# DIVING INTO THE
# MACHINE ROOM

Lecture 2

MAL2, Spring 2025

1

# DIVING INTO THE
# MACHINE ROOM

- How training a neural network works
- Activation functions
- Faster optimizers
- Learning rate scheduling
- Regularization
- General suggestions
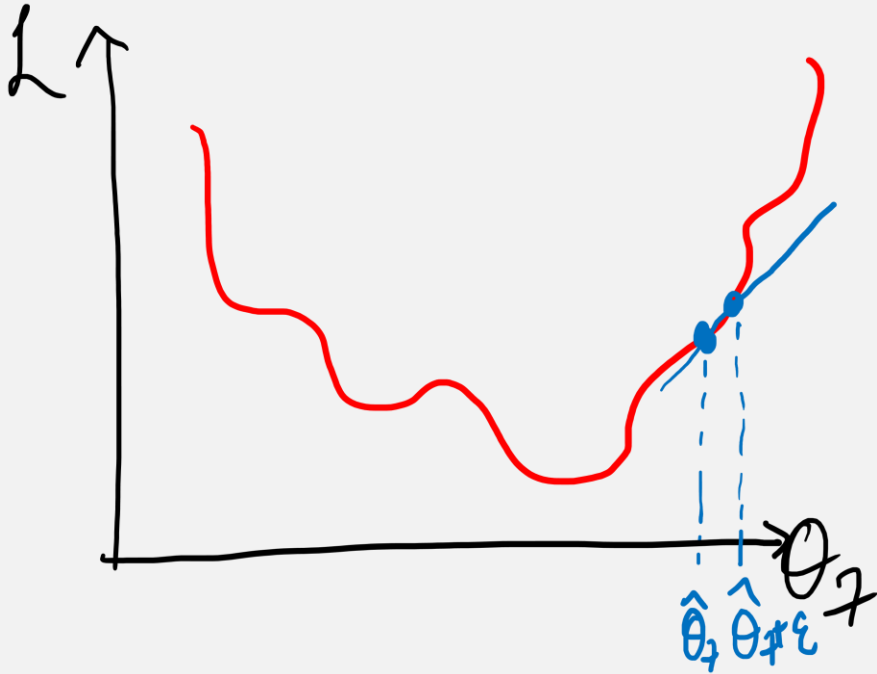
# GRADIENT DESCENT

gradient

weights & biases

$$\nabla \mathcal{L}(\Theta) = \left( \frac{\partial \mathcal{L}}{\partial \Theta_0}, \frac{\partial \mathcal{L}}{\partial \Theta_1}, \frac{\partial \mathcal{L}}{\partial \Theta_2}, \ldots, \frac{\partial \mathcal{L}}{\partial \Theta_n} \right)$$

← how do we calculate this?

loss function

$$\Theta_{new} = \Theta_{old} - \eta \nabla \mathcal{L}(\Theta)$$

repeat until $\nabla \mathcal{L}(\Theta) \approx 0$

# GRADIENT DESCENT – IDEA #1



$$\frac{\partial \mathcal{L}}{\partial \theta_7} \approx \frac{\mathcal{L}(\theta_7 + \varepsilon) - \mathcal{L}(\theta_7)}{\varepsilon}$$

INTRACTABLE

This procedure must then be repeated for all parameters at every single training step … and there may be **hundreds of thousands** of parameters!

# GRADIENT DESCENT – IDEA #2

input    hidden    output

Use the
chain rule

9 parameters
$L$ is optimized
in 9D space

$$L = MSE \sim (\hat{y} - y)^2$$

we want $\frac{\partial L}{\partial w}$
for all parameters
in the model

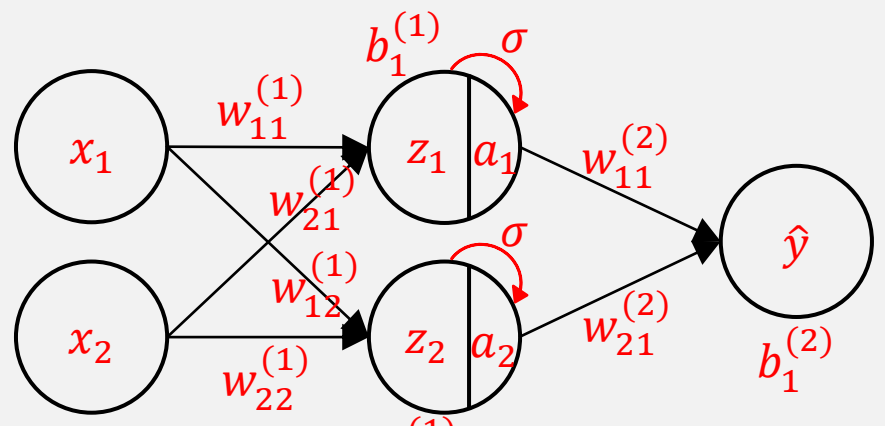$$\frac{\partial L}{\partial w_{21}^{(1)}}$$

# GRADIENT DESCENT – IDEA #2

$$\frac{\partial L}{\partial w_{21}^{(1)}}$$

$$L = (\hat{y} - y)^2$$



$$\frac{\partial L}{\partial w_{21}^{(1)}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a_1} \frac{\partial a_1}{\partial z_1} \frac{\partial z_1}{\partial w_{21}^{(1)}}$$

$$2(\hat{y}-y) \quad \cdot w_{11}^{(2)} \cdot \sigma' \cdot x_2$$

$$\hat{y} = w_{11}^{(2)} \cdot a_1 + w_{21}^{(2)} \cdot a_2 + b_1^{(2)}$$
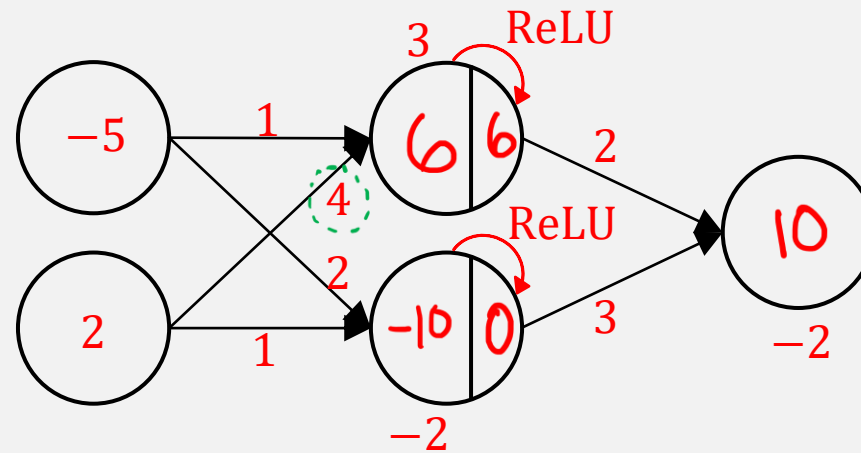
$$\sigma(z_1) \qquad \qquad \sigma(z_2)$$

