# Neural Networks: Foundations to Generative AI

Multi-layer Perceptron (MLP) +
Working with Keras API

# This Week's Agenda

1. What's in the News?

2. Building Blocks of NNs

   • Tensors (and relevant mathematical operations)
   • Activation Functions
   • Loss Functions
   • Backpropagation: Derivatives, Gradients & the Chain Rule (with examples)
   • Optimizers

3. Building a Linear Classifier

   • Overview of Keras and Tensorflow.
   • Implementing a linear classifier in Keras (now that we know the components).

# What's In the News?

INNOVATION & AI > TECHNOLOGY > DEVELOPER TOOLS

## Introducing Agentic Vision in Gemini 3 Flash

Jan 27, 2026
5 min read

Agentic Vision, a new capability in Gemini 3 Flash, combines visual reasoning with code execution to ground answers in visual evidence.

**Rohan Doshi**
Product Manager, Google DeepMind

Read AI-generated summary

Share

Agentic Vision ◎

Gemini 3 Flash

---

The Mashable 101 | Creator Hub | Tech | Science | Life | Social Good | Entertainment | Deals | Shopping

Home > Tech

## OpenAI says its mystery AI wearable is on track for 2026 as AI earbuds rumors spread

Are ChatGPT-powered earbuds in the works at OpenAI?

By Timothy Beck Werth on January 20, 2026

OpenAI

---

CNET
YOUR GUIDE TO A BETTER FUTURE

News | AI | Tech | Home | Entertainment | Wellness
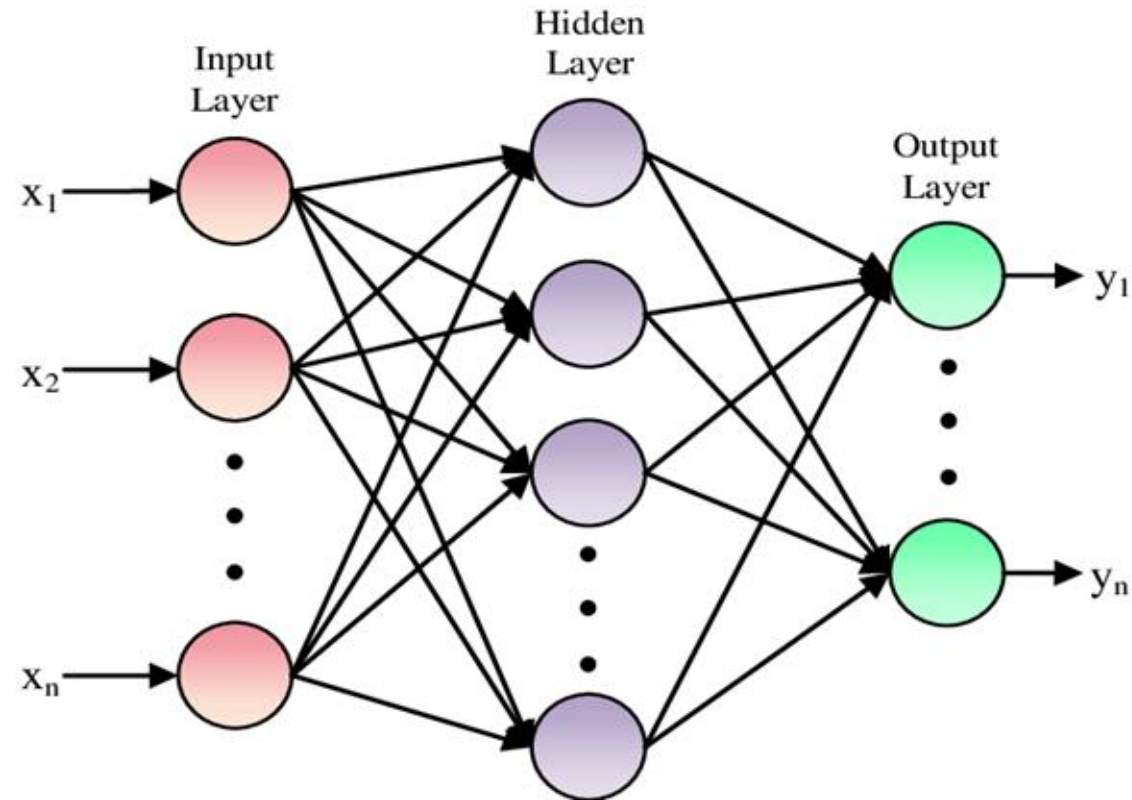
ai atlas > Tech > Services & Software > AI

