Practical Machine Learning Dr. Ashish Tendulkar Department of Computer Science and Engineering Indian Institute of Technology, Madras

Lecture - 29 Boosted Trees

You have already studied how to use tf.estimator API for training linear classifiers. We will use tf.estimator API to train a boosted tree classifiers, which are quite popular in the machine learning community for modeling structured data.

(Refer Slide Time: 00:37)

+ Text & Copy to Drive	Connect 👻 🖌 Editing
Copyright 2019 The TensorFlow Autors.	↑ ↓ ⊕ / ∎ NP
[] Licensed under the Apache License, Version 2.0 (the 'License');	
How to train Boosted Trees models in TensorFlow	
1 Vinne an Tensor Brance Brance Scoole Calabi	
This tutorial is an end-to-end walkthrough of training a Gradient Boosting model using decision trees with the tf.estimator API. Boosted the most popular and effective machine learning approaches for both regression and classification. It is an ensemble technique that combit several (think 10s, 100s or even 1000s) tree models.	Trees models are among ines the predictions from
Boosted Trees models are popular with many machine learning practitioners as they can achieve impressive performance with minimal $J_{\rm PP}$	perparameter tuning.
Load the titanic dataset	
You will be using the titanic dataset, where the (rather morbid) goal is to predict passenger survival, given characteristics such as gender, and	ge, class, etc.
[] fromfuture import absolute_import, division, print_function, unicode_literals	E
import mungy as np import mandas as pd from EPython.display import clear_output from antiplatilis import point as plt	
	All Room

The boosted tree model is an ensemble technique that combines predictions from tens and hundreds of trees. Another good point about boosted tree models is that they achieve the impressive performance on structured data with minimal hyper parameter tuning. Hyper parameter tuning is one of the big problems with neural networks. And we mostly need to use specialized services for hyper parameter tuning in neural networks because there are just to many hyper parameters.

(Refer Slide Time: 01:24)



So, in this exercise we will use titanic dataset where our goal is to predict the probability of passenger surviving. As we have already seen this titanic dataset in the context of logistic regression, we will not spend a lot of time on exploration of titanic dataset.

Let us download the titanic dataset. The titanic dataset already has a training and evaluation split. So, we read training and evaluation data in pandas dataframes and pop the label column from the dataframe to create another dataframe for labels. So, we have dftrain and dfeval as data frames containing features. And y_train and y_eval containing labels corresponding to train and the evaluation dataset.

(Refer Slide Time: 02:56)



Let us install let us install TensorFlow 2.0 and set a random seed to 123. This helps to make sure that across multiple runs we get the same results.

(Refer Slide Time: 03:13)

		ainst the eval set dfava	and v eval	to learn.					
For tra	aining you will use the fo	allowing features:	in, and y_cruit.						
			Feature Name	De	scription				
			sex	Gender of passe	nger				
	set consists of a training set and a ffraia and y_train are the trainin The model is tested against the eval ning you will use the following featu		age	Age of passenge	fr				
			n_siblings_spouses	# siblings and pa	artners aboard				Editin
			parch	# of parents and	children aboard				
		*	fare	Fare passenger p	paid.				
			class	Passenger's clas	is on ship				
			deck	Which deck pass	senger was on				
			embark_town	Which town pass	senger embarked f	om		Contraction of the second	
Ехр	lore the data		alone	If passenger was	s alone			70	
Let's f	first preview some of the	e data and create summ	ary statistics on the	training set.					-
[]	dftrain.head()								

We know that titanic dataset has about nine features that describes each passenger based on their gender, age, their class, their embarkation town and some other features.

(Refer Slide Time: 03:33)



We first convert the features into feature columns. Feature columns work with all TensorFlow estimators and their purpose is to define features for modeling. Gradient boosting estimator can utilize both numeric and categorical features. So, what we will do is, we will convert the categorical features into one hot encoding.

(Refer Slide Time: 04:03)



For every categorical column we first obtain the vocabulary and then we pass the vocabulary to categorical column with vocabulary list. And the output of that is fed into the indicator column, which converts the categorical column or each value in the categorical column into one hot categorical column.

Numerical features are fairly easy to handle. And we simply use numerical columns of feature column to represent each of the numerical features. We combine both categorical features and numerical features into feature_columns.

(Refer Slide Time: 05:19)



Let us look at the dense feature representation of the feature columns. And you can see that all the features are shown over here. Where there are some numerical features and categorical features are converted into one hot encoding as you can see it here. Next we need to create an input function. This will specify how data will be read into our model for both training and inference.

We will use the from_tensor_slices() to create the dataset. Here the dataset is in memory and stored in pandas dataframe. So, after creating the dataset we shuffle the dataset where we set the buffer size to the number of examples. And we cycle through the dataset as many times as we need. Here the n_epoch is set to none.

