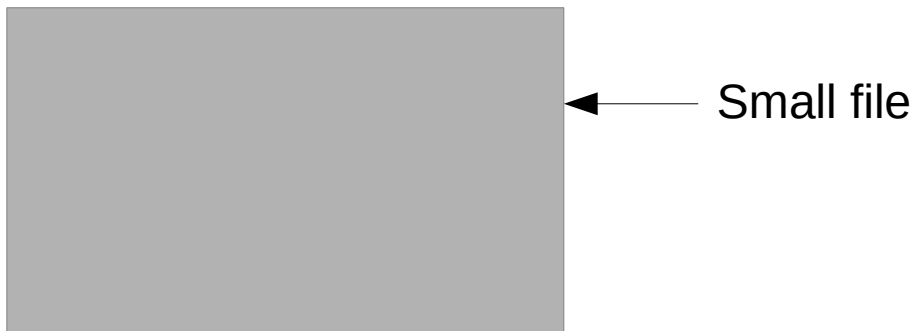




Data Management

Examining the Data

- Look at some random inputs
- Look at the largest and smallest inputs (by file size)
 - These are the most likely to be weird or corrupted



- Manually solve the task

Splitting

- Training set – used by SGD to learn model parameters
- Validation set – used by us to learn hyperparameters
- Test set – used once to measure final model performance

Overfitting

- Goal: Learn to classify *unseen* images
- Optimization objective: Learn to classify *training set* images

Training set (60-80% of data)
Model parameters overfit by SGD

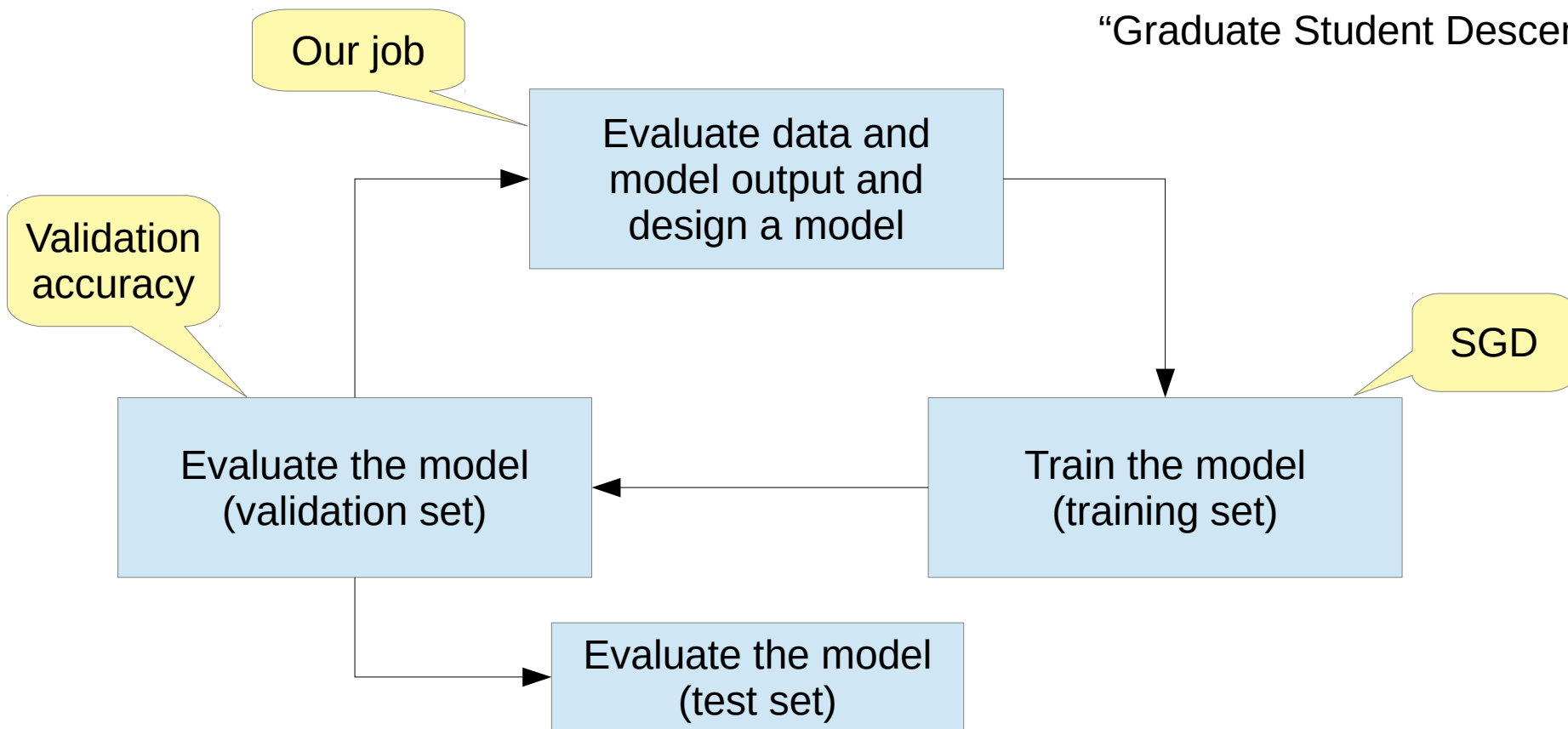
Validation (10-20%)
Hyperparameters overfit by us