$$z_1 = w_{11}^{(1)} x_1 + w_{21}^{(1)} x_2 + b_1^{(1)} \qquad z_2 = w_{12}^{(1)} x_1 + w_{22}^{(1)} x_2 + b_2^{(1)}$$
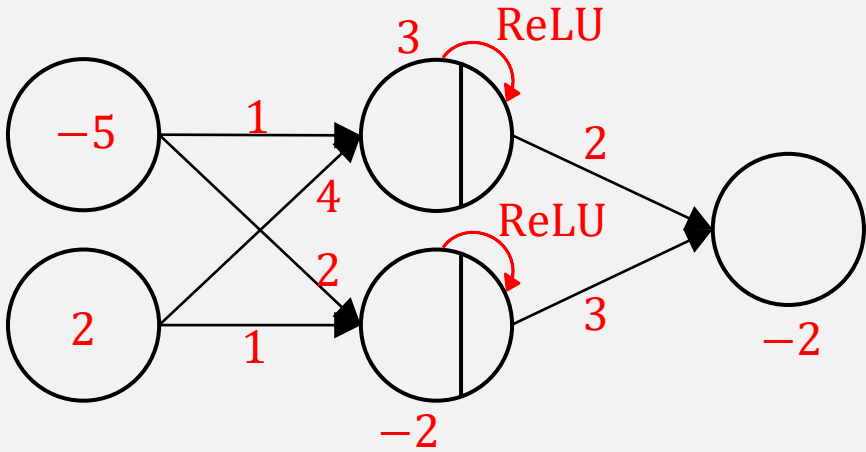
7

# THE BACKPROPAGATION ALGORITHM



$y = 8$

$\hat{y} = 10$

$$\nabla L = \begin{pmatrix} \partial L/\partial w_{11}^{(1)} \\ \partial L/\partial w_{21}^{(1)} \\ \partial L/\partial b_{1}^{(1)} \\ \partial L/\partial w_{12}^{(1)} \\ \partial L/\partial w_{22}^{(1)} \\ \partial L/\partial b_{2}^{(1)} \\ \partial L/\partial w_{11}^{(2)} \\ \partial L/\partial w_{21}^{(2)} \\ \partial L/\partial b_{1}^{(2)} \end{pmatrix} = \begin{pmatrix} -40 \\ 16 \\ 8 \\ 0 \\ 0 \\ 0 \\ 24 \\ 0 \\ 4 \end{pmatrix}$$

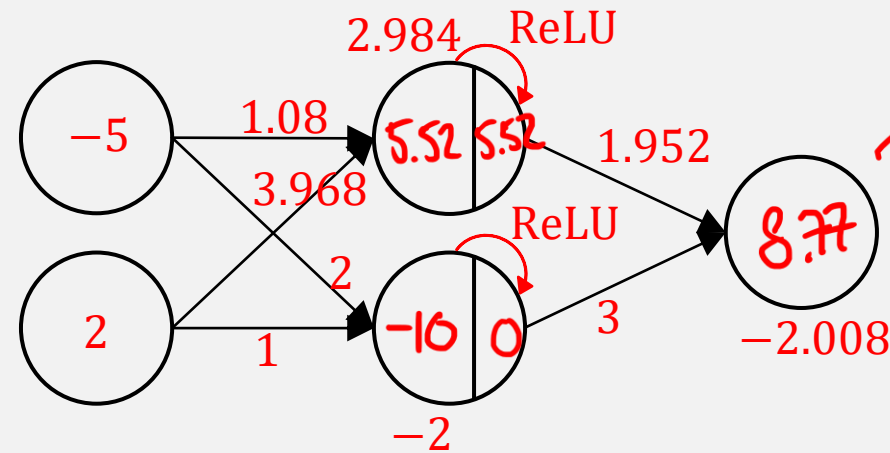$$\frac{\partial L}{\partial w_{21}^{(1)}} = 2(10-8) \cdot 2 \cdot 1 \cdot 2 = 16$$

$$\nabla L = \begin{pmatrix} \partial L / \partial w_{11}^{(1)} \\ \partial L / \partial w_{21}^{(1)} \\ \partial L / \partial b_1^{(1)} \\ \partial L / \partial w_{12}^{(1)} \\ \partial L / \partial w_{22}^{(1)} \\ \partial L / \partial b_2^{(1)} \\ \partial L / \partial w_{11}^{(2)} \\ \partial L / \partial w_{21}^{(2)} \\ \partial L / \partial b_1^{(2)} \end{pmatrix} = \begin{pmatrix} -40 \\ 16 \\ 8 \\ 0 \\ 0 \\ 0 \\ 24 \\ 0 \\ 4 \end{pmatrix}$$

$$\begin{pmatrix} w_{11}^{(1)} \\ w_{21}^{(1)} \\ b_1^{(1)} \\ w_{12}^{(1)} \\ w_{22}^{(1)} \\ b_2^{(1)} \\ w_{11}^{(2)} \\ w_{21}^{(2)} \\ b_1^{(2)} \end{pmatrix} = \begin{pmatrix} 1 \\ 4 \\ 3 \\ 2 \\ 1 \\ -2 \\ 2 \\ 3 \\ -2 \end{pmatrix} - 0.002 \begin{pmatrix} -40 \\ 16 \\ 8 \\ 0 \\ 0 \\ 0 \\ 24 \\ 0 \\ 4 \end{pmatrix} = \begin{pmatrix} 1.08 \\ 3.968 \\ 2.984 \\ 2 \\ 1 \\ -2 \\ 1.952 \\ 3 \\ -2.008 \end{pmatrix}$$