## Apple's Next Big Wearable Could Be an AirTag-Sized AI Pin

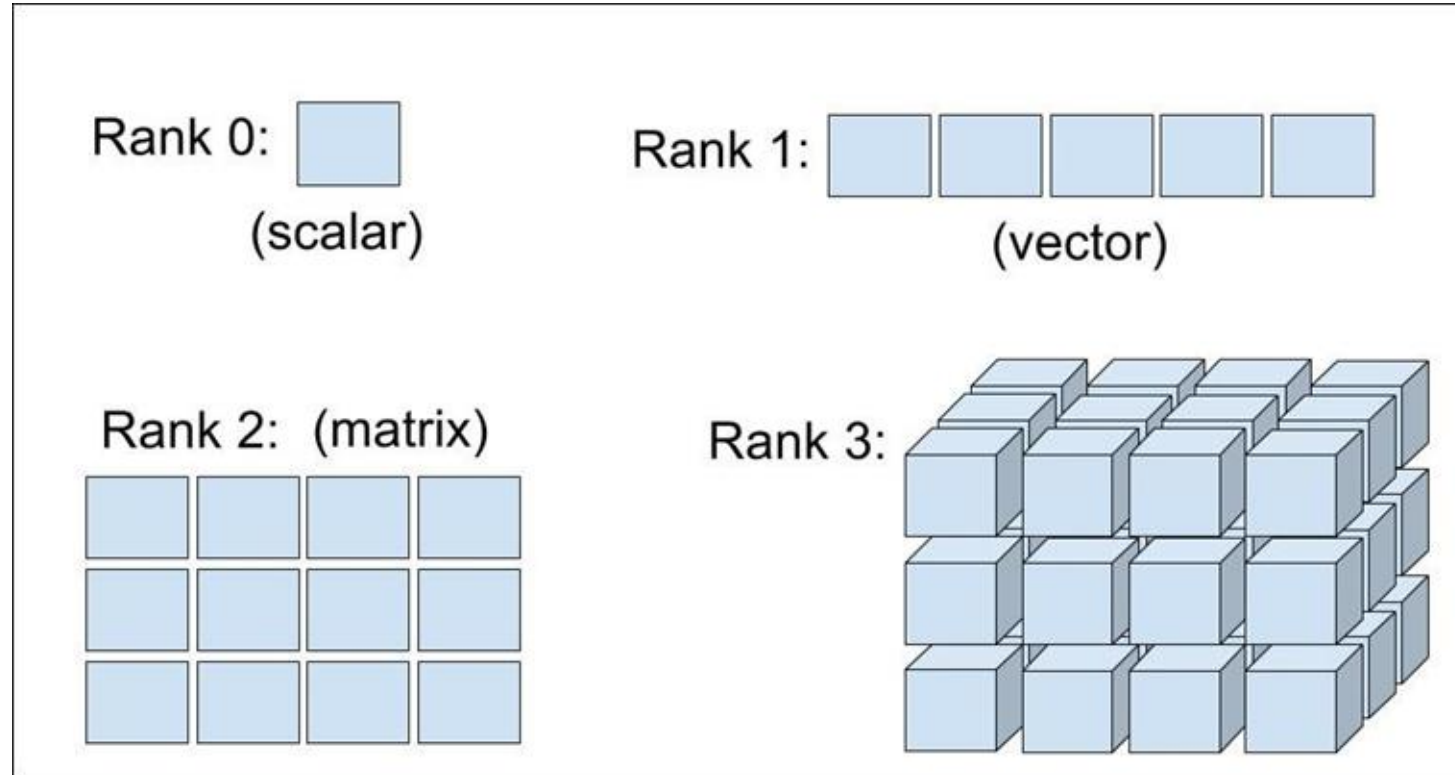Apple is reportedly working on a small wearable device with an embedded microphone and cameras.

