

These are simple examples of how to use GGML for inferencing. The first example uses convolutional neural network (CNN), the second one uses fully connected neural network.

- ## MNIST with CNN (使用CNN的MNIST)

- 这个实现在MNIST测试集上达到了大约99%的正确率

Setup the Python environemt and build the examples according to the main README.

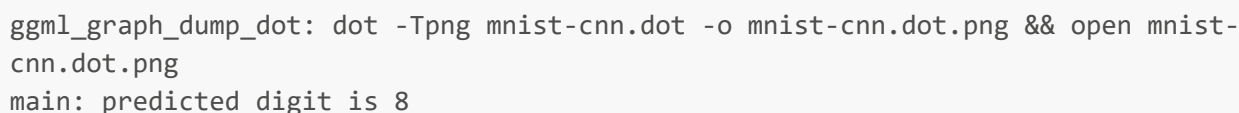
- Use the `mnist-cnn.py` script to train the model and convert it to GGUF format:

- ```
$ python3 ../examples/mnist/mnist-cnn.py train mnist-cnn-model
...
Keras model saved to 'mnist-cnn-model'
```

```
$ python3 ../examples/mnist/mnist-cnn.py convert mnist-cnn-model
...
Model converted and saved to 'mnist-cnn-model.gguf'
```

```
$./bin/mnist-cnn mnist-cnn-model.gguf ../examples/mnist/models/mnist/t10k-
images.idx3-ubyte
main: loaded model in 5.17 ms


```



The computational graph illustrates the data flow and operations for a neural network layer. It starts with an input node (x) of shape [128, 784] and a kernel node (k) of shape [2, 3, 3]. These are combined via a convolution operation (conv\_2d(x)) to produce a node of shape [0 [26, 26, 32]]. This node is then processed by a unary operation (unary(x)) to produce a node of shape [2 [26, 26, 32]]. A bias1 node (b1) of shape [26, 26] is added to this result via an arc 1 operation. The resulting node of shape [1 [26, 26, 32]] is then processed by another unary operation (unary(x)) to produce a node of shape [2 [11, 11, 64]]. A bias2 node (b2) of shape [11, 11] is added to this result via an arc 1 operation. The resulting node of shape [5 [11, 11, 64]] is then processed by a unary operation (unary(x)) to produce a node of shape [6 [11, 11, 64]]. This node is then processed by a permutation operation (permuted(x)) to produce a node of shape [8 [64, 5, 5]]. This node is then processed by a permutation operation (permuted(x)) to produce a node of shape [9 [64, 5, 5]]. This node is then processed by a permutation operation (permuted(x)) to produce a node of shape [10 [1600, 1]]. A dense w node (w) of shape [1600, 10] is added to this result via an arc 0 operation. The resulting node of shape [11 [16, 1]] is then processed by a permutation operation (permuted(x)) to produce a node of shape [12 [10, 1]]. A dense b node (b) of shape [10, 1] is added to this result via an arc 1 operation. The final output is produced by a softmax operation (softmax(x)).

A fully connected layer + relu, followed by a fully connected layer + softmax.

A Google Colab notebook for training a simple two-layer network to recognize digits is located [here](#). You can use this to save a pytorch model to be converted to qqml format.

- GGML "format" is whatever you choose for efficient loading. In our case, we just save the hyperparameters used plus the model weights and biases. Run `convert-h5-to-ggml.py` to convert your pytorch model. The

output format is (GGML“格式”是您选择的任何格式，以便高效加载。在我们的例子中，我们只保存使用的超参数以及模型权重和偏差。运行 `convert-h5-to-ggml.py` 来转换您的 `pytorch` 模型。输出格式为 )：

- magic constant (int32)
- repeated list of tensors
- number of dimensions of tensor (int32)
- tensor dimension (int32 repeated)
- values of tensor (int32)

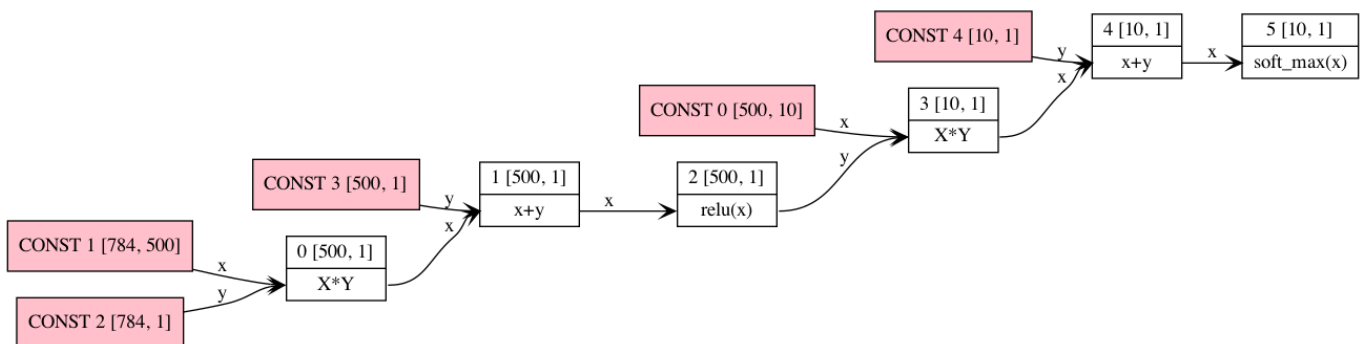
Run `convert-h5-to-ggml.py mnist_model.state_dict` where `mnist_model.state_dict` is the saved `pytorch` model from the Google Colab. For quickstart, it is included in the `mnist/models` directory.

```
mkdir -p models/mnist
python3 ../examples/mnist/convert-h5-to-ggml.py
../examples/mnist/models/mnist/mnist_model.state_dict
```

## Running the example

```
./bin/mnist ./models/mnist/ggml-model-f32.bin ../examples/mnist/models/mnist/t10k-images.idx3-ubyte
```

Computation graph:



## Web demo

The example can be compiled with Emscripten like this:

```
cd examples/mnist
emcc -I../include -I../include/ggml -I../examples ../src/ggml.c
../src/ggml-quant.c main.cpp -o web/mnist.js -s
EXPORTED_FUNCTIONS=['["_wasm_eval","_wasm_random_digit","_malloc","_free"]'] -s
EXPORTED_RUNTIME_METHODS=['["ccall"]'] -s ALLOW_MEMORY_GROWTH=1 --preload-file
models/mnist
```

Online demo: <https://mnist.ggerganov.com>