

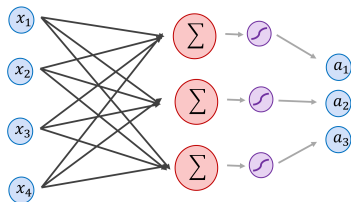
11-695: AI Engineering

Customizing NNs III

LTI/SCS

Spring 2020

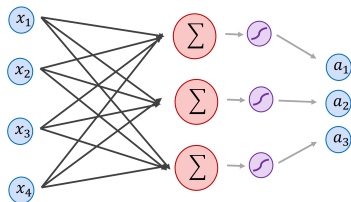
- 1 Choosing Activations
- 2 Callbacks
- 3 Save and Restore
- 4 Customization with Transfer Learning



- Recall: every NN has this form of operations

$$\hat{\mathbf{y}} = (\phi_n \circ \mathbf{f}_{W_n} \circ \phi_{n-1} \circ \mathbf{f}_{W_{n-1}} \dots \phi_1 \circ \mathbf{f}_{W_1})(\mathbf{X})$$

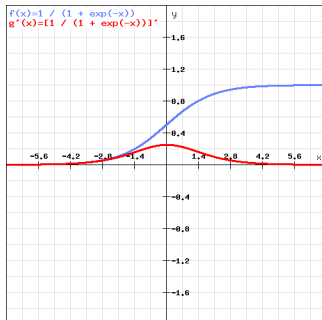
- Two basic operations
 - Linear: $o_i = \mathbf{f}_{W_i}(a_{i-1}) = W_i^T a_{i-1} + b_i$
 - Nonlinear (by activation functions): $a_i = \phi(o_i)$



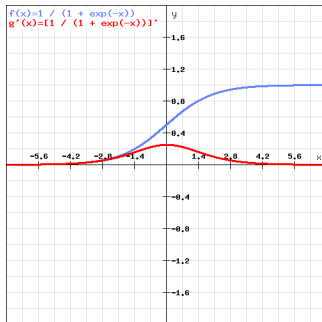
- Recall: every NN has this form of operations

$$\hat{\mathbf{y}} = (\phi_n \circ \mathbf{f}_{W_n} \circ \phi_{n-1} \circ \mathbf{f}_{W_{n-1}} \dots \phi_1 \circ \mathbf{f}_{W_1})(\mathbf{X})$$

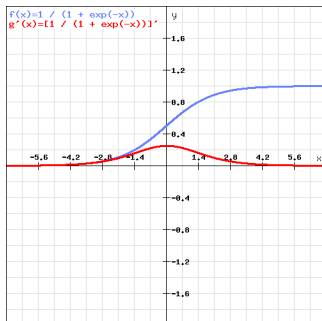
- Two basic operations
 - Linear: $o_i = \mathbf{f}_{W_i}(a_{i-1}) = W_i^T a_{i-1} + b_i$
 - Nonlinear (by activation functions): $a_i = \phi(o_i)$
- To model non linearity, activations need be *non-linear*.



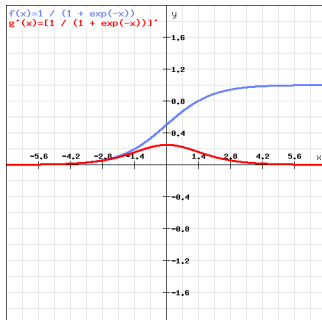
- $\sigma(x) = 1/(1 + \exp(-x))$ and $\sigma'(x) = \sigma(x)(1 - \sigma(x))$



- $\sigma(x) = 1/(1 + \exp(-x))$ and $\sigma'(x) = \sigma(x)(1 - \sigma(x))$
- Both are always positive

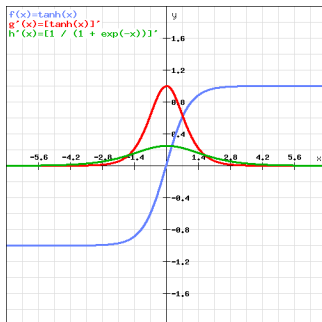


- $\sigma(x) = 1/(1 + \exp(-x))$ and $\sigma'(x) = \sigma(x)(1 - \sigma(x))$
- Both are always positive
- Squash unbounded values to $[0, 1]$, popular as logistic function



