

I was trying to clean some of my code so I can add more models. However, this can never happen without the awesome kernels from other talented Kagglers. Forgive me if I missed any.

- CLR from: <https://www.kaggle.com/hireme/fun-api-keras-f1-metric-cyclical-learning-rate/code>
- Based on SRK's kernel: <https://www.kaggle.com/sudalairakumar/a-look-at-different-embeddings>
- Vladimir Demidov's 2DCNN textClassifier: <https://www.kaggle.com/yekeno1/2dnn-textclassifier>
- Attention layer from Khoi Nguyen: <https://www.kaggle.com/suicaokhoailang/lstm-attention-baseline-0-652-lb>
- LSTM model from Strideradu: <https://www.kaggle.com/strideradu/word2vec-and-gensim-go-go-go>
- <https://www.kaggle.com/danofar/different-embeddings-with-attention-fork>
- <https://www.kaggle.com/ryanzhang/tfidf-naivebayes-logreg-baseline>
- Borrowed some idea from this model: <https://www.kaggle.com/c/jigsaw-toxic-comment-classification-challenge/discussion/52644>
- Sentence length seems to a good feature: <https://www.kaggle.com/thebrownking20/analyzing-quora-for-the-insinceres>

Some new things here:

- Take average of embeddings (Unweighted DME) instead of blending predictions: <https://arxiv.org/pdf/1804.07983.pdf>
- The original paper of this idea comes from: Frustratingly Easy Meta-Embedding – Computing Meta-Embeddings by Averaging Source Word Embeddings
- Modified the code to choose best threshold
- Robust method for blending weights: sort the val score and give the final weight

Some thoughts:

- Although I polished a kernel on Transformer, I will not use it
- Too much randomness in CuDNN. You may get different results by just running this kernel
- Blending rocks

```
In [ ]: # This Python 3 environment comes with many helpful analytics libraries installed
# It is defined by the kaggle/python docker image: https://github.com/kaggle/docker-python
# For example, here's several helpful packages to load in

import numpy as np # linear algebra
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)

# Input data files are available in the "../input/" directory.
# For example, running this (by clicking run or pressing Shift+Enter) will list the files in the input
directory

import os
print(os.listdir("../input"))

# Any results you write to the current directory are saved as output.
```

```
In [ ]: ## some config values
embed_size = 300 # how big is each word vector
max_features = 95000 # how many unique words to use (i.e num rows in embedding vector)
maxLen = 70 # max number of words in a question to use
```

Load packages and data

```
In [ ]: import os
import time
import numpy as np # linear algebra
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
from tqdm import tqdm
import math
from sklearn.model_selection import train_test_split
from sklearn import metrics
from sklearn.model_selection import GridSearchCV, StratifiedKFold
from sklearn.metrics import f1_score, roc_auc_score

from keras.preprocessing.text import Tokenizer
from keras.preprocessing.sequence import pad_sequences
from keras.layers import Dense, Input, CuDNNLSTM, Embedding, Dropout, Activation, CuDNNLSTM, Conv1D
from keras.layers import Bidirectional, GlobalMaxPooling1D, GlobalMaxPooling2D, GlobalAveragePooling1D
from keras.layers import Input, Embedding, Dense, Conv2D, MaxPool2D, concatenate
from keras.layers import Reshape, Flatten, Concatenate, Dropout, SpatialDropout1D
from keras.optimizers import Adam
from keras.models import Model
from keras import backend as K
from keras.engine.topology import Layer
from keras import initializers, regularizers, constraints, optimizers, layers
from keras.layers import concatenate
from keras.callbacks import *
```

```
In [ ]: def load_and_prec():
    train_df = pd.read_csv("../input/train.csv")
    test_df = pd.read_csv("../input/test.csv")
    print("Train shape : ",train_df.shape)
    print("Test shape : ",test_df.shape)

    ## fill up the missing values
    train_X = train_df["question_text"].fillna("#")
    test_X = test_df["question_text"].fillna("#")

    ## Tokenize the sentences
    tokenizer = Tokenizer(num_words=max_features)
    tokenizer.fit_on_texts(list(train_X))
    train_X = tokenizer.texts_to_sequences(train_X)
    test_X = tokenizer.texts_to_sequences(test_X)

    ## Pad the sentences
    train_X = pad_sequences(train_X, maxlen=maxlen)
    test_X = pad_sequences(test_X, maxlen=maxlen)

    ## Get the target values
    train_y = train_df["target"].values

    #shuffling the data
    np.random.seed(2018)
    trn_idx = np.random.permutation(len(train_X))

    train_X = train_X[trn_idx]
    train_y = train_y[trn_idx]

    return train_X, test_X, train_y, tokenizer.word_index
```

