



Intro to Neural Nets

Session 4: Model Fitting

Session Agenda

General Workflow

- Step 1: Get your model to overfit on training data (always possible).
- Step 2: Get your model to fit to validation data (this is quite exploratory)
- Step 3: Maximize out of sample performance by mitigating / delaying overfitting in training data.

Techniques to Mitigate Overfitting

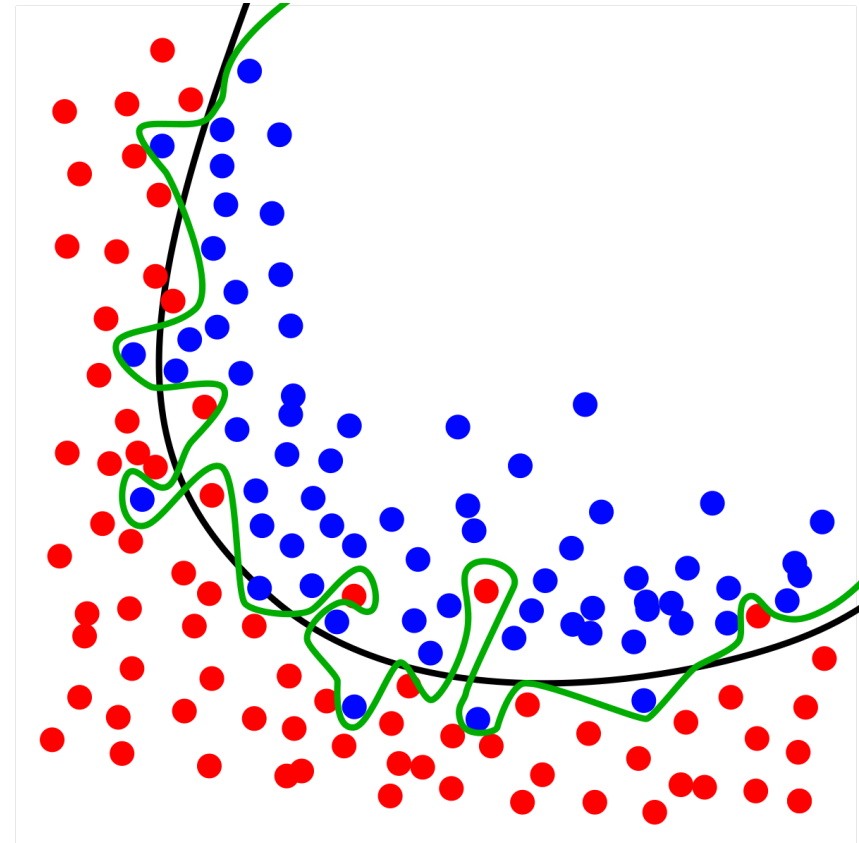
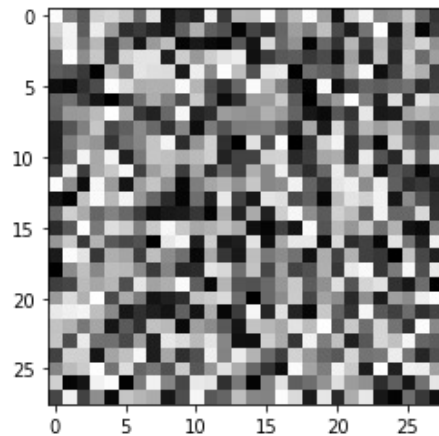
- Early stopping (we've seen this already, and we will do it regularly).
- Regularizing or constraining weights
- Dropout layer
- Adjust batch size
- Inject noise
- Get better / more data

Some Example Neural Nets

Overfitting

A Neural Network Can Easily Overfit...

- Let's look at an extreme case...