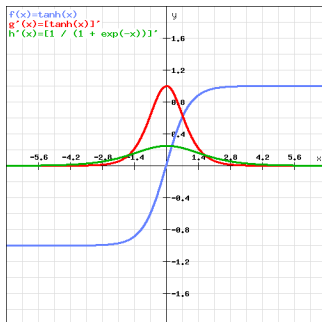
- $\sigma(x) = 1/(1 + \exp(-x))$ and $\sigma'(x) = \sigma(x)(1 - \sigma(x))$
- Both are always positive
- Squash unbounded values to $[0, 1]$, popular as logistic function
- Gradients vanishing problem

Tanh (Hyperbolic Tangent)



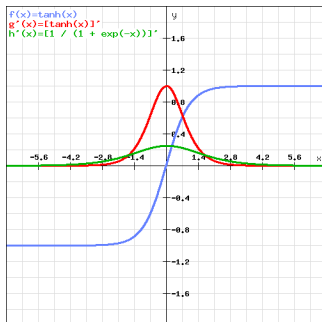
- $f(x) = \tanh(x)$ and $f'(x) = 1 - f^2(x)$

¹<http://yann.lecun.com/exdb/publis/pdf/lecun-98b.pdf>



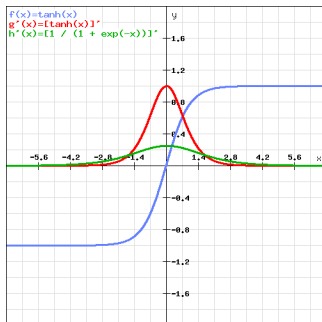
- $f(x) = \tanh(x)$ and $f'(x) = 1 - f^2(x)$
- Squash unbounded values to $[-1, 1]$, also used as logistic function

¹<http://yann.lecun.com/exdb/publis/pdf/lecun-98b.pdf>



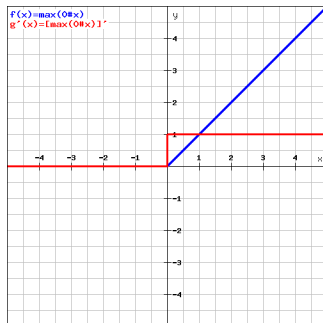
- $f(x) = \tanh(x)$ and $f'(x) = 1 - f^2(x)$
- Squash unbounded values to $[-1, 1]$, also used as logistic function
- Symmetric around origin, preferred over sigmoid¹

¹<http://yann.lecun.com/exdb/publis/pdf/lecun-98b.pdf>

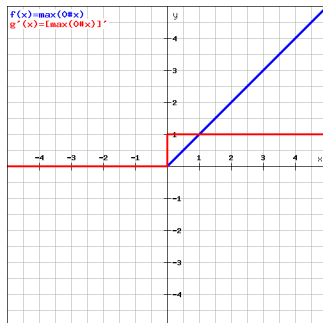


- $f(x) = \tanh(x)$ and $f'(x) = 1 - f^2(x)$
- Squash unbounded values to $[-1, 1]$, also used as logistic function
- Symmetric around origin, preferred over sigmoid¹
- Also, gradients vanishing problem

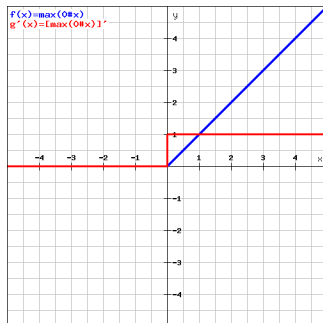
¹<http://yann.lecun.com/exdb/publis/pdf/lecun-98b.pdf>



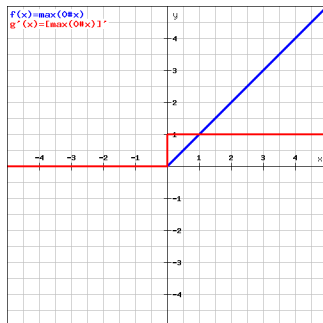
- $f(x) = x^+ = \max(0, x)$ and $f'(x) = \mathbf{1}_{x \geq 0}(x)$