Load embeddings

```
In [ ]: def load_glove(word_index):
    EMBEDDING_FILE = '../input/embeddings/glove.840B.300d/glove.840B.300d.txt'
    def get_coefs(word,*arr): return word, np.asarray(arr, dtype='float32')
    embeddings_index = dict(get_coefs(*o.split(" ")) for o in open(EMBEDDING_FILE))

    all_embs = np.stack(embeddings_index.values())
    emb_mean,emb_std = all_embs.mean(), all_embs.std()
    embed_size = all_embs.shape[1]

    # word index = tokenizer.word_index
    nb_words = min(max_features, len(word_index))
    embedding_matrix = np.random.normal(emb_mean, emb_std, (nb_words, embed_size))
    for word, i in word_index.items():
        if i >= max_features: continue
        embedding_vector = embeddings_index.get(word)
        if embedding_vector is not None: embedding_matrix[i] = embedding_vector

    return embedding_matrix

def load_fasttext(word_index):
    EMBEDDING_FILE = '../input/embeddings/wiki-news-300d-1M/wiki-news-300d-1M.vec'
    def get_coefs(word,*arr): return word, np.asarray(arr, dtype='float32')
    embeddings_index = dict(get_coefs(*o.split(" ")) for o in open(EMBEDDING_FILE) if len(o)>100)

    all_embs = np.stack(embeddings_index.values())
    emb_mean,emb_std = all_embs.mean(), all_embs.std()
    embed_size = all_embs.shape[1]

    # word index = tokenizer.word_index
    nb_words = min(max_features, len(word_index))
    embedding_matrix = np.random.normal(emb_mean, emb_std, (nb_words, embed_size))
    for word, i in word_index.items():
        if i >= max_features: continue
        embedding_vector = embeddings_index.get(word)
        if embedding_vector is not None: embedding_matrix[i] = embedding_vector

    return embedding_matrix

def load_para(word_index):
    EMBEDDING_FILE = '../input/embeddings/paragram_300_sl999/paragram_300_sl999.txt'
    def get_coefs(word,*arr): return word, np.asarray(arr, dtype='float32')
    embeddings_index = dict(get_coefs(*o.split(" ")) for o in open(EMBEDDING_FILE, encoding='utf8', err
ors='ignore') if len(o)>100)

    all_embs = np.stack(embeddings_index.values())
    emb_mean,emb_std = all_embs.mean(), all_embs.std()
    embed_size = all_embs.shape[1]

    # word index = tokenizer.word_index
    nb_words = min(max_features, len(word_index))
    embedding_matrix = np.random.normal(emb_mean, emb_std, (nb_words, embed_size))
    for word, i in word_index.items():
        if i >= max_features: continue
        embedding_vector = embeddings_index.get(word)
        if embedding_vector is not None: embedding_matrix[i] = embedding_vector

    return embedding_matrix
```