Regularization 1: Kernel Regularizer

Regularizing Weights Means Less Entropy

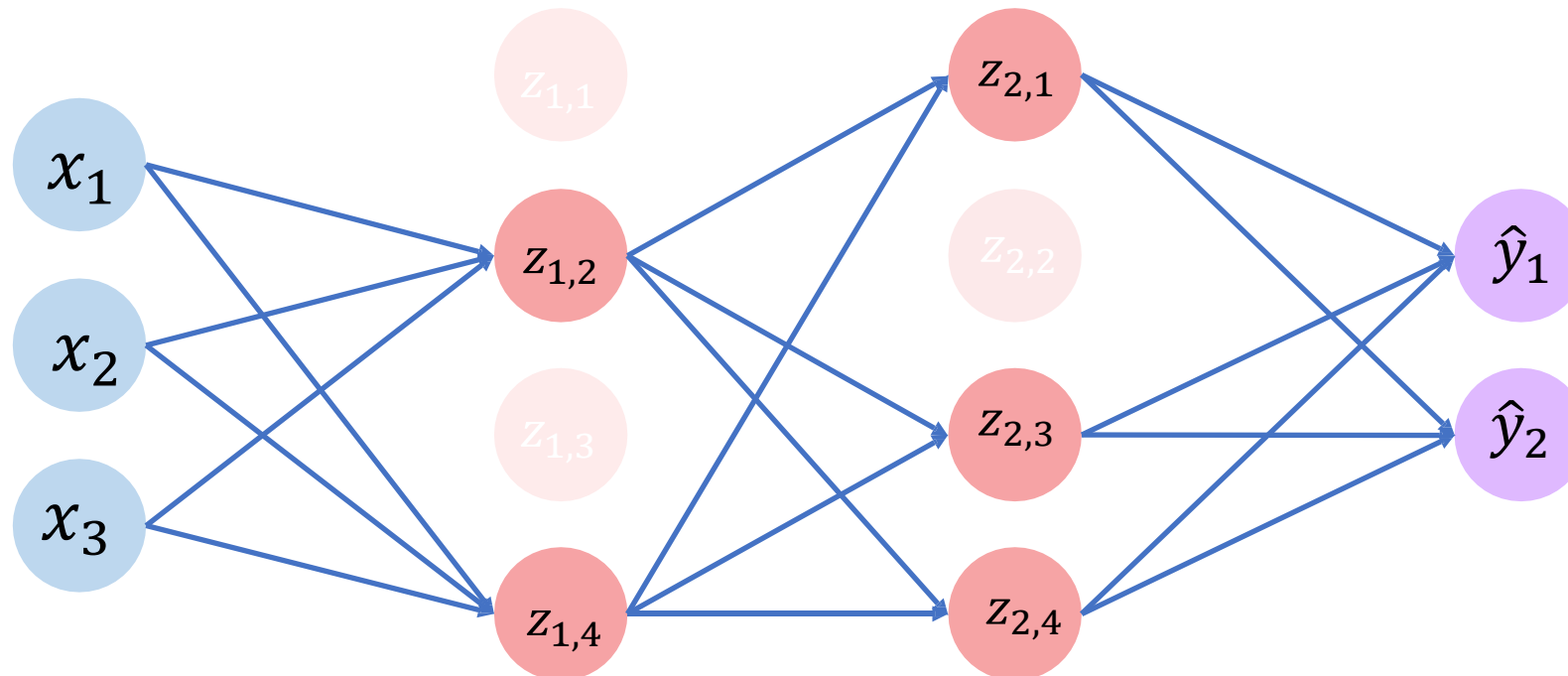
- Regularizing a layer's weights means the weights are updated less as their collective magnitude (e.g., sum) gets larger. Both L1 or L2 norms can be used here applied.
- This approach is typically used to improve the validation performance of smaller networks.

Listing 5.13 Adding L2 weight regularization to the model

```
1 from tensorflow.keras import regularizers
2 model = keras.Sequential([
3     layers.Dense(16,
4                   kernel_regularizer=regularizers.l2(0.002),
5                   activation="relu"),
6     layers.Dense(16,
7                   kernel_regularizer=regularizers.l2(0.002),
8                   activation="relu"),
9     layers.Dense(1, activation="sigmoid")
10 ])
11 model.compile(optimizer="rmsprop",
12               loss="binary_crossentropy",
13               metrics=["accuracy"])
14 history_l2_reg = model.fit(
15     train_data, train_labels,
16     epochs=20, batch_size=512, validation_split=0.4)
```