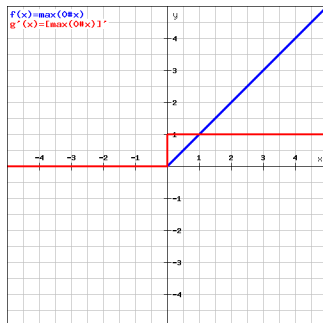
- $f(x) = x^+ = \max(0, x)$ and $f'(x) = \mathbf{1}_{x \geq 0}(x)$
- Both are *monotonic*, $f(x)$ is half-unbounded



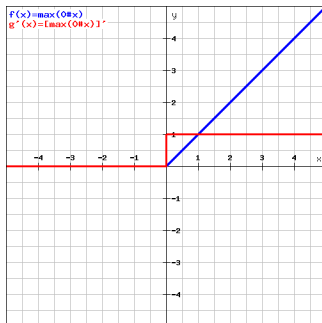
- $f(x) = x^+ = \max(0, x)$ and $f'(x) = \mathbf{1}_{x \geq 0}(x)$
- Both are *monotonic*, $f(x)$ is half-unbounded
- Inexpensive computation



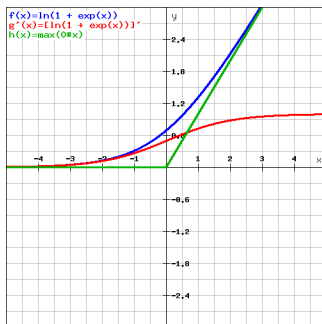
- $f(x) = x^+ = \max(0, x)$ and $f'(x) = \mathbf{1}_{x \geq 0}(x)$
- Both are *monotonic*, $f(x)$ is half-unbounded
- Inexpensive computation
- Avoid activating *all* neurons: sparse representations



- $f(x) = x^+ = \max(0, x)$ and $f'(x) = \mathbf{1}_{x \geq 0}(x)$
- Both are *monotonic*, $f(x)$ is half-unbounded
- Inexpensive computation
- Avoid activating *all* neurons: sparse representations
- *Hate* negative values,

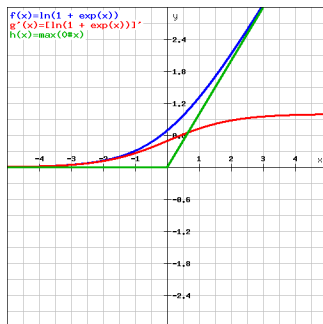


- $f(x) = x^+ = \max(0, x)$ and $f'(x) = \mathbf{1}_{x \geq 0}(x)$
- Both are *monotonic*, $f(x)$ is half-unbounded
- Inexpensive computation
- Avoid activating *all* neurons: sparse representations
- *Hate* negative values, but is most used by far



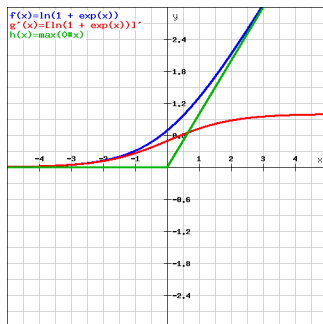
- $f(x) = \ln(1 + \exp(x))$ and $f'(x) = \sigma(x)$

²<http://proceedings.mlr.press/v15/glorot11a/glorot11a.pdf>



- $f(x) = \ln(1 + \exp(x))$ and $f'(x) = \sigma(x)$
- Smooth approximation of ReLU, *a.k.a* SmoothReLU

