# Backpropagation
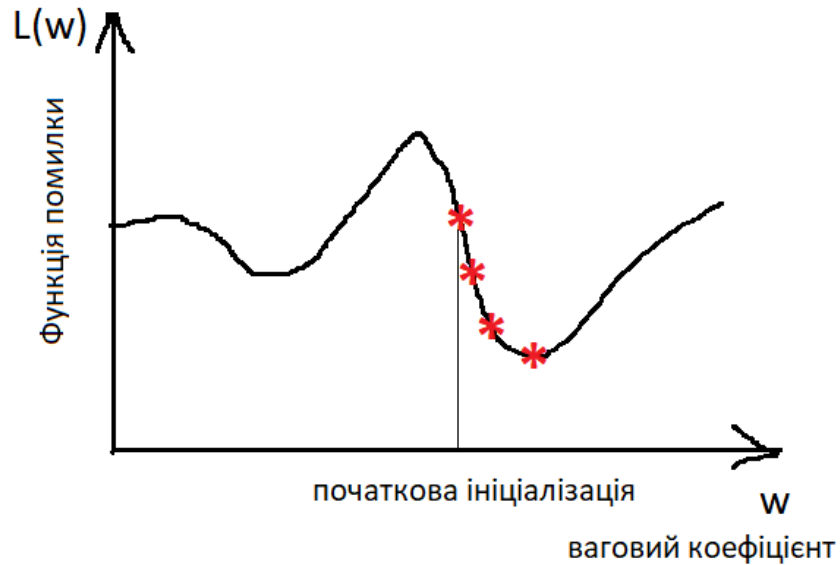
# Loss Function

$$L = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i)$$
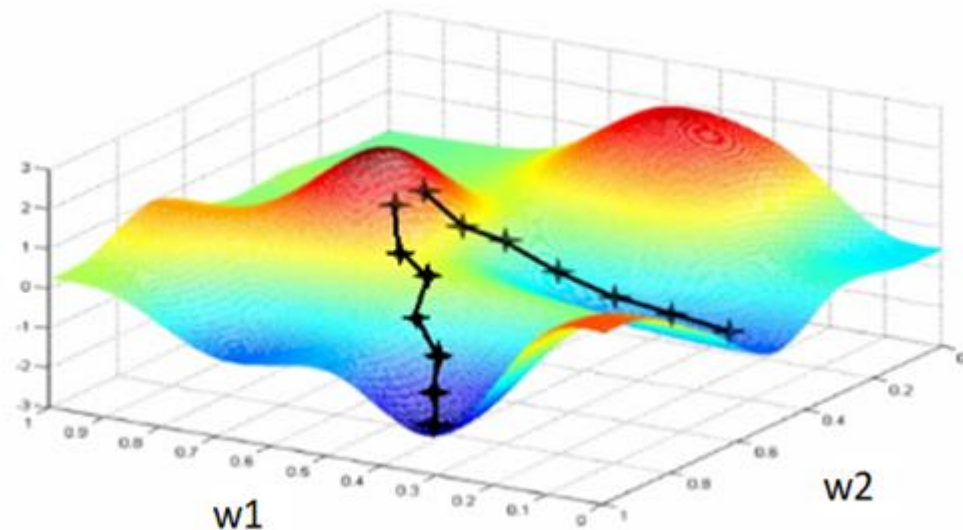
# Градієнтний спуск (Gradient descent)
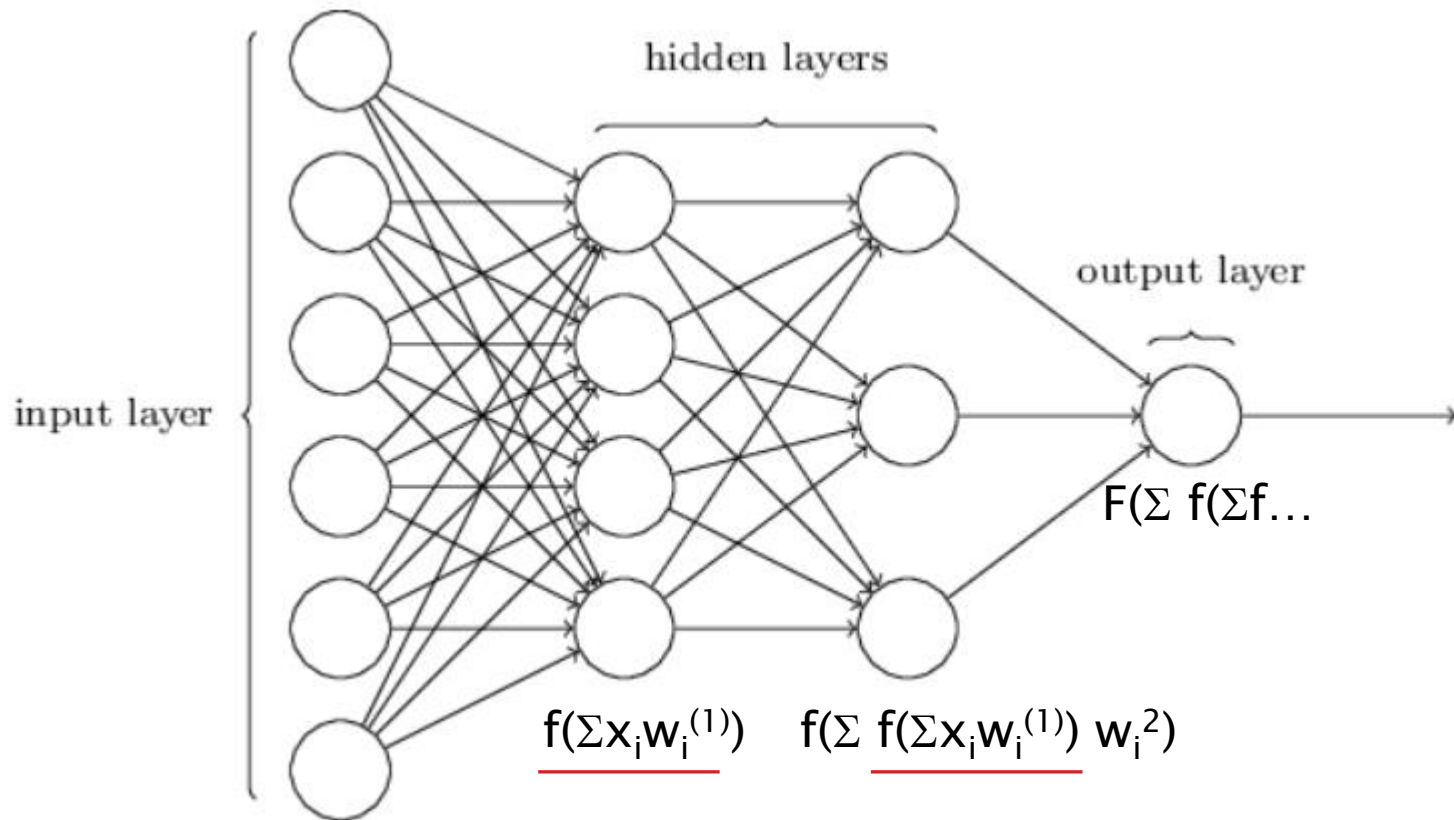
$$\vec{w} = \vec{w} - \eta \overrightarrow{\nabla_w} L$$

$$\vec{b} = \vec{b} - \eta \vec{\nabla}_b L$$



L(w)

Функція помилки

початкова ініціалізація

w

ваговий коефіцієнт

L(w1,w2)

w1

w2

input layer

hidden layers

output layer

$F(\Sigma\, f(\Sigma f\ldots$

$f(\Sigma x_i w_i^{(1)})$     $f(\Sigma\, f(\Sigma x_i w_i^{(1)})\, w_i^2)$

1986p: Д.Румельхарт, Дж.Хінтон, Р.Вільямс розробляють обчислювально ефективний алгоритм навчання нейромереж – метод зворотного поширення помилки.
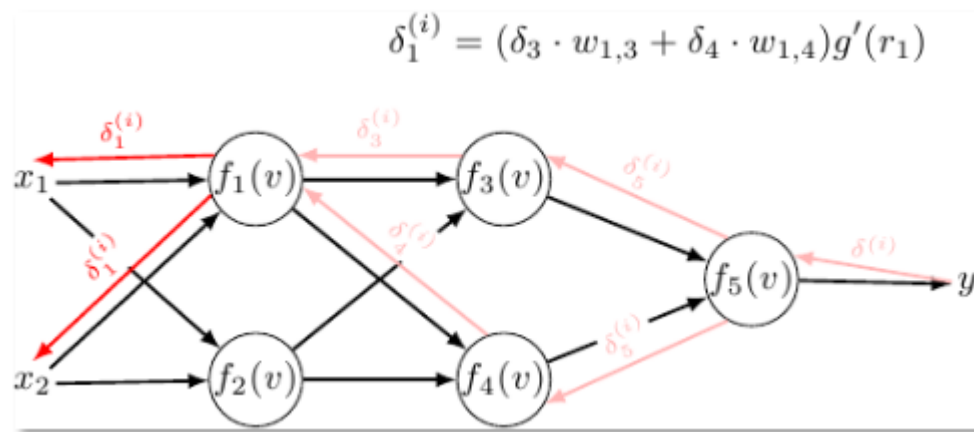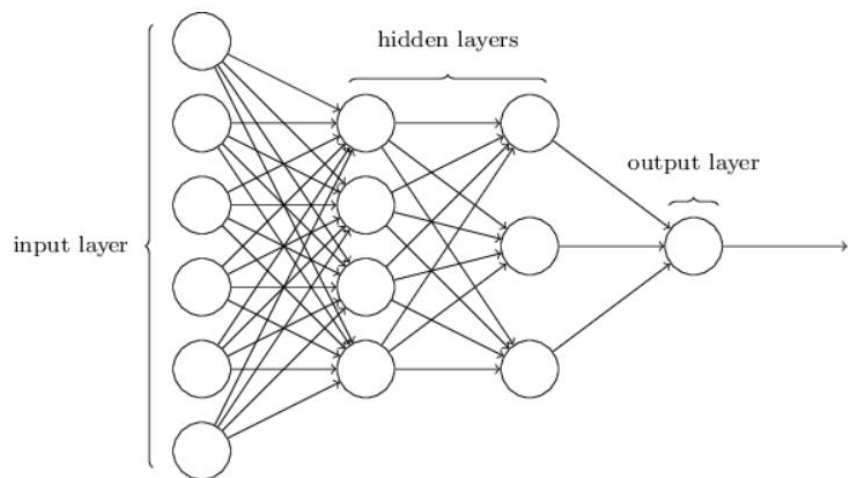


Джефрі Хінтон

Девід Румельхарт

Рональд Вільямс



Learning Internal Representations
by Error Propagation

D. E. RUMELHART, G. E. HINTON, and R. J. WILLIAMS

**THE PROBLEM**

We now have a rather good understanding of simple two-layer associative networks in which a set of input patterns arriving at an input layer are mapped directly to a set of output patterns at an output layer. Such networks have no *hidden* units. They involve only *input* and *output* units. In these cases there is no *internal representation*. The coding provided by the external world must suffice. These networks have proved useful in a wide variety of applications (cf. Chapters 2, 17, and 18). Perhaps the essential character of such networks is that they map similar input patterns to similar output patterns. This is what allows these networks to make reasonable generalizations and perform reasonably on patterns that have never before been presented. The similarity of patterns in a PDP system is determined by their overlap. The overlap in such networks is determined outside the learning system itself—by whatever produces the patterns.

