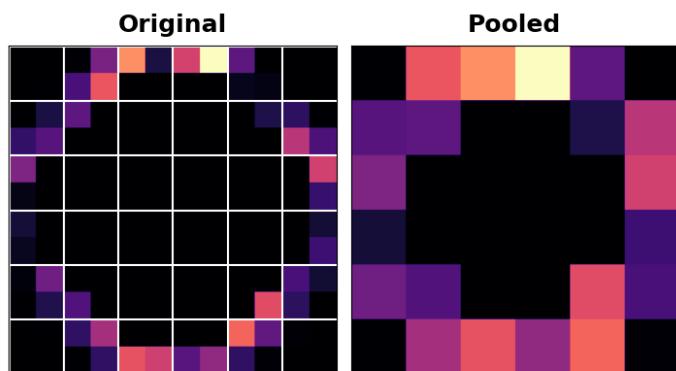
A `MaxPool2D` layer is much like a `Conv2D` layer, except that it uses a simple maximum function instead of a kernel, with the `pool_size` parameter analogous to `kernel_size`. A `MaxPool2D` layer doesn't have any trainable weights like a convolutional layer does in its kernel, however.



Notice that after applying the ReLU function (`Detect`) the feature map ends up with a lot of "dead space," that is, large areas containing only 0's (the black areas in the image). Having to carry these 0 activations through the entire network would increase the size of the model without adding much useful

information. Instead, we would like to *condense* the feature map to retain only the most useful part -- the feature itself.

This in fact is what **maximum pooling** does. Max pooling takes a patch of activations in the original feature map and replaces them with the maximum activation in that patch.



When applied after the ReLU activation, it has the effect of "intensifying" features. The pooling step increases the proportion of active pixels to zero pixels.

Example - Apply Maximum Pooling

```
import tensorflow as tf
import matplotlib.pyplot as plt
import warnings

plt.rc('figure', autolayout=True)
plt.rc('axes', labelweight='bold', labelsize='large',
      titleweight='bold', titlesize=18, titlepad=10)
plt.rc('image', cmap='magma')
warnings.filterwarnings("ignore") # to clean up output cells

# Read image
image_path = '../input/computer-vision-resources/car_feature.jpg'
image = tf.io.read_file(image_path)
image = tf.io.decode_jpeg(image)

# Define kernel
kernel = tf.constant([
    [-1, -1, -1],
    [-1,  8, -1],
    [-1, -1, -1],
```

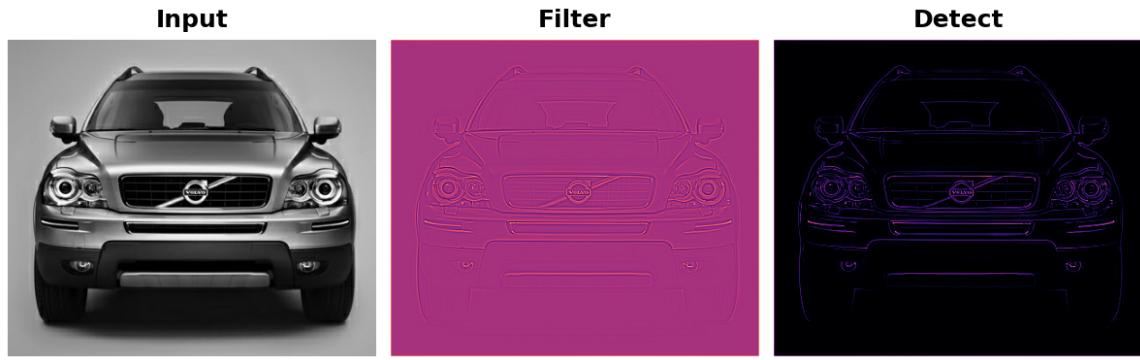
```
[], dtype=tf.float32)

# Reformat for batch compatibility.
image = tf.image.convert_image_dtype(image, dtype=tf.float32)
image = tf.expand_dims(image, axis=0)
kernel = tf.reshape(kernel, [*kernel.shape, 1, 1])

# Filter step
image_filter = tf.nn.conv2d(
    input=image,
    filters=kernel,
    # we'll talk about these two in the next lesson!
    strides=1,
    padding='SAME'
)

# Detect step
image_detect = tf.nn.relu(image_filter)

# Show what we have so far
plt.figure(figsize=(12, 6))
plt.subplot(131)
plt.imshow(tf.squeeze(image), cmap='gray')
plt.axis('off')
plt.title('Input')
plt.subplot(132)
plt.imshow(tf.squeeze(image_filter))
plt.axis('off')
plt.title('Filter')
plt.subplot(133)
plt.imshow(tf.squeeze(image_detect))
plt.axis('off')
plt.title('Detect')
plt.show();
```



We'll use another one of the functions in `tf.nn` to apply the pooling step, `tf.nn.pool`. This is a Python function that does the same thing as the `MaxPool2D` layer you use when model building, but, being a simple function, is easier to use directly.

```
import tensorflow as tf

image_condense = tf.nn.pool(
    input=image_detect, # image in the Detect step above
    window_shape=(2, 2),
    pooling_type='MAX',
    strides=(2, 2),
    padding='SAME',
)

plt.figure(figsize=(6, 6))
plt.imshow(tf.squeeze(image_condense))
plt.axis('off')
plt.show();
```