²<http://proceedings.mlr.press/v15/glorot11a/glorot11a.pdf>

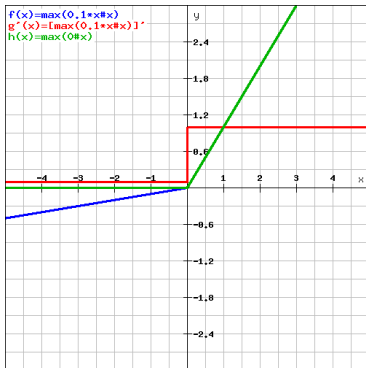


- $f(x) = \ln(1 + \exp(x))$ and $f'(x) = \sigma(x)$
- Smooth approximation of ReLU, *a.k.a* SmoothReLU
- Avoid zero saturation

²<http://proceedings.mlr.press/v15/glorot11a/glorot11a.pdf>

Neuron	MNIST	CIFAR10	NISTP	NORB
<i>With unsupervised pre-training</i>				
Rectifier	1.20%	49.96%	32.86%	16.46%
Tanh	1.16%	50.79%	35.89%	17.66%
Softplus	1.17%	49.52%	33.27%	19.19%
<i>Without unsupervised pre-training</i>				
Rectifier	1.43%	50.86%	32.64%	16.40%
Tanh	1.57%	52.62%	36.46%	19.29%
Softplus	1.77%	53.20%	35.48%	17.68%

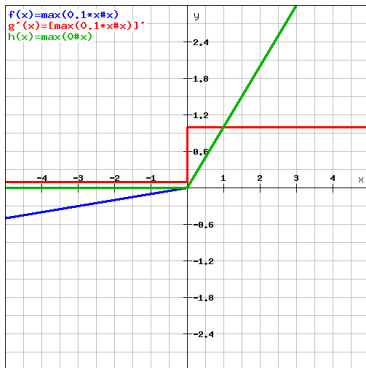
- $f(x) = \ln(1 + \exp(x))$ and $f'(x) = \sigma(x)$
- Smooth approximation of ReLU, *a.k.a* SmoothReLU
- Avoid zero saturation
- But is worse, and more expensive



- $f(x) = \max(0.01x, x)$

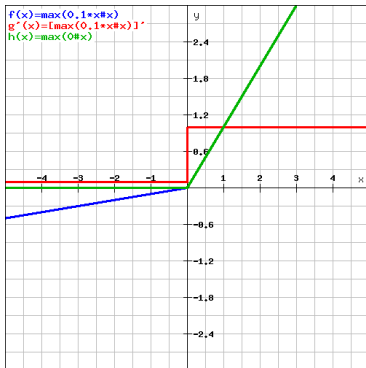
³https://ai.stanford.edu/~amaas/papers/relu_hybrid_icml2013_final.pdf

Leaky ReLU³: an(other) attempt



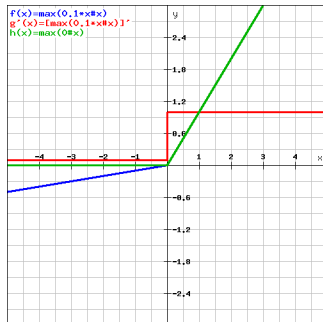
- $f(x) = \max(0.01x, x)$
- Fully unbounded

³https://ai.stanford.edu/~amaas/papers/relu_hybrid_icml2013_final.pdf



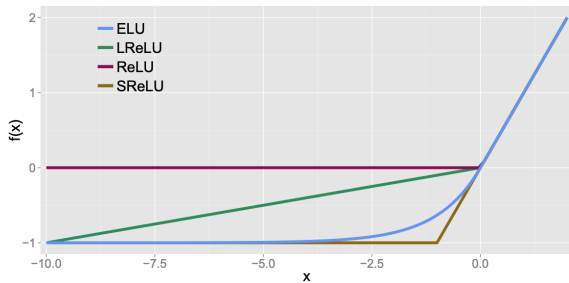
- $f(x) = \max(0.01x, x)$
- Fully unbounded
- Perform “nearly identically to standard rectifier DNNs”