Attention layer

```
In [ ]: # https://www.kaggle.com/suicaokhoailang/lstm-attention-baseline-0-652-lb

class Attention(Layer):
    def __init__(self, step_dim,
                  W_regularizer=None, b_regularizer=None,
                  W_constraint=None, b_constraint=None,
                  bias=True, **kwargs):
        self.supports_masking = True
        self.init = initializers.get('glorot_uniform')

        self.W_regularizer = regularizers.get(W_regularizer)
        self.b_regularizer = regularizers.get(b_regularizer)

        self.W_constraint = constraints.get(W_constraint)
        self.b_constraint = constraints.get(b_constraint)

        self.bias = bias
        self.step_dim = step_dim
        self.features_dim = 0
        super(Attention, self).__init__(**kwargs)

    def build(self, input_shape):
        assert len(input_shape) == 3

        self.W = self.add_weight((input_shape[-1],),
                                initializer=self.init,
                                name='%i_W'.format(self.name),
                                regularizer=self.W_regularizer,
                                constraint=self.W_constraint)
        self.features_dim = input_shape[-1]

        if self.bias:
            self.b = self.add_weight((input_shape[-1],),
                                    initializer='zero',
                                    name='%i_b'.format(self.name),
                                    regularizer=self.b_regularizer,
                                    constraint=self.b_constraint)
        else:
            self.b = None

        self.built = True

    def compute_mask(self, input, input_mask=None):
        return None

    def call(self, x, mask=None):
        features_dim = self.features_dim
        step_dim = self.step_dim

        eij = K.reshape(K.dot(K.reshape(x, (-1, features_dim)),
                               K.reshape(self.W, (features_dim, 1))), (-1, step_dim))

        if self.bias:
            eij += self.b

        eij = K.tanh(eij)

        a = K.exp(eij)

        if mask is not None:
            a *= K.cast(mask, K.floatx())

        a /= K.cast(K.sum(a, axis=1, keepdims=True) + K.epsilon(), K.floatx())

        a = K.expand_dims(a)
        weighted_input = x * a
        return K.sum(weighted_input, axis=1)

    def compute_output_shape(self, input_shape):
        return input_shape[0], self.features_dim
```