The constraint that similar input patterns lead to similar outputs can lead to an inability of the system to learn certain mappings from input to output. Whenever the representation provided by the outside world is such that the similarity structure of the input and output patterns are very different, a network without internal representations (i.e., a

$$\delta_1^{(i)} = (\delta_3 \cdot w_{1,3} + \delta_4 \cdot w_{1,4})g'(r_1)$$
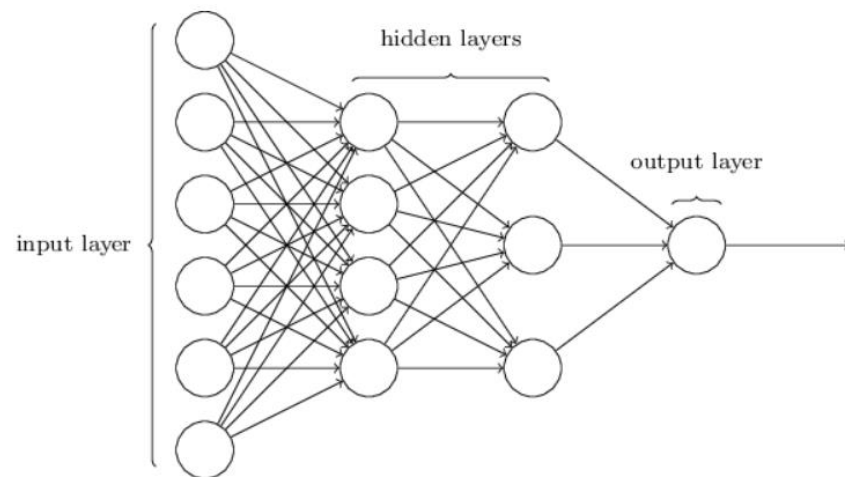
forward pass

прямий прохід

обчислюємо значення нейронів та функцію втрат

backward pass

зворотній прохід

обчислюємо похідні та корегуємо вагові коефіцієнти

# Chain rule

$f(g(x))$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial g} \cdot \frac{\partial g}{\partial x}$$

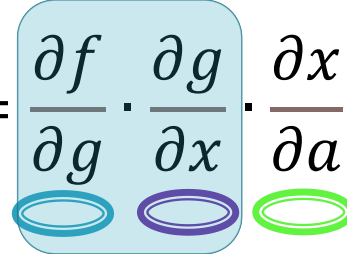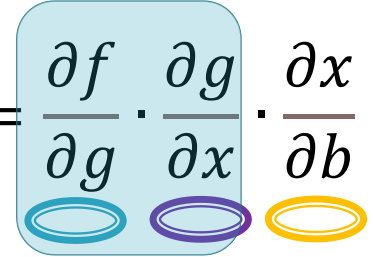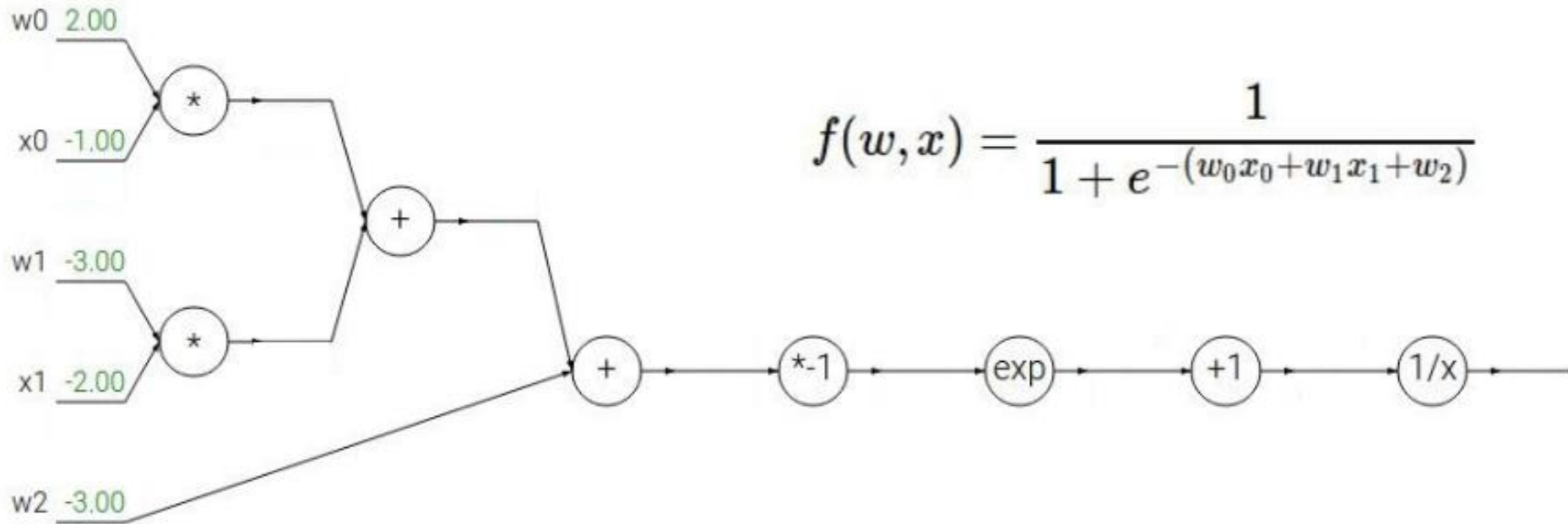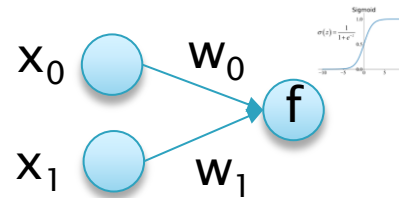upstream gradient            local gradient

$f(g(x(a)))$

$$\frac{\partial f}{\partial a} = \frac{\partial f}{\partial g} \cdot \frac{\partial g}{\partial x} \cdot \frac{\partial x}{\partial a} = \frac{\partial f}{\partial x} \cdot \frac{\partial x}{\partial a}$$
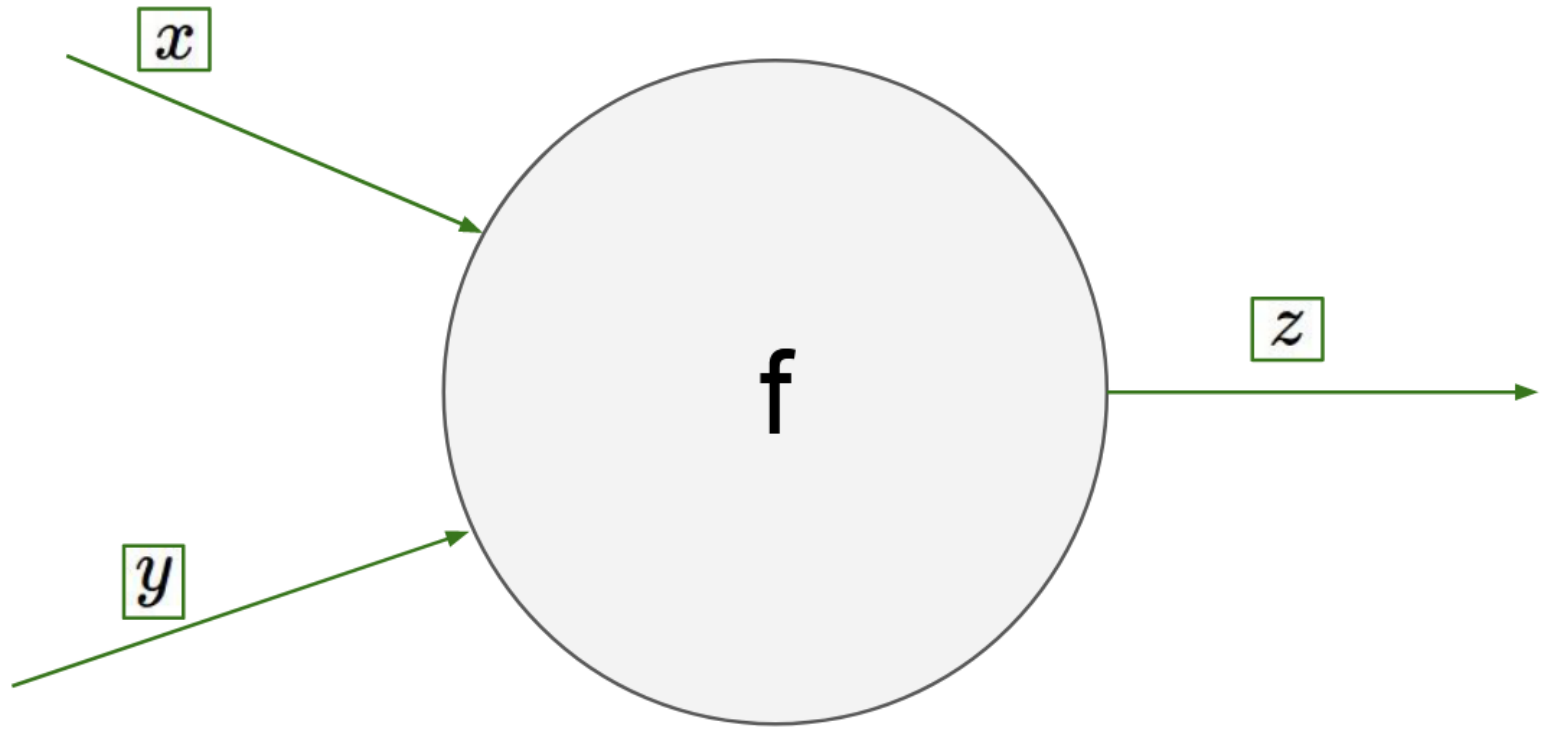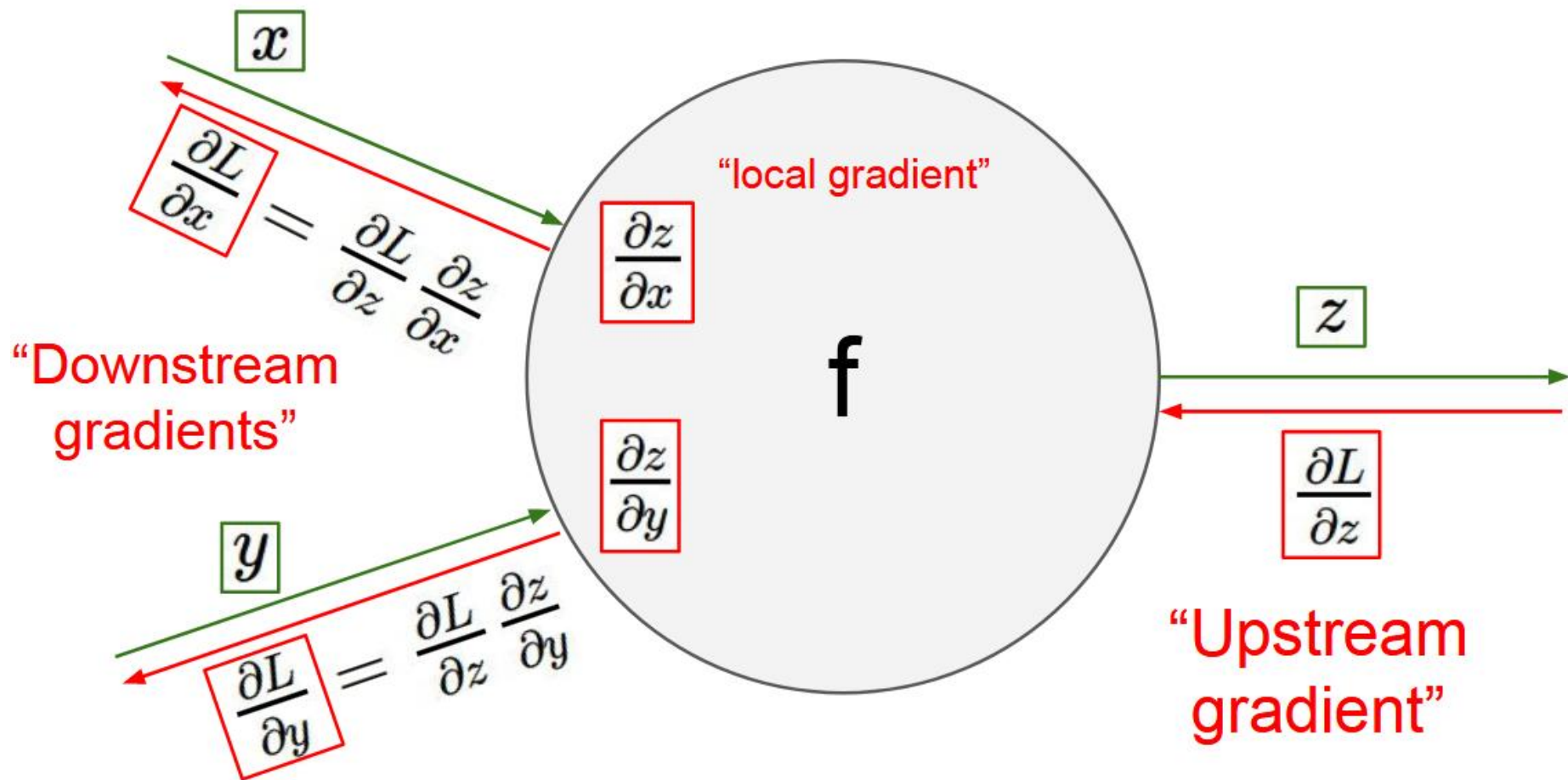
$f(g(x(a,b)))$

