

Lecture 3A

Lightning Introduction to Keras/TF

Training a DL Model for a Structured Data Problem



15.S04: Hands-on Deep Learning
Spring 2024
Farias, Ramakrishnan

(Recap) Summary of overall training flow

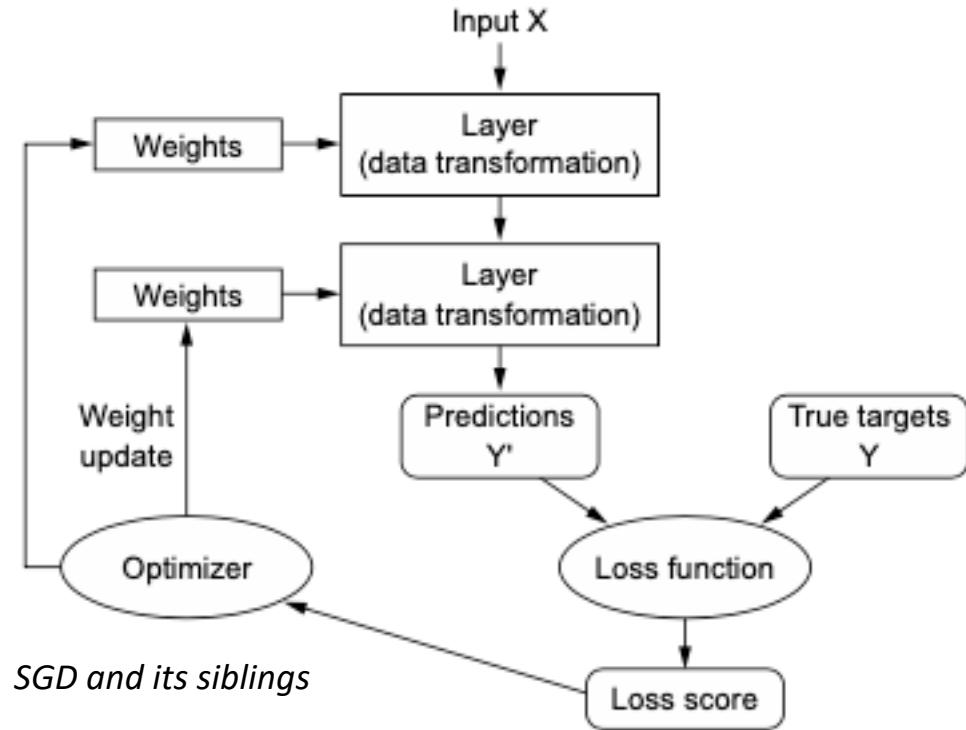


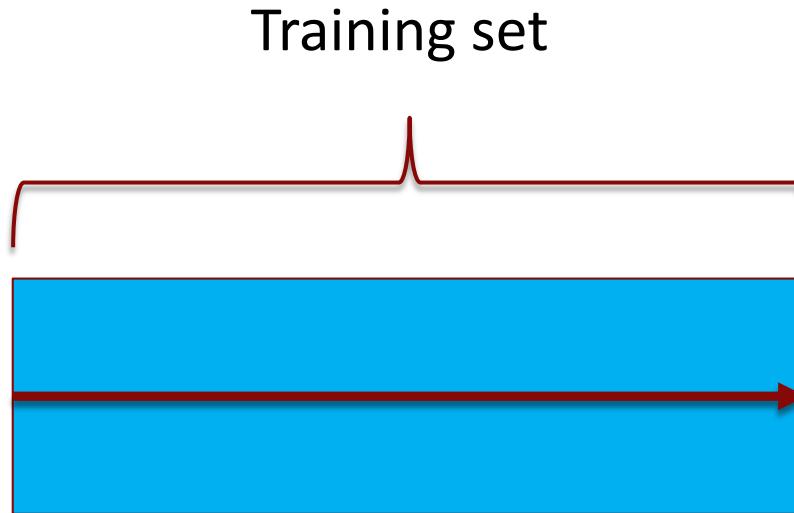
Figure 2.26 Relationship between the network, layers, loss function, and optimizer

(Recap) Gradient Descent vs Stochastic Gradient Descent

- At each iteration, use **all** data points to calculate the gradient of the loss function
- At each iteration, **randomly choose just a few** of the data points and use only these to compute the gradient of the loss function

Epochs and Batches

What is an epoch?



An epoch is one **pass** through the full training set.

But this plays out differently for Gradient Descent vs *Stochastic* Gradient Descent.

An epoch in Gradient Descent

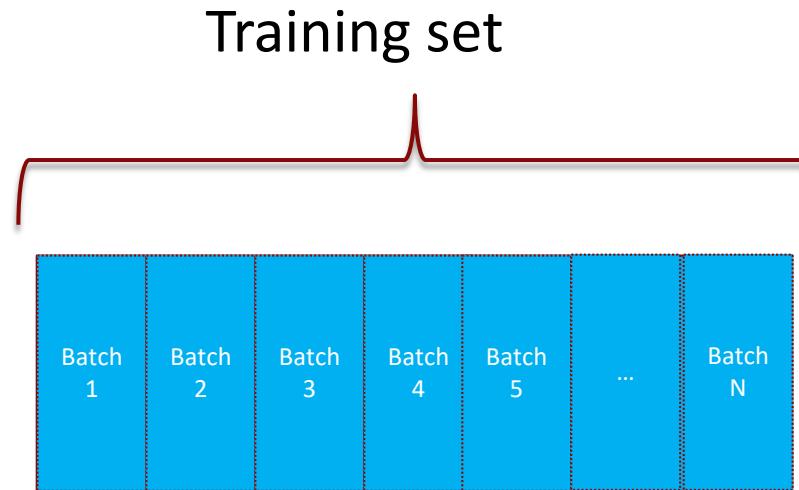


- We run every training sample through the network to get the predictions
- We calculate the gradient of the loss
- We update the parameters