³https://ai.stanford.edu/~amaas/papers/relu_hybrid_icml2013_final.pdf



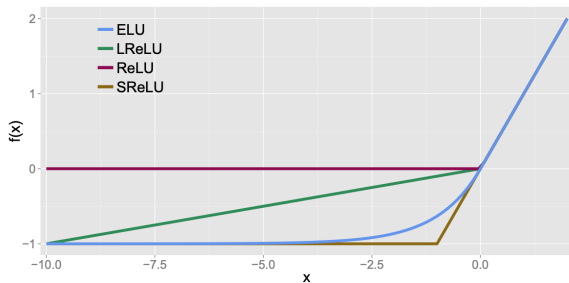
- $f(x) = \max(\alpha x, x)$ for small α
- Also called PReLU
- A generalized version of ReLU

⁴<https://arxiv.org/pdf/1502.01852.pdf>



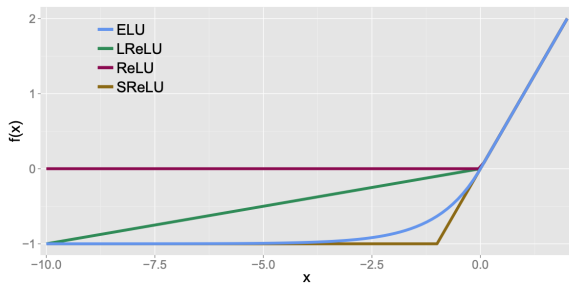
- $f(x) = \alpha(\exp(x) - 1)$ for negative x else x , for small α

⁵<https://arxiv.org/pdf/1511.07289.pdf>



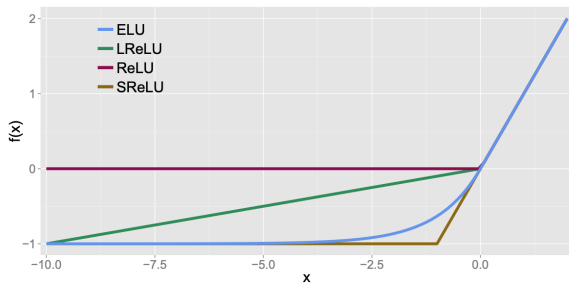
- $f(x) = \alpha(\exp(x) - 1)$ for negative x else x , for small α
- $f'(x) = f(x) + \alpha$ for negative x else 1

⁵<https://arxiv.org/pdf/1511.07289.pdf>



- $f(x) = \alpha(\exp(x) - 1)$ for negative x else x , for small α
- $f'(x) = f(x) + \alpha$ for negative x else 1
- Show better performance than previous rectifiers in many cases

⁵<https://arxiv.org/pdf/1511.07289.pdf>



- $f(x) = \alpha(\exp(x) - 1)$ for negative x else x , for small α
- $f'(x) = f(x) + \alpha$ for negative x else 1
- Show better performance than previous rectifiers in many cases
- Bring mean of activations to zero thus correct bias shifts

⁵<https://arxiv.org/pdf/1511.07289.pdf>

- API: `tf.keras.activations`

- 1 Choosing Activations
- 2 Callbacks
- 3 Save and Restore
- 4 Customization with Transfer Learning

callback (example from wikipedia)

```
1 def get_square(val):
2     """The callback."""
3     return val ** 2
4
5 def caller(func, val):
6     return func(val)
7
8 caller(get_square, 5) # 25
```

- A function which is passed as a parameter
- Usually work in asynchronous mode
- General usage involves triggering some thing after finishing a job.

callback_keras

```
1 from tf.keras.callbacks import Callback
2
3 class MyCallback(Callback):
4     ... # more later
5
6 my_callback = MyCallback()
7
8 model.fit(x_train, y_train,
9           batch_size=64,
10          steps_per_epoch=5,
11          epochs=3,
12          verbose=0,
13          callbacks=[my_callback])
```

- Is usually passed through training (`fit` function)

early_stopping

```
1 callback = tf.keras.callbacks.EarlyStopping(monitor='val_loss',
2                                             min_delta=0.001,
3                                             patience=3)
4
5 model.fit(x_train, y_train,
6         batch_size=64,
7         epochs=30,
8         validation_data=(x_val, y_val)),
9         validation_freq=1,
10        callbacks=[callback]) # note: always a list of Callbacks
```

- Usual usage: stop training after some rounds of no improvement
- Usually we monitor `val_loss` or `val_accuracy`