**Omar Gallaga**
Jan. 22, 2026 6:24 a.m. PT
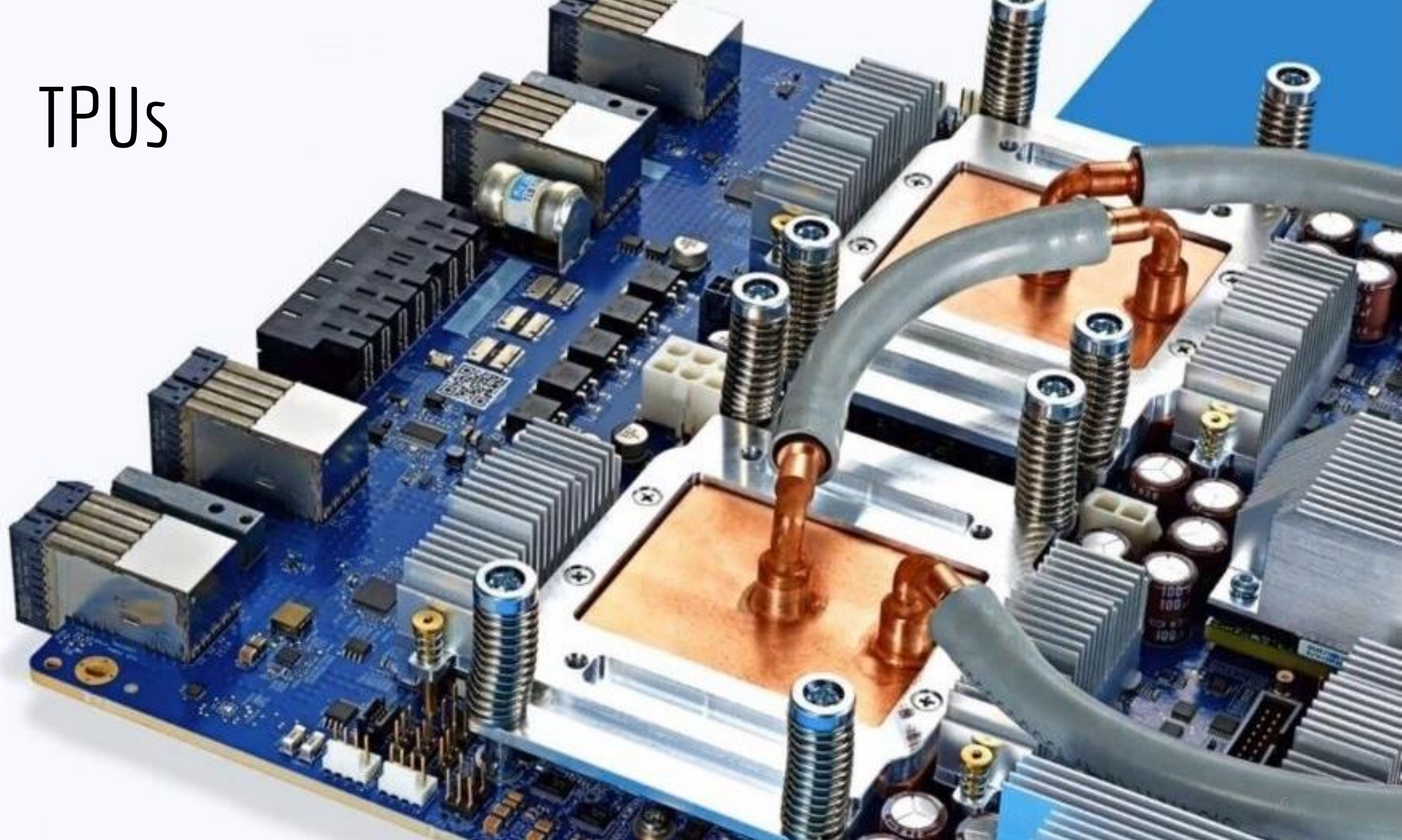
