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

Lecture - 18 Underfitting and Overfitting

Welcome to the next module of our course in this module, we will understand Underfitting and Overfitting through the code. You know that overfitting happens when our model has excess capacity to memorize the entire training data. So, what happens if we look at the learning curves we observe that the training error and validation error both reduce to begin with.

After a point the training error goes down, but validation error starts climbing up, if we are seeing that those kinds of learning curves we infer that the model is suffering from overfitting. On the other hand if your model is so simple that it does not have enough capacity to learn the model or to learn the patterns in the training data, then we then our model is suffered from underfitting.

In case of underfitting both training and test error are high. In this lab we will use IMDB movie review dataset. To demonstrate underfitting and overfitting we will initially build a baseline model, then we will build a model to underfit the data and overfit the data. We will first build a baseline model and then we will build a couple of models, so that our model underfits and overfits to the training data. Later in the lab we will demonstrate some of the techniques that help us to overcome the underfitting and overfitting problems.

(Refer Slide Time: 02:02)



Let us first connect to the colab collabrant time, install and import tensorflow 2.0, numpy and matplotlib.

(Refer Slide Time: 02:44)

ο O overfit-and-unde	rfit.ipynb 🗎		GD SHARE
File Edit View Insert	Runtime Tools Help		026596/9609.48
ODE E TENT + Propulations along Requirement alrea Requirement alrea	The form of the second	Wat) (0.0463) (
10			
Download the IMD Rather than using an ember used to demonstrate when	B dataset king as in the previous notebook, here we will multi-hot encode overfitting occurs, and how to fight it.	the sentences. This model will quickly overfit to the training set. It will b	N
Multi-hot-encoding our lists discretional vector that work Juozos + 10000	means turning them into vectors of 0s and 1s. Concretely, this v id be all-zeros except for indices 3 and 5, which would be ones.	could mean for instance turning the sequence $\left[1,\ S\right]$ into a 10.000-	
NPTEScain_data, train	labels), (test_data, test_labels) = keras.datasets.imdb	J.load_data(num_words+NUH_WORDS)	
def multi_hot_sequ	ences(sequences, dimension):		

So, you can see that we have successfully installed tensorflow 2.0 through the print command.

(Refer Slide Time: 03:04)



Next step is to setup the training data. IMDB dataset has movie reviews and each of the movie reviews is tagged with a label 0 or 1, 0 meaning the movie review is a negative and 1 meaning the movie review is positive.

So, it is a problem of identifying whether the movie review is positive or not. IMDB dataset is available in the keras.datasets. So, we do not have to write a lot of code to load the IMDB dataset. We can easily load IMDB datasets and we will be using a multi-hot encoding, where we turn the words into vectors of 0 and 1 in 10000 dimensional vector space.

So, if the word is present at a specific position that we will see one over there. Concretely, if we have a word whose index is 3 and 5 present in the in the in the document it gets converted into a 10000 dimensional vector, where all the entries are 0s except for indices 3 and 5.

(Refer Slide Time: 04:55)



So, let us try to understand that every review is converted into a multi-hot encoding, where we have we have a vector with 10000 positions. So, initially we convert each movie review into a sequence of word indices. So, the IMDB dataset is maintained as a sequence of numbers in the keras, we convert each review into a multi-hot encoding.

That means, if there are two words with indices 3 and 5 present in a review only the position corresponding to those words will have 1 in the list, rest of the other elements will have 0 in their list. So, that is what that is how we convert the review into multi-hot encoding. So, you can as you can see in the code we use a multi-hot encoding with 10000 words, we use the load_data command to a load training and test data along with their labels and then we define a function for converting the sequence into multi-hot sequences.

So, this particular function takes into takes sequence and dimension as argument and it essentially defines or it essentially tries to fill a vector of dimension 10000. We first convert the training data into multi-hot encoding and test data into multi-hot encoding using multi-hot sequences function and we are going to use the number of dimensions to be 10000.

(Refer Slide Time: 07:18)



So, let us run this and convert training and test data into multi-hot encoding. Now you can see that the data is being downloaded and getting loaded into memory and the data is converted into multi-hot encoding here. Let us look at some of the resulting multi-hot vectors.

(Refer Slide Time: 07:38)

141		/	DITING	
0	Journelausing Galla From <u>million://icenses.appdatedat.com/temportausit/icenses/appdatedat.com/temportausit</u>] - 0s Qua/step			
Let's can s	ook at one of the resulting multi-hot vectors. The word indices are sorted by frequency, so it is expected that there are more 1-values near index zero, as we e in this plot:			
0	plt.plot(traim_data(0))			
0	[(matplotlib.lines.Line2D at 0x79644bc7lcc0)]			
	88- 86- 84- 82-			

So, here we have plotted the multi-hot vector for the first training example and we can see that the index where the word is present is 1 everything else is 0. So, you can see that there are a lot of words from the initial index are present in this particular example and fewer words from letter indices are present here.

(Refer Slide Time: 08:07)

	le model using only Dens-	 covy to serve e layers as a baseline, then create smaller and larger versions, and compare the 	m.	/ EDITING
Create a baseli	e model			
<pre>baseline_</pre>	del = kerss.Sequential t_thape is only repu alger.Dens(15, activ ayers.Dens(1, activat del.compile(optimizer- loss=bin metrics=[del.summary()	<pre>i([ref have so that '_ismeary' upris, ref have so that '_ismeary' upris, tions ('book', issue that ''smear', tions 'theo'); ''sden', ''sden', ''sden', ''sden', ''stery_rossentropy'])</pre>		1
[] baseline_M	story = baseline_model	.fititrain_data, train_lambia, epoChasia, harCh_sizewil2, validatin_data(test_data, test_labels), verboseal)		

Now, that data is loaded let us try to build the model. We will first build a basic model where we will use where we will use a neural network with two hidden layers each layer having 16 units each. So, let us look at how this neural network looks like.

(Refer Slide Time: 08:26)



In base model we have neural network which has two layers of 16 units and then there is a single output unit and we have a multi-hot encoded vectors. So, the number of inputs is 10000. So, each of the 10000 inputs are fed into the first unilayer and the second unilayer and

in the output layer which gives us the probability of the review belonging or review being a positive review. So, let us see how to set this up in the code.