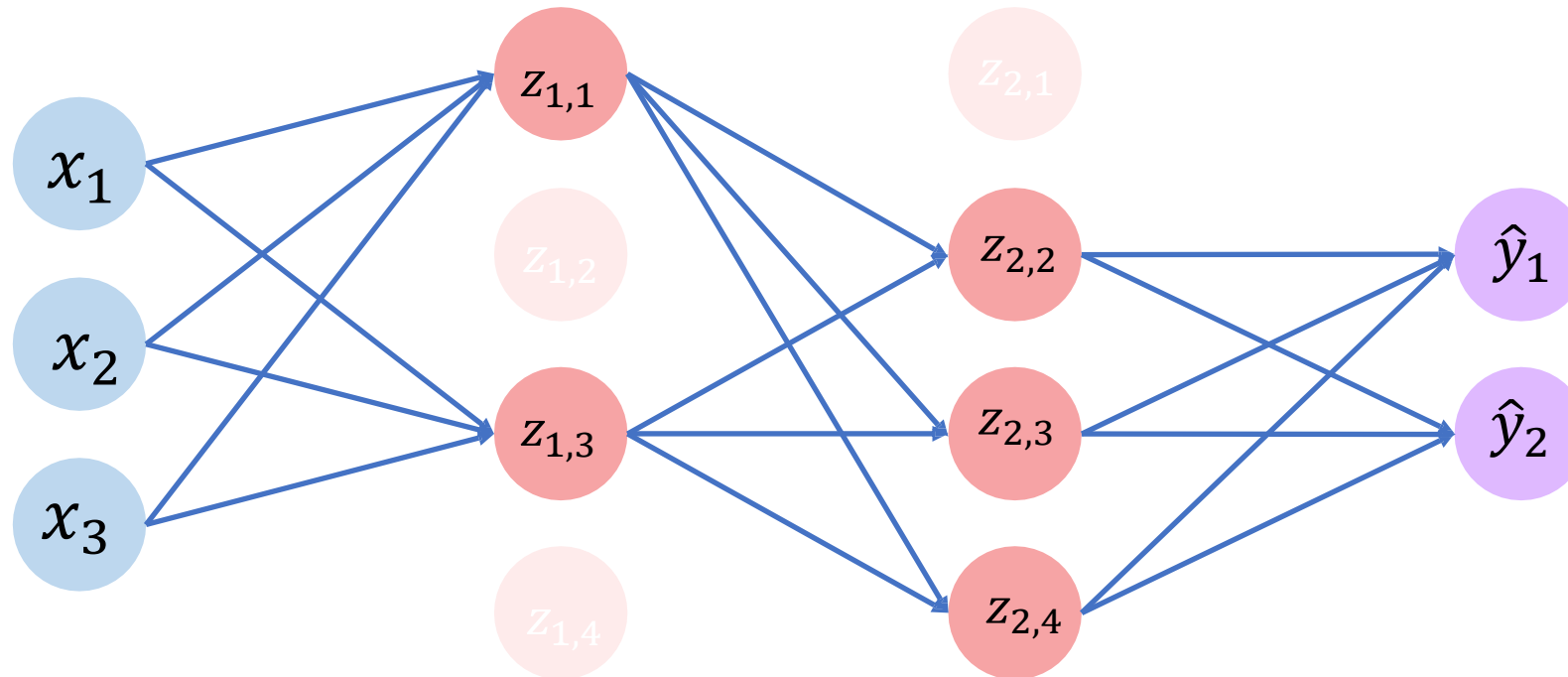
Regularization 2: Dropout

- During training, at each iteration, randomly set some activations to 0
 - Typically 'drop' 50% of activations in layer



Regularization 2: Dropout

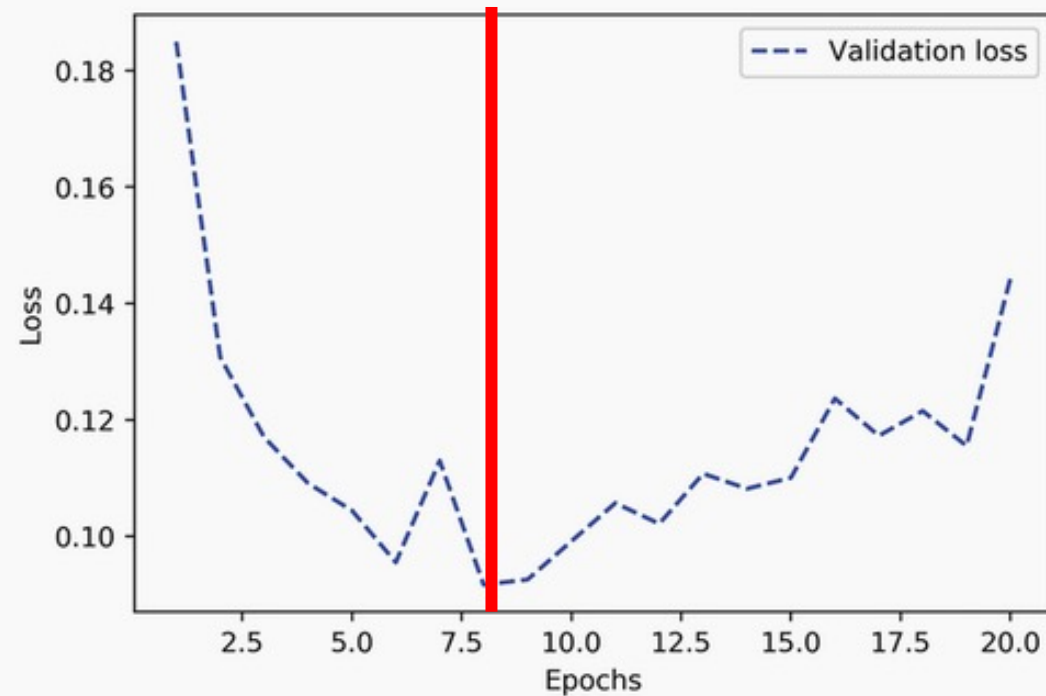
- During training, at each iteration, randomly set some activations to 0
 - Typically 'drop' 50% of activations in layer