9

2.984   ReLU

−5   1.08   5.52 | 5.52   1.952   → closer to 8

3.968

2

2   1   −10 | 0   ReLU   3   8.77   now repeat

−2.008

−2

$$
\begin{pmatrix} w_{11}^{(1)} \\ w_{21}^{(1)} \\ b_1^{(1)} \\ w_{12}^{(1)} \\ w_{22}^{(1)} \\ b_2^{(1)} \\ w_{11}^{(2)} \\ w_{21}^{(2)} \\ b_1^{(2)} \end{pmatrix} = \begin{pmatrix} 1 \\ 4 \\ 3 \\ 2 \\ 1 \\ -2 \\ 2 \\ 3 \\ -2 \end{pmatrix} - 0.002 \begin{pmatrix} -40 \\ 16 \\ 8 \\ 0 \\ 0 \\ 0 \\ 24 \\ 0 \\ 4 \end{pmatrix} = \begin{pmatrix} 1.08 \\ 3.968 \\ 2.984 \\ 2 \\ 1 \\ -2 \\ 1.952 \\ 3 \\ -2.008 \end{pmatrix}
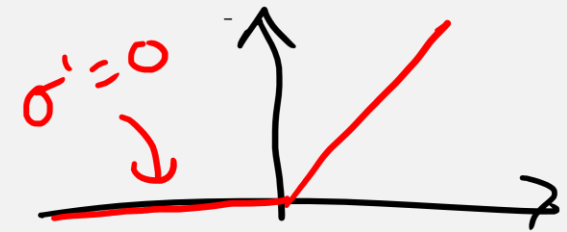$$

# THE BACKPROPAGATION ALGORITHM

Backpropagate through the network to
calculate all partial derivatives using the chain rule

**Backward pass**

Use the partial derivatives to
perform a gradient descent step

**Gradient descent**

**Calculate the loss**

Use the predictions to
calculate the loss

**Forward pass**

Data flows through the network
and predictions are recorded

$$\frac{\partial L}{\partial w_{12}^{(1)}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a_2} \frac{\partial a_2}{\partial z_2} \frac{\partial z_2}{\partial w_{12}^{(1)}} = 2(\hat{y} - y)w_{21}^{(2)}\sigma' x_1 = \boxed{0}$$

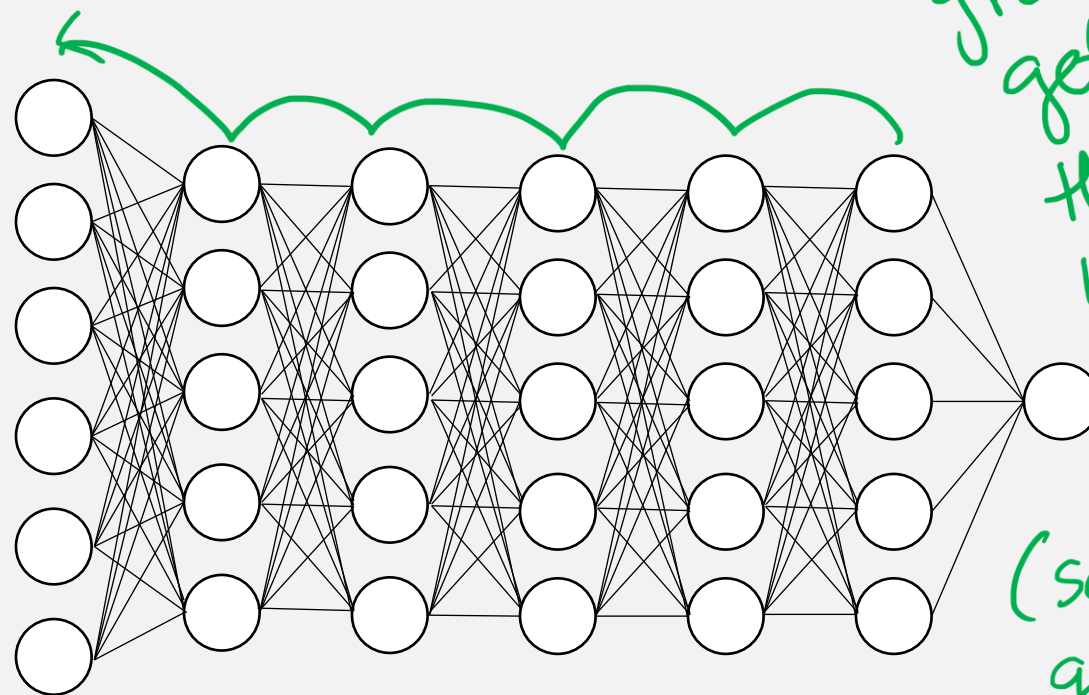$$\nabla L = \begin{pmatrix} \partial L/\partial w_{11}^{(1)} \\ \partial L/\partial w_{21}^{(1)} \\ \partial L/\partial b_1^{(1)} \\ \boxed{\partial L/\partial w_{12}^{(1)}} \\ \partial L/\partial w_{22}^{(1)} \\ \partial L/\partial b_2^{(1)} \\ \partial L/\partial w_{11}^{(2)} \\ \partial L/\partial w_{21}^{(2)} \\ \partial L/\partial b_1^{(2)} \end{pmatrix} = \begin{pmatrix} -40 \\ 16 \\ 8 \\ \boxed{0} \\ 0 \\ 0 \\ 24 \\ 0 \\ 4 \end{pmatrix}$$

*— why are these numbers 0?*

*derivative of ReLU*

*σ' = 0*

*"dying ReLU"*

*keeps outputting zero, killing the neuron*

*one of many problems we face →*

# VANISHING & EXPLODING GRADIENTS



gradients tend to
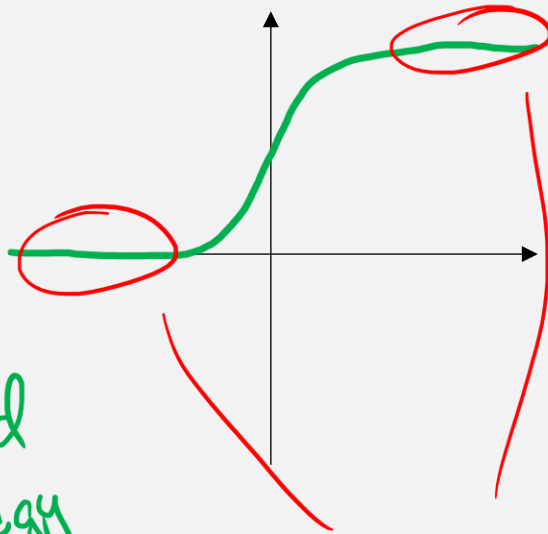get smaller & smaller,
the lower layers
hardly get trained