So, we can cycle through dataset as many times as we need. And we do not really use batching because all the data is in memory. And hence we have dataset.batch operation where

the batch size is the number of examples. So, this input_function returns it to dataset, which will be consumed during training and inference time.

So, we make the input functions for both training and evaluation. At a time of evaluation, we set shuffling to false and number of epochs to one. As the number of epoch is one, this makes sure that we iterate through the evaluation data only once.

(Refer Slide Time: 07:24)

✓ RAM ____ ✓ Editing Copy to Drive set = dataset reneat(r NPTEL Desn't use batching. h(NUM_EXAMPLES) Finalning and evaluation input functions. train_input_fn = make_input_fn(dftrain, y_train) eval input fn = make input fn(dfeval, v_eval, shuffle=False, n e - Train and evaluate the model Below you will do the following steps Initialize the model, specifying the features and hyperparameters.
 Feed the training data to the model using the train_input_fn and train the model using the train function 3. You will assess model performance using the evaluation set-in this example, the dfeval DataFrame. You will verify that the predictions match the labels from the y_eval array. Before training a Boosted Trees model, let's first train a linear classifier (logistic regression model). It is best practice to start with simpler model to establish a ↑↓©¢∎: linear_est = tf.estimator.LinearClassifier(feature_columns) # Train model.
linear_est.train(train_input_fn, max_steps=100) # Evaluation.
result = linear_est.evaluate(eval_input_fn)
clear_output()
print(pd.Series(result)) 0.765152 accuracy accuracy_baseline 9.625 0.832844 auc_precision_recall 0.789631 average_loss global_step 0.478908

Before building a boosted tree model let us build a logistics regression classifier to establish a baseline for this problem. So, we define logistic regression classifier with tf.estimator.LinearClassifier. And supply feature_columns as an argument.

We then train the classifier by specifying the input function and we train for maximum_steps. We evaluate the model with the evaluation_input_fn. And finally, we print the result of the evaluation.

(Refer Slide Time: 08:28)



You can see that we got accuracy of 76 percent. The baseline accuracy 62 percent and we got precision of 70 percent and recall of 64 percent. Now, let us train a boosted tree classifier for boosted trees there are TensorFlow supports boosted tree regressor and boosted tree classifier. Here since we are interested in predicting the survival or non survivals.

Here since our object is to predict whether passenger survives or not we are going to use boosted tree classifier. In BoostedTreeClassifier, we specify the feature_columns and the number of batches. We specify the maximum number of steps the model will stop training once the specified number of trees are built. And we evaluate by supplying the evaluation function.

Let us train the boosted tree classifier and check the output. So, after training the boosted tree classifier, we see that it achieves accuracy of 82 percent, which is 6 percentage points higher than the linear classifier. It achieves precision of 78 percent and recall of 73 percent, which is also higher than the linear classifier.

(Refer Slide Time: 10:09)

