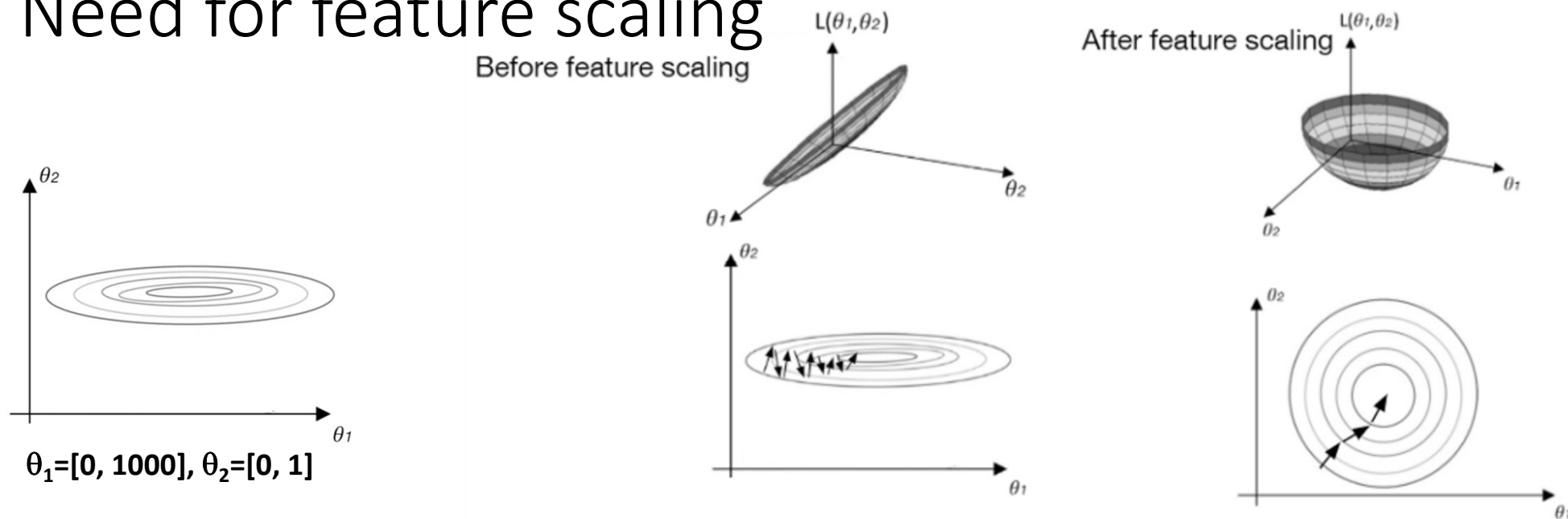


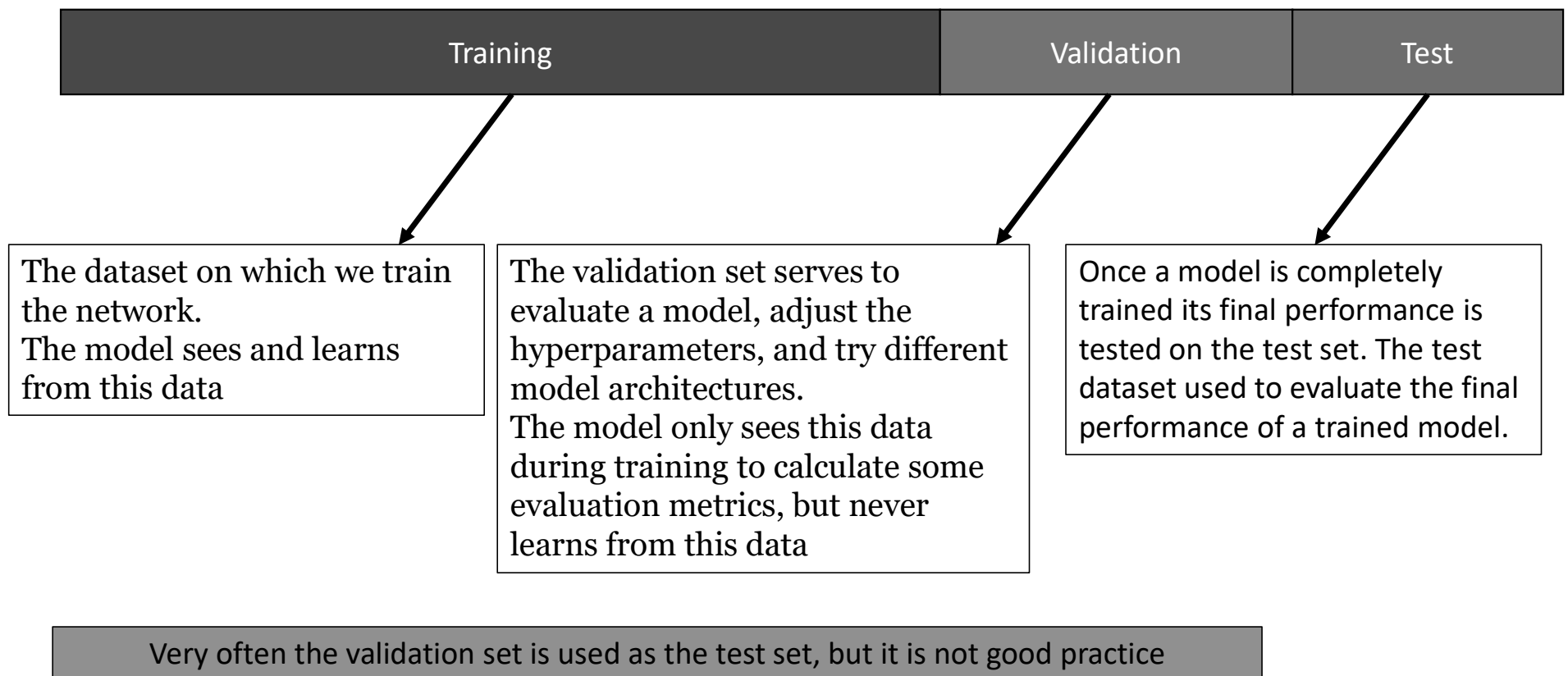
Need for feature scaling



- The GD optimizer would take bigger steps in the direction of smaller range feature, while making only little steps in the direction of higher range feature.
- The gradient would oscillate a lot, which would cause the model to take a longer time to find the global minimum or In the worst case, the minimum could not be found at all.