(sometimes they
get bigger & bigger,
then training diverges)

# DIVING INTO THE
## MACHINE ROOM

- How training a neural network works
- **Activation functions**
- Faster optimizers
- Learning rate scheduling
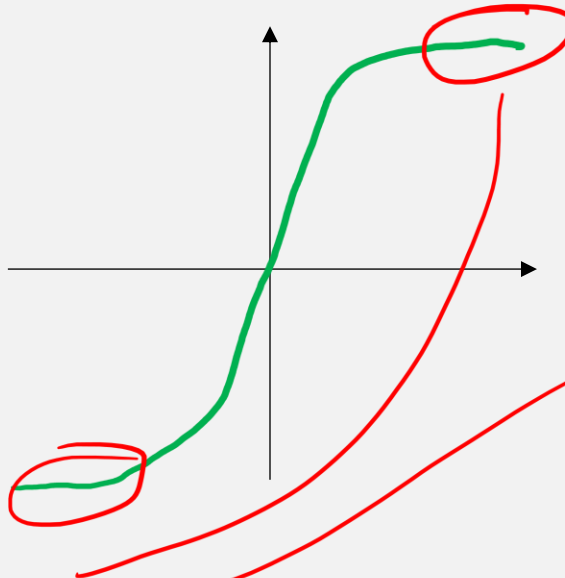- Regularization
- General suggestions
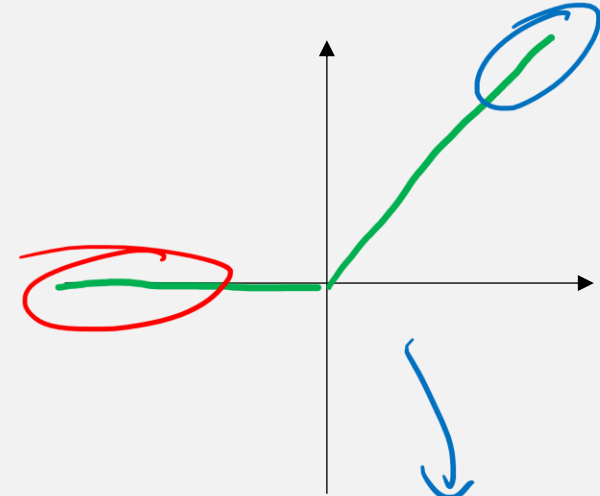
# ACTIVATION FUNCTIONS

**sigmoid**

**tanh**

**ReLU**

inspired
by biology

vanishing/dead gradients

doesn't
sturate
here

works best
for shallow
networks

# BETTER ACTIVATION FUNCTIONS

expential
ₐ linear unit

**Leaky ReLU**



not
differentiable

neurons don't die
(they go in
a coma)

$\max(\alpha z, z)$
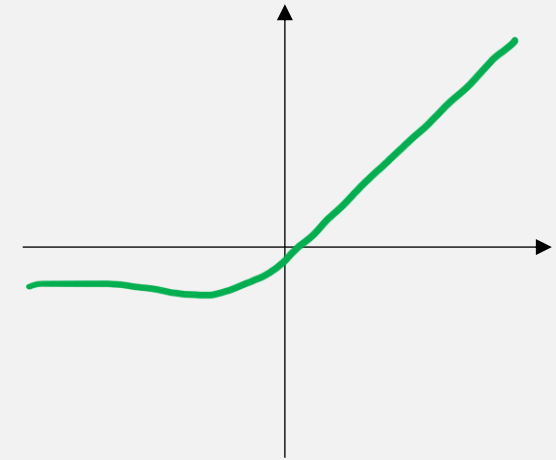
$\alpha$ hyperparameter
$0 < \alpha < 1$

**ELU**



$\begin{cases} \alpha(e^z - 1) & z < 0 \\ z & z \geq 0 \end{cases}$

outperforms all
variants of ReLU

**Swish**



$z \cdot \text{sigmoid}(z)$
outperforms everything
for deep NNs

16

# RECOMMENDATIONS

Try ReLU for shallow networks and Swish for deep networks

```
layer = Dense(100, activation="relu", kernel_initializer="he_normal")
layer = Dense(100, activation="swish", kernel_initializer="he_normal")
```

how the values of weights
and biases are initialized
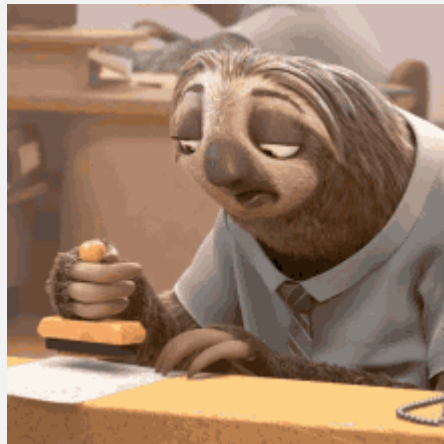
*for reference*

| Initialization method | Activation function |
|---|---|
| Glorot (default) | tanh, sigmoid, softmax |
| He | ReLU, Leaky ReLU, ELU, GELU, Swish, Mish |
| LeCun | SELU |