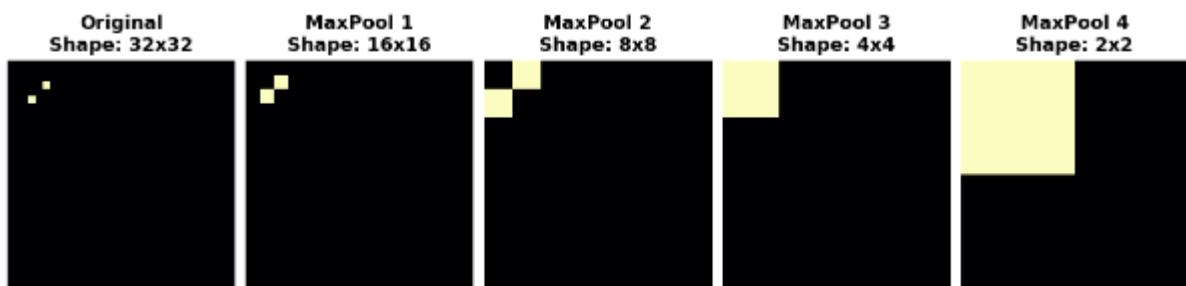


Translation Invariance

We called the zero-pixels "unimportant". Does this mean they carry no information at all? In fact, the zero-pixels carry *positional information*. The blank space still positions the feature within the image.

When `MaxPool2D` removes some of these pixels, it removes some of the positional information in the feature map. This gives a convnet a property called **translation invariance**. This means that a convnet with maximum pooling will tend not to distinguish features by their *location* in the image. ("Translation" is the mathematical word for changing the position of something without rotating it or changing its shape or size.)

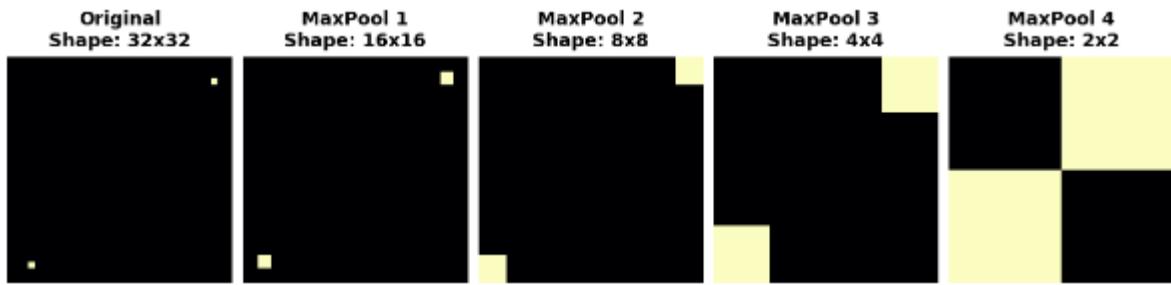
Watch what happens when we repeatedly apply maximum pooling to the following feature map.



The two dots in the original image became indistinguishable after repeated pooling. In other words, pooling destroyed some of their positional information. Since the network can no longer distinguish between them in the feature maps,

it can't distinguish them in the original image either: it has become *invariant* to that difference in position.

In fact, pooling only creates translation invariance in a network *over small distances*, as with the two dots in the image. Features that begin far apart will remain distinct after pooling; only *some* of the positional information was lost, but not all of it.



This invariance to small differences in the positions of features is a nice property for an image classifier to have. Just because of differences in perspective or framing, the same kind of feature might be positioned in various parts of the original image, but we would still like for the classifier to recognize that they are the same. Because this invariance is *built into* the network, we can get away with using much less data for training: we no longer have to teach it to ignore that difference. This gives convolutional networks a big efficiency advantage over a network with only dense layers.

The Sliding Window

```
import numpy as np
from itertools import product
from skimage import draw, transform

def circle(size, val=None, r_shrink=0):
    circle = np.zeros([size[0]+1, size[1]+1])
    rr, cc = draw.circle_perimeter(
        size[0]//2, size[1]//2,
        radius=size[0]//2 - r_shrink,
        shape=[size[0]+1, size[1]+1],
    )
    if val is None:
        circle[rr, cc] = np.random.uniform(size=circle.shape)
```

```

        else:
            circle[rr, cc] = val
    circle = transform.resize(circle, size, order=0)
    return circle

def show_kernel(kernel, label=True, digits=None, text_size=28
    # Format kernel
    kernel = np.array(kernel)
    if digits is not None:
        kernel = kernel.round(digits)

    # Plot kernel
    cmap = plt.get_cmap('Blues_r')
    plt.imshow(kernel, cmap=cmap)
    rows, cols = kernel.shape
    thresh = (kernel.max() + kernel.min()) / 2
    # Optionally, add value labels
    if label:
        for i, j in product(range(rows), range(cols)):
            val = kernel[i, j]
            color = cmap(0) if val > thresh else cmap(255)
            plt.text(j, i, val,
                      color=color, size=text_size,
                      horizontalalignment='center', verticalal
    plt.xticks([])
    plt.yticks([])

def show_extraction(image,
                    kernel,
                    conv_stride=1,
                    conv_padding='valid',
                    activation='relu',
                    pool_size=2,
                    pool_stride=2,
                    pool_padding='same',
                    figsize=(10, 10),
                    subplot_shape=(2, 2),
                    ops=['Input', 'Filter', 'Detect', 'Condens

```

```

        gamma=1.0):
# Create Layers
model = tf.keras.Sequential([
    tf.keras.layers.Conv2D(
        filters=1,
        kernel_size=kernel.shape,
        strides=conv_stride,
        padding=conv_padding,
        use_bias=False,
        input_shape=image.shape,
    ),
    tf.keras.layers.Activation(activation),
    tf.keras.layers.MaxPool2D(
        pool_size=pool_size,
        strides=pool_stride,
        padding=pool_padding,
    ),
])
layer_filter, layer_detect, layer_condense = model.layers
kernel = tf.reshape(kernel, [*kernel.shape, 1, 1])
layer_filter.set_weights([kernel])

# Format for TF
image = tf.expand_dims(image, axis=0)
image = tf.image.convert_image_dtype(image, dtype=tf.floa

# Extract Feature
image_filter = layer_filter(image)
image_detect = layer_detect(image_filter)
image_condense = layer_condense(image_detect)

images = {}
if 'Input' in ops:
    images.update({'Input': (image, 1.0)})
if 'Filter' in ops:
    images.update({'Filter': (image_filter, 1.0)})
if 'Detect' in ops:

```

```

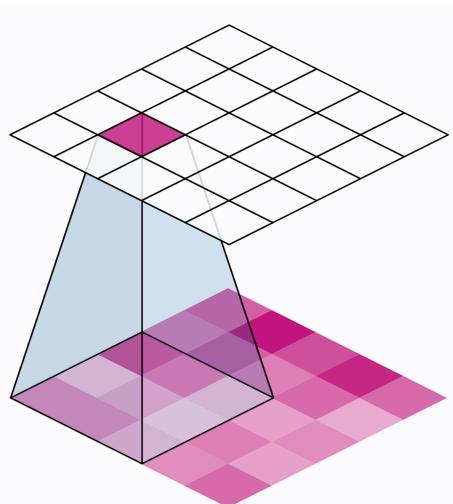
        images.update({'Detect': (image_detect, gamma)})
if 'Condense' in ops:
    images.update({'Condense': (image_condense, gamma)})

# Plot
plt.figure(figsize=figsize)
for i, title in enumerate(ops):
    image, gamma = images[title]
    plt.subplot(*subplot_shape, i+1)
    plt.imshow(tf.image.adjust_gamma(tf.squeeze(image), gamma))
    plt.axis('off')
    plt.title(title)

```

1. *filter* with a **convolution** layer
2. *detect* with **ReLU** activation
3. *condense* with a **maximum pooling** layer

The convolution and pooling operations share a common feature: they are both performed over a **sliding window**. With convolution, this "window" is given by the dimensions of the kernel, the parameter `kernel_size`. With pooling, it is the pooling window, given by `pool_size`.



There are two additional parameters affecting both **convolution** and **pooling** layers -- these are the `strides` of the window and whether to use `padding` at the image edges. The `strides` parameter says how far the window should move at

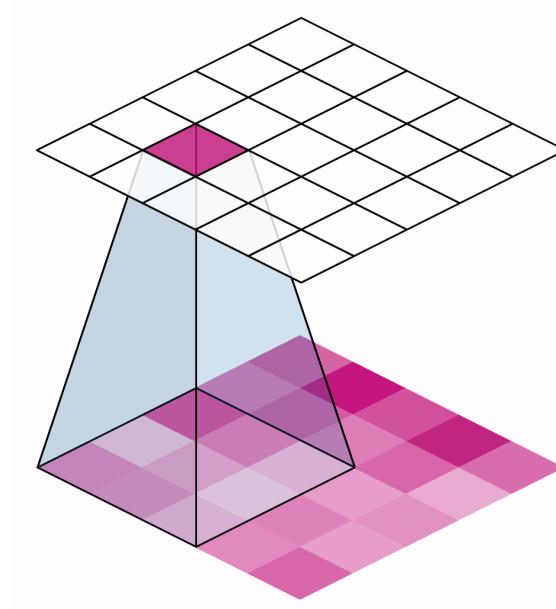
each step, and the `padding` parameter describes how we handle the pixels at the edges of the input.

```
from tensorflow import keras
from tensorflow.keras import layers

model = keras.Sequential([
    layers.Conv2D(filters=64,
                  kernel_size=3,
                  strides=1,
                  padding='same',
                  activation='relu'),
    layers.MaxPool2D(pool_size=2,
                     strides=1,
                     padding='same')
])
```

Stride

The distance the window moves at each step is called the **stride**. We need to specify the stride in both dimensions of the image: one for moving left to right and one for moving top to bottom. This animation shows `strides=(2, 2)`, a movement of 2 pixels each step.



What effect does the stride have? Whenever the stride in either direction is greater than 1, the sliding window will skip over some of the pixels in the input at each step.

Because we want high-quality features to use for classification, convolutional layers will most often have `strides=(1, 1)`. Increasing the stride means that we miss out on potentially valuable information in our summary. Maximum pooling layers, however, will almost always have stride values greater than 1, like `(2, 2)` or `(3, 3)`, but not larger than the window itself.

Finally, note that when the value of the `strides` is the same number in both directions, you only need to set that number; for instance, instead of `strides=(2, 2)`, you could use `strides=2` for the parameter setting.

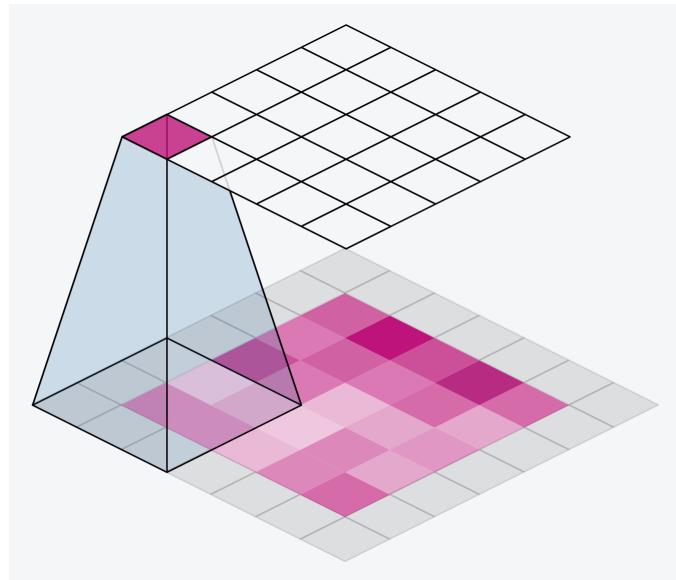
Padding

When performing the sliding window computation, there is a question as to what to do at the boundaries of the input. Staying entirely inside the input image means the window will never sit squarely over these boundary pixels like it does for every other pixel in the input. Since we aren't treating all the pixels exactly the same, could there be a problem?