So, we build a keras.Sequential model, where we are going to stack layers one over the other. We are going to use two dense layer each with 16 units with activation of relu. The output layer has one unit with sigmoid activation. We use adam as an optimizer and binary cross entropy as a loss because we are trying to solve a binary classification problem and we will track accuracy and binary cross entropy both during the training and this particular command will summarize the model.

(Refer Slide Time: 10:23)

0	baseline_model.summary(3		
θ	Model: "sequential"			
	Layer (type)	Output Shape	Param #	
	dense (Dense)	(None, 16)	160016	
	dense_1 (Dense)	(None, 16)	272	
0.	Total params: 100,101 Trainable params: 100, Non-trainable params: baselime_history = base	305 0 line_model.fit(train_data,		
-0		train_labels, epochsw30, betch_sizew31 validation_da verbosew2)	l, tae(test_dats, text_labels),	

The model summary is a useful way to understand what is going on in the model and to check whether model is setup as per our expectation. So, you can see that there are there are three layers in the model. So, the first hidden layer has 16 units, the second one also has 16 units and then the output layer has one unit.

Each unit from the first hidden layer has 10,001 parameters each because there are 10000 inputs plus a bias uni. This adds up to 160016 parameters. Similarly, we have 272 parameters for the second hidden layer and 17 parameters for the output layer.

So, after setting the model let us fit the model with the training data and validate it on the test data. So, we will train the model with 20 epochs with a batch size of 512 and we are going to

use adam optimizer for training the model and use binary cross entropy as a loss that will be optimized.

(Refer Slide Time: 14:40)

CODE D TEXT + CELL	CELL & COPY TO DRIVE	BAM TO PER PERING
Deseline_history + basel	ine_model.fit(train_data, train_labels, epochas0, batch_sizedil, validation_data=(test_data, test_labels), verbosex))	THE LOCAL
*** flag parsing goes to sto 0068781266816 deprecatio g: ch has the same broadcas validate on 25000 sample	ferr. n.py:323] From /usr/local/lib/python3.6/dist-packages/tu st rule as np.where ts	ensorflow/python/ops/math_grad.py:1250: add_dispatch_support.cl
0.4888 - accuracy: 0.80	075 - binary_crossentropy: 0.4888 - val_loss: 0.3356 - v	al_accuracy: 0.8770 - val_binary_crossentropy: 0.8356
0.2499 - accuracy: 0.93	113 - binary_crossentropy: 0.2499 - val_loss: 0.2851 - val_	al_accuracy: 0.8878 - val_binary_crossentropy: 0.2851
0.1834 - accuracy: 0.93	156 - binary_crossentropy: 0.1834 - val_loss: 0.2903 - va	al_accuracy: 0.8846 - val_binary_crossentropy: 0.2903
0.1490 - accuracy: 0.94	488 - binary_crossentropy: 0.1490 - val_loss: 0.3211 - va	al_accuracy: 0.8743 - val_binary_crossentropy: 0.3211
0.1230 - accuracy: 0.99	92 - binary_crossentropy: 0.1230 - val_loss: 0.3343 - v	al_accuracy: 0.8748 - val_binary_crossentropy: 0.3343
0.1023 - accuracy: 0.96	578 - binary_crossentropy: 0.1023 - val_loss: 0.3627 - va	al_accuracy: 0.8722 - val_binary_crossentropy: 0.3627
0.0058 - accuracy: 0.93	736 - binary_crossentropy: 0.0858 - val_loss: 0.3989 - va	al_accuracy: 0.8668 - val_binary_crossentropy: 0.3989
0.0732 - accuracy: 0.93	791 - binary_crossentropy: 0.0732 - val_loss: 0.4357 - va	al_accuracy: 0.8629 - val_binary_crossentropy: 0.4357
(100599 - accuracy: 0.90	846 - binary_crossentropy: 0.0599 - val_loss: 0.4759 - va	al_accuracy: 0.8628 - val_binary_crossentropy: 0.4759
0517 + accuracy: 0.90	174 - binary_crossentropy: 0.0517 - val_loss: 0.5173 - vo	al_accuracy: 0.8592 - val_binary_crossentropy: 0.5173
MPTEL	Md. blann concentration & API3	al annual a 1967 - us biana, contrastante à 5444

So, we can see that as we train the model the cross entropy loss is reducing, the validation loss is also reducing but after third layer it seems that the validation loss is going up while training loss is still going down and we can see the same thing happening for the cross entropy loss it initially went down then it started going up.

(Refer Slide Time: 15:08)



This points to some kind of an overfitting problem. So, we see that after 20 epochs our accuracy is 1 we got a perfect classifier. In validation, we only get 85 percent accuracy. Our job is to build a model that works well on the future data and not the one that works well under training data. So, there is a problem which we need to fix.

So, one way to one way to address the problem of overfitting is by reducing the number of parameters in the model. In neural network it is easy to reduce the complexity - we can either reduce the number of units in each of the layer or we altogether remove a layer that will result in lesser number of parameters. Another way to reduce the overfitting is by getting more of training data. If training data is not available then we have regularization techniques that we can use. In regularization, we can use either L1 or L2 regularization or a dropout regularization that is used in the context of neural networks.

In case of L1 and L2 regularization we add penalties that are proportional to the weights of the model. In case of L1 regularization, we add penalty that is proportional to the sum of absolute values of the parameter. In case of in case of L2 regularization, we had a penalty that is proportional to the sum of square of the value of the parameters.

In case of dropout we decide to randomly drop certain nodes from the hidden layer or input layer of a neural network. We normally said dropout between 20 percent to 50 percent that means 20 percent to 50 percent nodes will be randomly dropped out in each of the iteration in the neural network training. The dropout is only used only used while training, so in the same

manner even regularizations or dropout are used during training in the test we simply apply the test data on the model.

