There have been some issues regarding the correlation between CV and leaderboard scores in this competition. Every top-scoring public kernel has a much lower CV score than leaderboard score. It has also been very frustrating to tune a model to optimal CV score only to discover that the score on the Leaderboard is abysmal. In this kernel I am going to address this issue and propose a framework for robust local validation. The preprocessing and model architecture have stayed mostly the same as in my previous kernel. I'll also write about the impact of seeds on the score. Again, we'll start with standard imports. **Imports** In []: # standard imports import time import random import os from IPython.display import display import numpy as np import pandas as pd from scipy import stats import warnings # pytorch imports import torch import torch.nn as nn import torch.utils.data # imports for preprocessing the questions from keras.preprocessing.text import Tokenizer from keras.preprocessing.sequence import pad sequences # cross validation and metrics from sklearn.model_selection import StratifiedKFold from sklearn.metrics import f1 score # progress bars from tqdm import tqdm tqdm.pandas() # visualization import matplotlib.pyplot as plt import seaborn as sns warnings.filterwarnings("ignore", message="F-score is ill-defined and being set to 0.0 due to no predic ted samples.") %matplotlib inline Loading the data In []: train df = pd.read csv("../input/train.csv") test df = pd.read csv("../input/test.csv") print('Train data dimension: ', train df.shape) display(train df.head()) print('Test data dimension: ', test df.shape) display(test df.head()) In []: enable_local_test = True if enable local test: n test = len(test df)train df, local test df = (train df.iloc[:-n test].reset index(drop=True), train_df.iloc[-n_test:].reset_index(drop=True)) else: local test df = pd.DataFrame([[None, None, 0], [None, None, 0]], columns=['qid', 'question text', 'target']) n test = 2Here, we create a dataframe I call local test df. We will pretend that this dataframe is the actual test dataframe. The only difference: We know the labels for this one! So we do not have to blindly submit our model and pray for a good score, but can instead tune the score we achieve on this test dataframe. I have set the size of the local test dataframe to 4 times the size of the public test dataframe. That is a reasonable size (~200k rows) but more or less arbitrary. So we are going to test our model on two sets now: First, the regular test set, and second, a local test set which we know the labels for. Overall, the procedure is: split the data in a train and local test set perform CV on the train set when tuning the model: evaluate the predictions on the local test set · when submitting: predict the samples in the true test set make the size of the local test set 0 for best performance (set enable local test to False) We actually perform CV on the train side of a regular train / test split. It would be ideal to wrap the whole thing in another K-Fold cross validation procedure. That is, however, not feasible regarding computing power on my local machine and in the kaggle kernels. **Utility functions** In []: def seed everything(seed=1234): random.seed(seed) os.environ['PYTHONHASHSEED'] = str(seed) np.random.seed(seed) torch.manual seed(seed) torch.cuda.manual seed(seed) torch.backends.cudnn.deterministic = True seed everything() In []: | def threshold search(y true, y proba): best threshold = 0best score = 0 for threshold in tqdm([i * 0.01 for i in range(100)], disable=True): 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 []: def sigmoid(x):**return** 1 / (1 + np.exp(-x)) **Processing input** In []: embed size = 300 max features = 95000maxlen = 70puncts = [',', '.', '"', ':', ')', '(', '-', '!', '?', '|', ';', "'", '\$', '&', '/', '[', ']', '>', '%' In []: , '=', '#', '*', '+', '\\', '•', '~', '@', '£', 'Â', '<mark>'</mark>', '½', 'à', '...', ˙`', ˙*', '″', '-', '•', 'â', '▶', '-', '¢', '²', '¬', '░', '¶', '↑', '±', '¿', '▼', '=', '¦', '║', '░', ':', '½', '⊕', '▼', '∗', '\"', 'Ã', '.', '\', '∞', '・', ') ', '↓', '、', '│', ' (', '»', ', ', '♪', '╩', '╚', '³', '・', '╦', '╣', '╔', '╗', '━', '╩', $"\varnothing"$, "1", $"\leq"$, $"\ddagger"$, $"\sqrt"$,] def clean text(x): x = str(x)for punct in puncts: x = x.replace(punct, f' {punct} ') return x In []: for df in [train df, test df, local test df]: df["question_text"] = df["question_text"].str.lower() df["question_text"] = df["question_text"].apply(lambda x: clean_text(x)) df["question text"].fillna(" ## ", inplace=True) x train = train df["question text"].values x test = test df["question text"].values x_test_local = local_test_df["question_text"].values tokenizer = Tokenizer(num_words=max_features) tokenizer.fit_on_texts(list(x_train) + list(x_test_local)) x_train = tokenizer.texts_to_sequences(x_train) x test = tokenizer.texts to sequences(x test) x_test_local = tokenizer.texts_to_sequences(x_test_local) x_train = pad_sequences(x_train, maxlen=maxlen) x test = pad sequences(x test, maxlen=maxlen) x_test_local = pad_sequences(x_test_local, maxlen=maxlen) y train = train df['target'].values y test = local test df['target'].values Creating the embeddings matrix In []: def load glove(word index, max features): EMBEDDING FILE = '../input/embeddings/glove.840B.300d/glove.840B.300d.txt' def get coefs(word, *arr): return word, np.asarray(arr, dtype='float32')[:300] embeddings index = [] for o in tqdm(open(EMBEDDING FILE)): embeddings_index.append(get_coefs(*o.split(" "))) except Exception as e: print(e) embeddings_index = dict(embeddings_index) 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) + 1) 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, max features): 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')[:300] embeddings_index = [] for o in tqdm(open(EMBEDDING_FILE)): **if** len(o) <= 100: continue try: coefs = get coefs(*o.split(" ")) assert len(coefs[1]) == 300 embeddings index.append(coefs) except Exception as e: print(e) embeddings index = dict(embeddings index) all embs = np.stack(embeddings index.values()) emb mean,emb std = all embs.mean(), all embs.std() embed_size = all_embs.shape[1] nb words = min(max features, len(word index) + 1) 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, max features): EMBEDDING FILE = '../input/embeddings/paragram 300 s1999/paragram 300 s1999.txt' def get coefs(word, *arr): return word, np.asarray(arr, dtype='float32') embeddings index = [] for o in tqdm(open(EMBEDDING FILE, encoding="utf8", errors='ignore')): **if** len(o) <= 100: continue try: coefs = get coefs(*o.split(" ")) assert len(coefs[1]) == 300 embeddings index.append(coefs) except Exception as e: print(e) embeddings index = dict(embeddings index) 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) + 1) 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 In []: seed everything() glove embeddings = load glove(tokenizer.word index, max features) paragram embeddings = load para(tokenizer.word index, max features) embedding_matrix = np.mean([glove_embeddings, paragram_embeddings], axis=0) np.shape(embedding matrix) **Defining the model** The only thing I changed about the model is the size of the LSTM and GRU. They had 40 hidden units previously and 60 now. I also changed the number of K-Fold splits to 4. In []: splits = list(StratifiedKFold(n splits=4, shuffle=True, random state=10).split(x train, y train)) In []: | class Attention(nn.Module): def init (self, feature dim, step dim, bias=True, **kwargs): super(Attention, self). init (**kwargs) self.supports masking = True self.bias = bias self.feature dim = feature dim self.step dim = step dim self.features dim = 0weight = torch.zeros(feature dim, 1) nn.init.xavier uniform (weight) self.weight = nn.Parameter(weight) if bias: self.b = nn.Parameter(torch.zeros(step_dim)) def forward(self, x, mask=None): feature_dim = self.feature_dim step dim = self.step dim eij = torch.mm(x.contiguous().view(-1, feature dim), self.weight).view(-1, step dim) if self.bias: eij = eij + self.b eij = torch.tanh(eij) a = torch.exp(eij) if mask is not None: a = a * maska = a / torch.sum(a, 1, keepdim=True) + 1e-10 weighted input = x * torch.unsqueeze(a, -1)return torch.sum(weighted input, 1) class SpatialDropout(nn.Dropout2d): def forward(self, x): x = x.unsqueeze(2)# (N, T, 1, K) x = x.permute(0, 3, 2, 1) # (N, K, 1, T)x = super(SpatialDropout, self).forward(x) # (N, K, 1, T), some features are maskedx = x.permute(0, 3, 2, 1) # (N, T, 1, K)x = x.squeeze(2) # (N, T, K)return x In []: class NeuralNet(nn.Module): def __init__(self): super(NeuralNet, self). init () hidden size = 60self.embedding = nn.Embedding(max features, embed size) self.embedding.weight = nn.Parameter(torch.tensor(embedding matrix, dtype=torch.float32)) self.embedding.weight.requires_grad = False self.embedding dropout = SpatialDropout(0.1) self.lstm = nn.LSTM(embed size, hidden size, bidirectional=True, batch first=True) self.gru = nn.GRU(hidden size * 2, hidden size, bidirectional=True, batch first=True) self.lstm attention = Attention(hidden size * 2, maxlen) self.gru attention = Attention(hidden size * 2, maxlen) self.linear = nn.Linear(480, 16) self.relu = nn.ReLU() self.dropout = nn.Dropout(0.1) self.out = nn.Linear(16, 1)def forward(self, x): h = mbedding = self.embedding(x)h embedding = self.embedding dropout(h embedding) h_lstm, _ = self.lstm(h embedding) h_gru, _ = self.gru(h_lstm) h lstm atten = self.lstm attention(h lstm) h gru atten = self.gru attention(h gru) avg pool = torch.mean(h gru, 1) max_pool, _ = torch.max(h gru, 1) conc = torch.cat((h lstm atten, h gru atten, avg pool, max pool), 1) conc = self.relu(self.linear(conc)) conc = self.dropout(conc) out = self.out(conc) return out **Training** Regarding the training procedure, we use Cyclic LR with 5 epochs. I also made a separate function (train model) to train the model because we are going to use it multiple times. In []: batch size = 512 n = 5In []: | class CyclicLR(object): def init (self, optimizer, base lr=1e-3, max lr=6e-3, step size=2000, factor=0.6, min lr=1e-4, mode='triangular', gamma=1., scale fn=None, scale mode='cycle', last batch iteration=-1): if not isinstance(optimizer, torch.optim.Optimizer): raise TypeError('{} is not an Optimizer'.format(type(optimizer). name)) self.optimizer = optimizer if isinstance(base lr, list) or isinstance(base lr, tuple): if len(base lr) != len(optimizer.param groups): raise ValueError("expected {} base_lr, got {}".format(len(optimizer.param_groups), len(base_lr))) self.base lrs = list(base lr) else: self.base_lrs = [base_lr] * len(optimizer.param_groups) if isinstance(max lr, list) or isinstance(max lr, tuple): if len(max lr) != len(optimizer.param groups): raise ValueError("expected {} max lr, got {}".format(len(optimizer.param_groups), len(max_lr))) self.max lrs = list(max lr) self.max_lrs = [max_lr] * len(optimizer.param_groups) self.step size = step size if mode not in ['triangular', 'triangular2', 'exp range'] \ and scale fn is None: raise ValueError('mode is invalid and scale fn is None') self.mode = mode self.gamma = gamma if scale fn is None: if self.mode == 'triangular': self.scale fn = self. triangular scale fn self.scale mode = 'cycle' elif self.mode == 'triangular2': self.scale_fn = self._triangular2_scale_fn self.scale mode = 'cycle' elif self.mode == 'exp range': self.scale fn = self. exp range scale fn self.scale mode = 'iterations' else: self.scale fn = scale fnself.scale mode = scale mode self.batch step(last batch iteration + 1) self.last batch iteration = last batch iteration self.last loss = np.inf self.min lr = min lr self.factor = factor def batch_step(self, batch_iteration=None): if batch iteration is None: batch iteration = self.last batch iteration + 1 self.last_batch_iteration = batch iteration for param_group, lr in zip(self.optimizer.param_groups, self.get_lr()): param group['lr'] = lr def step(self, loss): if loss > self.last_loss: self.base lrs = [max(lr * self.factor, self.min lr) for lr in self.base lrs] self.max lrs = [max(lr * self.factor, self.min lr) for lr in self.max lrs] def _triangular_scale_fn(self, x): return 1. def triangular2 scale fn(self, x): return 1 / (2. ** (x - 1)) def exp range scale fn(self, x): return self.gamma**(x) def get lr(self): step size = float(self.step size) cycle = np.floor(1 + self.last_batch_iteration / (2 * step_size)) x = np.abs(self.last_batch_iteration / step_size - 2 * cycle + 1) lrs = []param lrs = zip(self.optimizer.param groups, self.base lrs, self.max lrs) for param_group, base_lr, max_lr in param_lrs: base_height = $(max_lr - base_lr) * np.maximum(0, (1 - x))$ if self.scale mode == 'cycle': lr = base_lr + base_height * self.scale_fn(cycle) lr = base lr + base height * self.scale fn(self.last batch iteration) lrs.append(lr) return lrs In []: def train_model(model, x_train, y_train, x_val, y_val, validate=True): optimizer = torch.optim.Adam(model.parameters()) step size = 300 