What the convolution does with these boundary values is determined by its `padding` parameter. In TensorFlow, you have two choices: either `padding='same'` or `padding='valid'`.

When we set `padding='valid'`, the convolution window will stay entirely inside the input. **The drawback is that the output shrinks (loses pixels)**, and shrinks more for larger kernels. This will limit the number of layers the network can contain, especially when inputs are small in size.

The alternative is to use `padding='same'`. The trick here is to **pad** the input with 0's around its borders, using just enough 0's to make the size of the output the *same* as the size of the input. This can have the effect however of diluting the influence of pixels at the borders. The animation below shows a sliding window with `'same'` padding.

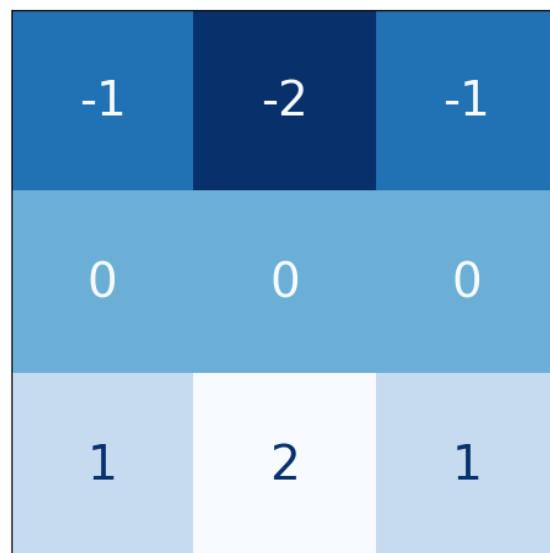


Example - Exploring Sliding Windows

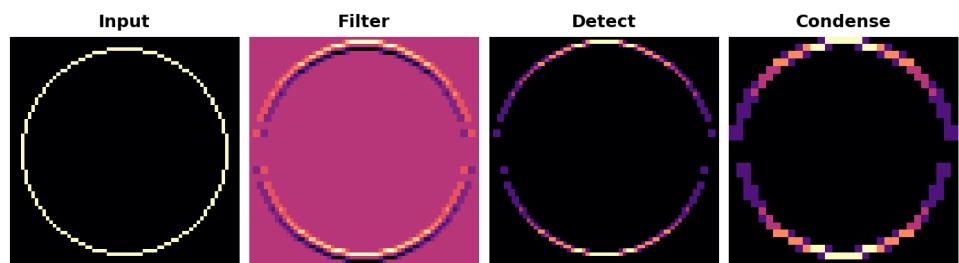
```
import tensorflow as tf
import matplotlib.pyplot as plt

plt.rc('figure', autolayout=True)
plt.rc('axes', labelweight='bold', labelsize='large',
      titleweight='bold', titlesize=18, titlepad=10)
plt.rc('image', cmap='magma')

image = circle([64, 64], val=1.0, r_shrink=3)
image = tf.reshape(image, [*image.shape, 1])
# Bottom sobel
kernel = tf.constant(
    [[-1, -2, -1],
     [0, 0, 0],
     [1, 2, 1]],
)
show_kernel(kernel)
```

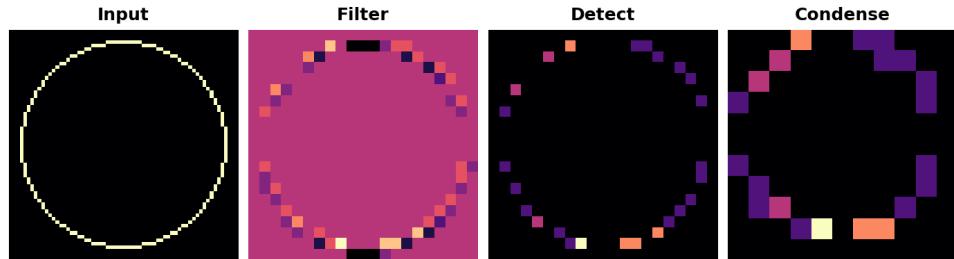


```
show_extraction(  
    image, kernel,  
    conv_stride=1,  
    pool_size=2,  
    pool_stride=2,  
  
    subplot_shape=(1, 4),  
    figsize=(14, 6),  
)
```



```
show_extraction(  
    image, kernel,  
    conv_stride=3,  
    pool_size=2,  
    pool_stride=2,  
  
    subplot_shape=(1, 4),
```

```
    figsize=(14, 6),  
)
```



This seems to reduce the quality of the feature extracted. Our input circle is rather "finely detailed," being only 1 pixel wide. A convolution with strides of 3 is too coarse to produce a good feature map from it.

Sometimes, a model will use a convolution with a larger stride in its initial layer. This will usually be coupled with a larger kernel as well. The ResNet50 model, for instance, uses 7×7 kernels with strides of 2 in its first layer. This seems to accelerate the production of large-scale features without the sacrifice of too much information from the input.

```
# Setup feedback system  
from learntools.core import binder  
binder.bind(globals())  
from learntools.computer_vision.ex4 import *  
  
import tensorflow as tf  
import matplotlib.pyplot as plt  
import learntools.computer_vision.visiontools as visiontools  
  
plt.rc('figure', autolayout=True)  
plt.rc('axes', labelweight='bold', labelsize='large',  
      titleweight='bold', titlesize=18, titlepad=10)  
plt.rc('image', cmap='magma')
```

```
from learntools.computer_vision.visiontools import edge, blur
```

```

image_dir = '../input/computer-vision-resources/'
circle_64 = tf.expand_dims(circle([64, 64], val=1.0, r_shrink=4), axis=-1)
kaggle_k = visiontools.read_image(image_dir + str('k.jpg'), color=False)
car = visiontools.read_image(image_dir + str('car_illus.jpg'))
car = tf.image.resize(car, size=[200, 200])
images = [(circle_64, "circle_64"), (kaggle_k, "kaggle_k"), (car, "car")]

plt.figure(figsize=(14, 4))
for i, (img, title) in enumerate(images):
    plt.subplot(1, len(images), i+1)
    plt.imshow(tf.squeeze(img))
    plt.axis('off')
    plt.title(title)
plt.show();

kernels = [(edge, "edge"), (blur, "blur"), (bottom_sobel, "bottom_sobel"),
            (emboss, "emboss"), (sharpen, "sharpen")]
plt.figure(figsize=(14, 4))
for i, (krn, title) in enumerate(kernels):
    plt.subplot(1, len(kernels), i+1)
    visiontools.show_kernel(krn, digits=2, text_size=20)
    plt.title(title)
plt.show()