⁶https://www.tensorflow.org/api_docs/python/tf/keras/callbacks/EarlyStopping

model_checkpoint

```
1 callbacks = [  
2     tf.keras.callbacks.ModelCheckpoint(  
3         filepath='mymodel_{epoch}',  
4         save_best_only=True,  
5         monitor='val_loss', # only if this improves  
6         save_weights_only=True,  
7         save_freq='epoch',  
8         verbose=1)  
9 ]  
10 model.fit(x_train, y_train,  
11         batch_size=64,  
12         epochs=30,  
13         #validation_data=(x_val, y_val)),  
14         #validation_freq=1,  
15         validation_split=0.3,  
16         callbacks=callbacks)
```

- Usual usage: save after one or a few epochs

⁷ https://www.tensorflow.org/api_docs/python/tf/keras/callbacks/ModelCheckpoint

tensorboard

```
1 callbacks = [  
2     tf.keras.callbacks.TensorBoard(log_dir='./logs',  
3                                     histogram_freq=1,  
4                                     write_graph=True,  
5                                     write_images=False)  
6 ]  
7 model.fit(x_train, y_train,  
8           batch_size=64,  
9           epochs=30,  
10          validation_split=0.3,  
11          callbacks=callbacks  
12          )
```

- Usual usage:
 - Monitor training and evaluations
 - Visualize histograms of activations
 - Visualize embeddings
- To visualize: `tensorboard --logdir=./logs`

⁸https://www.tensorflow.org/api_docs/python/tf/keras/callbacks/TensorBoard

custom_callback

```
1 class CustomCallback(tf.keras.callbacks.Callback):
2     def __init__(self, new_data):
3         super(CustomCallback, self).__init__()
4         self.data = new_data
5
6     # before/after fit()|eval()|predict()
7     def on_(train|test|predict)_(begin|end)(self, logs=None):
8         ...
9     # before/after each batch processing
10    def on_(train|test|predict)_batch_(begin|end)(self, batch, logs=None):
11        ...
12    # before/after each epoch
13    def on_epoch_(begin|end)(self, epoch, logs=None):
14        ...
```

- Basically you can customize however you want
- *E.g.* Val/Test after n batches, periodically sample some results to check, ...

logs_dict

```
1 class SomeCallback(tf.keras.callbacks.Callback):
2     def on_train_begin(self, logs={}):
3         self.losses = []
4         self accuracies = []
5     def on_batch_end(self, batch, logs={}):
6         self.losses.append(logs.get('loss'))
7         self accuracies.append(logs.get('accuracy'))
8         print('For batch {}, loss {:.2f}, acc {:.2f}.'.
9             .format(batch, logs['loss'], logs['accuracy']))
10
11 model.compile(optimizer='adam',
12               loss='categorical_crossentropy',
13               metrics=['accuracy'], # will go into logs
14               callbacks=[SomeCallback()])
15 # logs dict look like so while on training:
16 {'batch': 1000, 'size': 32, 'loss': 0.3528049, 'accuracy': 0.8058816}
```

- Store “the loss value, and all the metrics at the end of a batch or epoch”

- List of `tf.keras` predefined: [callbacks](#)
- Tutorial: [keras custom callbacks](#)
- Tutorial: [keras train and evaluate](#)
- Source [code](#)

- 1 Choosing Activations
- 2 Callbacks
- 3 Save and Restore**
- 4 Customization with Transfer Learning

`tf_save_restore`

```
1 model1 = Model1()
2 model2 = Model2()
3 opt = tf.train.AdamOptimizer()
4
5 # define a checkpoint object wrapping all objects in K-V type
6 checkpoint = tf.train.Checkpoint(model1=model1, model2=model2, opt=opt)
7 checkpoint_prefix = './ckpt' # directory to save
8
9 # save
10 checkpoint.save(checkpoint_prefix) # ckpt-1.index, ckpt-1.data-00000-of-00001, ...
11
12 # normally restore the latest one
13 checkpoint.restore(tf.train.latest_checkpoint(checkpoint_path))
```