scheduler = CyclicLR(optimizer, base lr=0.001, max lr=0.003, step size=step size, mode='exp range', qamma = 0.99994) train = torch.utils.data.TensorDataset(x train, y train) valid = torch.utils.data.TensorDataset(x val, y val) train_loader = torch.utils.data.DataLoader(train, batch_size=batch_size, shuffle=True) valid loader = torch.utils.data.DataLoader(valid, batch size=batch size, shuffle=False) loss fn = torch.nn.BCEWithLogitsLoss(reduction='mean').cuda() best_score = -np.inf for epoch in range(n epochs): start time = time.time() model.train() avg loss = 0. for x batch, y batch in tqdm(train loader, disable=True): y_pred = model(x_batch) scheduler.batch step() loss = loss fn(y pred, y batch) optimizer.zero grad() loss.backward() optimizer.step() avg loss += loss.item() / len(train loader) model.eval() valid preds = np.zeros((x val fold.size(0))) if validate: avg val loss = 0. for i, (x_batch, y_batch) in enumerate(valid_loader): y pred = model(x batch).detach() avg_val_loss += loss_fn(y_pred, y_batch).item() / len(valid_loader) valid preds[i * batch size:(i+1) * batch size] = sigmoid(y pred.cpu().numpy())[:, 0] search_result = threshold_search(y_val.cpu().numpy(), valid_preds) val_f1, val_threshold = search_result['f1'], search_result['threshold'] elapsed_time = time.time() - start_time print('Epoch $\{\}/\{\}\$ \t loss= $\{:.4f\}\$ \t val loss= $\{:.4f\}\$ \t time ={:.2f}s'.format(epoch + 1, n_epochs, avg_loss, avg_val_loss, val_f1, val_threshold, elapsed_time)) else: elapsed time = time.time() - start time print('Epoch {}/{} \t loss={:.4f} \t time={:.2f}s'.format(epoch + 1, n_epochs, avg_loss, elapsed time)) valid_preds = np.zeros((x_val_fold.size(0))) avg val loss = 0. for i, (x batch, y batch) in enumerate(valid loader): y_pred = model(x_batch).detach() avg_val_loss += loss_fn(y_pred, y_batch).item() / len(valid_loader) valid preds[i * batch size:(i+1) * batch size] = sigmoid(y pred.cpu().numpy())[:, 0] print('Validation loss: ', avg val loss) test preds = np.zeros((len(test loader.dataset))) for i, (x_batch,) in enumerate(test_loader): y pred = model(x batch).detach() test_preds[i * batch_size:(i+1) * batch_size] = sigmoid(y_pred.cpu().numpy())[:, 0] test preds local = np.zeros((len(test local loader.dataset))) for i, (x batch,) in enumerate(test local loader): y pred = model(x batch).detach() test_preds_local[i * batch_size:(i+1) * batch_size] = sigmoid(y_pred.cpu().numpy())[:, 0] return valid preds, test preds, test preds local In []: seed = 6017 In []: | x test cuda = torch.tensor(x test, dtype=torch.long).cuda() test = torch.utils.data.TensorDataset(x test cuda) test loader = torch.utils.data.DataLoader(test, batch size=batch size, shuffle=False) x test local cuda = torch.tensor(x test local, dtype=torch.long).cuda() test local = torch.utils.data.TensorDataset(x test local cuda) test_local_loader = torch.utils.data.DataLoader(test_local, batch_size=batch_size, shuffle=False) In []: | train preds = np.zeros(len(train df)) test preds = np.zeros((len(test_df), len(splits))) test preds local = np.zeros((n test, len(splits))) for i, (train idx, valid idx) in enumerate(splits): x train fold = torch.tensor(x train[train idx], dtype=torch.long).cuda() y_train_fold = torch.tensor(y_train[train_idx, np.newaxis], dtype=torch.float32).cuda() x_val_fold = torch.tensor(x_train[valid_idx], dtype=torch.long).cuda() y val fold = torch.tensor(y train[valid idx, np.newaxis], dtype=torch.float32).cuda() train = torch.utils.data.TensorDataset(x train fold, y train fold) valid = torch.utils.data.TensorDataset(x val fold, y val fold) train loader = torch.utils.data.DataLoader(train, batch size=batch size, shuffle=True) valid loader = torch.utils.data.DataLoader(valid, batch size=batch size, shuffle=False) print(f'Fold {i + 1}') seed everything(seed + i) model = NeuralNet() model.cuda() valid_preds_fold, test_preds_fold, test_preds_local_fold = train_model(model, x train fold, y train fold, x_val_fold, y_val_fold, validate=True) train preds[valid idx] = valid preds fold test preds[:, i] = test preds fold test_preds_local[:, i] = test_preds_local_fold **Evaluation** In []: search result = threshold search(y train, train preds) search_result Here we see our very low CV f1 score. But another metric that I have not seen in public kernels yet is the correlation between the test