```

imports various functions (`edge`, `blur`, `bottom_sobel`, `emboss`, `sharpen`, `circle`) from the module `visiontools` provided by `learntools.computer_vision`. These functions are used for performing different image processing operations.

`circle_64` generates a circular mask of size 64×64 pixels using the `circle` function from the imported module. The circular mask is created with a value of 1.0 and a radius shrink factor of 4. The `expand_dims` function is used to add an extra dimension to the mask to indicate that it has one channel (grayscale). This is done to match the channel dimensions of other images.

`kaggle_k` reads an image named 'k.jpg' from the specified image directory using the `read_image` function. The image is read as grayscale (1 channel) and stored in the variable `kaggle_k`.

read an image named '`car_illus.jpg`' from the image directory and store it in the variable `car`. The image is initially read as grayscale (1 channel). Subsequently,

the `tf.image.resize` function is used to resize the image to a size of 200×200 pixels.

creates a list of tuples named `images`. Each tuple contains an image along with its associated title. The images include the circular mask (`circle_64`), the '`kaggle_k`' image, and the resized '`car`' image.



edge	blur	bottom_sobel	emboss	sharpen
-1 -1 -1 -1 8 -1 -1 -1 -1	0.06 0.12 0.06 0.12 0.25 0.12 0.06 0.12 0.06	-1 -2 -1 0 0 0 1 2 1	-2 -1 0 -1 1 1 0 1 2	0 -1 0 -1 5 -1 0 -1 0

To choose one to experiment with, just enter it's name in the appropriate place below. Then, set the parameters for the window computation. Try out some different combinations and see what they do!

```
image = circle_64
kernel = bottom_sobel

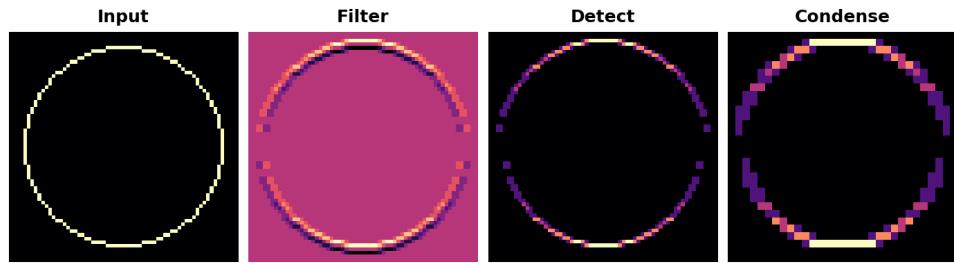
visiontools.show_extraction(
    image, kernel,
    conv_stride=1,
    conv_padding='valid',
    pool_size=2,
    pool_stride=2,
    pool_padding='same',

    subplot_shape=(1, 4),
```

```

    figsize=(14, 6),
)

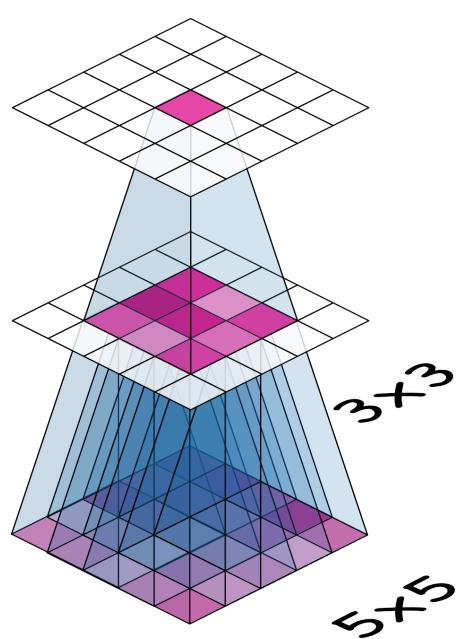
```



Trace back all the connections from some neuron and eventually you reach the input image. All of the input pixels a neuron is connected to is that neuron's **receptive field**. The receptive field just tells you which parts of the input image a neuron receives information from.

As we've seen, if your first layer is a convolution with 3×3 kernels, then each neuron in that layer gets input from a 3×3 patch of pixels (except maybe at the border).

What happens if you add another convolutional layer with $3 \times 3 \times 3$ kernels? Consider this next illustration:



Now trace back the connections from the neuron at top and you can see that it's connected to a 5×5 patch of pixels in the input (the bottom layer): each neuron in the 3×3 patch in the middle layer is connected to a 3×3 input patch, but they overlap in a 5×5 patch. So that neuron at top has a 5×5 receptive field.

Custom Convnets

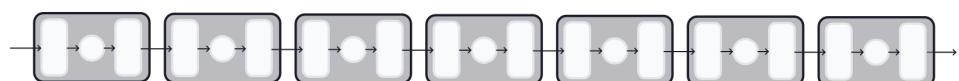
Simple to Refined

A single round of feature extraction can only extract relatively simple features from an image, things like simple lines or contrasts. These are too simple to solve most classification problems. Instead, convnets will repeat this extraction over and over, so that the features become more complex and refined as they travel deeper into the network.