$$w \leftarrow w - \alpha \frac{d\text{Loss}(w)}{dw}$$

This is done just **once** at the end of the epoch

An epoch in Stochastic Gradient Descent

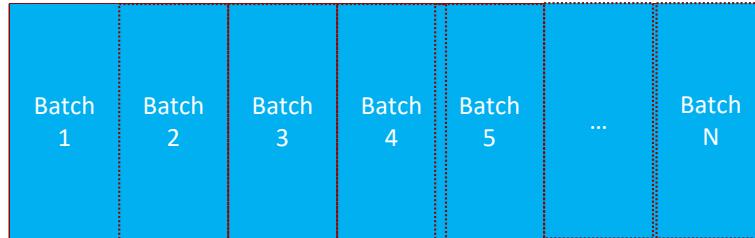


*But when we do Stochastic Gradient Descent (SGD), we process the data in **minibatches***, one after the other*

*we will refer to minibatches as batches from now on for simplicity

An epoch in Stochastic Gradient Descent

Training set



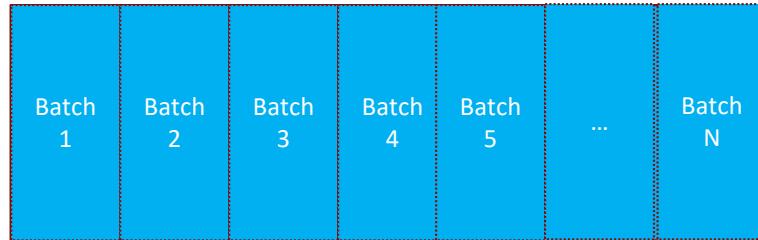
For each batch:

- We run the training samples in that batch through the network to get predictions
- We calculate the gradient of the loss
- We update the parameters

$$w \leftarrow w - \alpha \frac{d\text{Loss}(w)}{dw}$$

An epoch in Stochastic Gradient Descent

Training set



$$w \leftarrow w - \alpha \frac{d\text{Loss}(w)}{dw}$$

$$w \leftarrow w - \alpha \frac{d\text{Loss}(w)}{dw}$$

$$w \leftarrow w - \alpha \frac{d\text{Loss}(w)}{dw}$$

...

$$w \leftarrow w - \alpha \frac{d\text{Loss}(w)}{dw}$$

How many batches in an epoch when we do SGD?

of batches in one epoch = (Training set size / Batch size) rounded up

For Neural Heart Disease Model:

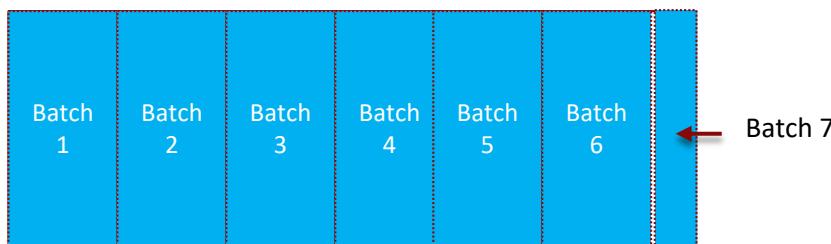
Training set size = 194

Batch size = 32

of batches in one epoch = (194/32) rounded up = 7

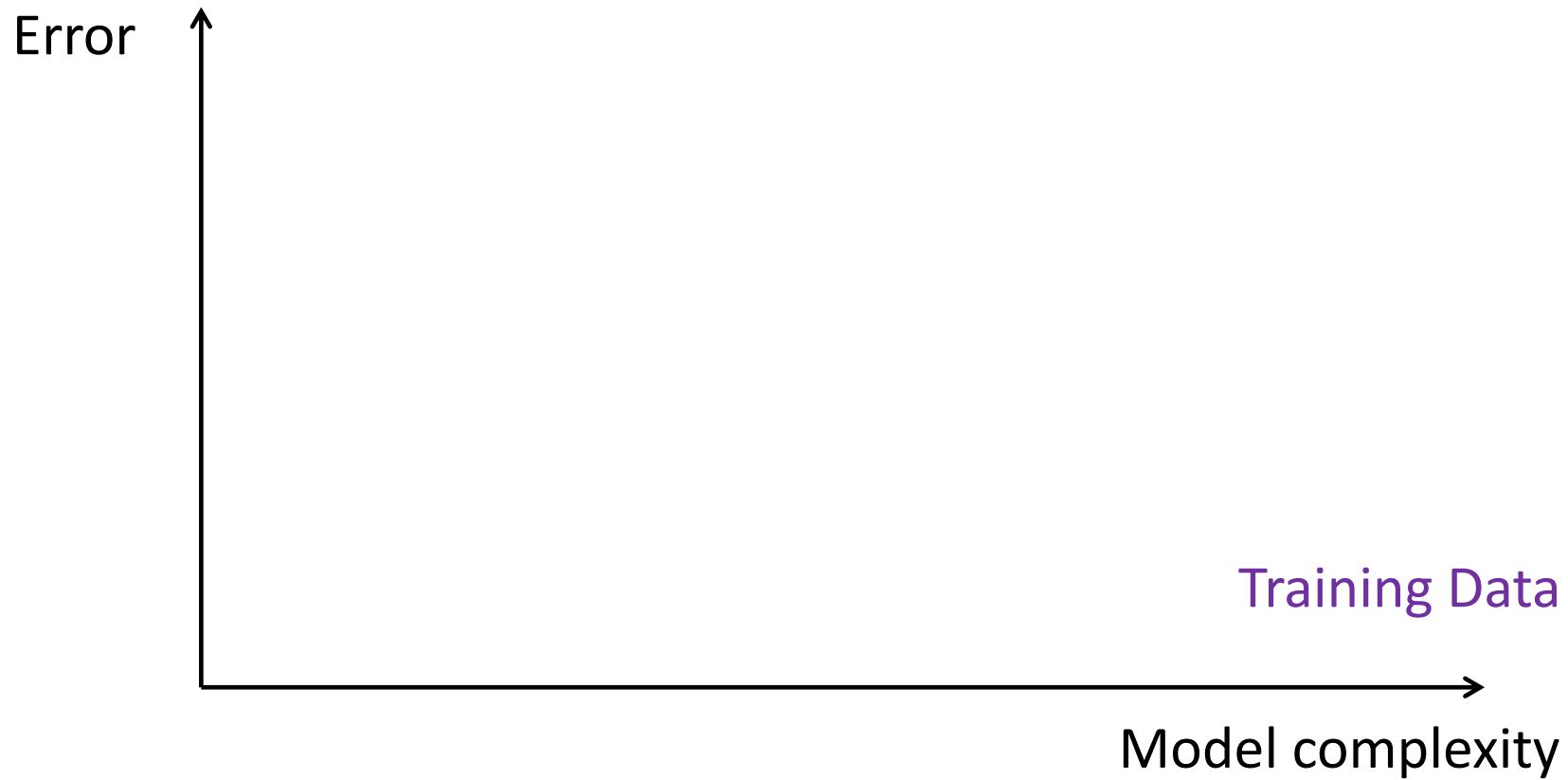
The first 6 batches have 32 samples each, and the 7th batch has the last 2 samples.