→ *needs*

*to match activation functions*

17

## DIVING INTO THE
# MACHINE ROOM

- How training a neural network works

- Activation functions

- **Faster optimizers**

- Learning rate scheduling

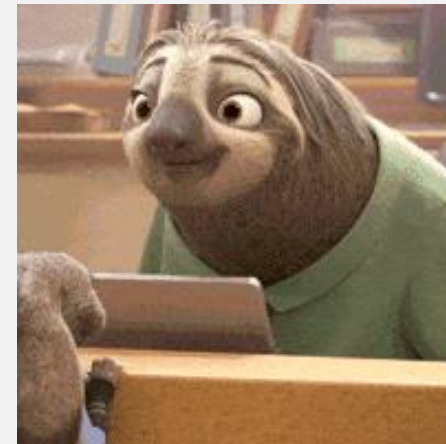- Regularization

- General suggestions

# FASTER OPTIMIZERS

#GradientDescent
UpdatingParameters

#GradientDescent
FinishingAnEpoch

#GradientDescent
WhenItConverges

The point is: Gradient descent can be painfully slow

# FASTER OPTIMIZERS

**SGD**
Stochastic gradient descent

+ momentum

+ adaptive learning rate

**Momentum**

**AdaGrad**
Adaptive gradient

+ Nesterov trick

+ scaling decay

**NAG**
Nesterov accelerated gradients

**RMSProp**
Root mean squared propagation

**Nadam**
Nesterov + Adam

**Adam**
Adaptive moment estimation

+ weight decay

**AdamW**
Adam + weight decay

# MOMENTUM

"the step we just took probably wasn't a terrible idea"
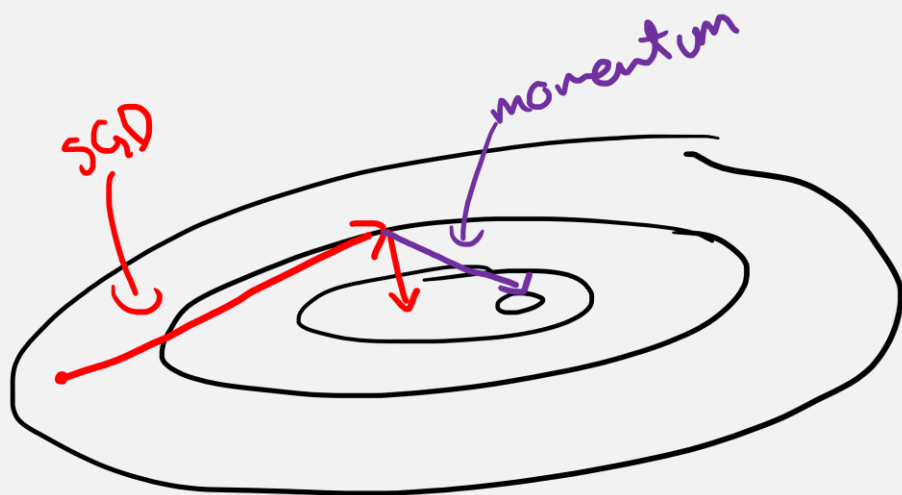
Gradient descent

$$\theta \leftarrow \theta - \eta \nabla L(\theta)$$

Momentum

momentum

$$m \leftarrow \beta m - \eta \nabla L(\theta)$$

$$\theta \leftarrow \theta + m$$



SGD

momentum

```
optimizer = SGD(learning_rate=0.001, momentum=0.9)
```

# THE NESTEROV TRICK

The trick works because $m$ generally points in the right direction, so the gradient here is slightly more accurate
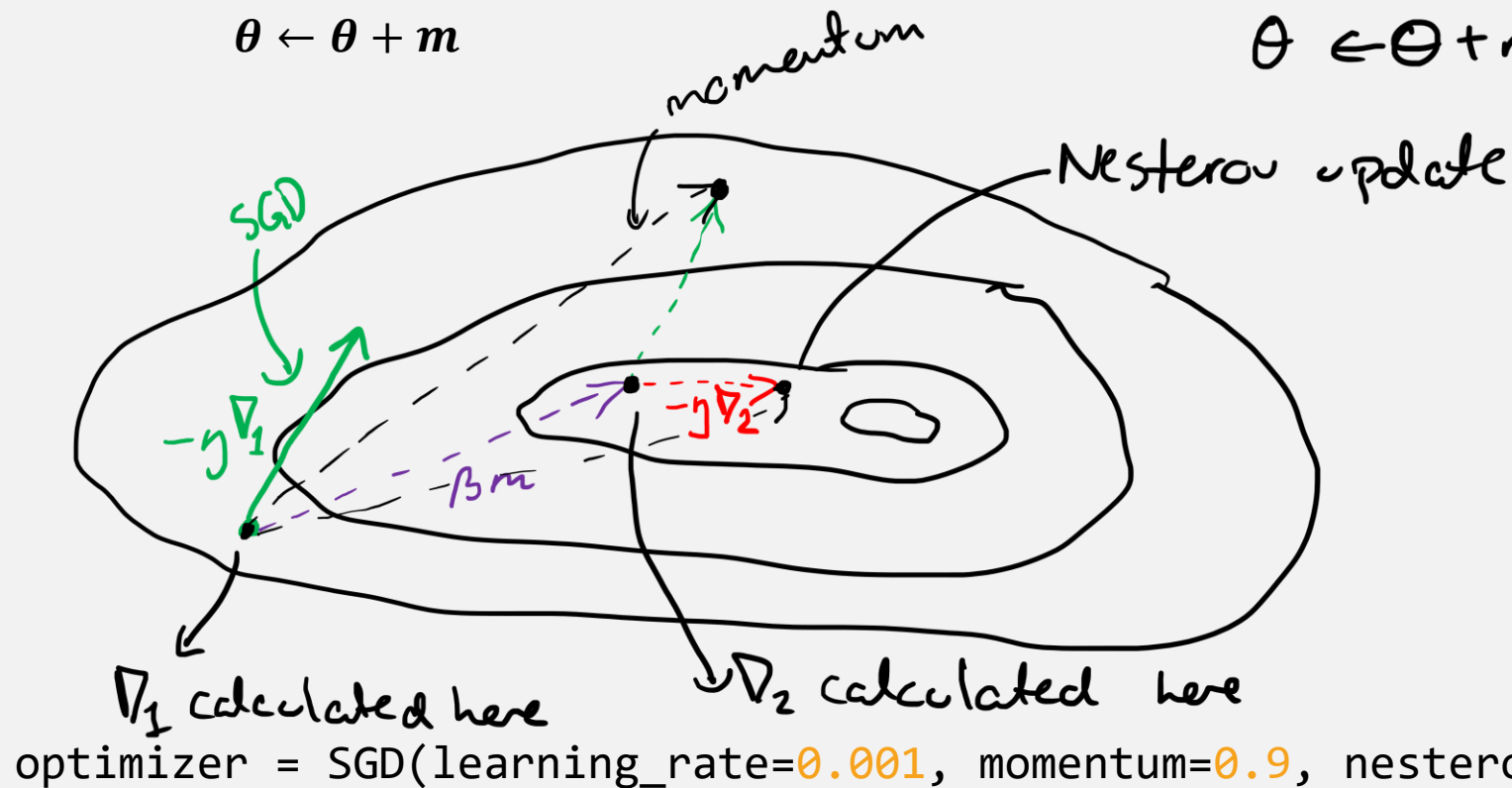
**Momentum**