$$\frac{\partial f}{\partial a} = \frac{\partial f}{\partial g} \cdot \frac{\partial g}{\partial x} \cdot \frac{\partial x}{\partial a}$$

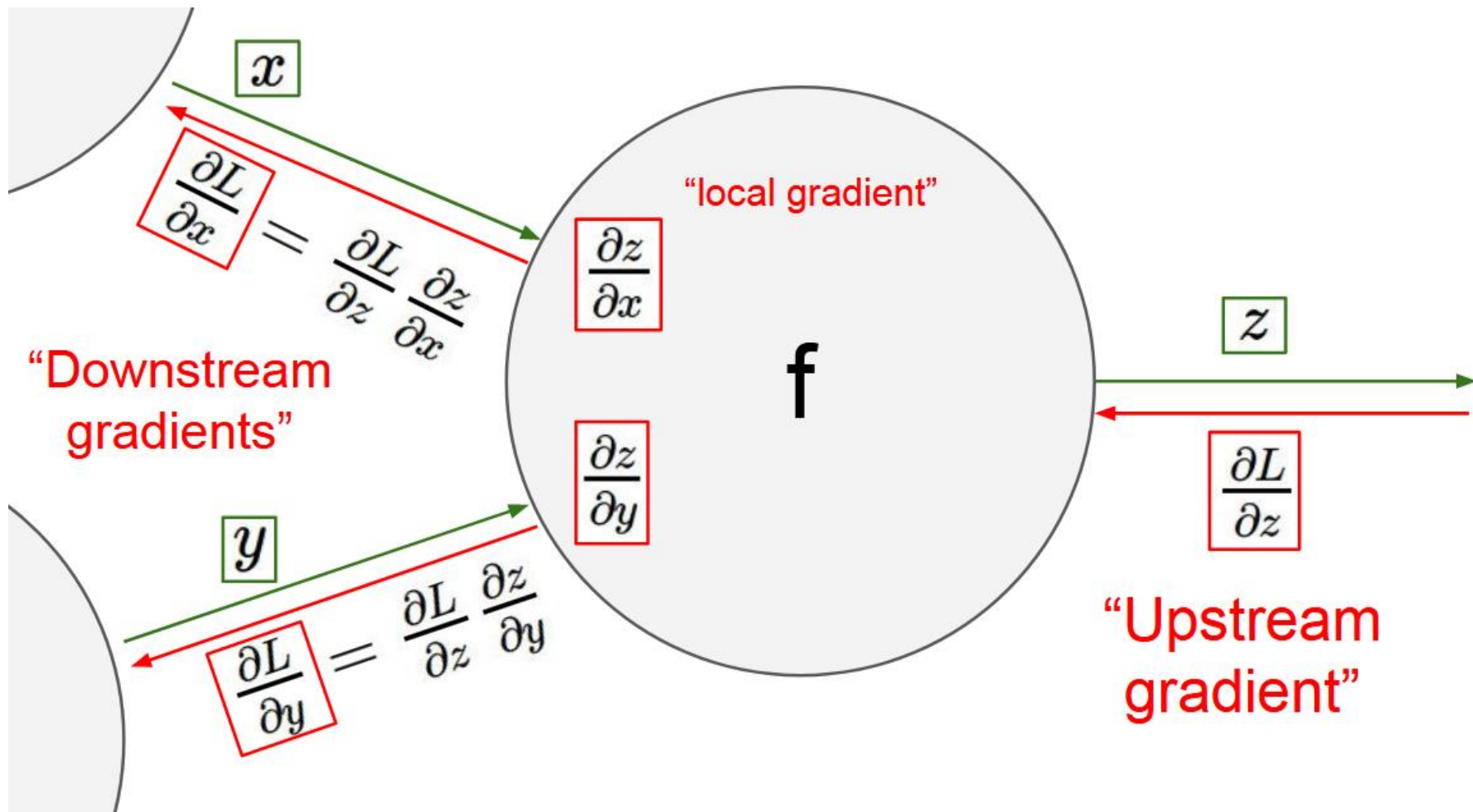$$\frac{\partial f}{\partial b} = \frac{\partial f}{\partial g} \cdot \frac{\partial g}{\partial x} \cdot \frac{\partial x}{\partial b}$$

# Computational graph


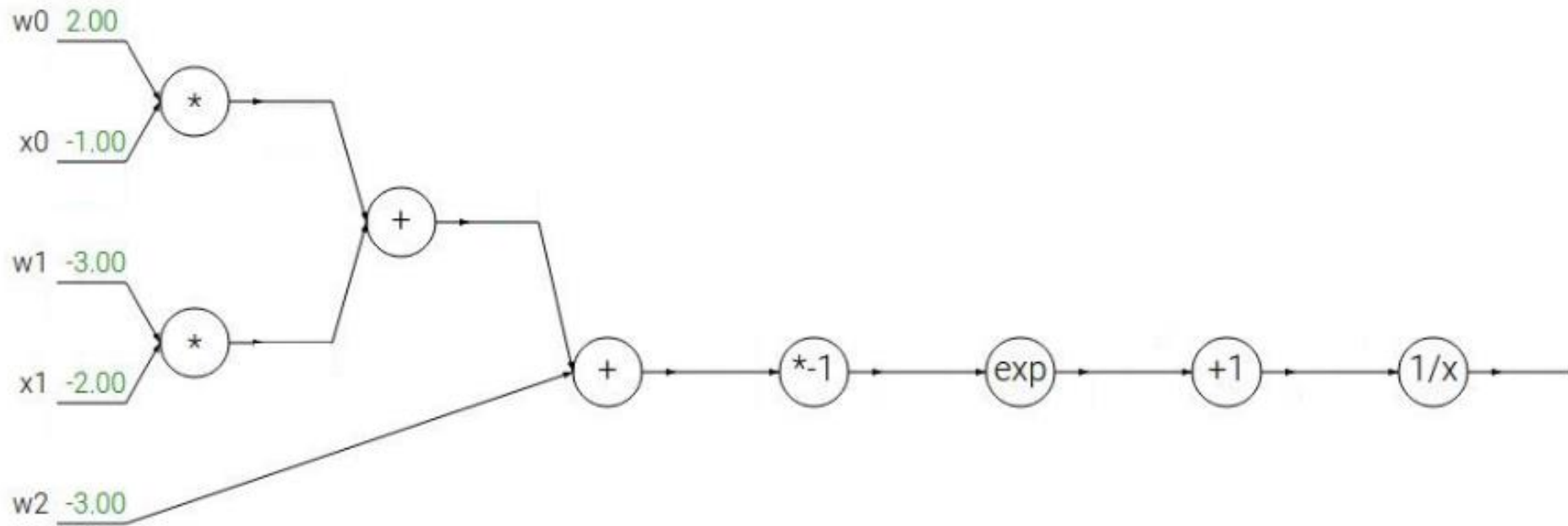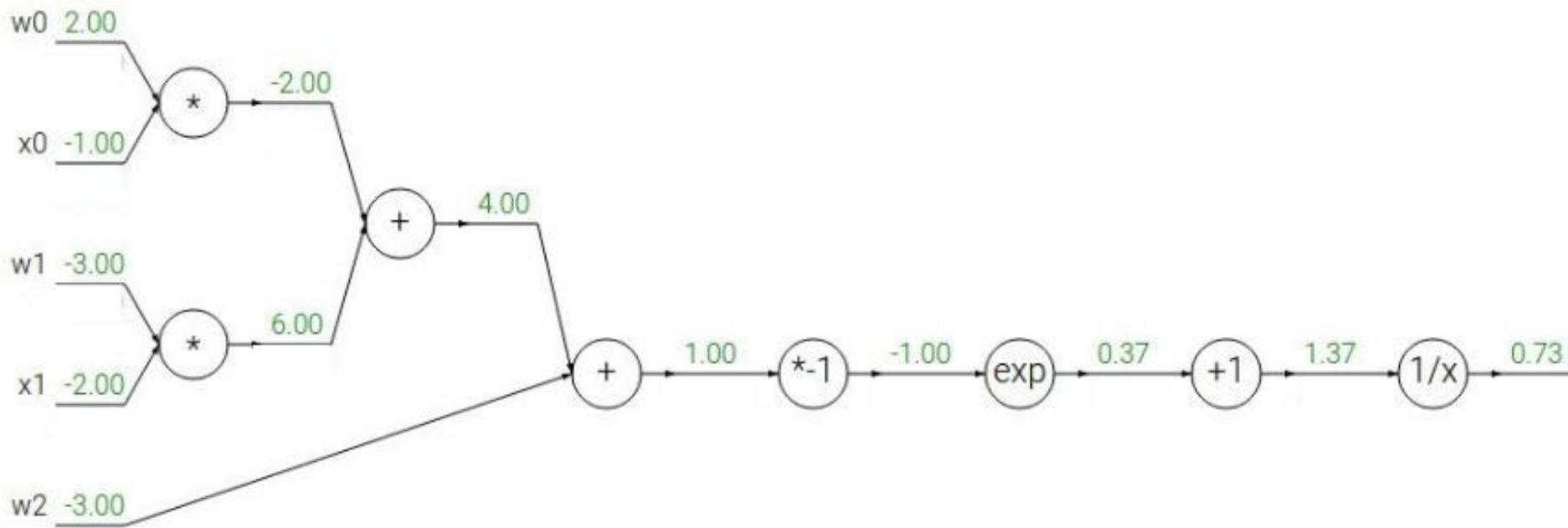
$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