$$32 * 6 + 2 = 194$$

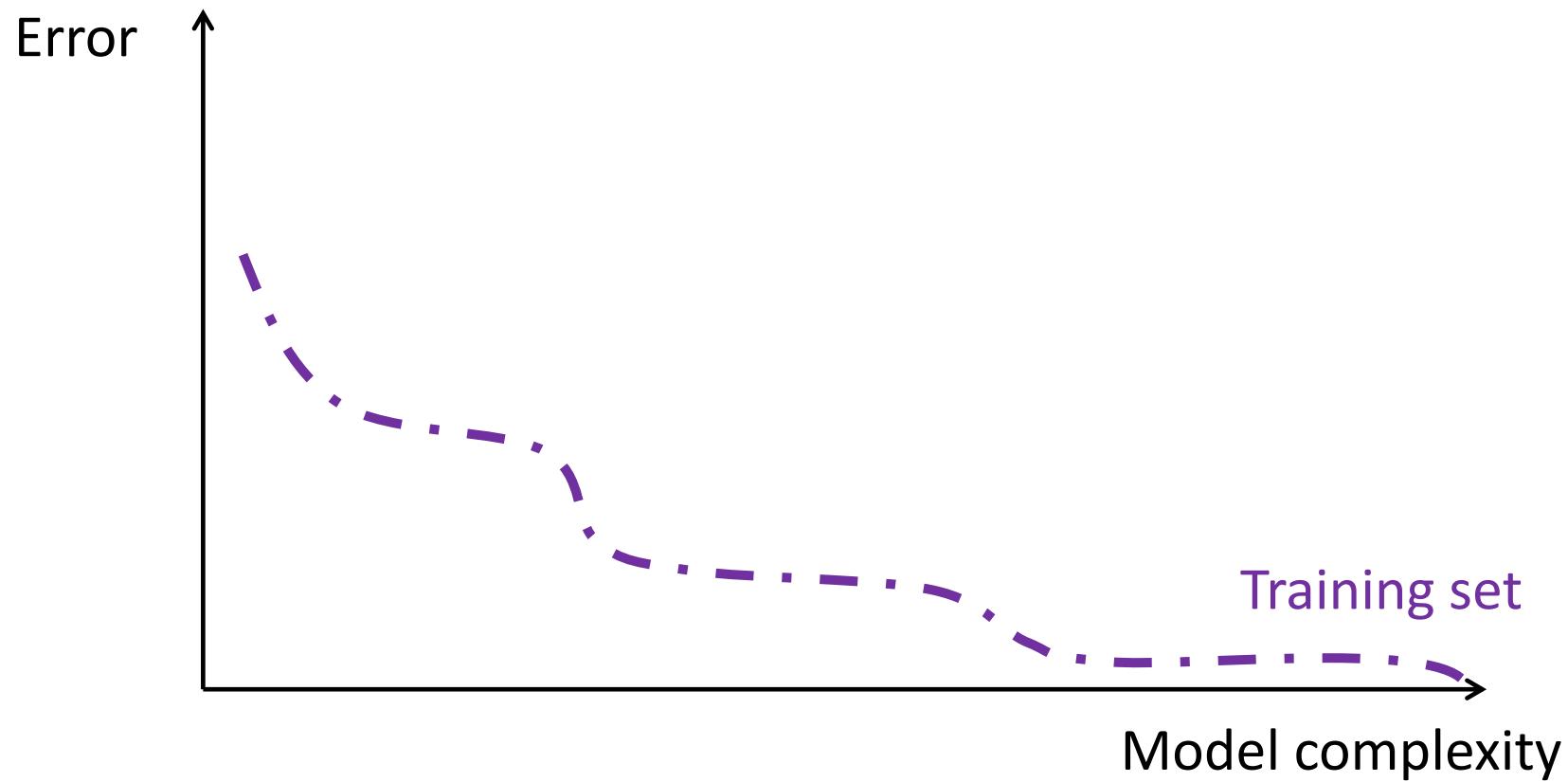


Overfitting and Regularization

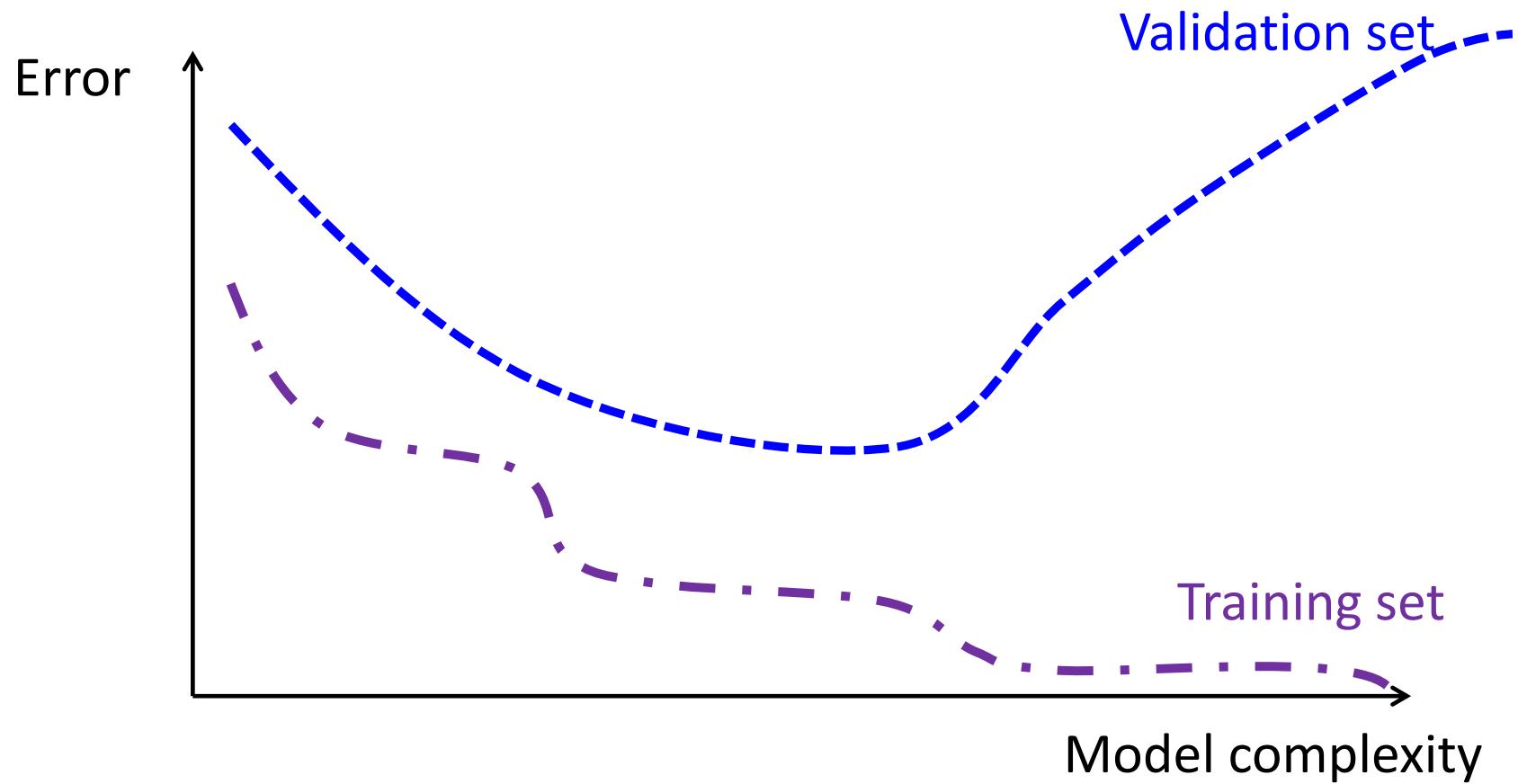
Recall Underfitting vs. Overfitting



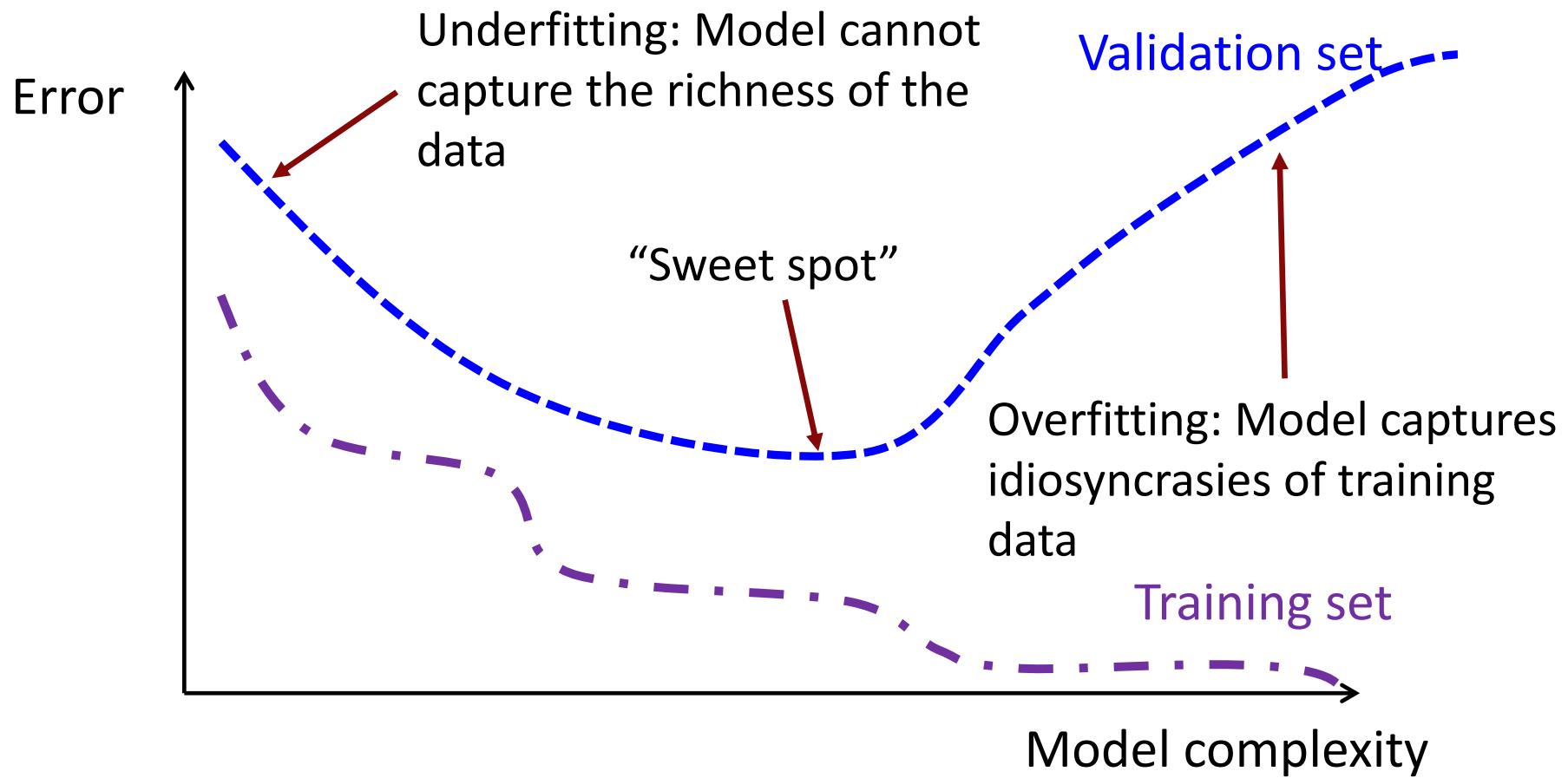
Recall Underfitting vs. Overfitting