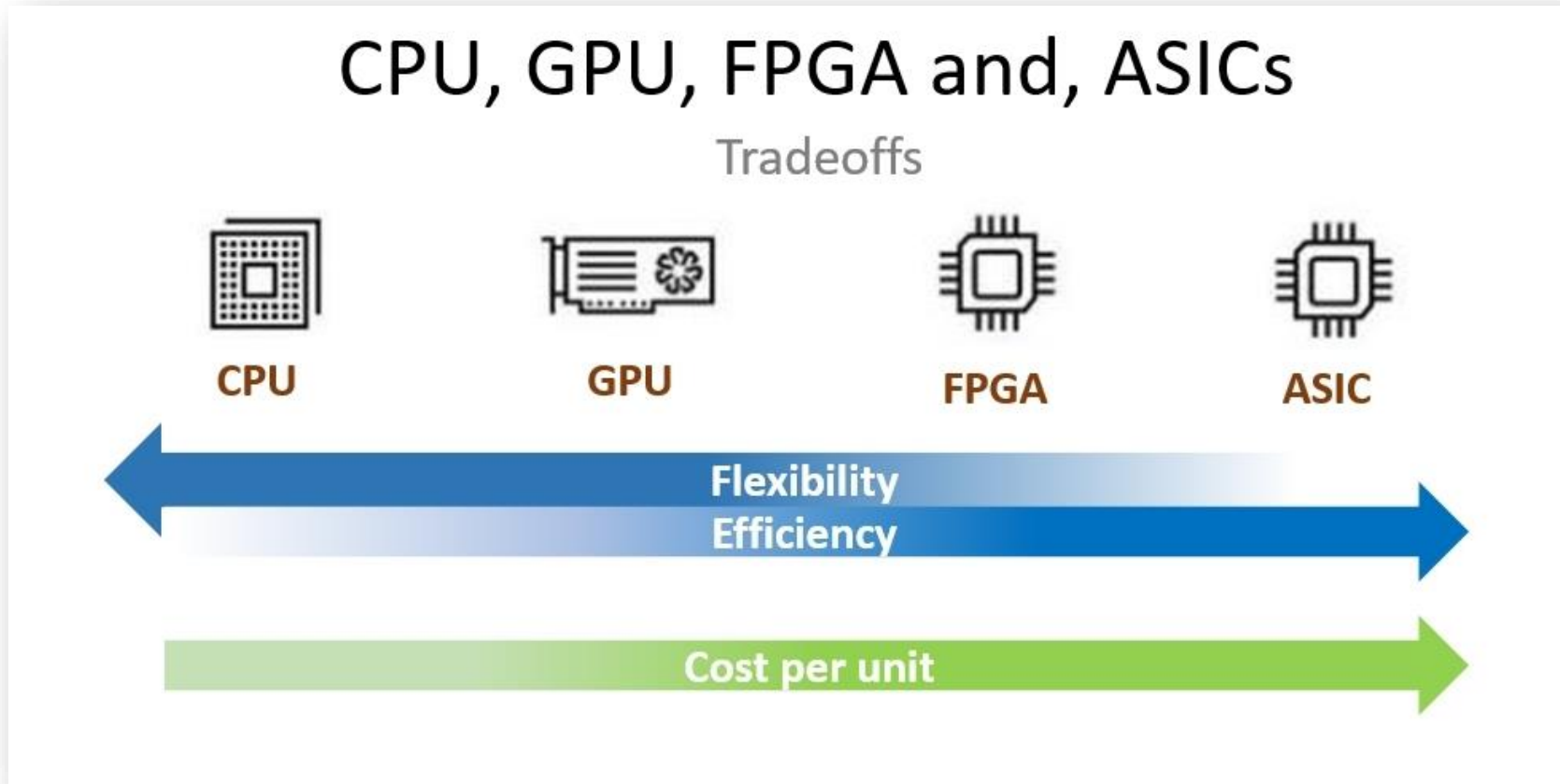2 min read