$x$

$y$

$z$

**f**

$x$

$\dfrac{\partial L}{\partial x} = \dfrac{\partial L}{\partial z}\dfrac{\partial z}{\partial x}$

"Downstream gradients"

$y$

$\dfrac{\partial L}{\partial y} = \dfrac{\partial L}{\partial z}\dfrac{\partial z}{\partial y}$

"local gradient"

$\dfrac{\partial z}{\partial x}$

f

$\dfrac{\partial z}{\partial y}$

$z$

$\dfrac{\partial L}{\partial z}$
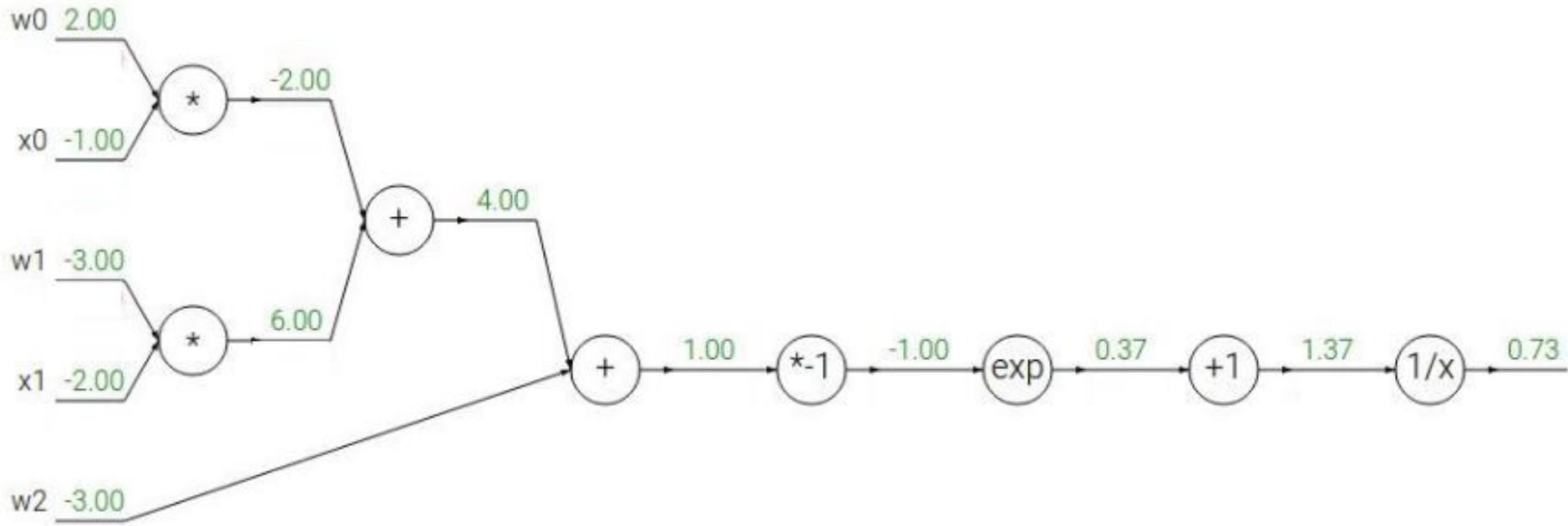
"Upstream gradient"

# Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

# Another example:

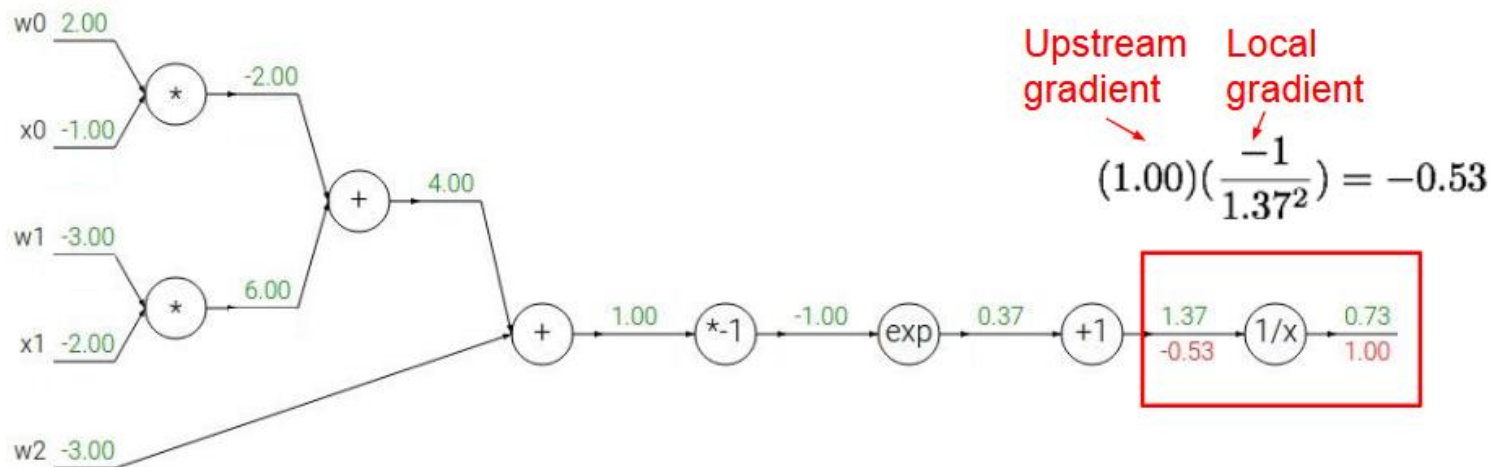$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

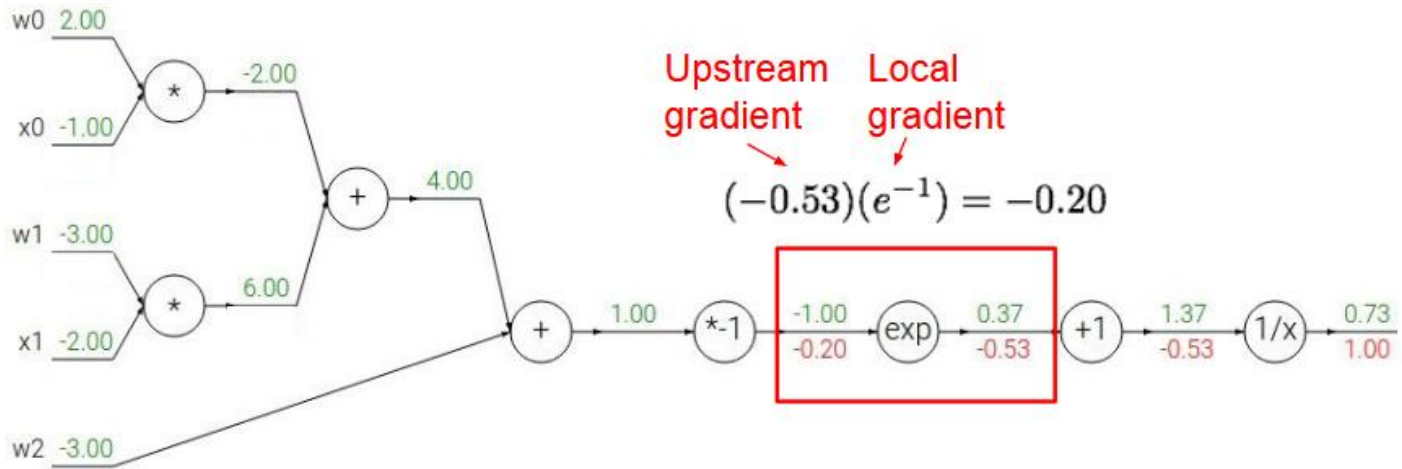# Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x \qquad \bigg| \qquad f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a \qquad \bigg| \qquad f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

# Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



Upstream gradient    Local gradient

$$(1.00)(\frac{-1}{1.37^2}) = -0.53$$

$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x \qquad \Bigg| \qquad f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a \qquad \Bigg| \qquad f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

# Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



Upstream gradient · Local gradient

$$(-0.53)(1) = -0.53$$

$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

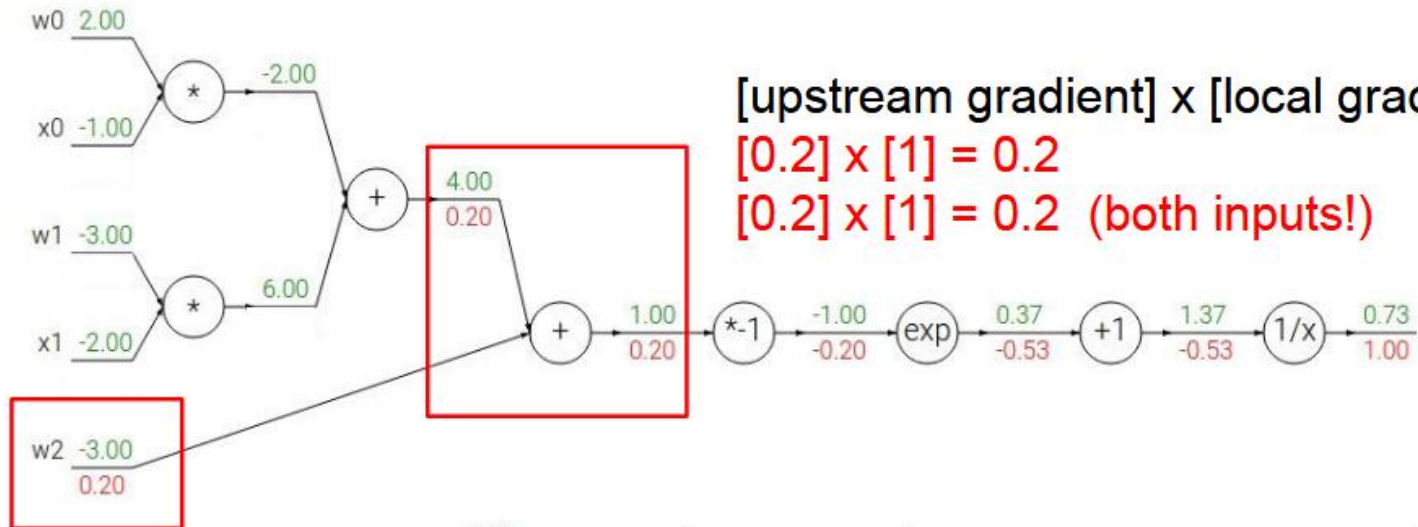# Another example:

$$f(w,x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



Upstream gradient   Local gradient

$$(-0.53)(e^{-1}) = -0.20$$

| | | |
|---|---|---|
| $f(x) = e^x$ | $\rightarrow$ | $\frac{df}{dx} = e^x$ |
| $f_a(x) = ax$ | $\rightarrow$ | $\frac{df}{dx} = a$ |

| | | |
|---|---|---|
| $f(x) = \frac{1}{x}$ | $\rightarrow$ | $\frac{df}{dx} = -1/x^2$ |
| $f_c(x) = c + x$ | $\rightarrow$ | $\frac{df}{dx} = 1$ |

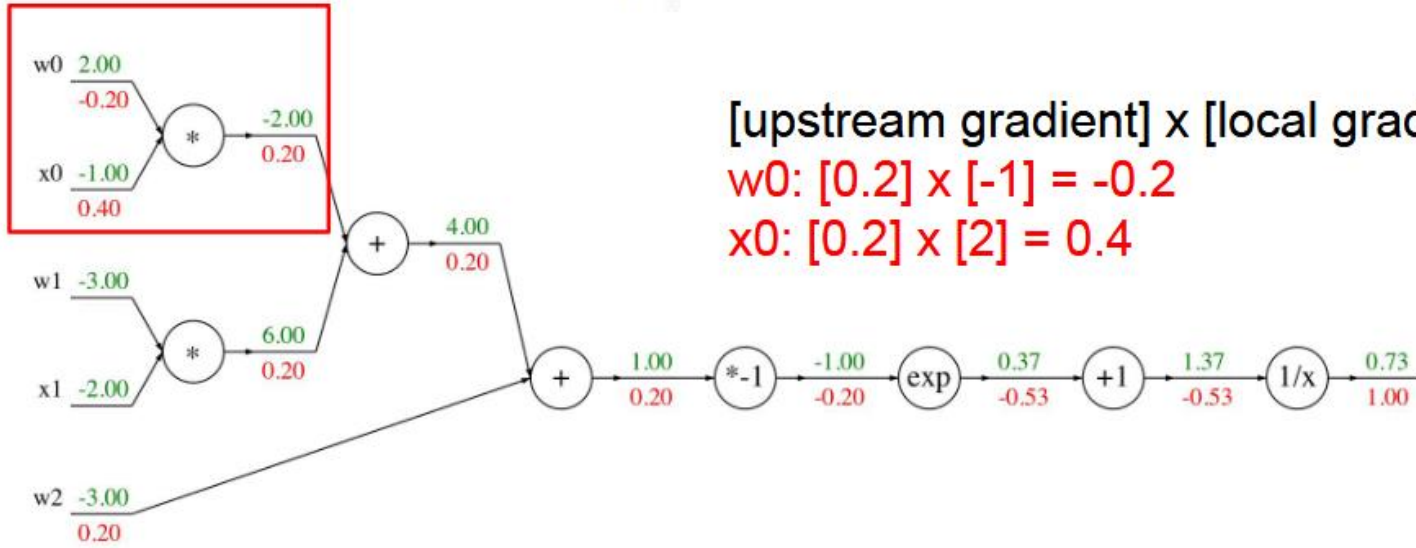# Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



Upstream gradient    Local gradient

$$(-0.20)(-1) = 0.20$$

$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

# Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



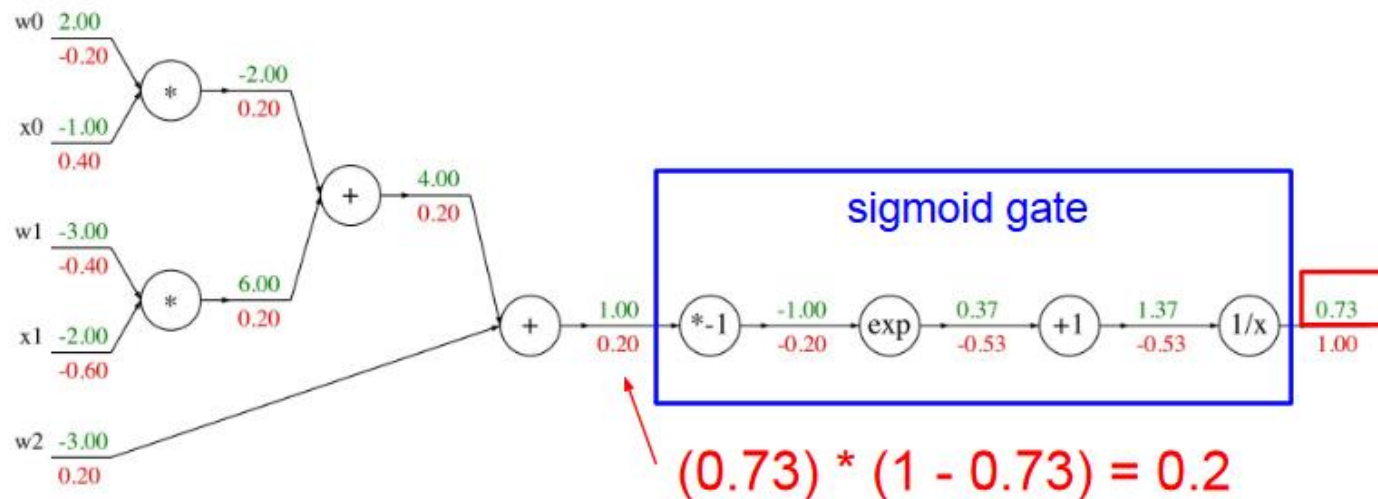[upstream gradient] x [local gradient]
[0.2] x [1] = 0.2
[0.2] x [1] = 0.2  (both inputs!)

$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

# Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



[upstream gradient] x [local gradient]
w0: [0.2] x [-1] = -0.2
x0: [0.2] x [2] = 0.4

$f(x) = e^x$ $\rightarrow$ $\dfrac{df}{dx} = e^x$

$f_a(x) = ax$ $\rightarrow$ $\dfrac{df}{dx} = a$

$f(x) = \dfrac{1}{x}$ $\rightarrow$ $\dfrac{df}{dx} = -1/x^2$

$f_c(x) = c + x$ $\rightarrow$ $\dfrac{df}{dx} = 1$

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$
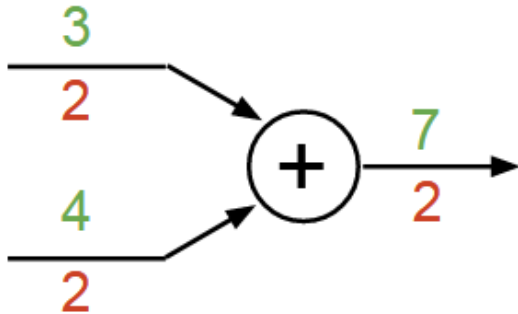
$$\boxed{\sigma(x) = \frac{1}{1 + e^{-x}}}$$ **sigmoid function**

$$\frac{d\sigma(x)}{dx} = \frac{e^{-x}}{(1 + e^{-x})^2} = \left(\frac{1 + e^{-x} - 1}{1 + e^{-x}}\right) \left(\frac{1}{1 + e^{-x}}\right) = (1 - \sigma(x))\,\sigma(x)$$



sigmoid gate

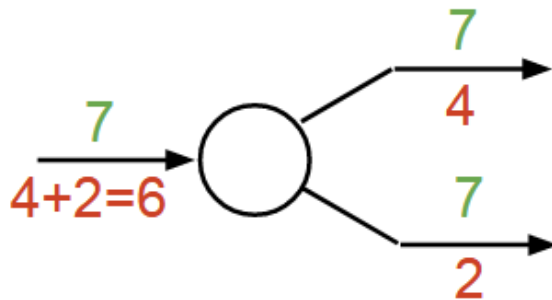(0.73) * (1 - 0.73) = 0.2

# Patterns in gradient flow

**add** gate: gradient distributor



**mul** gate: "swap multiplier"



**copy** gate: gradient adder
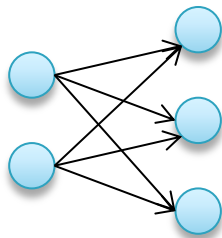


**max** gate: gradient router

$$W = \begin{pmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{pmatrix}$$

$$\frac{\partial L}{\partial W} - ? \qquad \frac{\partial L}{\partial X} - ?$$

$$X = \begin{pmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{pmatrix}$$



$$Y = \begin{pmatrix} y_{11} & y_{12} & y_{13} \\ y_{21} & y_{22} & y_{23} \end{pmatrix}$$