$$m \leftarrow \beta m - \eta \nabla L(\theta)$$

$$\theta \leftarrow \theta + m$$

Nesterov

$$m \leftarrow \beta m - \eta \nabla L (\theta + \beta m)$$
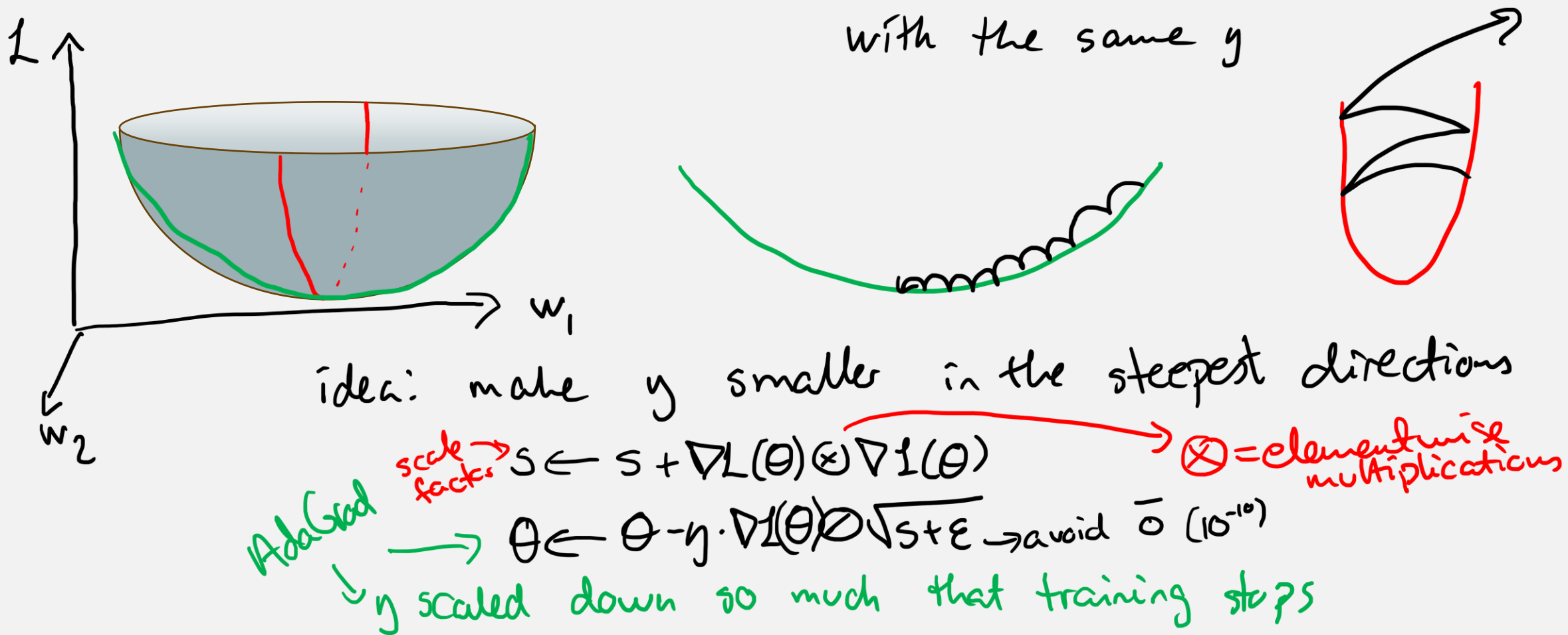
$$\theta \leftarrow \theta + m$$

Nesterov update

momentum

SGD

$-\eta \nabla_1$

$\beta m$

$-\eta \nabla_2$

$\nabla_1$ calculated here

$\nabla_2$ calculated here

```
optimizer = SGD(learning_rate=0.001, momentum=0.9, nesterov=True)
```

22

# SCALING DECAY

add decay to s so it doesn't explode

**AdaGrad**

$$s \leftarrow s + \nabla L(\boldsymbol{\theta}) \otimes \nabla L(\boldsymbol{\theta})$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \nabla L(\boldsymbol{\theta}) \oslash \sqrt{s + \varepsilon}$$ $\longrightarrow$ same thing