Recall Underfitting vs. Overfitting



Recall Underfitting vs. Overfitting

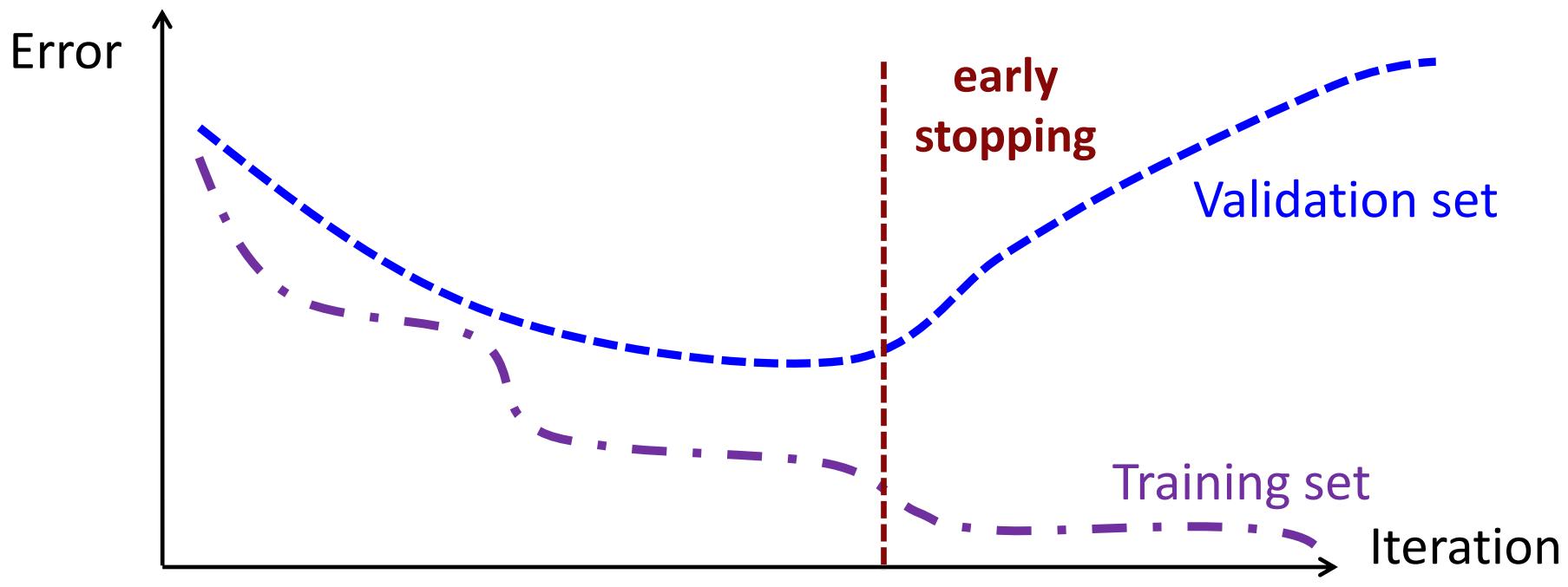


Overfitting in Neural Networks

- To learn smart representations of complex, unstructured data, the NN needs to have large “capacity” i.e., many layers and many neurons in each layer
- But this raises the likelihood of overfitting so we need to add *regularization*
- Several regularization methods have been developed to address this problem