$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial Y}\frac{\partial Y}{\partial W} \qquad \frac{\partial L}{\partial X} = \frac{\partial L}{\partial Y}\frac{\partial Y}{\partial X}$$

$$\frac{\partial L}{\partial W} = \begin{pmatrix} \dfrac{\partial L}{\partial w_{11}} & \dfrac{\partial L}{\partial w_{12}} & \dfrac{\partial L}{\partial w_{13}} \\[2mm] \dfrac{\partial L}{\partial w_{21}} & \dfrac{\partial L}{\partial w_{22}} & \dfrac{\partial L}{\partial w_{23}} \end{pmatrix} \qquad \frac{\partial L}{\partial X} = \begin{pmatrix} \dfrac{\partial L}{\partial x_{11}} & \dfrac{\partial L}{\partial x_{12}} \\[2mm] \dfrac{\partial L}{\partial x_{21}} & \dfrac{\partial L}{\partial x_{22}} \end{pmatrix} \qquad \frac{\partial L}{\partial Y} = \begin{pmatrix} \dfrac{\partial L}{\partial y_{11}} & \dfrac{\partial L}{\partial y_{12}} & \dfrac{\partial L}{\partial y_{13}} \\[2mm] \dfrac{\partial L}{\partial y_{21}} & \dfrac{\partial L}{\partial y_{22}} & \dfrac{\partial L}{\partial y_{23}} \end{pmatrix}$$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial Y}\frac{\partial Y}{\partial w_{11}} = \begin{pmatrix} \dfrac{\partial L}{\partial y_{11}} & \dfrac{\partial L}{\partial y_{12}} & \dfrac{\partial L}{\partial y_{13}} \\[2mm] \dfrac{\partial L}{\partial y_{21}} & \dfrac{\partial L}{\partial y_{22}} & \dfrac{\partial L}{\partial y_{23}} \end{pmatrix} \begin{pmatrix} x_{11} & 0 & 0 \\ x_{21} & 0 & 0 \end{pmatrix} = \frac{\partial L}{\partial y_{11}}x_{11} + \frac{\partial L}{\partial y_{21}}x_{21}$$

$$\frac{\partial Y}{\partial w_{11}} = \begin{pmatrix} x_{11} & 0 & 0 \\ x_{21} & 0 & 0 \end{pmatrix} \quad \frac{\partial Y}{\partial w_{12}} = \begin{pmatrix} 0 & x_{11} & 0 \\ 0 & x_{21} & 0 \end{pmatrix} \quad \frac{\partial Y}{\partial w_{13}} = \begin{pmatrix} 0 & 0 & x_{11} \\ 0 & 0 & x_{21} \end{pmatrix}$$
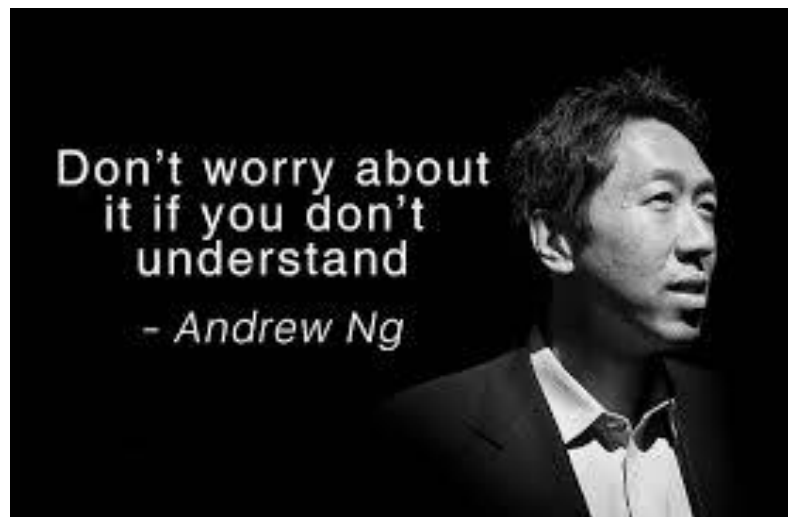
$$\frac{\partial Y}{\partial w_{21}} = \begin{pmatrix} x_{12} & 0 & 0 \\ x_{22} & 0 & 0 \end{pmatrix} \quad \frac{\partial Y}{\partial w_{22}} = \begin{pmatrix} 0 & x_{12} & 0 \\ 0 & x_{22} & 0 \end{pmatrix} \quad \frac{\partial Y}{\partial w_{23}} = \begin{pmatrix} 0 & 0 & x_{12} \\ 0 & 0 & x_{22} \end{pmatrix}$$

$$\frac{\partial L}{\partial W} = \begin{pmatrix} \frac{\partial L}{\partial w_{11}} & \frac{\partial L}{\partial w_{12}} & \frac{\partial L}{\partial w_{13}} \\ \frac{\partial L}{\partial w_{21}} & \frac{\partial L}{\partial w_{22}} & \frac{\partial L}{\partial w_{23}} \end{pmatrix} = \begin{pmatrix} (\frac{\partial L}{\partial y_{11}}x_{11} + \frac{\partial L}{\partial y_{21}}x_{21}) & (\frac{\partial L}{\partial y_{12}}x_{11} + \frac{\partial L}{\partial y_{22}}x_{21}) & (\frac{\partial L}{\partial y_{13}}x_{11} + \frac{\partial L}{\partial y_{23}}x_{21}) \\ (\frac{\partial L}{\partial y_{11}}x_{12} + \frac{\partial L}{\partial y_{21}}x_{22}) & (\frac{\partial L}{\partial y_{12}}x_{12} + \frac{\partial L}{\partial y_{22}}x_{22}) & (\frac{\partial L}{\partial y_{13}}x_{12} + \frac{\partial L}{\partial y_{23}}x_{22}) \end{pmatrix}$$

$$= \begin{pmatrix} x_{11} & x_{21} \\ x_{12} & x_{22} \end{pmatrix} \begin{pmatrix} \frac{\partial L}{\partial y_{11}} & \frac{\partial L}{\partial y_{12}} & \frac{\partial L}{\partial y_{13}} \\ \frac{\partial L}{\partial y_{21}} & \frac{\partial L}{\partial y_{22}} & \frac{\partial L}{\partial y_{23}} \end{pmatrix} = X^T \frac{\partial L}{\partial Y}$$

$$\boxed{\frac{\partial L}{\partial W} = X^T \frac{\partial L}{\partial Y}}$$

$$\boxed{\frac{\partial L}{\partial X} = \frac{\partial L}{\partial Y} W^T}$$

- http://cs231n.stanford.edu/handouts/linear-backprop.pdf
- http://cs231n.stanford.edu/handouts/derivatives.pdf
- http://cs231n.stanford.edu/slides/2021/lecture_4.pdf
- https://cs231n.github.io/optimization-2/

# Model

Sequential API

Functional API

# Functional model

Input    Layers    Model

```python
input = Input(shape=(28,28))
x = Flatten()(input)
x = Dense(128, activation="relu")(x)
predictions = Dense(10, activation="softmax")(x)


func_model = Model(inputs=input, outputs=predictions)
```
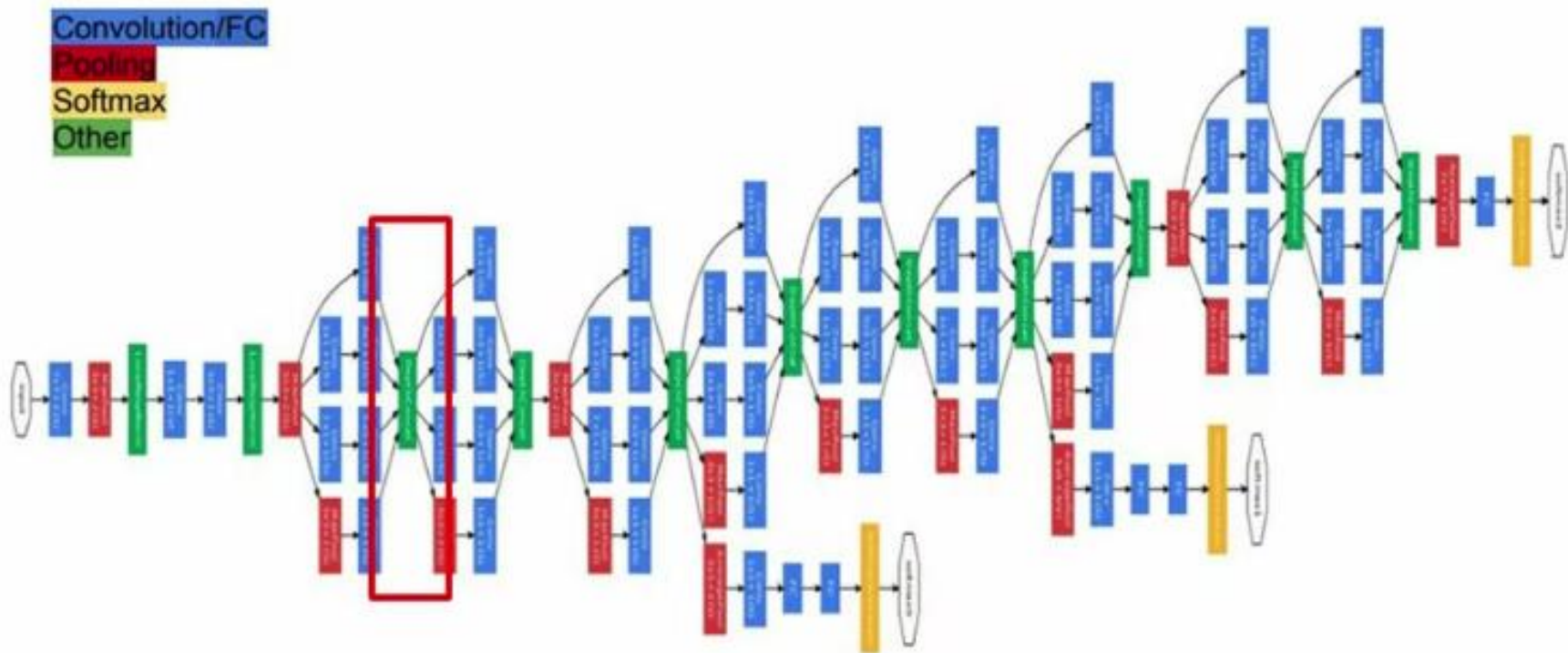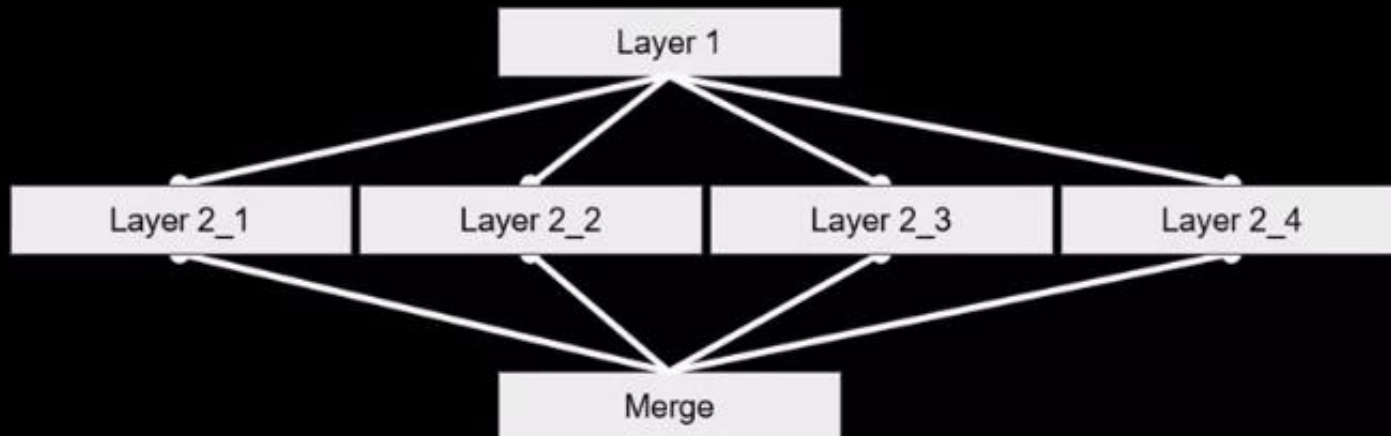
* 
```python
func_model = Model(inputs=[input1, input2], outputs=[output1, output2])
```

```python
dense = layers.Dense(64, activation="relu")
x = dense(inputs)
```

```python
model = keras.Model(inputs=inputs, outputs=outputs, name="mnist_model")
```