- Using `tf.train.Checkpoint` to wrap any object you have
- Can use validation score to save the best latest checkpoint

keras_manual_save_weights_only

```
1 model = create_model() # including compile()
2
3 # train
4 model.fit(x, y, ...)
5
6 # save
7 model.save_weights('./checkpoints/my_checkpoint')
8
9 # load later
10 # sometimes, create a new exact same architecture then load weights
11 model.load_weights('./checkpoints/my_checkpoint')
```

- Need to have an exact same architecture to load weights
- Save some space

`keras_manual_save_whole`

```
1 model = create_model() # including compile()
2
3 # train
4 model.fit(x, y, ...)
5
6 # save
7 model.save('./checkpoints/my_checkpoint') # whole model
8 # or
9 model.save('./checkpoints/my_checkpoint.h5') # force to HDF5 format
10
11 # load
12 new_model = tf.keras.models.load_model('my_model.h5')
```

- Whole model means: architecture, weights and optimizer
- Restore is simpler: no need to have a scaffold to load weights into

model_checkpoint

```
1 callbacks = [  
2     tf.keras.callbacks.ModelCheckpoint(  
3         filepath='mymodel_{epoch}',  
4         save_best_only=True,  
5         monitor='val_loss', # only if this improves  
6         save_weights_only=True,  
7         save_freq='epoch',  
8         verbose=1)  
9 ]  
10 model.fit(x_train, y_train,  
11         batch_size=64,  
12         epochs=30,  
13         #validation_data=(x_val, y_val)),  
14         #validation_freq=1,  
15         validation_split=0.3,  
16         callbacks=callbacks)
```

- Usual usage: save after one or a few epochs, or save best models

model_checkpoint

```
1 callbacks = [  
2     tf.keras.callbacks.ModelCheckpoint(  
3         filepath='mymodel_{epoch}',  
4         save_best_only=True,  
5         monitor='val_loss', # only if this improves  
6         save_weights_only=True,  
7         save_freq='epoch',  
8         verbose=1)  
9 ]  
10 model.fit(x_train, y_train,  
11         batch_size=64,  
12         epochs=30,  
13         #validation_data=(x_val, y_val)),  
14         #validation_freq=1,  
15         validation_split=0.3,  
16         callbacks=callbacks)
```

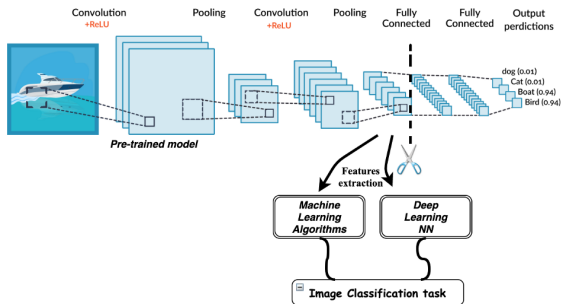
- Usual usage: save after one or a few epochs, or save best models
- Same options: save only weights (`model.save_weights`) or the whole model (`model.save`)

tf_keras_restore

```
1 # create an exact same model
2 model = Model() # must be the same structure
3
4 # restore
5 model.load_weights("mymodel.best.hdf5")
6
7 # compile and eval
8 model.compile(loss=..., optimizer=..., metrics=[...])
9 model.evaluate(x_test, y_test, verbose=1)
```

- Same use cases: restore from disk to evaluate a model or fine tune a pretrained model
- Same restore as manual operations above

- 1 Choosing Activations
- 2 Callbacks
- 3 Save and Restore
- 4 Customization with Transfer Learning



- General use case: fine tuning with new yet similar data
- General techniques: re-use (ideally most) part of known knowledge into new data
- Implementation: re-use (ideally most) architecture and weights

keras_load_source

```
1 # create an exact same model and restore weights
2 source_model = Model() # must be the same structure
3 source_model.load_weights("mymodel.best.hdf5")
4 source_model.compile(loss=..., optimizer=..., metrics=[...])
5
6 # or load the whole model as a whole
7 source_model = tf.keras.models.load_model('my_model.h5')
8
9 # customize
10 ... # soon later
```