Regularization strategy: *Early Stopping*

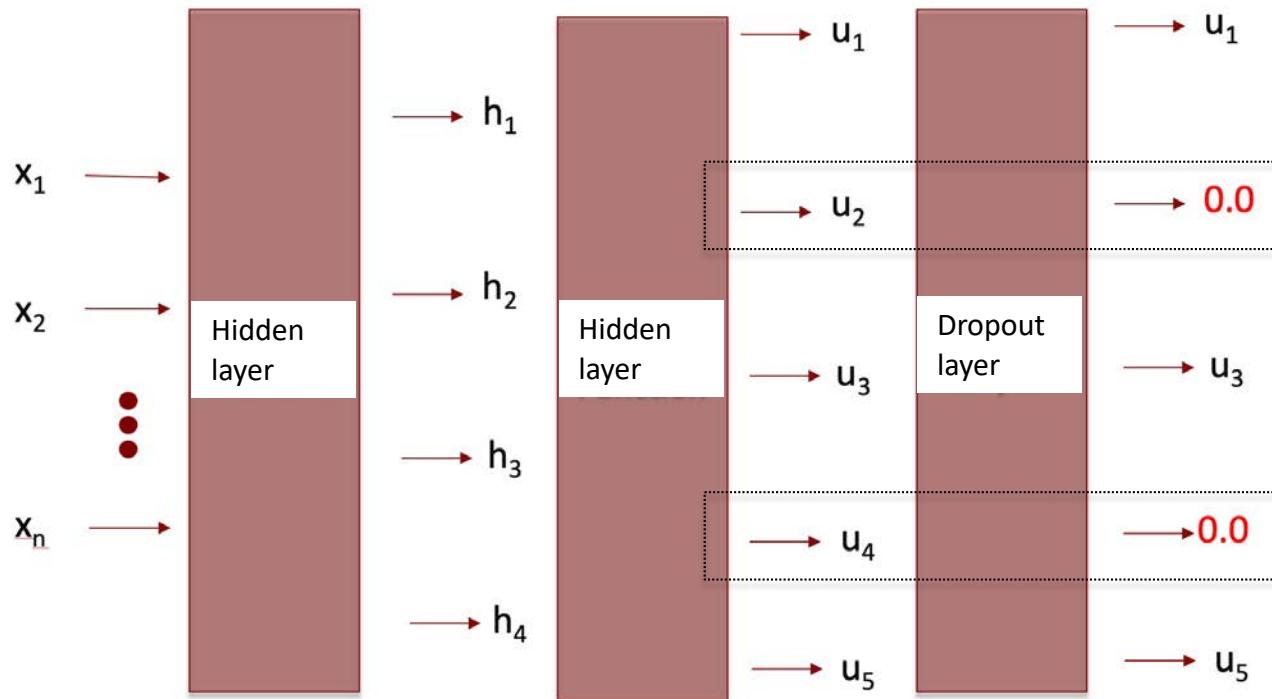
Stop the training early before the training loss is minimized by monitoring the loss on a validation dataset.



We will cover this in Lecture 4

Regularization strategy: *Dropout*

Randomly zero out the output from some of the nodes (typically 50% of the nodes) in a hidden layer (implemented as a “dropout layer” in Keras)



Summary: Creating and training a DNN from scratch

- We get the data ready
- We design i.e., “lay out” the network
 - Choose the number of hidden layers and the number of ‘neurons’ in each layer
 - Pick the right output layer based on the type of the output (more on this shortly)
- We pick
 - An appropriate loss function based on the type of the output (more on this shortly)
 - An optimizer from the many SGD flavors that are available and a “good” learning rate
- We decide on a regularization strategy
- We set things up in Keras/Tensorflow and start training!

Lightning Intro to Tensorflow/Keras

What's a Tensor?

What's a Tensor?

Tensor of rank 0 (Scalar)

42

What's a Tensor?

Tensor of rank 0 (Scalar)

42

Tensor of rank 1 (aka Vector)

(42, 23.4, 11.2)

What's a Tensor?

Tensor of rank 0 (Scalar)

42

Tensor of rank 1 (aka Vector)

(42, 23.4, 11.2)

Tensor of rank 2 (aka Matrix)

1.1	1.2	1.3	1.4	1.5	1.6	1.7	1.8	1.11	1.12	1.13	1.14	1.15	1.16	1.17	1.18	1.19	1.20	1.21	1.22	1.23	1.24	1.25	1.26	1.27	1.28
1.2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1.3	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1.4	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1.5	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1.6	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1.7	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1.8	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1.11	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1.12	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1.13	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1.14	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1.15	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1.16	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1.17	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1.18	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1.19	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1.20	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1.21	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1.22	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1.23	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1.24	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1.25	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1.26	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1.27	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1.28	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

Image credit: fast.ai

What's a Tensor?

Tensor of rank 0 (Scalar)

42

Tensor of rank 1 (aka Vector)

(42, 23.4, 11.2)

Tensor of rank 2 (aka Matrix)

	[,1]	[,2]	[,3]	[,4]	[,5]	[,6]	[,7]	[,8]	[,9]	[,10]	[,11]	[,12]	[,13]	[,14]	[,15]	[,16]	[,17]	[,18]	[,19]	[,20]	[,21]	[,22]	[,23]	[,24]	[,25]	[,26]	[,27]	[,28]		
[,1]	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
[,2]	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
[,3]	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
[,4]	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
[,5]	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
[,6]	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
[,7]	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
[,8]	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
[,9]	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
[,10]	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
[,11]	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
[,12]	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
[,13]	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
[,14]	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
[,15]	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
[,16]	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
[,17]	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
[,18]	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
[,19]	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
[,20]	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
[,21]	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
[,22]	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
[,23]	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
[,24]	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
[,25]	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
[,26]	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
[,27]	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
[,28]	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