Early Stopping

Monitoring Validation Performance and then Manually Limiting Epoch Count

- This is the sole approach we have been taking thus far.



Adjust Batch Size

Deep Learning

Ian Goodfellow
Yoshua Bengio
Aaron Courville

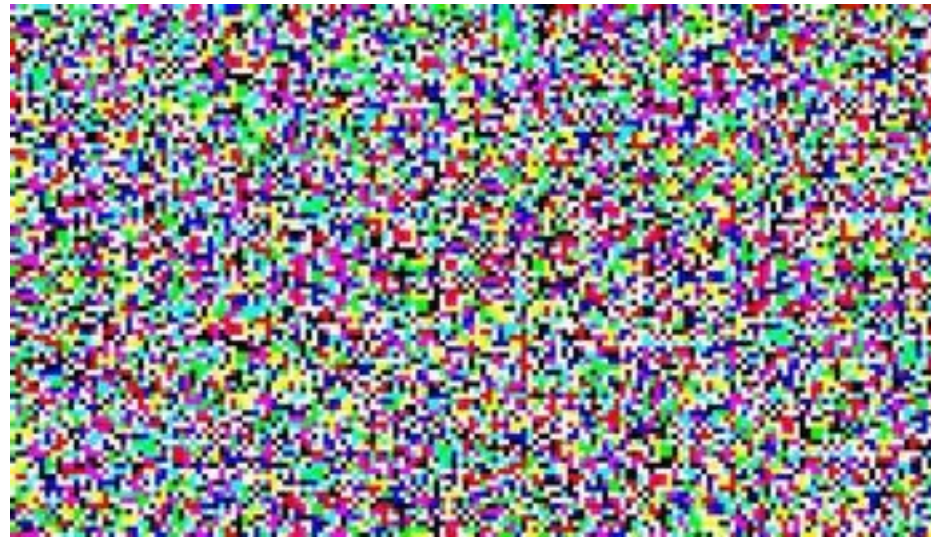
Minibatch sizes are generally driven by the following factors:

- Larger batches provide a more accurate estimate of the gradient, but with less than linear returns.
- Multicore architectures are usually underutilized by extremely small batches. This motivates using some absolute minimum batch size, below which there is no reduction in the time to process a minibatch.
- If all examples in the batch are to be processed in parallel (as is typically the case), then the amount of memory scales with the batch size. For many hardware setups this is the limiting factor in batch size.
- Some kinds of hardware achieve better runtime with specific sizes of arrays. Especially when using GPUs, it is common for power of 2 batch sizes to offer better runtime. Typical power of 2 batch sizes range from 32 to 256, with 16 sometimes being attempted for large models.
- Small batches can offer a regularizing effect ([Wilson and Martinez, 2003](#)), perhaps due to the noise they add to the learning process. Generalization error is often best for a batch size of 1. Training with such a small batch size might require a small learning rate to maintain stability because of the high variance in the estimate of the gradient. The total runtime can be very high as a result of the need to make more steps, both because of the reduced learning rate and because it takes more steps to observe the entire training set.