F1 score and CLR

```
In [ ]: # https://www.kaggle.com/hireme/fun-api-keras-f1-metric-cyclical-learning-rate/code

class CyclicalLR(Callback):
    """This callback implements a cyclical learning rate policy (CLR).
    The method cycles the learning rate between two boundaries with a
    constant frequency, as detailed in this paper (https://arxiv.org/abs/1506.01186).
    The amplitude of the cycle can be scaled on a per-iteration or
    per-cycle basis.
    This class has three built-in policies, as put forth in the paper.
    "triangular":
        A basic triangular cycle w/ no amplitude scaling.
    "triangular2":
        A basic triangular cycle that scales initial amplitude by half each cycle.
    "exp_range":
        A cycle that scales initial amplitude by gamma**(cycle iterations) at each
        iteration.
    For more detail, please see paper.

    # Example
    ``python
    clr = CyclicalLR(base_lr=0.001, max_lr=0.006,
                    step_size=2000., mode='triangular')
    model.fit(X_train, Y_train, callbacks=[clr])

    Class also supports custom scaling functions:
    ``python
    clr_fn = lambda x: 0.5*(1+np.sin(x*np.pi/2.))
    clr = CyclicalLR(base_lr=0.001, max_lr=0.006,
                    step_size=2000., scale_fn=clr_fn,
                    scale_mode='cycle')
    ... model.fit(X_train, Y_train, callbacks=[clr])

    # Arguments
        base_lr: initial learning rate which is the
            lower boundary in the cycle.
        max_lr: upper boundary in the cycle. Functionally,
            it defines the cycle amplitude (max_lr - base_lr).
            The lr at any cycle is the sum of base_lr
            and some scaling of the amplitude; therefore
            max_lr may not actually be reached depending on
            scaling function.
        step_size: number of training iterations per
            half cycle. Authors suggest setting step_size
            2-8 x training iterations in epoch.
        mode: one of {triangular, triangular2, exp_range}.
            Default 'triangular'.
            Values correspond to policies detailed above.
            If scale_fn is not None, this argument is ignored.
        gamma: constant in 'exp_range' scaling function:
            gamma**(cycle iterations)
        scale_fn: Custom scaling policy defined by a single
            argument lambda function, where
            0 <= scale_fn(x) <= 1 for all x >= 0.
            mode parameter is ignored
        scale_mode: {'cycle', 'iterations'}.
            Defines whether scale_fn is evaluated on
            cycle number or cycle iterations (training
            iterations since start of cycle). Default is 'cycle'.
    """

    def __init__(self, base_lr=0.001, max_lr=0.006, step_size=2000., mode='triangular',
                 gamma=1., scale_fn=None, scale_mode='cycle'):
        super(CyclicalLR, self).__init__()

        self.base_lr = base_lr
        self.max_lr = max_lr
        self.step_size = step_size
        self.mode = mode
        self.gamma = gamma
        if scale_fn == None:
            if self.mode == 'triangular':
                self.scale_fn = lambda x: 1.
            elif self.mode == 'triangular2':
                self.scale_fn = lambda x: 1/(2.**((x-1)
            elif self.mode == 'cycle':
                self.scale_fn = lambda x: 1.
            elif self.mode == 'exp_range':
                self.scale_fn = lambda x: gamma**(x)
            self.scale_mode = 'iterations'
        else:
            self.scale_fn = scale_fn
            self.scale_mode = scale_mode
        self.clr_iterations = 0.
        self.trn_iterations = 0.
        self.history = {}

        self.reset()

    def _reset(self, new_base_lr=None, new_max_lr=None,
               new_step_size=None):
        """Resets boundary iterations.
        Optional boundary/step size adjustment.
        """
        if new_base_lr != None:
            self.base_lr = new_base_lr
        if new_max_lr != None:
            self.max_lr = new_max_lr
        if new_step_size != None:
            self.step_size = new_step_size
            self.clr_iterations = 0.

    def clr(self):
        cycle = np.floor(1+self.clr_iterations/(2*self.step_size))
        x = np.abs(self.clr_iterations/self.step_size - 2*cycle + 1)
        if self.scale_mode == 'cycle':
            return self.base_lr + (self.max_lr-self.base_lr)*np.maximum(0, (1-x))*self.scale_fn(cycle)
        else:
            return self.base_lr + (self.max_lr-self.base_lr)*np.maximum(0, (1-x))*self.scale_fn(self.clr_iterations)

    def on_train_begin(self, logs={}):
        logs = logs or {}

        if self.clr_iterations == 0:
            K.set_value(self.model.optimizer.lr, self.base_lr)
        else:
            K.set_value(self.model.optimizer.lr, self.clr())

    def on_batch_end(self, epoch, logs=None):
        logs = logs or {}
        self.trn_iterations += 1
        self.clr_iterations += 1

        self.history.setdefault('lr', []).append(K.get_value(self.model.optimizer.lr))
        self.history.setdefault('iterations', []).append(self.trn_iterations)

        for k, v in logs.items():
            self.history.setdefault(k, []).append(v)

        K.set_value(self.model.optimizer.lr, self.clr())

def f1(y_true, y_pred):
    """
    metric from here
    https://stackoverflow.com/questions/43547402/how-to-calculate-f1-macro-in-keras
    """
    def recall(y_true, y_pred):
        """Recall metric.

        Only computes a batch-wise average of recall.

        Computes the recall, a metric for multi-label classification of
        how many relevant items are selected.
        """
        true_positives = K.sum(K.round(K.clip(y_true * y_pred, 0, 1)))
        possible_positives = K.sum(K.round(K.clip(y_true, 0, 1)))
        recall = true_positives / (possible_positives + K.epsilon())
        return recall

    def precision(y_true, y_pred):
        """Precision metric.

        Only computes a batch-wise average of precision.

        Computes the precision, a metric for multi-label classification of
        how many selected items are relevant.
        """
        true_positives = K.sum(K.round(K.clip(y_true * y_pred, 0, 1)))
        predicted_positives = K.sum(K.round(K.clip(y_pred, 0, 1)))
        precision = true_positives / (predicted_positives + K.epsilon())
        return precision
    precision = precision(y_true, y_pred)
    recall = recall(y_true, y_pred)
    return 2*((precision*recall)/(precision+recall+K.epsilon()))
```