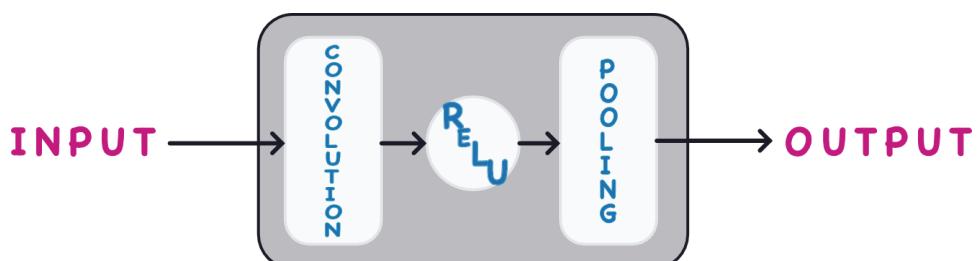


Convolutional Blocks

It does this by passing them through long chains of **convolutional blocks** which perform this extraction.



These convolutional blocks are stacks of `Conv2D` and `MaxPool2D` layers, whose role in feature extraction we learned about in the last few lessons.



Each block represents a round of extraction, and by composing these blocks the convnet can combine and recombine the features produced, growing them

and shaping them to better fit the problem at hand. The deep structure of modern convnets is what allows this sophisticated feature engineering and has been largely responsible for their superior performance.

Example - Design a Convnet

Step 1 - Load Data

```
import os, warnings
import matplotlib.pyplot as plt
from matplotlib import gridspec

import numpy as np
import tensorflow as tf
from tensorflow.keras.preprocessing import image_dataset_from_directory

# Reproducability
def set_seed(seed=31415):
    np.random.seed(seed)
    tf.random.set_seed(seed)
    os.environ['PYTHONHASHSEED'] = str(seed)
    os.environ['TF_DETERMINISTIC_OPS'] = '1'
set_seed()

# Set Matplotlib defaults
plt.rc('figure', autolayout=True)
plt.rc('axes', labelweight='bold', labelsize='large',
       titleweight='bold', titlesize=18, titlepad=10)
plt.rc('image', cmap='magma')
warnings.filterwarnings("ignore") # to clean up output cells

# Load training and validation sets
ds_train_ = image_dataset_from_directory(
    '../input/car-or-truck/train',
    labels='inferred',
    label_mode='binary',
    image_size=[128, 128],
    interpolation='nearest',
    batch_size=64,
```

```

        shuffle=True,
    )
ds_valid_ = image_dataset_from_directory(
    '../input/car-or-truck/valid',
    labels='inferred',
    label_mode='binary',
    image_size=[128, 128],
    interpolation='nearest',
    batch_size=64,
    shuffle=False,
)

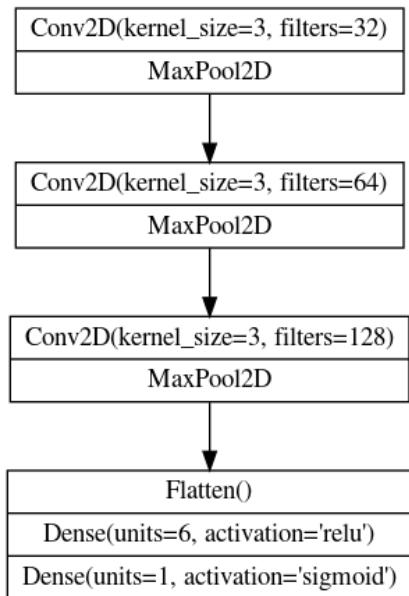
# Data Pipeline
def convert_to_float(image, label):
    image = tf.image.convert_image_dtype(image, dtype=tf.float32)
    return image, label

AUTOTUNE = tf.data.experimental.AUTOTUNE
ds_train = (
    ds_train_
    .map(convert_to_float)
    .cache()
    .prefetch(buffer_size=AUTOTUNE)
)
ds_valid = (
    ds_valid_
    .map(convert_to_float)
    .cache()
    .prefetch(buffer_size=AUTOTUNE)
)

```

Found 5117 files belonging to 2 classes.
 Found 5051 files belonging to 2 classes.

Step 2 - Define Model



Now we'll define the model. See how our model consists of three blocks of `Conv2D` and `MaxPool2D` layers (the base) followed by a head of `Dense` layers. We can translate this diagram more or less directly into a Keras `Sequential` model just by filling in the appropriate parameters

```

from tensorflow import keras
from tensorflow.keras import layers

model = keras.Sequential([
    # First Convolutional Block
    layers.Conv2D(filters=32, kernel_size=5, activation="relu",
                  input_shape=[128, 128, 3]),
    layers.MaxPool2D(),

    # Second Convolutional Block
    layers.Conv2D(filters=64, kernel_size=3, activation="relu"),
    layers.MaxPool2D(),

    # Third Convolutional Block
    layers.Conv2D(filters=128, kernel_size=3, activation="relu"),
    layers.MaxPool2D(),

    # Classifier Head
    layers.Flatten(),
    layers.Dense(units=6, activation="relu"),
    layers.Dense(units=1, activation="sigmoid"),
])