$\rho \sim 0.9$

RMSProp

$$s \leftarrow \rho s + (1-\rho) \nabla L(\theta) \otimes \nabla L(\theta)$$

# WEIGHT DECAY

At each training step, multiply
all weights by, say 0.99

$\rightarrow$ built-in regularization

# RECOMMENDATIONS



SGD
Stochastic gradient descent

+ momentum

+ adaptive learning rate

Momentum

AdaGrad
Adaptive gradient

+ Nesterov trick

+ scaling decay

NAG
Nesterov accelerated gradients

RMSProp
Root mean squared propagation

Nadam
Nesterov + Adam

Adam
Adaptive moment estimation

+ weight decay

AdamW
Adam + weight decay

## DIVING INTO THE
# MACHINE ROOM

- How training a neural network works

- Activation functions

- Faster optimizers

- **Learning rate scheduling**
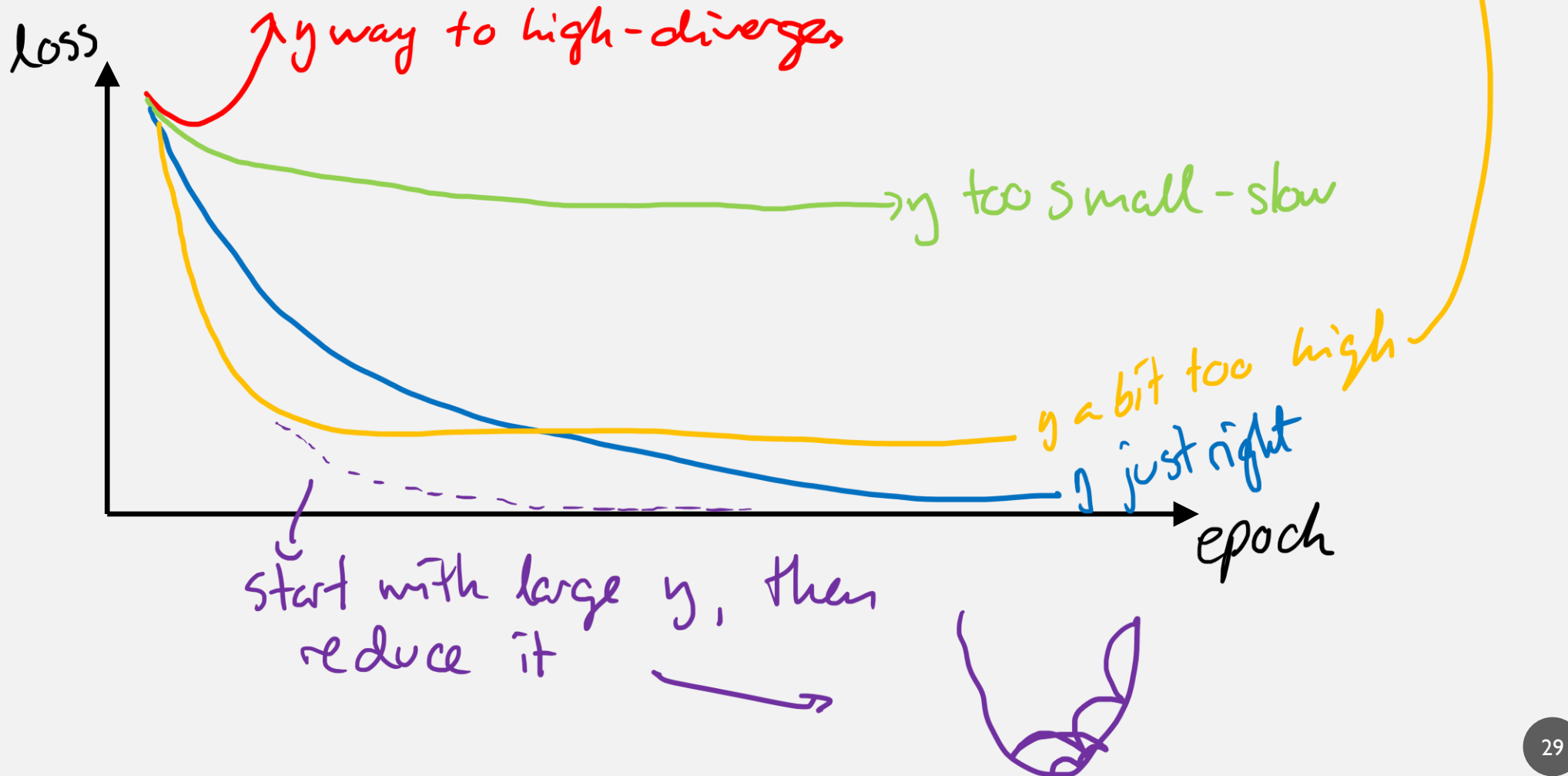
- Regularization

- General suggestions

# LEARNING RATE SCHEDULING

loss

y way to high - diverges

y too small - slow

y a bit too high

y just right

epoch

start with large y, then reduce it

# LEARNING RATE SCHEDULING

## Power

$$\eta(t) = \frac{\eta_0}{1 + t/s}$$

iteration #

after s steps: $\frac{\eta_0}{2}$

2s $\quad \frac{\eta_0}{3}$

3s $\quad \frac{\eta_0}{4}$

## Exponential

$$\eta(t) = \eta_0 \, 0.1^{t/s}$$

reduces by a factor of 10 every s steps

## Piecewise constant

$\eta = 0.1$ for 5 epochs

$\eta = 0.01$ for 10

0.001 for 50

. . . . .

## Performance

measure valid. error every N steps, divide $\eta$ by $\lambda$ when error stops dropping

All of them work pretty well – performance scheduling is a good default

# CODING A LEARNING RATE SCHEDULE

```python
#power scheduling
optimizer = SGD(learning_rate=0.01, decay=1e-4)

#exponential scheduling
def exp_decay_fn(epoch):
    return 0.01 * 0.1 ** (epoch/20)

lr_scheduler = tf.keras.callbacks.LearningRateScheduler(exp_decay_fn)
history = model.fit([...], callbacks = [lr_scheduler])

#piecewise constant scheduling
def piece_fn(epoch):
    if epoch < 5:
        return 0.01
    elif epoch < 15:
        [...]
        and so on - use LearningRateScheduler callback with piece_fn

#performance scheduling
lr_scheduler = tf.keras.callbacks.ReduceLROnPlateau(factor = 0.5, patience = 5)
history = model.fit([...], callbacks = [lr_scheduler])
```

DIVING INTO THE
# MACHINE ROOM

- How training a neural network works
- Activation functions
- Faster optimizers
- Learning rate scheduling
- **Regularization**
- General suggestions

# L1 AND L2 REGULARIZATION
*tools to avoid overfitting*

penalty for large coefficients

$$\hookrightarrow \quad \alpha \sum_i \theta_i^2 \quad (L2)$$
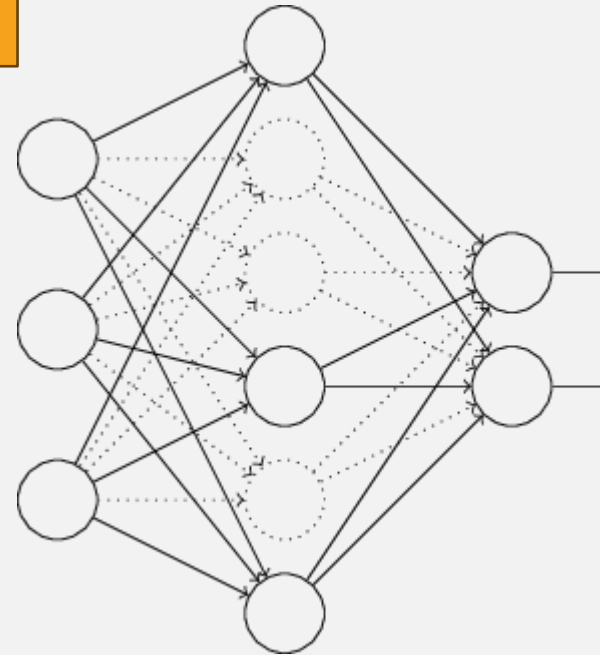
$$\alpha \sum_i |\theta_i| \quad (L1)$$