predictions of each fold (in this case the predictions of the local test set). In []: pd.DataFrame(test_preds_local).corr() That is astonishingly low! I am used to seeing correlations of > 99% here. Remember that the model architecture of each of these predictions is exactly the same! The only difference is some of the training data and the seed used to initialize the parameters. Because we have very low correlations it makes sense that, when stacking the predictions of each fold, the score gets much higher. In []: f1 score(y test, test preds local.mean(axis=1) > search result['threshold']) And now we see a score that is about the same as what could be expected on the leaderboard. So the reason why models that have a high CV score often score badly on the leaderboard is that they seem to have a higher correlation between folds than models with a lower CV score. So the challenge we face with neural networks in this competition is finding the perfect tradeoff between CV score and correlation between the predictions of each fold. So when tuning a model, it makes little sense to only track the change in CV score. We have to tune models on a local test set in order to get a valid estimate of how well it will perform on the leaderboard. In []: submission = test df[['qid']].copy() submission['prediction'] = test preds.mean(axis=1) > search result['threshold'] submission.to_csv('submission.csv', index=False) A note on seeds You might have noticed the line declaring the random seed to a cryptic value of 6017 above. That is because I hyperparameter-tuned the random seed. That might sound horrifying but, in my opinion, it makes sense in this competition. The problem when tuning the seed without a local test set is that you are bound to overfit to the public test set which will be exchanged in stage 2. However, if we tune the seed on a cross-validation of local test sets, we do not have this risk. And the seed does make a huge difference. Not only on the public leaderboard but also on the local test set that is close to the size of the test set used in stage 2. I said in my first kernel that evaluating the model multiple times will not be necessary anymore because PyTorch behaves deterministically. But I have to correct that statement: PyTorch does behave deterministically, but that only means that we can run the model with one fixed seed and get the same result. That solves the problem of reproducability. But it does not change impact of the seed on the score. If you change some parameter of the model and have an unlucky seed, you might believe that the change was bad. But it could just have been the seed. The code below wraps the regular K-Fold CV in a K-Fold CV for the local test set. Seeds are selected randomly. Every seed takes about 1 hour to evaluate on my machine (GTX 1080 TI). if enable local test: x train full = np.concatenate([x train, x test local]) y train full = np.concatenate((y train, y test)) x_train_full = x_train y_train_full = y_train seed_splits = list(StratifiedKFold(n_splits=4, shuffle=True, random_state=11).split(x_train, y_tr ain)) seeds = list(np.random.randint(0, 100000, 8)) train scores = [] test_scores = [] for seed i, seed in enumerate(seeds): seed test scores = [] seed_train_scores = [] for i, (train index, val index) in enumerate (seed splits): seed x train = x train full[train index] seed_y_train = y_train_full[train_index] seed x test = x train full[val index] seed_y_test = y_train_full[val index] splits = list(StratifiedKFold(n splits=4, shuffle=True, random state=10).split(seed x tra in, seed y train)) train preds = np.zeros(len(seed x train)) test preds local = np.zeros((len(seed x test), len(splits))) x test local cuda = torch.tensor(seed x test, dtype=torch.long).cuda() test local = torch.utils.data.TensorDataset(x test local cuda) test local loader = torch.utils.data.DataLoader(test local, batch size=batch size, shuffl e=False) print(f'Seed Fold {i + 1}\n') for i, (train idx, valid idx) in enumerate(splits): x train fold = torch.tensor(seed x train[train idx], dtype=torch.long).cuda() y_train_fold = torch.tensor(seed_y_train[train_idx, np.newaxis], dtype=torch.float32) .cuda() x val fold = torch.tensor(seed x train[valid idx], dtype=torch.long).cuda() y val fold = torch.tensor(seed y train[valid idx, np.newaxis], dtype=torch.float32).c