```

```
])  
model.summary()
```

```
Model: "sequential"  
-----  
Layer (type)          Output Shape         Param #  
=====  
conv2d (Conv2D)       (None, 128, 128, 32)  2432  
  
max_pooling2d (MaxPooling2D) (None, 64, 64, 32)  0  
)  
  
conv2d_1 (Conv2D)      (None, 64, 64, 64)    18496  
  
max_pooling2d_1 (MaxPooling 2D) (None, 32, 32, 64)  0  
  
conv2d_2 (Conv2D)      (None, 32, 32, 128)   73856  
  
max_pooling2d_2 (MaxPooling 2D) (None, 16, 16, 128)  0  
  
flatten (Flatten)      (None, 32768)        0  
  
dense (Dense)          (None, 6)            196614  
  
dense_1 (Dense)        (None, 1)            7  
  
=====  
Total params: 291,405  
Trainable params: 291,405  
Non-trainable params: 0  
-----
```

Step 3 - Train

```
model.compile(  
    optimizer=tf.keras.optimizers.Adam(epsilon=0.01),  
    loss='binary_crossentropy',  
    metrics=['binary_accuracy'])  
  
history = model.fit(  
    ds_train,  
    validation_data=ds_valid,  
    epochs=40,  
    verbose=0,  
)
```

Notice in this definition is how the number of filters doubled block-by-block: 32, 64, 128. This is a common pattern. Since the `MaxPool2D` layer is reducing the size of the feature maps, we can afford to increase the quantity we create.

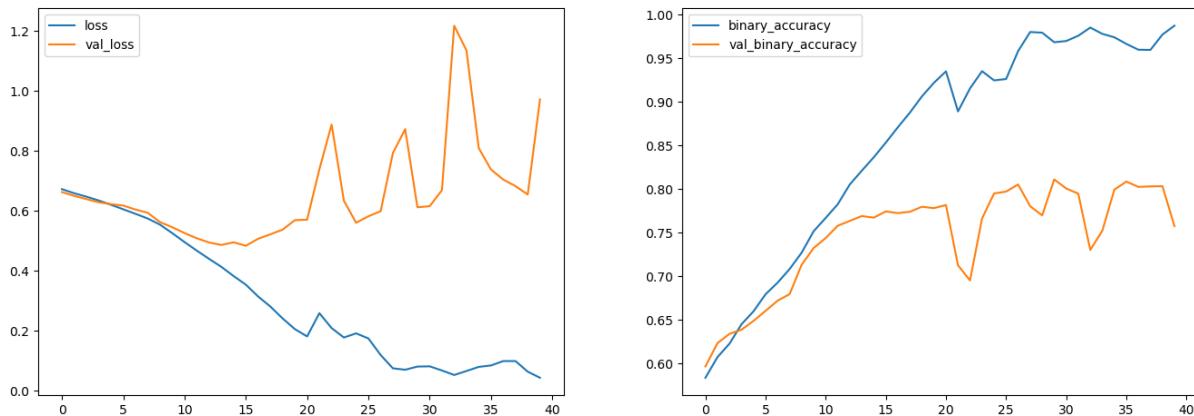
```
import pandas as pd

history_frame = pd.DataFrame(history.history)
history_frame.loc[:, ['loss', 'val_loss']].plot()
history_frame.loc[:, ['binary_accuracy', 'val_binary_accuracy']]
```

- `model.compile`: This method compiles the model, which means it configures the model for training. During compilation, you specify the optimizer, loss function, and metrics that will be used during the training process.
- `optimizer=tf.keras.optimizers.Adam(epsilon=0.01)`: The optimizer used for gradient descent is Adam, a popular optimization algorithm. The `epsilon` parameter here is an optional hyperparameter that helps prevent division by zero. It's set to 0.01 in this case.
- `loss='binary_crossentropy'`: This specifies the loss function to be used during training. `'binary_crossentropy'` is often used for binary classification problems. It measures the difference between the true labels and the predicted probabilities for each class.
- `metrics=['binary_accuracy']`: This defines the evaluation metric used during training and validation. In this case, it's `'binary_accuracy'`, which calculates the accuracy of binary classification predictions.
- `ds_train`: This is the training dataset (presumably a TensorFlow dataset) that contains the training samples and labels.
- `validation_data=ds_valid`: This is the validation dataset used to monitor the model's performance during training. It contains validation samples and labels. Monitoring the validation performance helps prevent overfitting.
- `epochs=40`: The number of epochs specifies how many times the entire training dataset will be used to update the model's weights. In this case, the model will be trained for 40 epochs.
- `verbose=0`: This parameter controls the verbosity during training. A value of 0 means silent mode, where no output will be printed during training.
- `pd.DataFrame(history.history)`: The `history` object contains various training metrics recorded after each epoch. This line creates a Pandas DataFrame from the history object, making it easier to work with the data.
- `history_frame.loc[:, ['loss', 'val_loss']].plot()`: This line plots the training loss and validation loss over epochs. It uses Pandas' `.loc` indexer to select the

'loss' and 'val_loss' columns from the DataFrame and then calls `.plot()` to create the plot.

- `history_frame.loc[:, ['binary_accuracy', 'val_binary_accuracy']].plot()`: Similarly, this line plots the binary accuracy and validation binary accuracy over epochs.



Data Augmentation

The Usefulness of Fake Data

The best way to improve the performance of a machine learning model is to train it on more data. The more examples the model has to learn from, the better it will be able to recognize which differences in images matter and which do not. More data helps the model to **generalize** better.

One easy way of getting more data is to use the data you already have. If we can transform the images in our dataset in ways that preserve the class, we can teach our classifier to **ignore those kinds of transformations**. For instance, whether a car is facing left or right in a photo doesn't change the fact that it is a *Car* and not a *Truck*. So, if we **augment** our training data with flipped images, our classifier will learn that "left or right" is a difference it should ignore.

And that's the whole idea behind data augmentation: add in some extra fake data that looks reasonably like the real data and your classifier will improve.