(Refer Slide Time: 17:51)

<pre>Lefs create a model with less hidden units to compare against the baseline model that we just created for smaller_model = teres.Sequential({ teres.Sequential({ isrs:Jeyr:Dense(*, activations 'relu', isrg(*, isrg(</pre>	Crea	ste a smaller model	
<pre>smaller_model + %eres.Sequential([</pre>	Lets	create a model with less hidden units to compare against the baseline model that we just created	
And train the model using the same data: [] smaller_history + smaller_model.fit(train_deta_ train_labels_	0	<pre>smalls=_model = kerns.Sequential(wars.layers.Dense(k_sclinations'relu', input_shapes(WAR_WORDS,)), kerns.layers.Dense(i, activations'signedd'))) smaller_model.compile(optimizers'adam', loss=%inery_crossentropy', matrians'sicromsty', %inery_crossentropy')) smaller_model.umery() </pre>	
[] smaller_history + smaller_model.fit(train_data, train_labels,	And	train the model using the same data:	
apobražu, hatto,iise512, valiation_data=(test_data, test_labels), verboase1		smalter_hittery = smalter_houst.fri(frain_dats, epoch=0; battot_size=0; validatio_data(test_data, test_labels), validatio_data(test_data, test_labels),	

So, in the first strategy let us create a smaller model here. So, instead of using 16 units we used 4 units, so naturally the number of parameters will go down that we just have 4 units. So, let us run and see how many parameters are there in this.

(Refer Slide Time: 18:09)

0	keras.layers.Dense(4, keras.layers.Dense(1,]) smaller_model.compile(opti loss metr	activations'relu'), activations'signoid') mizers'adam', s'binary_crossentropy', ics=['accuracy', 'binary_	(crossentropy"))	
0	<pre>smaller_model.summary() Model: "sequential_1"</pre>			
	Layer (type) dense_3 (Dense)	Output Shape (None, 4)	Param # 40004	
	dense_4 (Dense) dense_5 (Dense)	(None, 4) (None, 1)	20	
	Total params: 405029 Trainable params: 40,029 Non-trainable params: 0			
And	train the model using the same	data:		
0	amailer_history + smaller_	model.fit(train_data, train_labels, epochs+20,		

So, you can see that the number of parameters have come down at least four folds. So, just by reducing the number of units in each of the layer we got it down to 40000 parameters. Let us train this smaller model and see what happens.

(Refer Slide Time: 18:44)

SWAILER_DISTORY * SMAILER_BOOK	itititrain_sata, train_labels,		
	epochsal0, batch_size=512, validation_data=(test_data, test_labels), verbose=2)		
A Train on 25000 samples, valid	ate on 25000 samples		
Epoch 1/20 25000/25000 - 3s - loss: 0.64 Epoch 2/20	99 - accuracy: 0.7705 - binary_crossentropy: 0.6095	0 - val_loss: 0.5303 - val_accuracy: 0.8313 - val_bina	ry_c
25000/25000 - 3s - loss: 0.44 Epoch 3/20	37 - accuracy: 0.8691 - binary_crossentropy: 0.4437	<pre>' val_loss: 0.4030 - val_accuracy: 0.8675 - val_bina</pre>	ry_ci
25000/25000 - 3s - 1oss: 0.31 Epoch 4/20	32 - accuracy: 0.9007 - binary_crossentropy: 0.323	t - val_loss: 0.3291 - val_accuracy: 0.8795 - val_bina	ry_ci
25000/25000 - 3s - loss: 0.25 Epoch 5/20 25000/25000 - 3s - loss: 0.25	28 - accuracy: 0.9216 - binary_crossentropy: 0.2528 36 - accuracy: 0.9216 - binary_crossentropy: 0.2528	<pre>s - val_loss: 0.2991 - val_accuracy: 0.8848 - val_bina i - val_loss: 0.2977 - val_accuracy: 0.0066 - val_bina</pre>	ry_c
Epoch 6/20 25000/25000 - 3s - 1oss: 0.11	49 - accuracy: 0.9410 - binary crossentropy: 0.1154	<pre>+ val_toss: 0.2826 + val_accuracy: 0.8875 + val_bina</pre>	ry_c
Epoch 7/20 25000/25000 - 3s - loss: 0.10	i15 - accuracy: 0.9513 - binary_crossentropy: 0.161	i - val_loss: 0.2855 - val_accuracy: 0.8848 - val_bina	ry_c
Epoch 8/20 25000/25000 - 3s - loss: 0.14	28 - accuracy: 0.9572 - binary_crossentropy: 0.1429	8 - val_loss: 0.2914 - val_accuracy: 0.8822 - val_bina	ry_c
Epoch 9/20 25000/25000 - 3s - loss: 0.12 Epoch 10/20	69 - accuracy: 0.9627 - binary_crossentropy: 0.1265	0 - val_loss: 0.3029 - val_accuracy: 0.8784 - val_bina	ry_c
25000/25000 - 3s - loss: 0.11 Epoch 11/20	40 - accuracy: 0.9678 - binary_crossentropy; 0.1146	<pre># - val_loss: 0.3083 - val_accuracy: 0.8784 - val_bina</pre>	ry_c
1000/25000 - 3s - loss: 0.10	22 - accuracy: 0.9726 - binary_crossentropy: 0.1023	1 - val_loss: 0.3186 - val_accuracy: 0.8762 - val_bina	ry_c
Eoch 12/20 2000/25000 - 3s - loss: 0.05 100ch 13/20	24 - accuracy: 0.9766 - binary_crossentropy: 0.0924	4 - val_loss: 0.3294 - val_accuracy: 0.8751 - val_bina	ŋ
PEP" 29000/25000 - 3s - loss: 0.05	33 - accuracy: 0.9799 - binary crossentropy: 0.0833	+ val loss: 0.3430 - val accuracy: 0.8731 - val bina	ev e

So, you can see that even this smaller model seems to be overfitting the training accuracy is close to 1 but validation accuracy is quite low. That means, the model is probably memorizing the training data it has not enough capacity to memorize the training data. But going back you can see that the model started overfitting somewhere around, we can look at the validation loss is going down and epoch 7, it started overfitting where validation data validation loss started going up.

If we go back and check when our baseline model started overfitting our baseline model started overfitting right from the third epoch. So, you know by creating smaller model we are able to we are able to delay the overfitting of the model. So, if we use smaller model and if you stop training around fifth epoch we should still be fine.