- We can prevent this problem by feature scaling

<https://www.deeplearning-academy.com/p/ai-wiki-data-preprocessing>

Data splitting strategy

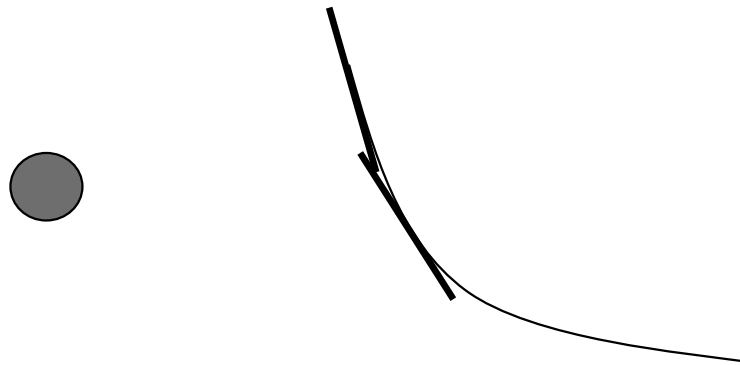


Data splitting strategy

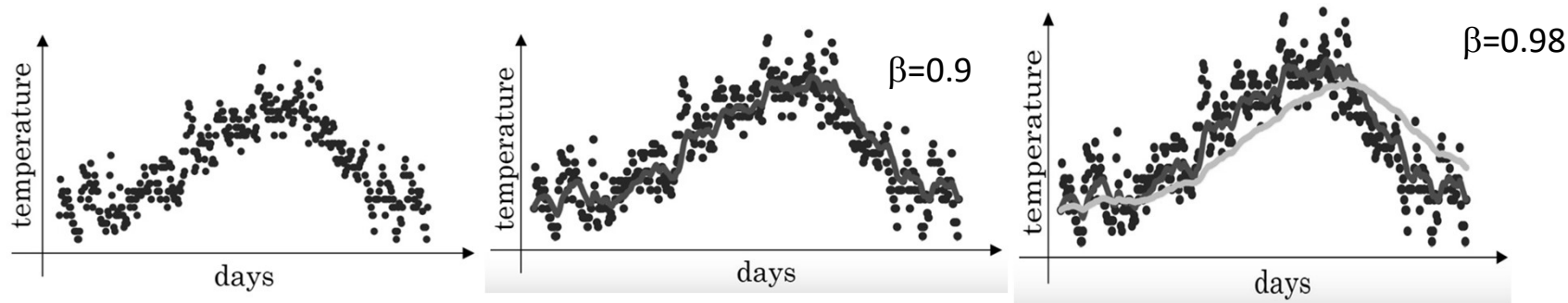
- Training size has to be larger
- 1000 to 50,000 70% Train, 20% validation and 10% for test
- >>50,000, 90% Train, 5% validation and 5% for test