never use with AdamW – use Adam instead

$\alpha$

```
layer = Dense([...], kernel_regularizer=tf.keras.regularizers.l2(0.01))
```
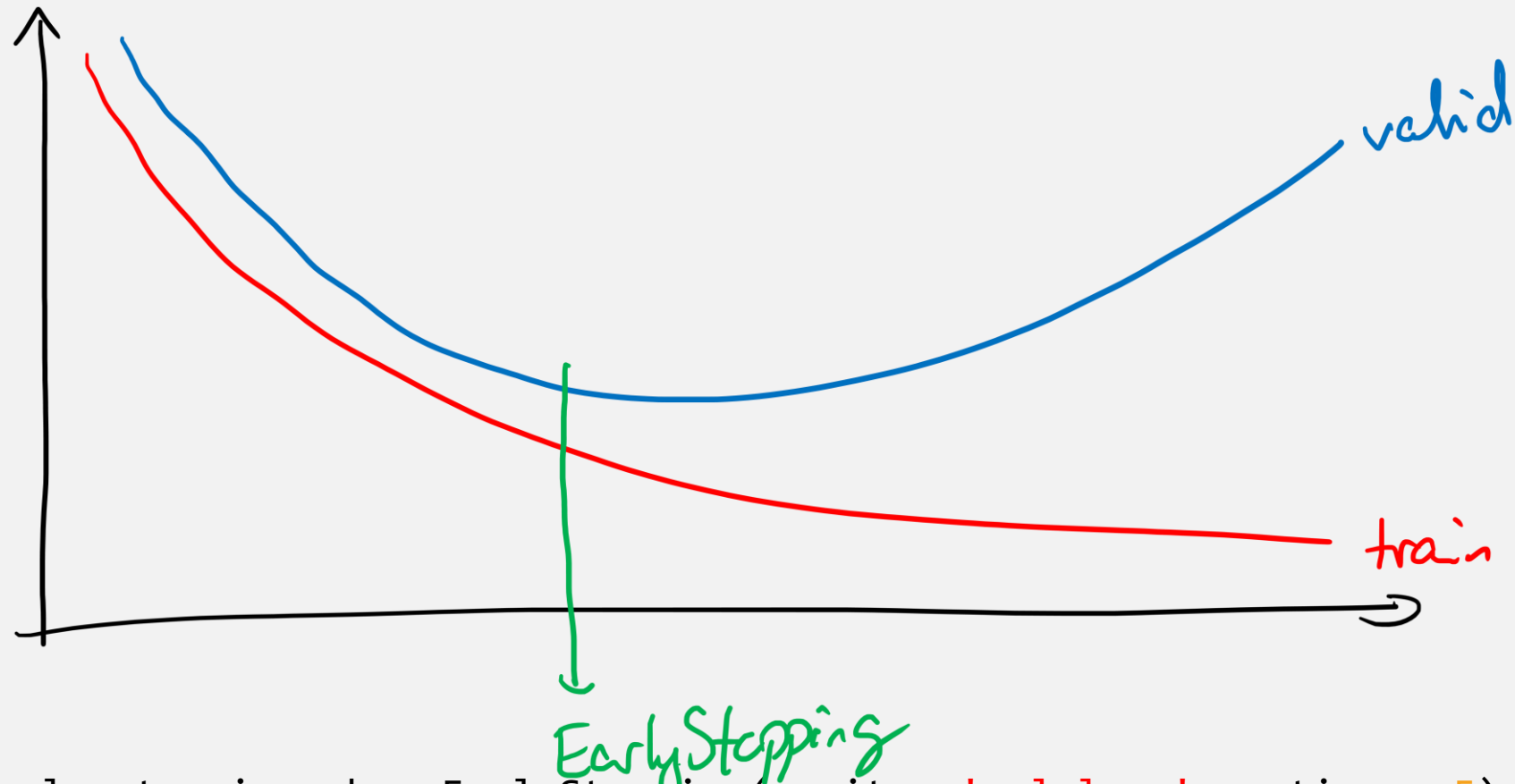
33

# DROPOUT REGULARIZATION

At every training step, every neuron has a probability $p$ of being entirely ignored!

So the model can't rely too much on any particular neuron, making it more robust

```
tf.keras.layers.Dropout(0.2)
```

# BUT THE BEST WAY TO PREVENT OVERFITTING IS USUALLY ...



valid

train

EarlyStopping

```
early_stopping_cb = EarlyStopping(monitor='val_loss', patience=5)
history = model.fit([...], callbacks = [early_stopping_cb])
```
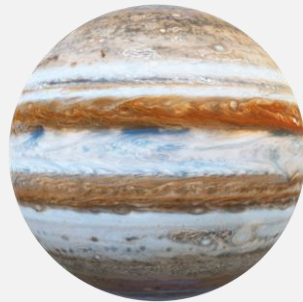
DIVING INTO THE
# MACHINE ROOM

- How training a neural network works
- Activation functions
- Faster optimizers
- Learning rate scheduling
- Regularization
- **General suggestions**

# GENERAL SUGGESTIONS

| Hyperparameter | Recommendations |
|---|---|
| Activation function | ReLU if shallow.<br>Swish if deep (Leaky ReLU as a faster alternative). |
| Optimizer | Start with Adam or AdamW, but switch to NAG if it doesn't work out. Never use SGD or AdaGrad. |
| Learning rate schedule | Performance scheduling is pretty good and requires very little hyperparameter tuning, but others may do better if you do it right. |
| Regularization | EarlyStopping and weight decay will usually do the trick, but look at others if you can't get rid of overfitting. Never use L1/L2 with Adam-type optimizers. |

# TWO TASKS

**Find your neural network from last week**



and make it better now that you are cleverer

**Scan this QR code**



and tell me about something you are still unsure about

You have 20 minutes