(Refer Slide Time: 20:09)

	100E E	TEXT	+ CELL	+ CELL	4 00	PY TO DRIVE						~	RAM	1	. /	EDITING	
0	Epoch 25000 Epoch	1 16/28 2/25000 - 1 17/20	3s - lo	s: 0.061	i + accur	acy: 0.987	75 - bi	inary_crossent	opy: 0.06	514 - væl_loss	: 0.3835	- val_accu	recy:	0.8687	- val_b	inery_cr	0
0	25000 Epoch	/25000 -	35 - 10	s: 0.054	4 - accur	acy: 0.996	03 - bi	inary_crossent	opy: 0.85	548 - val_loss	: 0.3996	- val_accu	racy:	0.8682	- val_b	inary_cr	0
	25000	/25000 -	35 - 10	5: 0.0494) - accur	acy: 0.991	16 - b;	inary_crossent	opy: 0.04	490 - val_lost	: 0.4130	• val_accu	racy:	0.8661	- val_b	inary_cr	.0
	Epoch 25888	19/28	11 - 10	5: 0.844	d - accur	acy: 0.991	10 - bi	inary crossent	opy: 0.04	440 - val loss	: 0.4272	- val accu	necv:	0.8654	- val b	inary co	
	Epoch	20/20		7.5													
		/25000 -	35 - 10	5: 0.039;	accur	NGY., U. 994	• b)	inary_crossent	opy: 0.03	192 - Val_lost	: 0,4423	- val_accu	racy:	0.8644	- Val_b	inary_cn	
Cres As ar	ate a b	igger mor	jei create ar	even large m would v	r model, a	sd see how o	quickly	it begins overfitti	g. Next, let'i	's add to this ber	ichmark a	setwork (hat)	has mu	uch more			
Crea As ar capa	ate a b in exerci acity, far bigge	igger mor ise, you can rmore than in_sode1 =	fel create a the probl	even large m would v pels. Segu	r model, a rarrant intial([td See how a	quickly	it begins overfitti	g. Next, let's	's add to this ber	ichmark a	setwork that I	has mi	uch more			
Cres As ar capa	ate a b in exerci acity, far bigge k k k 2)	igger mod ise, you can more than in_model + inras.layer ienas.layer	fel create a the probl s.Dense s.Dense s.Dense	even large en would w dels.Secu 512, acti 512, acti 1, activa	r model a variant intial([utions'n iions'iig	nd See how o du', input du'), sid')	quickly t_shape	it begins overfitti =(twp(uoRoS,)),	g. Next, let's	's add to this ber	ichmark a	setwork that I	has mi	uch more			
Crea As ar capa []	ate a b in exerci- acity, far bigge k k k 2) bigge	igger mot ise, you can more than in sodel = in sodel = in sodel = co in sodel = co	del create a the probl s. Dense s. Dense s. Dense s. Dense s. Dense s. Dense s. Dense	even large en would w dels. Segun 532, acti 532, acti 1, active timizere' pse 'binary tricse['e	r model, a variant: intial([rations'rig idae', '_crossen :curacy',	nd see how o nlu", insut llu", nois") ropy", binery_tro	quickly t_shape	it begins overfits =(%#(JKROS,)), opy*])	g. Next, let's	's add to this ber	ichmark a	setwork that h	has mu	uch more			

If we stop our training after fifth epoch we will not get the overfitting that we see after 20 epochs.

(Refer Slide Time: 20:10)