SGD

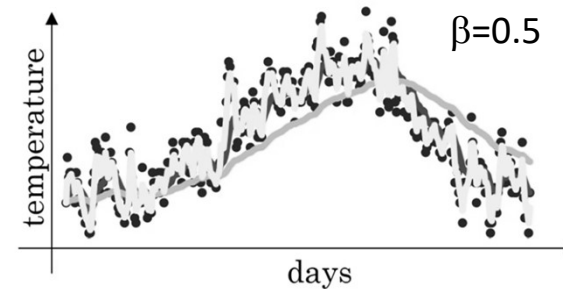
- 1) cannot adapt to changing gradient values
- 2) it cannot change the learning rate to accelerate/ decelerate the minimization based on the gradient values



Exponentially weighted moving average



$$\begin{aligned}v_0 &= 0 \\v_1 &= 0.9 v_0 + 0.1 \theta_1 \\v_2 &= 0.9 v_1 + 0.1 \theta_2 \\v_3 &= 0.9 v_2 + 0.1 \theta_3 \\&\vdots \\v_t &= 0.9 v_{t-1} + 0.1 \theta_t\end{aligned}$$



$$v_t = \beta v_{t-1} + (1 - \beta) \theta_t$$

$v_t \sim$ averaging over $1/(1-\beta)$ values

- deeplearning.ai [Andrew NG]

Exponentially weighted moving average

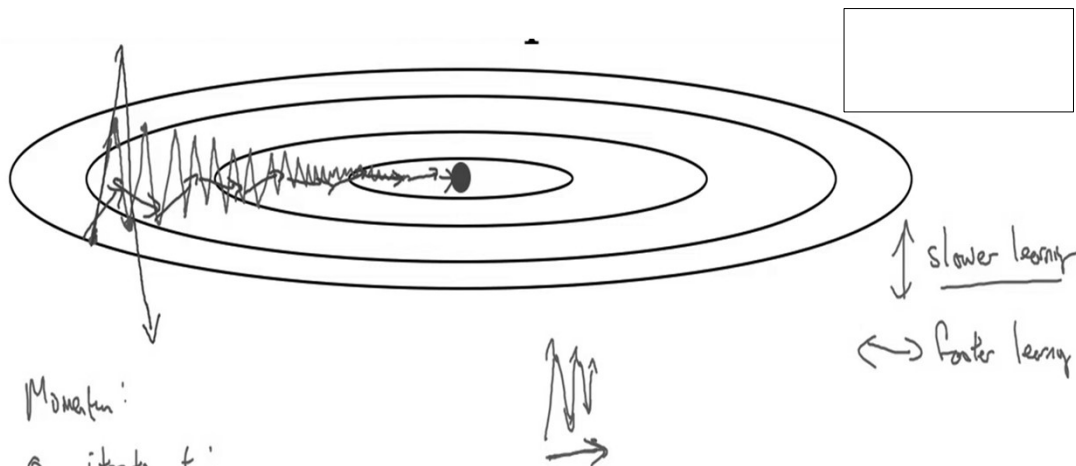
$$v_n = \sum_{k=0}^n (1 - \beta) \times \beta^k \theta_n$$

$$V_{10} = 0.1\theta_{10} + 0.1 \times (0.9) \theta_9 + 0.1 \times (0.9)^2 \theta_8 + 0.1 \times (0.9)^3 \theta_7 + \dots + 0.1 \times (0.9)^{10} \theta_1$$

$$v_t = \beta v_{t-1} + (1 - \beta) \theta_t$$

Advantage Less memory

Gradient descent with momentum



$$W = W - \alpha dw$$

On iteration t :

Compute dW, db on the current mini-batch

$$v_{dW} = \beta v_{dW} + (1 - \beta) dW$$

$$v_{db} = \beta v_{db} + (1 - \beta) db$$

$$W = W - \alpha v_{dW}, \quad b = b - \alpha v_{db}$$

Momentum:

On iteration t :

Compute dW, db on current mini-batch.