	e copy to Driv	3		Disk 🔳			'	Editi	ng
A	accuracy	0.825758							N
0	accuracy_baseline	0.625000							1
0	auc	0.872360							
	auc_precision_recall	0.857325							
	average_loss	0.411853							
	label/mean	0.375000							
	loss	0.411853							
	precision	0.784946							
	prediction/mean	0.382282							
	recall	0.737374							
	global_step	100.000000							
	dtype: float64								
collec	tion, of examples at once. E	arlier, the eval_input_fn is	efined using the entire evaluation set.						
					1	4 0	9	t ii	
0	pred_dicts = list(est.p	redict(eval_input_fn))			^ ·	4 0	9 (i i	
0	pred_dicts = list(est.p probs = pd.Spries([pred	redict(eval_input_fn)) ['probabilities'][1] for	ed in pred_dicts])		^ ·	+ 0	9 4	i i	
0	<pre>pred_dicts = list(est.p probs = pd.Spries([pred probs plot(kind='hist'.</pre>	redict(eval_input_fn)) ['probabilities'][1] for bins=20. title='nredicts	ed in pred_dicts])		^ ·	1	9 (1	
0	<pre>pred_dicts = list(est.p probs = pd.Skries([pred probs.plot(kind='hist', plt.show()</pre>	redict(eval_input_fn)) ['probabilities'][1] for bins=20, title='predicte	red in pred_dicts]) probabilities')		↑ \	4 0	9 4	1	
0	<pre>pred_dicts = list(est.p probs = pd.Spries([pred probs.plot(kind='hist', plt.show()</pre>	<pre>redict(eval_input_fn)) ['probabilities'][1] for bins=20, title='predicte</pre>	<pre>ved in pred_dicts]) probabilities')</pre>		^ •	1 0	9 4		
•	pred_dicts = list(est.p probs = pd.Skries([pred probs.plot(kind='hist', plt.show() pred	redict(eval_input_fn)) ('probabilities'][1] for bins=20, title='predicte icted probabilities	ved in pred_dicts]) probabilities')		^ `	4 0	9		
•	<pre>pred_dicts = list(est.p probs = pd.skries([pred probs.plot(kind='hist', plt.show()</pre>	redict(eval_input_fn)) ('probabilities'][1] for bins=20, title='predicte icted probabilities	ed in pred_dicts]) probabilities')		<u> </u>	↓ (94	1	
•	pred_dicts = list(est.p probs = pd.Spries([pred probs.plot(kind='hist', plt.show() 60	redict(eval_input_fn)) ('probabilities'][1] for bins=20, title='predicte icted probabilities	red in pred_dicts]) probabilities')		^ \	t (9 4		
•	pred_dicts = list(est.probs = pd.5kries([pred probs.plot(kind='hist', plt.show() 50	redict(eval input fn)) ('probabilities')[1] for bins=20, title='predicte cted probabilities	ed in pred_dicts]) probabilities')		^ \	t c	9 4		
•	pred_dicts = list(est.p probs = pd.Spries([pred probs.plot(kind='hist', plt.show() 50 - pred	redict(eval_input_fn)) ('probabilities'[[1] for bins=20, title='predicts icted probabilities	red in pred_dicts]) probabilities')		<u> </u>	¥ (9 4		
•	pred_dicts = list(est.p probs = pd.Spries([pred probs.plct(kind='hist', plt.show() 50 50 50 50	redict(eval input_fn)) ('probabilities'][1] for bins=20, title='predicte icted probabilities	ed in pred_dicts]) probabilities')		^ `	¥ (•		
•	pred_dicts = list(est.p probs = pd.Spries([pred probs.plot(kind='hist', plt.show()	redict(eval_input_fn)) ('probabilities')[i] for bins=20, title='predicte icted probabilities	red in pred_dictu]) probabilities')		• •	¥ (9		
•	pred_dits = list(est.p probs = pd.Spries([pred probs.plot(kind='hist', plt.show()	redict(eval_input_fn)) ['probabilities'][1] fon bins=20, title='predicte icted probabilities	ed in pred_dicts]) probabilities')		^ .	÷.	9		Ig
•	pred_dicts = list(est.probs = pd.dirts = list(est.probs = pd.dirts=(lpred probs.plot(kind='hist', pt.show() pred 00 production pred 00	redict(eval_input_fn)) ['probabilities'][1] for bins=20, title='predicte icted probabilities	red in pred_dicts]) probabilities')		^ ·	÷.	3		
•	pred_dicts = list(est,probs = pd.deries[[prob probs.plct(kind='hist', probs.plct(kind='hist', predb.pl	redict(eval_input_fn)) ['probabilities'][1] for bins-10, titles'predicte cted probabilities	ed in pred_dicts]) probabilities')		<u>^</u>	Ψ (3		
•	pred dicts = list(ast,p probs = pd.3priss([pred probs.plst(ind="hist", pred pl:show] 20 20 20 20 20 20 20 20 20 20 20 20 20	redict(eval_input_fn) [probabilities][1] for bins=20, title='predicte icted probabilities	<pre>eed in pred_dicts]) probabilities')</pre>		<u>^</u>	1 c	3		
•	pred dists = list(ert.pr prebl = pd.sprest(pred prebl.splct(ind="hist", pred prebl.splct(ind="hist", pred pred prebl.splct(ind="hist", pred pred prebl.splct(ind="hist", pred prebl.splct(ind="hist", pred prebl.splct(pred prebl.s	redict(weg_input_fn)) [rotabilities][i] for binsa00, titles/predicte (ted probabilities	<pre>ed in pred_dicts]) probabilities')</pre>		<u>^</u>	1 c	3		
•	pred dicts = list(est,p probs = pd.3priss()pred probs.plot(kind="hist", probs.plot() 00 00 00 00 00 00 00 00 00 00 00 00 00	redict(wa] input fn) [probabilities][1] for bins=20, title="predict cited probabilities	eed in pred_dicts]) probabilities')		<u>^</u>	1 c			

So, we can see that the model the boosted tree model performs better than the logistic regression model in this particular dataset. So, now, we can use the model to make predictions on a passenger from the evaluation set. TensorFlow models are optimized to make predictions in a batch or collection of examples all at once.

So, we give the eval_input_function and this eval_input_function is defined on the entire evaluation set. We look at the probabilities from prediction and plot it. So, you can see that there are a lot of passengers with probability of surviving as only 0.1. And there are few passengers who have got probability equals to 1.

(Refer Slide Time: 11:32)



Let us plot the ROC curve with a ROC with roc_curve from sklearn.metrics package. For ROC curve, we have to specify the actual labels and the probabilities. And it returns the false positive rate and true positive rate. You plot ROC curve with false positive rate on the x axis and y and true positive rate on the y axis. The ROC curve gives us a better idea about tradeoff between two positive rate and false positive rate.