[7] Zów/JSów - S Loss: 0.0400 - accuracy: 0.0916 - binary_crossentropy: 0.0400 - val_loss: 0.4100 - val_accuracy: 0.0461 - val_binary_crossentropy: 0.0400 - val_loss: 0.4100 - val_accuracy: 0.0461 - val_binary_crossentropy: 0.0400 - val_loss: 0.4100 - val_accuracy: 0.0014 - val_binary_crossentropy: 0.0400 - val_loss: 0.4272 - val_accuracy: 0.0014 - val_binary_crossentropy: 0.0400 - val_loss: 0.4272 - val_accuracy: 0.0014 - val_binary_crossentropy: 0.0012 - val_loss: 0.4272 - val_accuracy: 0.0014 - val_binary_crossentropy: 0.0012 - val_loss: 0.4272 - val_accuracy: 0.0014 - val_binary_crossentropy: 0.0012 - val_loss: 0.4272 - val_eccuracy: 0.0014 - val_binary_crossentropy: 0.0012 - val_loss: 0.4273 - val_eccuracy: 0.0014 - val_binary_crossentropy: 0.0012 - val_loss: 0.4273 - val_eccuracy: 0.0014 - val_binary_crossentropy: 0.0012 - val_loss: 0.4273 - val_eccuracy: 0.0014 - val_binary_crossentropy: 0.0012 - val_loss: 0.4273 - val_eccuracy: 0.0014 - val_binary_crossentropy: 0.0012 - val_loss: 0.4273 - val_eccuracy: 0.0014 - val_binary_crossentropy: 0.0012 - val_loss: 0.4273 - val_eccuracy: 0.0014 - val_binary_crossentropy: 0.0012 - val_loss: 0.4273 - val_eccuracy: 0.0014 - val_binary_crossentropy: 0.0012 - val_eccuracy: 0.0014 - val_binary_crossentropy: 0.0012 - val_binary_cro	<pre>[7] 2500/2500 - 3s - loss: 0.000 - accuracy: 0.0016 - binary_crossentropy: 0.0000 - val_loss: 0.410 - val_accuracy: 0.0011 - val_binary_cro topon 19/20 2500/2500 - 3s - loss: 0.0002 - accuracy: 0.0000 - binary_crossentropy: 0.0000 - val_loss: 0.4272 - val_eccuracy: 0.0004 - val_binary_cro 2000/25000 - 3s - loss: 0.0002 - accuracy: 0.0000 - binary_crossentropy: 0.00002 - val_loss: 0.4272 - val_eccuracy: 0.0004 - val_binary_cro 2000/25000 - 3s - loss: 0.0002 - accuracy: 0.0000 - binary_crossentropy: 0.00002 - val_loss: 0.4272 - val_eccuracy: 0.0004 - val_binary_cro 2000/25000 - 3s - loss: 0.0002 - accuracy: 0.0000 - binary_crossentropy: 0.0000 - val_loss: 0.4272 - val_eccuracy: 0.0004 - val_binary_cro C Create a bigger model As an exercise, you can create an even larger model and see how quickly it begins overfitting Next lefts add to this benchmark a network that has much more capacity for more than the problem would warrant: D bigger_model - kerss.model.sequential({</pre>		DOE 🖪 TEXT 🛊 CELL 🛊 CELL 🎄 COPY TO DRIVE	/ EDITING
Create a bigger model As an exercise, you can create an even larger model, and see how quickly it begins overfitting. Next, let's add to this benchmark a network that has much more capacity, far more than the problem would warrant: bigger_model = ker'ss.model.is.Sequential({ ter'ss.layers.benes(lat.iseuential({ ter'ss.l	Create a bigger model As an exercise, you can create an even larger model, and see how quickly it begins overfitting. Next, lefs add to this benchmark a network that has much more capacity, far more than the problem would warrant: Isgrer_model = kerss.model: Sequential({	[7]	25000/25000 - 3s - los: 0.0400 - accuracy: 0.0918 - binary_crossentropy: 0.0400 - val_tos: 0.4110 - val_eccuracy: 0.8661 - Epoch 30/20 25000/25000 - 3s - loss: 0.0440 - accuracy: 0.0930 - binary_crossentropy: 0.0440 - val_loss: 0.4272 - val_eccuracy: 0.8654 - Epoch 20/20 25000/25000 - 3s - loss: 0.0392 - accuracy: 0.0980 - binary_crossentropy: 0.0392 - val_loss: 0.4423 - val_eccuracy: 0.8644 -	val_binary_cros val_binary_cros val_binary_cros
As an exercise, you can create an even larger model, and see how quickly it begins overfitting. Next, let's add to this benchmark a network that has much more capacity, far more than the problem would warrant: bigger_model = keres.model.isevential({ terms.layers.benes(la, latticetions", "about_input_thatper(NAP_URDS,)), terms.layers.benes(l, activations", "about_input.terms)) bigger_model.summer()	As an exercise, you can create an even larger model, and see how quickly it begins overfitting. Next, lefs add to this benchmark a network that has much more capacity, far more than the problem would warrant: bigger_model = kerrs.woodin.fequential({ kerss.layers.Dense(18), inclustions "kla", input_thapse(NUM_LOROS,)), kerss.layers.Dense(1, activations "kla", input_thapse(NUM_LOROS,)), kerss.layers.Dense(1, activations "kla", input_thapse(NUM_LOROS,)), kerss.layers.Dense(1, activations "kla", input_thapse(NUM_LOROS,)), bigger_model.comple(optimizer 'stats', 'stats', 'prot_thapse(NUM_LOROS,)), bigger_model.comple(optimizer 'stats', 'stats', 'prot_thapse(NUM_LOROS,)), bigger_model.sumery() And, again, train the model using the same data:	Cre	· e a bigger model	
<pre>bigger_model = kerss.modelx.lsquartial({ kerss.lsquart.squartial({ kerss.lsquart.squartial({ kerss.lsquart.squartial({ kerss.lsquart.squartial({ kerss.lsquart.squartial({ kerss.lsquartial({ kerss.lsquartial(</pre>	<pre>bigger_model = kerss.models.lequential({ kerss.lupers.dense(BE, lictivation* value*, inspec=(NUPLUORDS,)), kerss.lupers.dense(1, activation* value*, kerss.lupers.dense(1, activation* value*, ligner_model.compl(cottisizer* kada*,</pre>	As a capa	exercise, you can create an even larger model, and see how quickly it begins overfitting. Next, let's add to this benchmark a network that has much more ony, far more than the problem would warrant:	
bigger_model.compile(optimizer='adam', loss-binary_crossentropy', metrics('accuracy','binary_crossentropy']) bigger_model.summary()	bigger_model.compile(motivitar-idade')	0	<pre>bigger_model = kerss.models.Sequential({ kers.lyser.Genes(32, jetistions_relut, input_shape=(NUM_LORDS,)), kerse.lyser.Genes(32, settistion=relut), kerse.lyser.Genes(32, settistion=relut), kerse.lyser.Genes(1, settistion='ifpoid')</pre>	1
bigger_model.summary()	bigger_model.sumery() And, spain, train the model using the same data:		bigger_model.compile(optimizer='adam', less=binary_crossentropy', metrics=['crossentropy'] binary_crossentropy'])	
	And, again, train the model using the same data:			
Alan Mare - Mare and Alafada day had bada		And	bigger, model using the same data:	

Let us go to the other extreme and create a bigger model and see how fast it overfits. We create a bigger model in neural network we can simply add more units in each of the layer or add more layer itself. So, that way we will have more capacity in the model or more parameters in the model and since we have more parameters in the model, if you do not have

enough training data the model tries to memorize the training data and gives us perfect output on our training set but it did not perform well on unseen data.

0	bigger_model.summary() karss.layers.Dense(52 karas.layers.Dense(52 karas.layers.Dense(52 karas.layers.Dense(1,)) bigger_model.compile(optim loss estri bigger_model.summary()	<pre>s.Sequential([, setivations'relu', inp , setivations'relu'), activations'signolo') vizers'adam', 'binary_crossentropy', cse['accuracy','binary_c'</pre>	vt_shape=(NUP_NCHCS,)), rossentropy'])	
0	Model: "sequential_2"			
	Layer (type) dense_6 (Dense)	Output Shape (None, 512)	Param #	
	dense_7 (Dense)	(None, 512)	262656	
	dense_8 (Dense) Total params: \$,383,681 Trainable params: \$,383, Non-trainable params: 0	(None, 1)	513	
And	again, train the model using the	same data:		

(Refer Slide Time: 20:58)

So, let us compile this model and you can see that now we have far more number of parameters. So, there are close to 5.3 million parameters as against 160k parameters that we saw in our baseline model.