Image credit: fast.ai

Tensor of rank 3 (aka “cube”)

	[,1]	[,2]	[,3]	[,4]	[,5]	[,6]	[,7]	[,8]	[,9]	[,10]	[,11]	[,12]	[,13]	[,14]	[,15]	[,16]	[,17]	[,18]	[,19]	[,20]	[,21]	[,22]	[,23]	[,24]	[,25]	[,26]	[,27]	[,28]		
[,1]	147	131	138	144	131	134	144	135	133	145	146	147	148	149	150	151	152	153	154	155	156	157	158	159	160	161	162	163		
[,2]	131	125	128	134	122	127	135	129	127	137	132	126	129	133	130	134	136	138	135	137	139	136	138	134	132	135	137	139		
[,3]	125	119	122	128	116	123	130	127	125	132	126	120	123	129	126	128	131	129	133	135	132	134	136	138	130	132	134	136	138	
[,4]	128	122	125	134	118	127	131	129	127	135	126	121	124	130	127	129	132	128	134	136	133	135	137	139	131	133	135	137	139	
[,5]	122	116	119	128	113	123	129	120	118	132	121	115	118	124	121	119	126	123	120	127	129	126	128	130	122	124	126	128	130	
[,6]	125	119	122	134	116	123	130	127	125	132	126	120	123	129	126	128	131	127	133	135	132	134	136	138	130	132	134	136	138	
[,7]	127	121	124	134	118	126	131	128	126	135	127	121	124	130	127	129	132	128	134	136	133	135	137	139	131	133	135	137	139	
[,8]	129	123	126	134	120	127	131	129	127	135	128	122	125	131	128	129	133	129	135	137	134	136	138	140	132	134	136	138	140	
[,9]	126	120	123	134	113	125	130	122	119	132	121	115	118	124	121	116	127	124	121	128	130	125	127	129	131	123	125	127	129	131
[,10]	123	117	120	134	110	122	129	118	115	132	119	113	116	123	119	117	126	123	120	127	129	126	128	130	114	116	118	120	122	

Can you give an example of a rank-4 tensor?

What's a Tensor?

See Chapter 2.2 of text for more detail

Tensorflow

Tensorflow (TF) is a library that provides

- Automatic calculation of the gradient of (complicated) loss functions

$$\nabla Loss(w) = \left[\frac{\partial Loss}{\partial w_1}, \frac{\partial Loss}{\partial w_2}, \dots, \frac{\partial Loss}{\partial w_n} \right]$$

Tensorflow

Tensorflow (TF) is a library that provides

- Automatic calculation of the gradient of (complicated) loss functions
- Library of state-of-the-art optimizers

Tensorflow

Tensorflow (TF) is a library that provides

- Automatic calculation of the gradient of (complicated) loss functions
- Library of state-of-the-art optimizers
- Automatic distribution of computational load across servers

Tensorflow

Tensorflow (TF) is a library that provides

- Automatic calculation of the gradient of (complicated) loss functions
- Library of state-of-the-art optimizers
- Automatic distribution of computational load across servers
- Automatic adaptation of code to work on parallel hardware (GPUs and TPUs)



Keras “sits on top of” Tensorflow ...

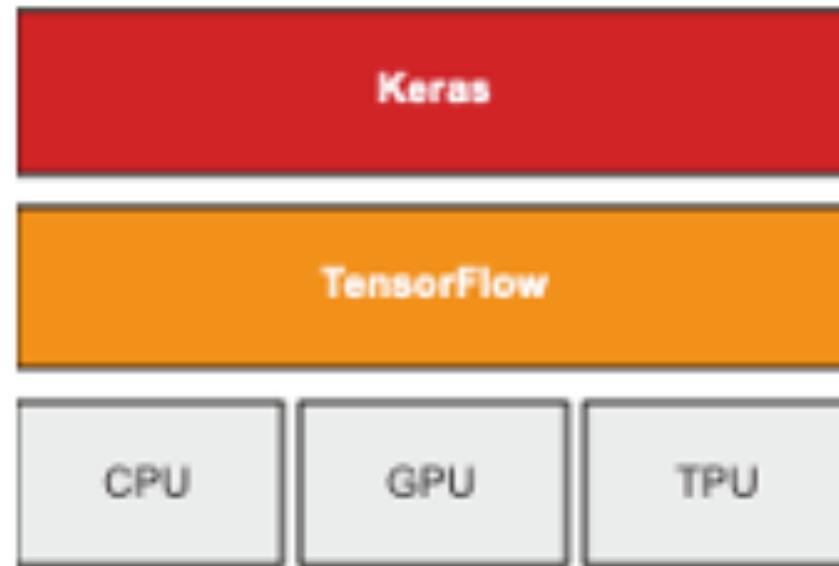


Image: Page 70 of textbook

... and provides “convenience” features



- Pre-defined **layers**
- Incredibly flexible ways to specify network **architectures**
- Easy ways to **preprocess** data
- Easy ways to **train** models and **report** metrics
- **Pre-trained models** you can download and customize