Test (10-20%)
Not used in training, not overfit

Be sure to split
randomly

Tuning

“Graduate Student Descent”





Initialization

Initialization – All Zero

Linear 1

$$x = ???$$

$$v_1 = W_1 x$$

$$v_1 = 0$$

$$g_1 = 0$$

$$g_1 = W_1^T g_2$$

$$g_2 = 0$$

$$\nabla_{W_1} \ell(v_3) = g_2 x^T$$

$$0$$

ReLU

$$v_2 = \max(v_1, 0)$$

$$v_2 = 0$$

$$g_2 = g_3 [v_1 > 0]$$

$$g_3 = 0$$

Linear 2

$$v_3 = W_2 v_2$$

$$v_3 = 0$$

$$g_3 = W_2^T g_o$$

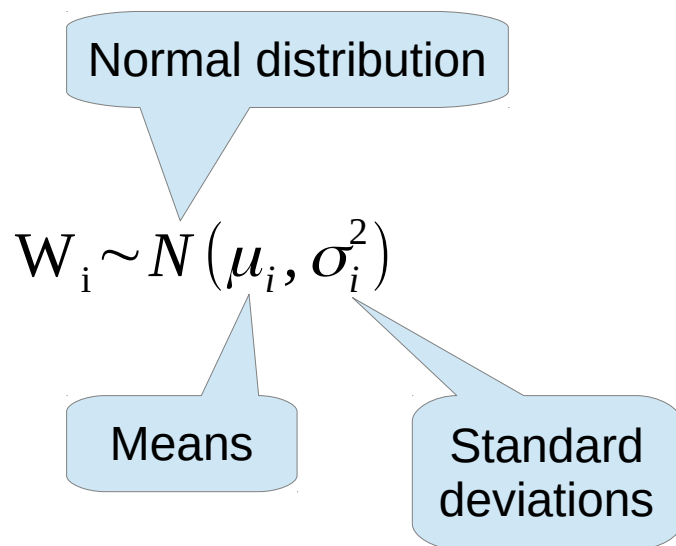
$$g_o = ???$$

$$\nabla_{W_2} \ell(v_3) = g_o v_2^T$$

$$0$$

$$g_o = \nabla_{v_3} \ell(v_3)$$

Initialization – Random



In practice, the mean is zero and we need to choose a standard deviation

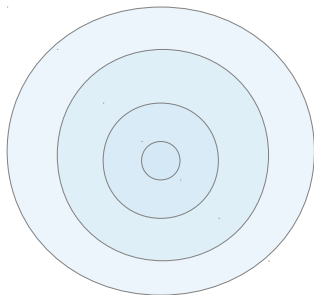
Initialization – Scaling

$$\mathbf{o} = \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}$$

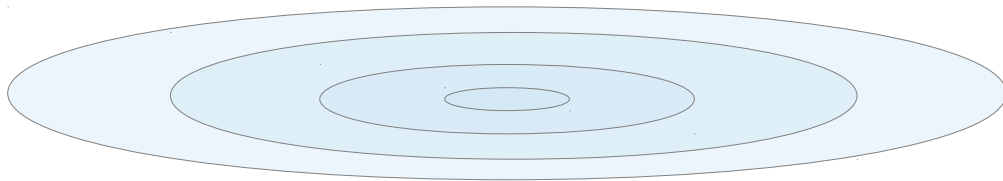
$$\nabla_{\mathbf{W}_1} \ell(\mathbf{o}) = \left(\mathbf{W}_2^T \left(\nabla_{\mathbf{o}} \ell(\mathbf{o}) \right) \right) \mathbf{x}^T$$

$$\nabla_{\mathbf{W}_2} \ell(\mathbf{o}) = \left(\nabla_{\mathbf{o}} \ell(\mathbf{o}) \right) \left(\mathbf{W}_1 \mathbf{x} \right)^T$$

$$\|\mathbf{W}_1\| \approx \|\mathbf{W}_2\|$$



$$\|\mathbf{W}_1\| \gg \|\mathbf{W}_2\|$$



Choosing a Standard Deviation

Math

(See the supplementary material for details)

Kaiming Initialization

- Choose either activations or gradients and keep the magnitude roughly constant:

- Activations: $W \sim N\left(0, \frac{2}{n_{\text{in}}}\right)$

- Gradients: $W \sim N\left(0, \frac{2}{n_{\text{out}}}\right)$

Xavier Initialization

- Keep both activation and gradient magnitudes roughly constant throughout the network

$$W \sim N\left(0, \frac{4}{n_{\text{in}} + n_{\text{out}}}\right)$$

In Practice

- The PyTorch defaults are usually good enough.
- The last layer can be initialized to zero.
 - If the previous layers are not zero, the last activation is not zero, so we still get gradients.