(Refer Slide Time: 21:12)

0	bigger_history + bigger_model.fit	(train_data, train_labels, epochsz0, batch_sizes512, validatio_data+(test_data, test_labels), verbose=2)	
	Train on 15000 samples, validat Epoch 1/20 25000/15000 - 16s - loss: 0.349 Epoch 1/20 25000/15000 - 17s - loss: 0.452 Epoch 1/20 25000/15000 - 17s - loss: 0.452 Epoch 1/20 25000/15000 - 16s - loss: 0.407 Epoch 1/20 15000/15000 - 16s - loss: 2.199 Epoch 1/2000 - 16s - loss: 1.366 Epoch 1/2000 - 16s - loss: 1.366 Epoch 1/2000 - 16s - loss: 1.366	<pre>x on 25000 samples i - accuracy: 0.8507 - binary_crossentropy: 0.340 6 - accuracy: 0.9452 - binary_crossentropy: 0.140 4 - accuracy: 0.9838 - binary_crossentropy: 0.052 8 - accuracy: 0.9987 - binary_crossentropy: 0.007 5e-04 - accuracy: 1.0000 - binary_crossentropy: 2 5e-04 - accuracy: 1.0000 - binary_crossentropy: 2 5e-04 - accuracy: 1.0000 - binary_crossentropy: 2 5e-04 - accuracy: 1.0000 - binary_crossentropy: 1</pre>	3 - val_loss: 0.1993 - val_accuracy: 0.8772 - val_binary_cr 6 - val_loss: 0.3335 - val_accuracy: 0.8726 - val_binary_cr 4 - val_loss: 0.4479 - val_accuracy: 0.8873 - val_binary_cr 8 - val_loss: 0.5914 - val_accuracy: 0.8876 - val_binary_cr 2015e-04 - val_loss: 0.6679 - val_accuracy: 0.8689 - val_b 2099e-04 - val_loss: 0.7101 - val_accuracy: 0.8692 - val_b 2.3656-04 - val_loss: 0.775 - val_accuracy: 0.8692 - val_b
_	¥		_

Let us try to train the model for 20 epochs with the same batch size of 512 you can see that after 4 epochs. So, the model right away started overfitting. So, you can see that the validation loss is increasing right after the second epoch. So, in the first epoch it was 0.29 it went up to 0.33 and then it went it kept going up all the way and at the same time you can see that the cross entropy loss or the training loss is coming down training accuracy is going up, but the validation accuracy is reducing. These are the signs that this is probably the model with quite a large capacity or excess capacity and it is proven to overfitting pretty quickly in the training cycle.

(Refer Slide Time: 22:25)



So, we can see that after training the model for 20th for 20 epochs you can see that the training loss is very small close to 0. But the validation loss is quite high and validation loss never reduced rather it started growing up right from the second epoch. So, this points to the fact that the model is overfitting pretty fast in the bigger model and we will compare how the size of the model affects the possibility of overfitting.

(Refer Slide Time: 22:58)



So, you can see that there are there are six lines in this particular plot. So, the solid lines are for the training loss and dotted lines are for the validation loss. So, all of them we will have will have the same number of iterations in one epoch and on the y axis we have got a binary cross entropy loss.

So, you can see that the green line corresponds to the bigger model the blue line corresponds to the baseline model whereas, the orange line corresponds to the smaller model. We will have to keep an eye on the dotted lines because those give us validation loss and we infer that there is a overfitting when validation loss starts going up while training loss keeps going down. So, for the smaller model, the validation and training both the losses are going down and around eighth iteration the validation loss started climbing up.

So, we can say that the smaller model started overfitting after 8 epochs. The baseline model started overfitting around third epoch. Whereas, in case of bigger model it overfit almost instantly right from the first iteration, training loss was coming down, but validation loss never came up. Hence, we see that the smaller the model the longer it will take to overfit.

(Refer Slide Time: 25:02)



Now, that we have visualized the overfitting of the model, how can we prevent overfitting? So, there are three strategies - (1) you can get more data overfitting can be prevented. But getting more data is not always an option because getting more data is can be costly sometimes. (2) We have is to add regularization or a penalty term to the loss function. So, there are multiple ways of performing regularization couple of them are based on L1 and L2 penalty.

So, L1 regularization adds penalty proportional to the sum of absolute values of the weight, whereas L2 regularization adds penalty proportional to sum of square of the weights of the parameters.

(Refer Slide Time: 26:03)



Let us see how to add L1 and L2 regularization to the neural network through tf.keras.regularizers (Refer Time: 26:32). So, along with each of the layer dense we add along with each of the dense layer while defining the model we add a kernel_regularizer argument. So, kernel_regularizer we can either use 11 or 12. So, here we are using 12 and we are also specifying what is the regularization rate.

(Refer Slide Time: 26:57)



The regularization term for the loss function is added. Here, the regularization penalty is 0.001 and we are using 12 regularizer here. So, we use 12 regularizers in both layers in order

to prevent overfitting. So, this is how you can define your 12 regularizer and then we can compile the model and train it you can also add 11 regularization in the similar manner.

(Refer Slide Time: 27:39)