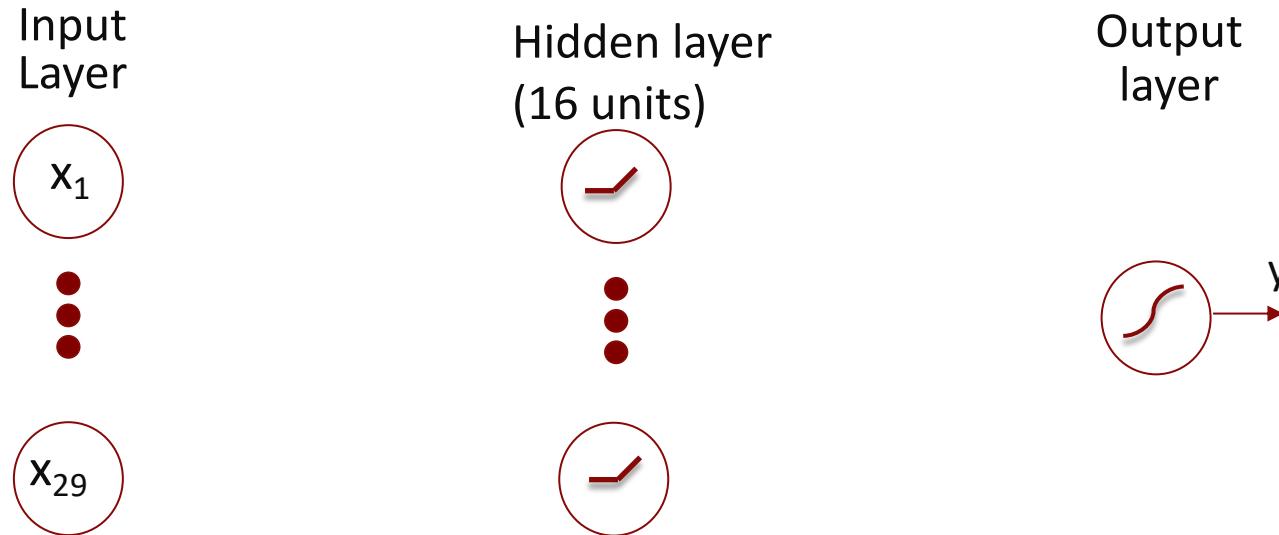
Keras APIs

- There are three broad ways to build DL models with Keras
 - Sequential
 - Functional API
 - Subclassing
- **We will almost exclusively use the Functional API.** The model we built for heart disease prediction is an example.
- Please read 7.2.2 of the textbook to understand in detail how the Keras Functional API works



Check out the wealth of introductory
and advanced material, with
accompanying colabs, at
tensorflow.org and keras.io

Let's revisit the Neural Model for Heart Disease Prediction we designed previously



```
input = keras.Input(shape=29)
h = keras.layers.Dense(16, activation="relu")(input)
output = keras.layers.Dense(1, activation="sigmoid")(h)
model = keras.Model(input, output)
```



Let's train this model!

Training Checklist

- We get the data ready (will cover in the colab)
- We design i.e., “lay out” the network **1 hidden layer with 16 ReLU neurons**
 - Choose ***the number of hidden layers*** and ***the number of ‘neurons’ in each layer***
 - Pick the ***right output layer*** based on the type of the output **Sigmoid**
- We pick
 - An appropriate ***loss function*** based on the type of the output _____
 - An ***optimizer from the many SGD flavors*** that are available
- We decide on a ***regularization strategy***
- We set things up in Keras/Tensorflow and start training!

Training Checklist

- We get the data ready (will cover in the colab)
- We design i.e., “lay out” the network **1 hidden layer with 16 ReLU neurons**
 - Choose ***the number of hidden layers*** and ***the number of ‘neurons’ in each layer***
 - Pick the ***right output layer*** based on the type of the output **Sigmoid**
- We pick
 - An appropriate ***loss function*** based on the type of the output **binary crossentropy**
 - An ***optimizer from the many SGD flavors*** that are available
- We decide on a ***regularization strategy***
- We set things up in Keras/Tensorflow and start training!

Training Checklist

- We get the data ready (will cover in the colab)
- We design i.e., “lay out” the network **1 hidden layer with 16 ReLU neurons**
 - Choose ***the number of hidden layers*** and ***the number of ‘neurons’ in each layer***
 - Pick the ***right output layer*** based on the type of the output **Sigmoid**
- We pick
 - An appropriate ***loss function*** based on the type of the output **binary crossentropy**
 - An ***optimizer from the many SGD flavors*** that are available “**adam**”
- We decide on a ***regularization strategy*** **Early stopping**
- We set things up in Keras/Tensorflow and start training!

Colab

Predicting Heart Disease

Before we start coding ...

- Don't worry if you don't understand every detail of what we will do in class.
- But go through the Colab notebooks carefully later, play around with the code and make sure you understand every line

Colab General Instructions

Step 1

Make your own copy of the notebook



Step 2

Request a GPU for your notebook*



Notebook settings

Hardware accelerator
GPU

To get the most out of Colab, avoid using a GPU unless you need one. [Learn more](#)

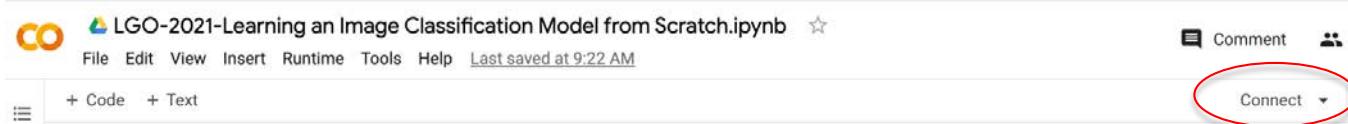
Omit code cell output when saving this notebook

Cancel

Save

Step 3

Start your notebook



You need to do steps 1 and 2 just the first time you use a notebook. From the second time onwards, jump to Step 3.

MIT OpenCourseWare
<https://ocw.mit.edu>

15.773 Hands-on Deep Learning

Spring 2024

For information about citing these materials or our Terms of Use, visit: <https://ocw.mit.edu/terms>.