Convolution/FC
Pooling
Softmax
Other

```python
inputs = keras.Input(shape=(32, 32, 3), name="img")
x = layers.Conv2D(32, 3, activation="relu")(inputs)
x = layers.Conv2D(64, 3, activation="relu")(x)
block_1_output = layers.MaxPooling2D(3)(x)

x = layers.Conv2D(64, 3, activation="relu", padding="same")(block_1_output)
x = layers.Conv2D(64, 3, activation="relu", padding="same")(x)
block_2_output = layers.add([x, block_1_output])
```

```python
encoder_input = keras.Input(shape=(28, 28, 1), name="img")
x = layers.Conv2D(16, 3, activation="relu")(encoder_input)
x = layers.Conv2D(32, 3, activation="relu")(x)
x = layers.MaxPooling2D(3)(x)
x = layers.Conv2D(32, 3, activation="relu")(x)
x = layers.Conv2D(16, 3, activation="relu")(x)
encoder_output = layers.GlobalMaxPooling2D()(x)

encoder = keras.Model(encoder_input, encoder_output, name="encoder")
encoder.summary()

x = layers.Reshape((4, 4, 1))(encoder_output)
x = layers.Conv2DTranspose(16, 3, activation="relu")(x)
x = layers.Conv2DTranspose(32, 3, activation="relu")(x)
x = layers.UpSampling2D(3)(x)
x = layers.Conv2DTranspose(16, 3, activation="relu")(x)
decoder_output = layers.Conv2DTranspose(1, 3, activation="relu")(x)

autoencoder = keras.Model(encoder_input, decoder_output, name="autoencoder")
autoencoder.summary()
```
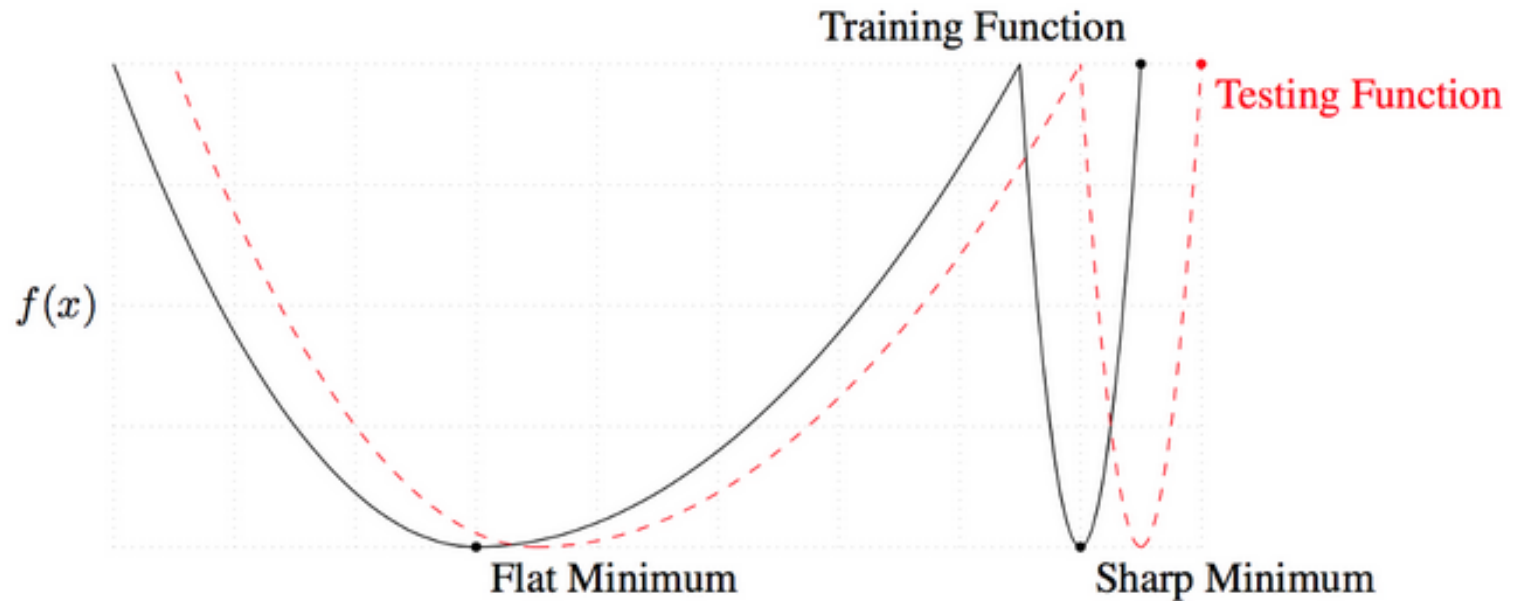
```python
encoder_input = keras.Input(shape=(28, 28, 1), name="original_img")
x = layers.Conv2D(16, 3, activation="relu")(encoder_input)
x = layers.Conv2D(32, 3, activation="relu")(x)
x = layers.MaxPooling2D(3)(x)
x = layers.Conv2D(32, 3, activation="relu")(x)
x = layers.Conv2D(16, 3, activation="relu")(x)
encoder_output = layers.GlobalMaxPooling2D()(x)

encoder = keras.Model(encoder_input, encoder_output, name="encoder")
encoder.summary()


decoder_input = keras.Input(shape=(16,), name="encoded_img")
x = layers.Reshape((4, 4, 1))(decoder_input)
x = layers.Conv2DTranspose(16, 3, activation="relu")(x)
x = layers.Conv2DTranspose(32, 3, activation="relu")(x)
x = layers.UpSampling2D(3)(x)
x = layers.Conv2DTranspose(16, 3, activation="relu")(x)
decoder_output = layers.Conv2DTranspose(1, 3, activation="relu")(x)

decoder = keras.Model(decoder_input, decoder_output, name="decoder")
decoder.summary()


autoencoder_input = keras.Input(shape=(28, 28, 1), name="img")
encoded_img = encoder(autoencoder_input)
decoded_img = decoder(encoded_img)
autoencoder = keras.Model(autoencoder_input, decoded_img, name="autoencoder")
autoencoder.summary()
```

# Ансамблі нейронних мереж

Narrow and wide optima. Flat minimum will produce similar loss during training and testing. Narrow loss, however, will give very different results during training and testing. In other words, wide minimum is more generalizable than narrow. Source.

https://towardsdatascience.com/stochastic-weight-averaging-a-new-way-to-get-state-of-the-art-results-in-deep-learning-c639ccf36a
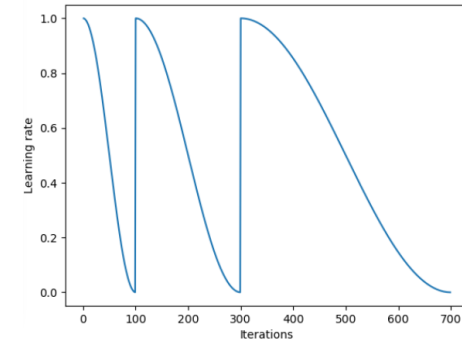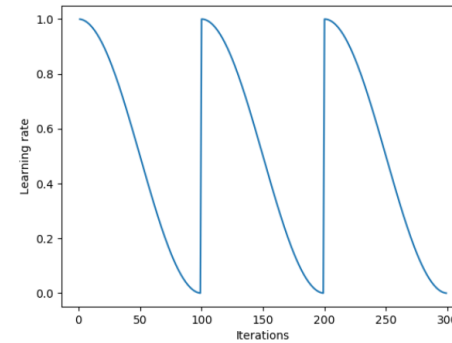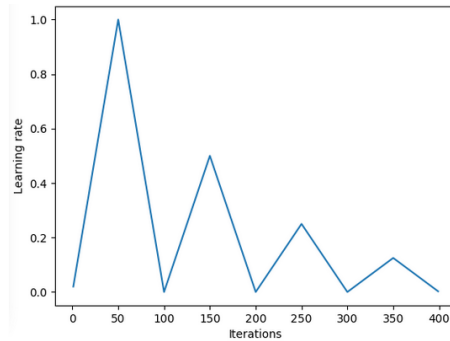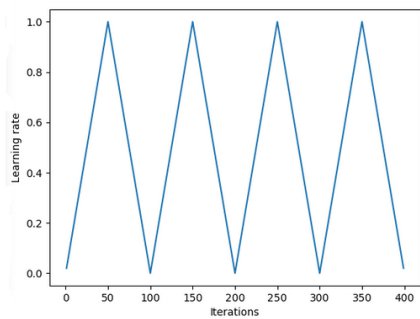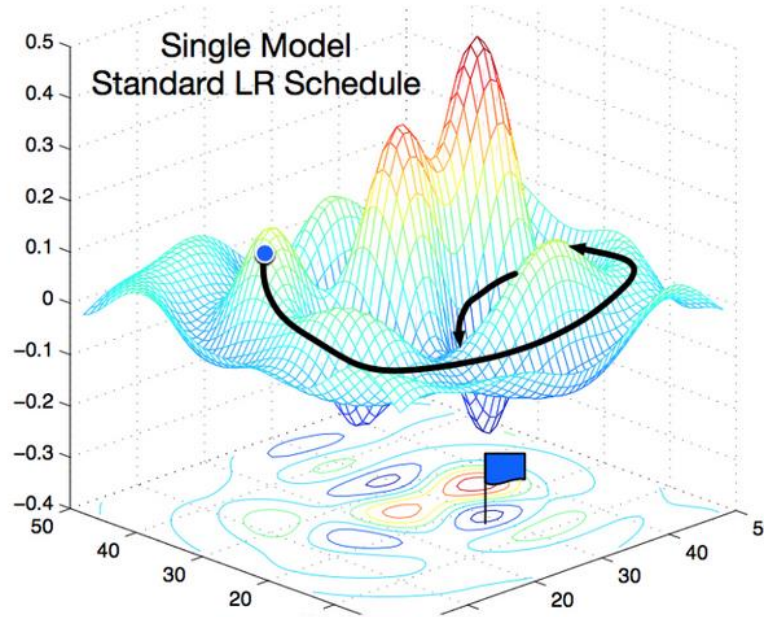
https://arxiv.org/pdf/1609.04836.pdf