	ODE 🖪 TEXT 🛊 CELL 🛊 CELL 💩 COPY TO DRIVE	V RAM MM - / EDITIN
0	validation_data+(test_data, test_labels), +erbota+2)	
Θ	Train on 25000 samples, validate on 25000 samples	
	Epoch 1/20 Menne / Fenne 14 - Jane A 5275 - Account A 2111 - Manne contractore	A 4661 - only langer & 2008 - only appropriate & 2008 - only blance -
	Each 1/38	strast - AstToost status - AstTococarl. status - AstTotus.AT
	25000/25000 - 3s - loss: 0.5012 - accuracy: 0.9063 - binary crossentropy:	0.2573 - val loss: 0.3332 - val accuracy: 0.8878 - val binary (
	Epoch 3/20	
	25000/25000 - 3s - loss: 0.2483 - accuracy: 0.9296 - binary_crossentropy:	0.1979 - val_loss: 0.3362 - val_accuracy: 0.8862 - val_binary_(
	Epoch 4/20	
	25000/25000 - 3s - loss: 0.2257 - accuracy: 0.9402 - binary_crossentropy:	0.1714 - val_loss: 0.1541 - val_accurecy: 0.8792 - val_binary_r
	Epoch 5/20	
	25000/25000 - 3s - loss: 0.2103 - accuracy: 0.9484 - binary_crossentropy:	0.1534 - val_loss: 0.3565 - val_accuracy: 0.8767 - val_binary_s
	16000/15000 - Is - Tess: 0.1076 - accuracy: 0.0515 - Minary consideration	A TIRE , wal loss: A 1858 , wal accuracy: A 8765 , wal binary -
	Each 7/28	artist , definite a sets , definited, artists , definited.
	25000/25000 - 3s - loss: 0.1895 - accuracy: 0.9563 - binary crossentropy:	0.1290 - val_loss: 0.4042 - val_accuracy: 0.8715 - val_binary_c
	Epoch 8/28	
	25000/25000 - 3s - loss: 0.1826 - accuracy: 0.9585 - binary_crossentropy:	8.1207 - val_loss: 0.4190 - val_accuracy: 0.8702 - val_binary_t
	Epoch 9/20	
	25000/25000 - 3s - loss: 0.1767 - accuracy: 0.9607 - binary_crossentropy:	0.1134 - val_loss: 0.4376 - val_accuracy: 0.8646 - val_binary_r
	Epoch 10/20	a size of a second size of a contract of and the second
	25000/25000 - 35 - 1055: 0.1714 - accuracy: 0.9037 - 01nary_crossentropy: Easth 11/10	0.1073 - Val_1055: 0.4473 - Val_accuracy: 0.8061 - Val_Dinary_
	Stana/Stana - 14 - Joss: B 1661 - accuracy: B 9667 - binary crossanteony:	A 1811 - val loss: A 4517 - val accuracy: A 8545 - val binary /
	Epoch 12/20	the second
a	25680/25888 - 3s - Joss: 0.1628 - accuracy: 0.9659 - binary_crossentropy:	0.0968 - val_loss: 0.4789 - val_accuracy: 0.8619 - val_binary_
11	Alech 13/20	
19	00/25008 - 3s - loss: 0.1578 - accuracy: 0.9705 - binary_crossentropy:	0.0912 - val_loss: 0.4912 + val_accuracy: 0.8614 - val_binary_t
N	Joch 14/20	
NPT	10000/25000 - 3s - loss: 0.1552 - accuracy: 0.9706 - binary_crossentropy:	0.0880 - val_loss: 0.5098 - val_accuracy: 0.8574 - val_binary_c
	Epoch 15/20	
	25000/25000 - 3t - 10th: 0.1548 - accuracy: 0.9706 - Dinary crossentrody:	0.0071 - val 1015: 0.5298 - val accuracy: 0.8589 - val binary (

(Refer Slide Time: 27:56)