Inject Noise

Just Like Smaller Batch Size, but Purposeful

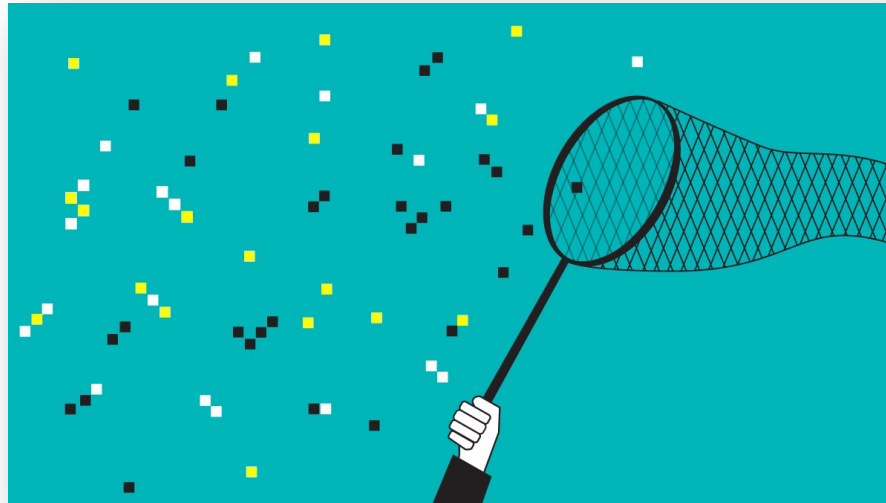
- One can manually jitter model weights at each iteration by ‘adding’ random noise.
- You can add noise to any model component, including the inputs, activations and outcome labels.



Get More Data!

More Training Examples

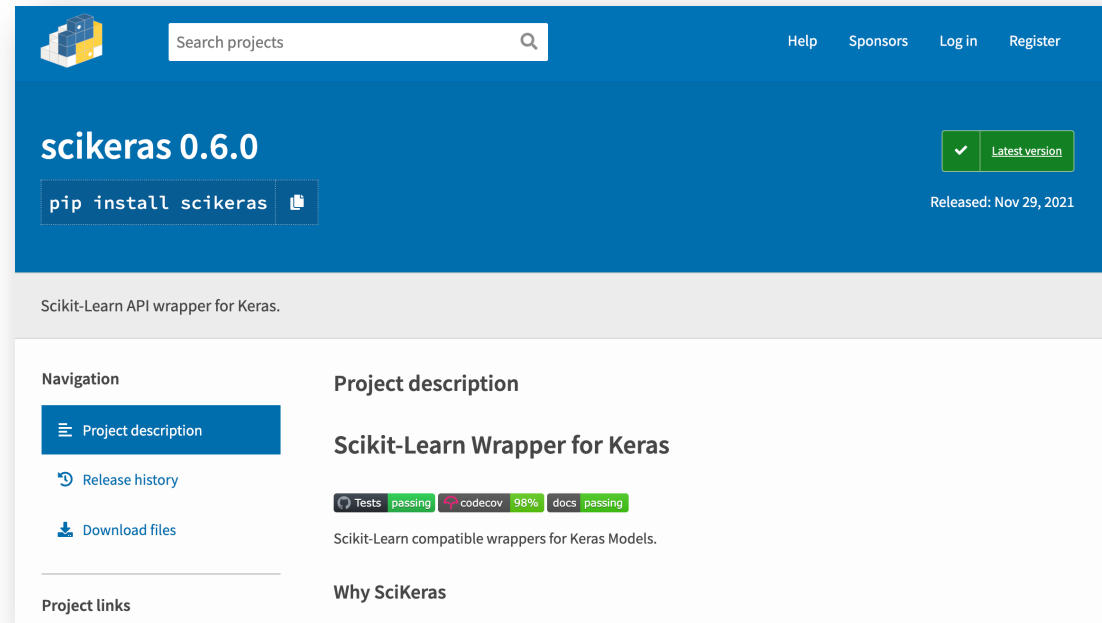
- Means you can have a bigger training data-set, which will presumably contain more information for the model to extract.
- This is the approach that often yields the best marginal returns, though it can also be most costly.



Scikit-Learn Wrapper for Keras

Facilitate Hyperparameter Tuning and Cross-validation of a Deep Net

- These wrappers can be used with the Sequential API, with two caveats:
- First, you need to formally specify the shape of the input layer.
- Second, you need to install scikeras in Google colab to use it.



Let's Walk Through a Few Examples