LSTM models

```
In [ ]: def model_lstm_atten(embedding_matrix):
    inp = Input(shape=(maxlen,))
    x = Embedding(max_features, embed_size, weights=[embedding_matrix], trainable=False)(inp)
    x = SpatialDropout1D(0.1)(x)
    x = Bidirectional(CuDNNLSTM(40, return_sequences=True))(x)
    y = Bidirectional(CuDNNLSTM(40, return_sequences=True))(x)

    atten_1 = Attention(maxlen)(x) # skip connect
    atten_2 = Attention(maxlen)(y)
    avg_pool = GlobalAveragePooling1D()(y)
    max_pool = GlobalMaxPooling1D()(y)

    conc = concatenate([atten_1, atten_2, avg_pool, max_pool])
    conc = Dense(16, activation='relu')(conc)
    conc = Dropout(0.1)(conc)
    outp = Dense(1, activation='sigmoid')(conc)

    model = Model(inputs=inp, outputs=outp)
    model.compile(loss='binary_crossentropy', optimizer='adam', metrics=[f1])

    return model
```

Train and predict

```
In [ ]: # https://www.kaggle.com/strideradu/word2vec-and-gensim-go-go-go
def train_pred(model, train_X, train_y, val_X, val_y, epochs=2, callback=None):
    for e in range(epochs):
        model.fit(train_X, train_y, batch_size=512, epochs=1, validation_data=(val_X, val_y), callbacks
        = callback, verbose=0)
        pred_val_y = model.predict([val_X], batch_size=1024, verbose=0)

        best_score = metrics.f1_score(val_y, (pred_val_y > 0.33).astype(int))
        print("Epoch: ", e, "- Val F1 Score: {:.4f}".format(best_score))

        pred_test_y = model.predict([test_X], batch_size=1024, verbose=0)
        print('=' * 60)
        return pred_val_y, pred_test_y, best_score
```

Main part: load, train, pred and blend

```
In [ ]: train_X, test_X, train_y, word_index = load_and_prec()
embedding_matrix_1 = load_glove(word_index)
embedding_matrix_2 = load_fasttext(word_index)
embedding_matrix_3 = load_para(word_index)
```

```
In [ ]: ## Simple average: http://aclweb.org/anthology/N18-2031

# We have presented an argument for averaging as
# a valid meta-embedding technique, and found experimental
# performance to be close to, or in some cases
# better than that of concatenation, with the
# additional benefit of reduced dimensionality

## Unweighted DME in https://arxiv.org/pdf/1804.07983.pdf

# "The downside of concatenating embeddings and
# giving that as input to an RNN encoder, however,
# is that the network then quickly becomes inefficient
# as we combine more and more embeddings."
```

```
embedding_matrix = np.mean([embedding_matrix_1, embedding_matrix_3], axis = 0)
np.shape(embedding_matrix)
```

```
In [ ]: # https://www.kaggle.com/ryanzhang/tfidf-naivebayes-logreg-baseline

def threshold_search(y_true, y_proba):
    best_threshold = 0
    best_score = 0
    for threshold in [i * 0.01 for i in range(100)]:
        score = f1_score(y_true=y_true, y_pred=y_proba > threshold)
        if score > best_score:
            best_threshold = threshold
            best_score = score
    search_result = {'threshold': best_threshold, 'f1': best_score}
    return search_result
```

```
In [ ]: DATA_SPLIT_SEED = 2018
clr = CyclicalLR(base_lr=0.001, max_lr=0.002,
                 step_size=3000., mode='exp_range',
                 gamma=0.9994)

train_meta = np.zeros(train_y.shape)
test_meta = np.zeros(test_X.shape(0))
splits = list(StratifiedKFold(n_splits=4, shuffle=True, random_state=DATA_SPLIT_SEED).split(train_X, tr
ain_y))

for idx, (train_idx, valid_idx) in enumerate(splits):
    X_train = train_X[train_idx]
    y_train = train_y[train_idx]
    X_val = train_X[valid_idx]
    y_val = train_y[valid_idx]
    model = model_lstm_atten(embedding_matrix)
    pred_val_y, pred_test_y, best_score = train_pred(model, X_train, y_train, X_val, y_val, epochs
    = 8, callback = [clr,])
    train_meta[valid_idx] = pred_val_y.reshape(-1) / len(splits)
    test_meta += pred_test_y.reshape(-1) / len(splits)
```

```
In [ ]: sub = pd.read_csv("../input/sample_submission.csv")
sub.prediction = test_meta > 0.33
sub.to_csv("submission.csv", index=False)
```

```
In [ ]: f1_score(y_true=train_y, y_pred=train_meta > 0.33)
```