		A compute course	1
		Otok Disk	100.000
1	D 25000/25000 - 3s - 10ss: 0.1548	- accuracy: 0.9706 - binary_crossentropy: 0.0871 - val_loss: 0.5298 - val_accuracy: 0.8589 -	val_binary_cro
0	Epoch 16/20		
	25000/25000 - 3s - loss: 0.1526	- accuracy: 0.9788 - binary_crossentropy: 0.0835 - val_loss: 0.5355 - val_accuracy: 0.8565 -	val_binary_cro
	25000/25000 - 3s - 10ss: 0.1486	- accuracy: 0.9731 - binary crossentrony: 0.8794 - val loss: 0.5464 - val accuracy: 0.8571 -	val binary cro
	Epoch 18/20		
	25000/25000 - 3s - loss: 0.1469	- accuracy: 0.9736 - binary_crossentropy: 0.0774 - val_loss: 0.5855 - val_accuracy: 0.8529 -	val_binary_cro
	Epoch 19/20 255552/255553 - 34 - 10441 - 0 1453	- accuracy: a 0716 - binary constantions: a 0701 - usl loss: 0 0701 - usl accuracy: 0 0017 -	ual bicard con
	Epoch 20/20 Epoch 20/20	• accuracy: 0.9/10 • Dinary_crossencropy: 0.0/91 • Vai_1055: 0.5/05 • Vai_eccuracy: 0.8557 •	var_onary_cro
	25000/25000 - 3s - loss: 0.1409	- accuracy: 0.9765 - binary_crossentropy: 0.8699 - val_loss: 0.5764 - val_accuracy: 0.8537 -	val_binary_cro
			and course
-			
1	2(0.001) means that every coefficient in th	e weight matrix of the layer will add 0.001 * weight_coefficient_value**2 to the total loss of the network. No	te
1	2(0.001) means that every coefficient in th hat because this penalty is only added at trai	e weight matrix of the layer will add 0.001 * weight_coefficient_value**2 to the total loss of the network. No uning time, the loss for this network will be much higher at training than at test time.	te
1 0 1	2(0.001) means that every coefficient in th hat because this penalty is only added at trai tere's the impact of our L2 regularization per	we weight matrix of the layer will add 0.001 * weight_coefficient_vvlue**2 to the total loss of the network. No uning time, the loss for this network will be much higher at training than at test time. nally,	te
1 8	2(0.001) means that every coefficient in th hat because this penalty is only added at trai kere's the impact of our L2 regularization per	e weight matrix of the layer will add 0.005 * weight_coefficient_value**2 to the total loss of the network. No ning time, the loss for this network will be much higher at training than at test time. nahy.	te
1 e H	2(0.001) means that every coefficient in th hat because this penalty is only added at trai tere's the impact of our L2 regularization per] plot_history([('baseline', baselji	ee weight mathic of the layer will add 0.005 * weight_coefficient_value**2 to the total loss of the network. No ning time, the loss for this network will be much higher at training than at test time. nally: ge_history),	te
1 8 H	2(8, 601) means that every coefficient in th hat because this penalty is only added at tra- iere's the impact of our L2 regularization per } plot_history(((taskline', baskli ('12', 12, mode_hist	we weight matrix of the layer will add 0.001 * weight_coefficient_value**2 to the total loss of the network. No uning time, the loss for this network will be much higher at training than at test time. nally: lag_history), tery)]	te
1 8 H	2(0.001) means that every coefficient in the because this penalty is only added at tra- tere's the impact of our L2 regularization per plot_history([('taseline', baseline') '12', 12_mode_hist	ex weight matrix of the layer will add 0.001 * weight_coefficient_value**2 to the total loss of the network. No using time, the loss for this network will be much higher at training than at test time. nally: ae,history), fory)])	te
1 8 H	<pre>2(0.001) means that every coefficient in th hat because this penalty is only added at trai evers the impact of our L2 regularization per) plot_history({('tessiline', bessili ('12', 12_mosel_hist 's you can see, the L2 regularized model has</pre>	en weight matrix of the layer will add 0.005 * weight_coefficient_value**2 to the total loss of the network. No aning time, the loss for this network will be much higher at training than at test time. nally: a_history), tory)) become much more resistant to overfitting than the baseline model, even though both models have the same num	te ber
1 8 H	2(0.001) means that every coefficient in the because this penalty is only added at tra- tere's the impact of our L2 regularization per plot_history([('baseline', baselin ('l2', l2_mode_hist d parameters.	we weight matrix of the layer will add 0.001 * weight_coefficient_vvlue**2 to the total loss of the network. No uning time, the loss for this network will be much higher at training than at test time. nally: se_fistory), site one much more resistant to overfitting than the baseline model, even though both models have the same num	te ber
1 8 H	2(0.001) means that every coefficient in th the because this penalty is only added at tra- tere's the impact of our L2 regularization per) plot_history([('baseline', baseline', ('l2', l2_model_hist is you can see, the L2 regularized model has (parameters.	ie weight matrix of the layer will add 0.001 * weight_coefficient_value**2 to the total loss of the network. No aning time, the loss for this network will be much higher at training than at test time. nally: se_fistory), sitecome much more resistant to overfitting than the baseline model, even though both models have the same num	te
1 8 H	2(0.491) means that every coefficient in th hat because this penalty is only added at trai evers the impact of our L2 regularization per) plot_history([('taskline', baskli ('12', 12,molei_hist d parameters.	ex weight matrix of the layer will add 0.001 * weight_coefficient_value**2 to the total loss of the network. No uning time, the loss for this network will be much higher at training than at test time. nally: ae_history), etry)]) s become much more resistant to overfitting than the baseline model, even though both models have the same num	ber
1 8 H	2(8.601) means that every coefficient in th the because this penalty is only added at tra- tere's the impact of our L2 regularization per) plot_histery([('tasellee', beall' ('ta', l2,model_hist ('parmeters.	ie weight matrix of the layer will add 0.001 * weight_coefficient_value**2 to the total loss of the network. No uning time, the loss for this network will be much higher at training than at test time. nally: aa_history), dag_history), become much more resistant to overfitting than the baseline model, even though both models have the same num	ber
	2(0.001) means that every coefficient in the because this penalty is only added at tra- tere's the impact of our L2 regularization per plot_history([['baseline', baseline', ['baseline', baseline', baseline', ['baseline', baseline', baseline', ['baseline', baseline', baseline', ['baseline', baseline', baseline', ['baseline', baseline', baseline', ['baseline', baseline', baseline', baseline', ['baseline', baseline', baseline', baseline', ['baseline', baseline', baseline', baseline', ['baseline', baseline', baseline', baseline', baseline', ['baseline', baseline', baseline', baseline', baseline', ['baseline', baseline', bas	ieweight matrix of the layer will add 0.001 * weight_coefficient_value**2 to the total loss of the network. No aning time, the loss for this network will be much higher at training than at test time. nally: se_fistory), s become much more resistant to overfitting than the baseline model, even though both models have the same num st commonly used regularization techniques for neural networks, developed by Hinton and his students at the Unive	ber

Now, let us compare how the training of regularized model compares with the baseline model.

(Refer Slide Time: 28:05)



Let us plot the loss with respect to the epoch. So, we can see that the blue model corresponds to the baseline model, whereas the orange model corresponds to the baseline model with L2 regularization. You can see that the blue line which is which is a baseline model started overfitting around third epoch, whereas the regularized model took slightly longer to start overfitting. We can see that the regularized model is resistant to overfitting for some more epochs than the unregularized model.

(Refer Slide Time: 28:43)



You can see the 12 regularized model has become much more resistant to overfitting than the baseline model, even though both models have the same number of parameters. The other way of adding regularization or other way of regularizing neural network is by adding dropouts. This results in randomly dropping out a number of output features in the layer during training, dropout is not applied at the test time. It is very important to note that at the test time, we do not drop any units, instead the layers output values are scaled down by a factor equal to a dropout rate.

So, let us see how to use dropout in the context of neural networks, keras has a dropout layer we can use a dropout layer with dropout rate as it is parameter.

(Refer Slide Time: 29:24)



So, you can see that this particular dropout will be applied to a layer that is specified just before that. So, this dropout will be applied to the first hidden layer and then this dropout of point five will be applied to the next hidden layer and after specifying the dropout we can again retrain the model and compare its performance with the original baseline model.

(Refer Slide Time: 30:08)



So, let us compare the baseline model with the with the dropout added. So, now compare this also with the L2 regularization, we can see that after adding L2 regularization the model marginally improved its resistance to overfitting. But after applying dropout we can see that there is a substantial improvement in the resistance to the dropout and the model overfits after few more epochs than the original model.

So, to recap we what we did is we actually build a model on IMDB dataset, the model that we build a baseline model on IMDB dataset then we build a very big model on IMDB dataset that was overfitting. In order to prevent overfitting we first reduce the capacity of the model by building a smaller model and we also added regularization to it, we studied L1 and L2 regularization and then dropout regularization. So, these are the common things that are done to prevent overfitting in the neural network.

The first strategy is to get the training data if its possible. If training data more training data is not possible, we can reduce the capacity of the model or add L1 or L2 regularization or dropout. Dropout is the most commonly used regularization technique in neural network. Apart from that we can also use data augmentation to generate more training data from the available data and batch normalization to prevent overfitting. Hope you understood underfitting and overfitting through this coding exercise and you had fun learning these concepts.