$$v_{dW} = \beta v_{dW} + (1 - \beta) dW$$

$$v_{db} = \beta v_{db} + (1 - \beta) db$$

$$v_{\theta} = \beta v_{\theta} + (1 - \beta) \theta_t$$

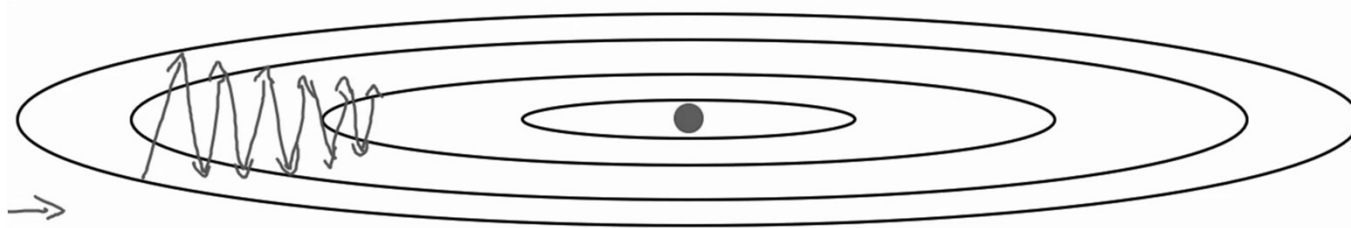
$$W = W - \alpha v_{dW}, \quad b = b - \alpha v_{db}$$

Hyperparameters: α, β

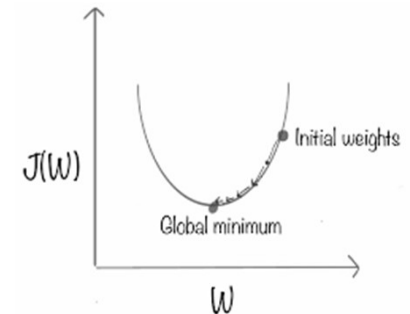
$$\beta = 0.9$$

- Large gradient \rightarrow smaller weight update

Adaptive learning rate



↑ slow
↔ fast



$$dw = dC_f/dw$$

On iteration t :

Compute dw, db on current mini-batch
 \nwarrow element-wise

$$S_{dw} = \beta S_{dw} + (1-\beta) \frac{dw^2}{2}$$

$$S_{db} = \beta S_{db} + (1-\beta) db^2$$

$$w := w - \alpha \frac{dw}{\sqrt{S_{dw}}}$$

$$b := b - \alpha \frac{db}{\sqrt{S_{db}}}$$

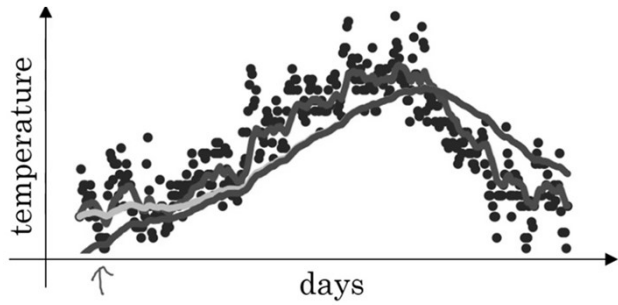
Large gradient \rightarrow smaller weight update

small gradient \rightarrow larger weight update

Adaptive learning rate

$$\frac{\alpha}{\sqrt{S}}$$

Bias correction in EWMA



$$v_t = \frac{\beta}{(1 - \beta^t)} v_{t-1} + \frac{(1 - \beta)}{(1 - \beta^t)} \theta_t$$

$$v_t = \beta v_{t-1} + (1 - \beta) \theta_t$$

ADAM optimizer (Adaptive Moment Estimate)

$$V_{dw}=0, S_{dw}=0. \quad V_{db}=0, S_{db}=0$$

On iteration t :

Compute dw, db using current mini-batch

$$V_{dw} = \beta_1 V_{dw} + (1-\beta_1) dw, \quad V_{db} = \beta_1 V_{db} + (1-\beta_1) db \quad \leftarrow \text{"momentum"} \beta_1$$

$$S_{dw} = \beta_2 S_{dw} + (1-\beta_2) dw^2, \quad S_{db} = \beta_2 S_{db} + (1-\beta_2) db^2 \quad \leftarrow \text{"RMSprop"} \beta_2$$