# Model Ensembles

- Same model, different initializations

- Top models discovered during cross-validation

- Different checkpoints of a single model

- Running average of parameters during training

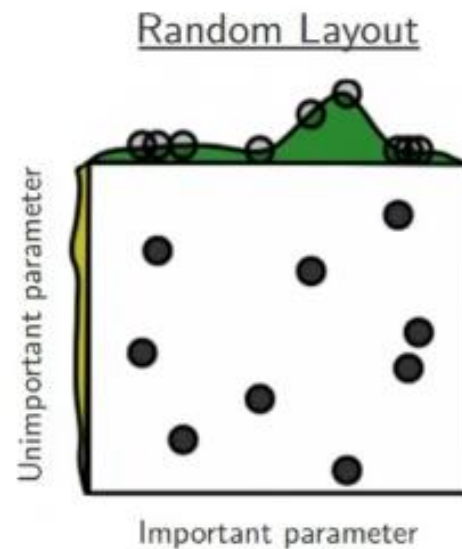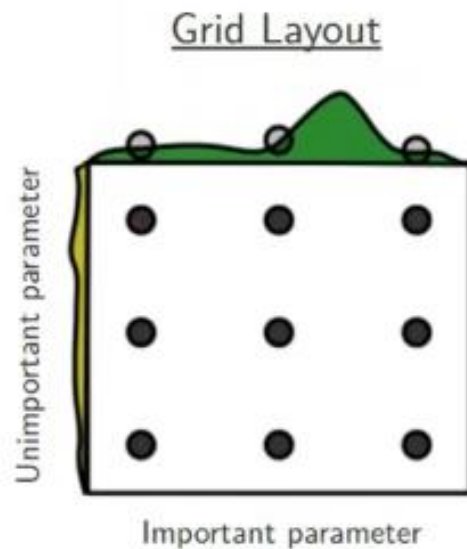$$w_{\text{SWA}} \leftarrow \frac{w_{\text{SWA}} \cdot n_{\text{models}} + w}{n_{\text{models}} + 1,} ,$$
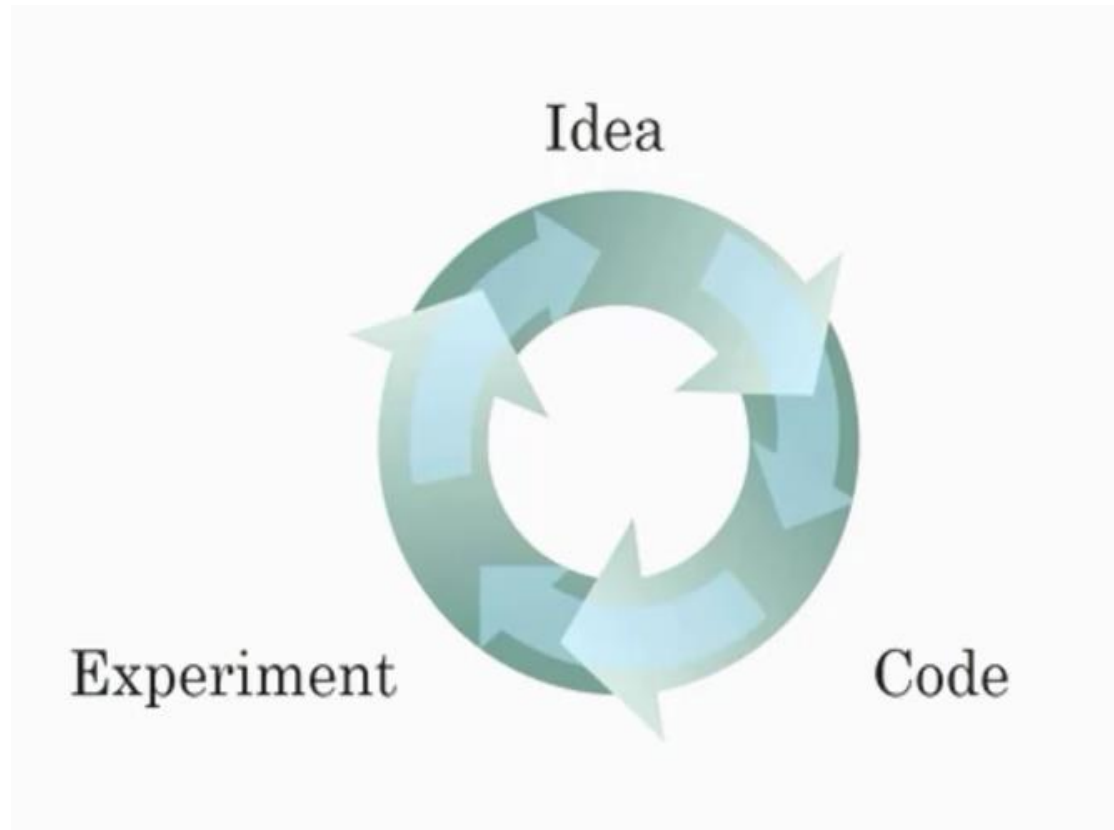
# Different checkpoints of a single model

# Ensembling

```python
inputs = keras.Input(shape=(128,))
y1 = model1(inputs)
y2 = model2(inputs)
y3 = model3(inputs)
outputs = layers.average([y1, y2, y3])
ensemble_model = keras.Model(inputs=inputs, outputs=outputs)
```

# Підбір гіперпараметрів



Grid Layout — Random Layout

☰ Filter

**ML basics with Keras** ⌃
  Basic image classification
  Basic text classification
  Text classification with TF Hub
  Regression
  Overfit and underfit
  Save and load
  **Tune hyperparameters with the Keras Tuner**
  More examples on keras.io ☒

Load and preprocess data    ⌄

ADVANCED

Customization    ⌄

Distributed training    ⌄

Vision    ⌄

Text    ⌄

Audio    ⌄

Structured data    ⌄

TensorFlow  >  Learn  >  TensorFlow Core  >  Tutorials

Was this helpful?   👍  👎

# Introduction to the Keras Tuner    🔖 ▾

🔵 Run in Google Colab        🐙 View source on GitHub        ⬇ Download notebook

## Overview

The Keras Tuner is a library that helps you pick the optimal set of hyperparameters for your TensorFlow program. The process of selecting the right set of hyperparameters for your machine learning (ML) application is called *hyperparameter tuning* or *hypertuning*.

Hyperparameters are the variables that govern the training process and the topology of an ML model. These variables remain constant over the training process and directly impact the performance of your ML program. Hyperparameters are of two types:

1. **Model hyperparameters** which influence model selection such as the number and width of hidden layers

2. **Algorithm hyperparameters** which influence the speed and quality of the learning algorithm such as the learning rate for Stochastic Gradient Descent (SGD) and the number of nearest neighbors for a k Nearest Neighbors (KNN) classifier

In this tutorial, you will use the Keras Tuner to perform hypertuning for an image classification application.

https://www.tensorflow.org/tutorials/keras/keras_tuner

## get_best_hyperparameters method

[source]

```
Tuner.get_best_hyperparameters(num_trials=1)
```

Returns the best hyperparameters, as determined by the objective.

This method can be used to reinstantiate the (untrained) best model found during the search process.

### Example

```
best_hp = tuner.get_best_hyperparameters()[0]
model = tuner.hypermodel.build(best_hp)
```

### Arguments

- **num_trials**: Optional number of HyperParameters objects to return.

### Returns

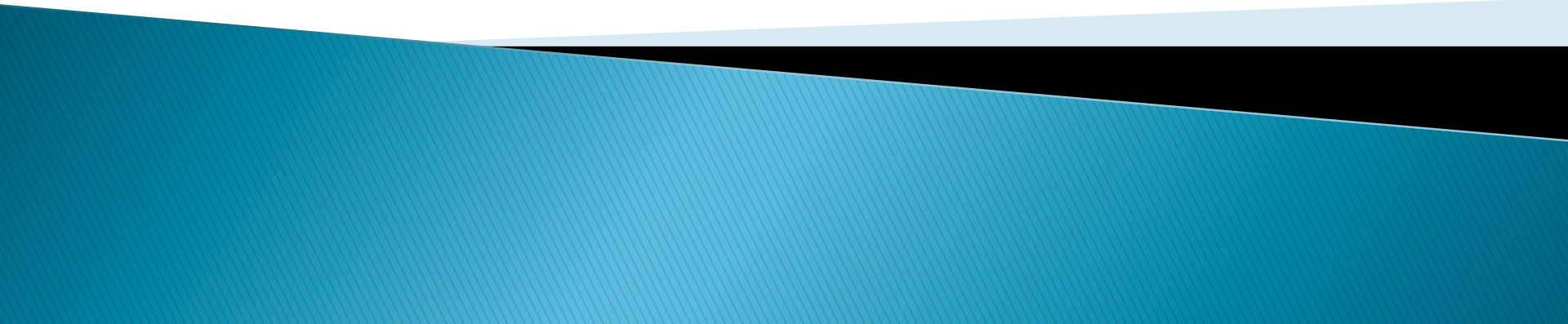List of HyperParameter objects sorted from the best to the worst.

---

## get_best_models method

[source]

```
Tuner.get_best_models(num_models=1)
```

Returns the best model(s), as determined by the tuner's objective.

# Save and load models

≡ Filter

TensorFlow tutorials
Quickstart for beginners
Quickstart for experts

**BEGINNER**

ML basics with Keras    ⌃
    Basic image classification
    Basic text classification
    Text classification with TF Hub
    Regression
    Overfit and underfit
    **Save and load**
    Tune hyperparameters with the Keras Tuner
    More examples on keras.io ↗

Load and preprocess data    ⌄

**ADVANCED**

Customization    ⌄

Was this helpful?    👍  👎

# Save and load models    🔖 ▾

**On this page** ⌄
Options
Setup
    Installs and imports
    Get an example dataset
    Define a model
Save checkpoints during training
    Checkpoint callback usage
    Checkpoint callback options
...

**CO** Run in Google Colab        View source on GitHub        ⬇ Download notebook

Model progress can be saved during and after training. This means a model can resume where it left off and avoid long training times. Saving also means you can share your model and others can recreate your work. When publishing research models and techniques, most machine learning practitioners share:

- code to create the model, and