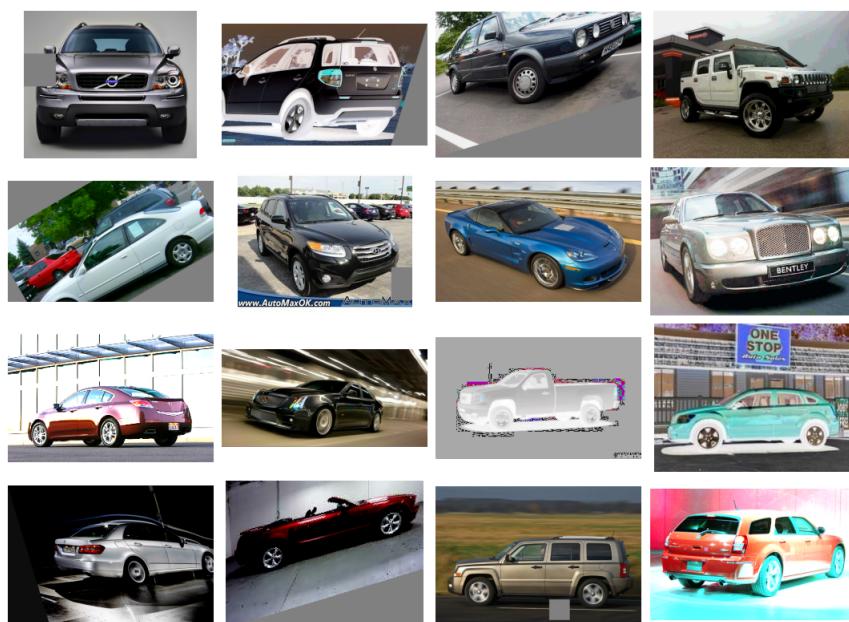
Using Data Augmentation

Typically, many kinds of transformation are used when augmenting a dataset. These might include rotating the image, adjusting the color or contrast, warping

the image, or many other things, usually applied in combination. Here is a sample of the different ways a single image might be transformed.



Data augmentation is usually done **online**, meaning, as the images are being fed into the network for training. Recall that training is usually done on mini-batches of data. This is what a batch of 16 images might look like when data augmentation is used.



Each time an image is used during training, a new random transformation is applied. This way, the model is always seeing something a little different than what it's seen before. This extra variance in the training data is what helps the model on new data.

It's important to remember though that not every transformation will be useful on a given problem. Most importantly, whatever transformations you use should not mix up the classes. If you were training a [digit recognizer](#), for instance, rotating images would mix up '9's and '6's. In the end, the best approach for finding good augmentations is the same as with most ML problems: try it and see!

Example - Training with Data Augmentation

Keras lets you augment your data in two ways. The first way is to include it in the data pipeline with a function like [ImageDataGenerator](#). The second way is to include it in the model definition by using Keras's [preprocessing layers](#). This is the approach that we'll take.

```
import os, warnings
import matplotlib.pyplot as plt
from matplotlib import gridspec

import numpy as np
import tensorflow as tf
from tensorflow.keras.preprocessing import image_dataset_from_directory

# Reproducability
def set_seed(seed=31415):
    np.random.seed(seed)
    tf.random.set_seed(seed)
    os.environ['PYTHONHASHSEED'] = str(seed)
    #os.environ['TF_DETERMINISTIC_OPS'] = '1'
set_seed()

# Set Matplotlib defaults
plt.rc('figure', autolayout=True)
plt.rc('axes', labelweight='bold', labelsize='large',
       titleweight='bold', titlesize=18, titlepad=10)
plt.rc('image', cmap='magma')
```

```

warnings.filterwarnings("ignore") # to clean up output cells

# Load training and validation sets
ds_train_ = image_dataset_from_directory(
    '../input/car-or-truck/train',
    labels='inferred',
    label_mode='binary',
    image_size=[128, 128],
    interpolation='nearest',
    batch_size=64,
    shuffle=True,
)
ds_valid_ = image_dataset_from_directory(
    '../input/car-or-truck/valid',
    labels='inferred',
    label_mode='binary',
    image_size=[128, 128],
    interpolation='nearest',
    batch_size=64,
    shuffle=False,
)

# Data Pipeline
def convert_to_float(image, label):
    image = tf.image.convert_image_dtype(image, dtype=tf.float32)
    return image, label

AUTOTUNE = tf.data.experimental.AUTOTUNE
ds_train = (
    ds_train_
    .map(convert_to_float)
    .cache()
    .prefetch(buffer_size=AUTOTUNE)
)
ds_valid = (
    ds_valid_
    .map(convert_to_float)
)

```

```
.cache()
.prefetch(buffer_size=AUTOTUNE)
)
```

Found 5117 files belonging to 2 classes.

Found 5051 files belonging to 2 classes.

Step 2 - Define Model

```
from tensorflow import keras
from tensorflow.keras import layers
from tensorflow.keras.layers.experimental import preprocessing

pretrained_base = tf.keras.models.load_model(
    '../input/cv-course-models/cv-course-models/vgg16-pretrained'
)
pretrained_base.trainable = False

model = keras.Sequential([
    # Preprocessing
    preprocessing.RandomFlip('horizontal'), # flip left-to-right
    preprocessing.RandomContrast(0.5), # contrast change by up to 50%
    # Base
    pretrained_base,
    # Head
    layers.Flatten(),
    layers.Dense(6, activation='relu'),
    layers.Dense(1, activation='sigmoid'),
])

```

Step 3 - Train and Evaluate

```
model.compile(
    optimizer='adam',
    loss='binary_crossentropy',
    metrics=['binary_accuracy'],
)
```

```
history = model.fit(  
    ds_train,  
    validation_data=ds_valid,  
    epochs=30,  
    verbose=0,  
)  
  
import pandas as pd  
  
history_frame = pd.DataFrame(history.history)  
  
history_frame.loc[:, ['loss', 'val_loss']].plot()  
history_frame.loc[:, ['binary_accuracy', 'val_binary_accuracy']]
```