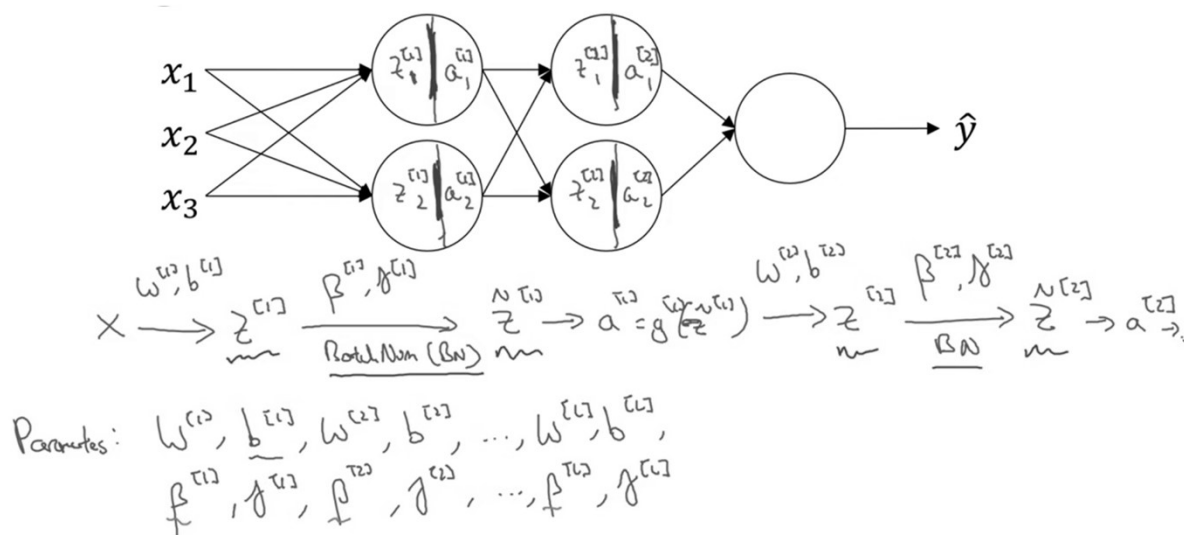
$$V_{dw}^{\text{corrected}} = V_{dw} / (1-\beta_1^t), \quad V_{db}^{\text{corrected}} = V_{db} / (1-\beta_1^t)$$

$$S_{dw}^{\text{corrected}} = S_{dw} / (1-\beta_2^t), \quad S_{db}^{\text{corrected}} = S_{db} / (1-\beta_2^t)$$

$$w := w - \alpha \frac{V_{dw}^{\text{corrected}}}{\sqrt{S_{dw}^{\text{corrected}} + \epsilon}}, \quad b := b - \alpha \frac{V_{db}^{\text{corrected}}}{\sqrt{S_{db}^{\text{corrected}} + \epsilon}}$$

$\alpha \rightarrow$ needs to be tuned
 $\beta_1 = 0.9 \rightarrow$ momentum
 $\beta_2 = 0.999 \rightarrow$ RMS

Batch normalization works



Given some intermediate values in NN $z^{(1)}, \dots, z^{(n)}$

$$\mu = \frac{1}{n} \sum_i z^{(i)}$$

$$\sigma^2 = \frac{1}{n} \sum_i (z^{(i)} - \mu)^2$$

$$z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$\hat{z}^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta$$

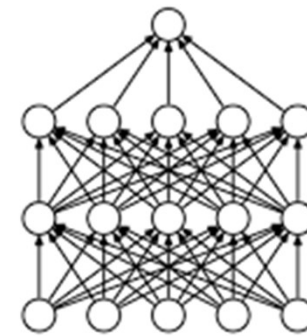
If $\sigma^2 = 0$, $\beta = \mu$
 then $\hat{z}^{(i)} = z^{(i)}$

learnable parameters of model.

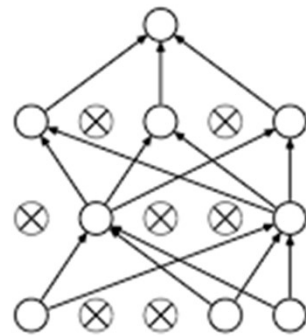
- internal covariance shift: The change in the distribution of inputs to layers in the network is referred to "internal covariate shift"
- Batch normalization accelerates training, in some cases by halving the epochs or better,
- provides some regularization, reducing generalization error

dropout

- → drop out 20% to 50% of the neurons randomly during training
- → all the neurons are activated during testing
- Weights are multiplied by the probability value (probability of drop out)
- Weights are updated for a simpler neural net
→ generalization capability of NN



(a) Standard Neural Net



(b) After applying dropout.

Best practices

- Normalize the features
 - Use Kaiming He initialization for wights
 - Use ReLu activation function for hidden layers
 - Use Softmax AF for output layer for classification, no activation for regression
 - Use RMSE cost function for regression
 - Use CrossEntropy (Cat/ Sp cat) for classification → large dataset with lot of class labels → use sparse categorical cross entropy
 - Use batch normalization
 - Use drop out or L2 regularization if over fit
 - Use ADAM optimizer
-
- With all these figured out generally, the architecture of the NNet is what needs to be hyper tuned for a larger range
 - Best architecture → NAS
 - Network pruning