uda() train = torch.utils.data.TensorDataset(x train fold, y train fold) valid = torch.utils.data.TensorDataset(x val fold, y val fold) train loader = torch.utils.data.DataLoader(train, batch size=batch size, shuffle=True valid loader = torch.utils.data.DataLoader(valid, batch size=batch size, shuffle=Fals print(f'Fold {i + 1}') seed everything(seed + i) model = NeuralNet() model.cuda() valid_preds_fold, test_preds_fold, test_preds_local_fold = train_model(model, x train fold, y train fold, x val fold, y_val_fold, va lidate=**True**) train preds[valid idx] = valid preds fold test_preds_local[:, i] = test_preds_local_fold train search result = threshold search(seed y train, train preds) seed train scores.append(train search result['f1']) test score = f1 score(seed y test, test preds local.mean(axis=1) > train search result['t hreshold']) seed test scores.append(test score) train score = np.mean(seed train scores) test score = np.mean(seed test scores) train scores.append(train score) test scores.append(test score) print('\ni={} \t seed={}'.format(seed i, seed, test score)) In []: # loading the results from my local machine here # you have to trust me on this ;) $\texttt{test scores} = [0.6894145809793863, \ 0.6904706309470233, \ 0.6905915253597362, \ 0.6908101789878276, \ 0.691033]$ 4464526553, 0.6916507797390641, 0.6903868185698696, 0.6908830283890897] train scores = [0.669555770620476, 0.6708382008438574, 0.6700974173065081, 0.6701065866112219, 0.670477 8141088164, 0.6708436318389969, 0.6705310002773053, 0.6710429366071224] seeds = [42853, 73399, 21152, 58237, 25688, 6017, 29547, 65803] Because seed tuning would exceed the runtime of kernels, I copied the results into this kernel. Now we can evaluate it. In []: eval df = pd.DataFrame() eval df['cv score'] = train scores eval df['local test score'] = test scores eval_df['seed'] = seeds eval df.head() In []: eval df.loc[[eval df['local test score'].idxmax()]] In []: plt.figure(figsize=(14, 14)) sns.violinplot(x='level_0', y=0, data=eval_df[['cv_score', 'local_test_score']].unstack().reset_index plt.title('Distribution of scores for different random seeds') plt.ylabel('') plt.xlabel('') plt.show() In []: eval df['local test score'].describe() The seed really does have a huge influence on the score. And the influence shown here is evalulated on the whole training set, it is surely even stronger on the ~ 50k rows in the public test set. And keep in mind that the statistic where calculated on a small sample of only 8 seeds. **Possible Shortcomings** The validation technique shown in this kernel only evaluates the model on a subset of the data. It would be ideal to wrap the procedure in another K-Fold cross-validation, but that is computationally hardly feasible (except for seed tuning where it is absolutely necessary). The model behaves differently when the data in the local test set is added to the training data (e. g. batches are shuffled differently). Thus, it is impossible to tune the seed on a model that is exactly the same as the one used for submitting. When tuning the model using the shown technique, you will tune it so that it behaves ideally with ~1M training samples. The best architecture for the model also changes when more training data is added (e. g. less need for regularization), so it might again not behave ideally when submitting. **Takeaway** The much higher scores on the leaderboard compared to CV scores are caused by a low correlation between folds of K-Fold CV. When tuning a model, you have to find the best tradeoff between CV score and correlation between folds. The seed is a valid hyperparameter to tune when not tuning it to the public LB. Because the seed has a huge influence on the score, the LB score of top public kernels is not a good indicator on how good the model architecture is. Although this is a kernels-only competition, local compute does matter a lot because you will likely not be able to achieve a good score on the leaderboard without tuning the seed. All of the points above are my current beliefs. I might be wrong about some of them. I'm looking forward to discussion in the comments. Thanks for reading! Note: Version 1 is the one scoring 0.696 on the Leaderboard. The only difference is that <code>enable_local_test</code> is set to <code>False</code>. Runtime is also only 70 Minutes.