- First is to load the source model with weights
- Techniques are the same as above for customized models

keras_load_predefined_source

```
1 # Create the base model from the pre-trained model MobileNet V2
2 source_model = tf.keras.applications.ResNet50(input_shape=None,
3                                               include_top=False,
4                                               weights='imagenet')
5
6 # customize
7 ... # soon later
```

- `tf.keras.applications` has a list of predefined networks
- Weights are automatically downloaded when we load
- `input_shape` is inferred if `None`, or explicitly assigned
- Top FC layer can be optionally excluded

⁹https://github.com/keras-team/keras-applications/tree/master/keras_applications

keras_add_layers

```
1 # add layers
2 global_average_layer = tf.keras.layers.GlobalAveragePooling2D()
3 prediction_layer = keras.layers.Dense(1)
4
5 # define new target model
6 target_model = tf.keras.Sequential([
7     source_model,
8     global_average_layer,
9     prediction_layer
10 ])
11
12 # compile
13 base_learning_rate = 0.0001
14 target_model.compile(optimizer=tf.keras.optimizers.RMSprop(lr=base_learning_rate),
15                     loss=tf.keras.losses.BinaryCrossentropy(),
16                     metrics=['accuracy'])
```

- Define new layer(s) and define a new model
- `keras.Sequential` accepts a model as a part

keras_add_layers_random

```
1 # define new target model
2 target_model = tf.keras.Sequential()
3
4 for layer in source_model.layers[:-10]: # exclude last 10 layers from copying
5     target_model.add(layer)
6
7 # add your new layers here
8 target_model.add(global_average_layer)
9 target_model.add(prediction_layer)
10
11 # compile
12 ...
```

- Define a new empty model and copy layers
- Weights should have been loaded before copy

keras_add_layers_random_two

```
1 # define new target model
2 target_model = tf.keras.Model()
3
4 taken = model.layers[-10].output
5
6 # add your new layers here
7 global_average_layer = tf.keras.layers.GlobalAveragePooling2D()
8 prediction_layer = keras.layers.Dense(1)
9
10 # make them attached sequentially
11 out = global_average_layer(taken)
12 out = prediction_layer(out)
13
14 # define new model with tf.keras.Model
15 target_model = tf.keras.Model(tf.keras.layers.Input(shape=(224,224,3)), out)
16
17 # compile
18 ...
```

- Take output of the point we want
- And attach new layers sequentially

freeze_source_weights

```
1 # freeze source weights which is part of target model
2 # should do before adding new layers
3 source_model.trainable = False
4
5 # define target_model
6 target_model = tf.keras.Sequential([
7     source_model,
8     global_average_layer,
9     prediction_layer
10 ])
11
12 # compile
13 target_model.compile(...)
14
15 # train new layers only
16 history = target_model.fit(x_train, y_train, ...) # we use history later
17
18 # eval
19 target_model.evaluate(x_test, y_test, ...)
```

- First, train new layers only

fine_tune

```
1 # un-freeze source weights which is part of target model
2 source_model.trainable = True
3
4 # Freeze all the layers before the 'fine_tune_at' layer
5 for layer in base_model.layers[:100]: # fine tune from 100th layer onwards
6     layer.trainable = False
7
8 # NOTE: compile again
9 target_model.compile(loss=tf.keras.losses.BinaryCrossentropy(),
10                      optimizer = tf.keras.optimizers.RMSprop(lr=base_learning_rate/10),
11                      metrics=['accuracy'])
12
13 # continue training
14 total_epochs = initial_epochs + fine_tune_epochs
15 history_fine = target_model.fit(train_batches,
16                                epochs=total_epochs,
17                                initial_epoch = history.epoch[-1],
18                                validation_data=validation_batches)
```

- Then, fine tune pretrained weights
- For deep CNN, should fine tunes more specialized layers only

- Tutorial: `tf.train.Checkpoint`
- Tutorial: `tf.keras-Save and Serialize`
- Tutorial: `tf.keras-Based Save and Restore`
- Tutorial: `SavedModel format`
- Tutorial: `Load pretrained model and add new layers`