# Recall: Multi-layer Perceptron (MLP)

# Tensors



Rank 0: (scalar)

Rank 1: (vector)

Rank 2: (matrix)

Rank 3:

*Question:* What sort of data (give an example) would be stored in a rank-3 tensor? How about a rank-4 tensor?

# TPUs

# An Aside: GPU vs. ASIC



CPU, GPU, FPGA and, ASICs

Tradeoffs

CPU    GPU    FPGA    ASIC

Flexibility
Efficiency

Cost per unit

# Forward Pass

# The Perceptron

# Neuron / Network Components



Inputs    Weights    Summation and Bias    Activation    Output

*Question: What rank tensor are x, w and b here?*
*What will the shape of y be?*
*What is the order of operations in a forward pass?*

# Neuron / Network Components



*Question:* Which of these values are constants?
Which are trainable parameters?

# Multiplication

$$y_1 = \varphi\,(\mathbf{x_1} \cdot \mathbf{w_1} + b_1)$$

$$\begin{bmatrix} a_1 & a_2 & a_3 \\ b_1 & b_2 & b_3 \\ c_1 & c_2 & c_3 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} a_1 x + a_2 y + a_3 z \\ b_1 x + b_2 y + b_3 z \\ c_1 x + c_2 y + c_3 z \end{bmatrix}$$

# Matrix Addition (Broadcast)

$$y_1 = \varphi\,(x_1 \cdot w_1 + b_1)$$

Shape of the Two Tensors Needs to Conform
- A + B will only work if A is cleanly divisible by B (or vice versa)

Sum Element-wise
- Replicate B until it matches A's dimensions, then perform element-wise addition.

We Use This for the Addition Step
- Add x*w and b (bias)



np.arange(3)+5

np.ones((3, 3))+np.arange(3)

np.arange(3).reshape((3, 1))+np.arange(3)

# Neuron / Network Components

# Activation Functions

$$y_1 = \varphi\,(x_1 \cdot w_1 + b_1)$$

### Sigmoid

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

(a)

### Tanh

$$\sigma(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

(b)

### ReLU

$$ReLU(z) = \begin{cases} z, z > 0 \\ 0, otherwise \end{cases}$$

(c)

### LeakyReLU(a=0.2)

$$LeakyReLU(z) = \begin{cases} z, z > 0 \\ az, otherwise \end{cases}$$

(d)

# Activation Functions

$$y_1 = \varphi\,(x_1 \cdot w_1 + b_1)$$

# Multi-Class, Single-Label

$$y_1 = \varphi \, (x_1 \cdot w_1 + b_1)$$

**Softmax (MLOGIT):**

We have D inputs (x's).
We have k outputs (classes).

So, W is a (D,k) matrix and X is a (D,1) matrix.

That means, A is a (k,1) matrix.

That means Y is also a (k,1) matrix.

$$A = W^T X,$$

$$Y = \text{softmax}(A),$$

$$Y_i = \frac{e^{A_i}}{\sum_{j=1}^{k} e^{A_j}}.$$

Output Layer

$y_1$

$y_n$

D

k

# Multi-Class, Single-Label

$$y_1 = \varphi \left( x_1 \cdot w_1 + b_1 \right)$$



© Gordon Burtch, 2026

# Multi-Class, Multi-Label

## Many Non-Exclusive Labels

- We would create a sigmoid output layer with one output for each class we are predicting.
- Train on all labels together.

# We Know Enough for a Forward Pass

## Calculate Output of Each Node Sequentially

$$y_1 = \varphi \left( x_1 \cdot w_{1,1} + x_2 \cdot w_{1,2} + \cdots + b_1 \right)$$

$$y_2 = \varphi \left( x_1 \cdot w_{2,1} + x_2 \cdot w_{2,2} + \cdots + b_2 \right)$$

...

## Eventually We Obtain Model's Predictions

Multi-Layer Perceptron (MLP) – Dense, Fully-connected, Feed-forward

# Calculate Loss

# Loss Functions

## Cross-Entropy / Log-Loss

$$H_p(q) = -\frac{1}{N} \sum_{i=1}^{N} y_i \cdot log(p(y_i)) + (1 - y_i) \cdot log(1 - p(y_i))$$

- Typical for binary outcomes. Value grows exponentially larger as the predicted probability moves away from the true 0,1 label.
- Multi-category outcomes have an analogous loss function known as categorical cross-entropy.

$$CE = -\sum_{i}^{C} t_i log(s_i)$$

## MAE / L1 Loss

$$MAE = \frac{\sum_{i=1}^{n} |\, y_i - \hat{y}_i \,|}{n}$$

- Typical for continuous outcomes. Errors are penalized homogenously, in magnitude and direction. This should look familiar!

## MSE / Quadratic / L2 Loss

$$MSE = \frac{\sum_{i=1}^{n} (y_i - \hat{y}_i)^2}{n}$$

- Typical for continuous outcomes, larger errors penalized exponentially more. This should look familiar!

# Backpropagation

# Derivative = "Rate" of Change



epsilon_y

y = f(x)

y

x

epsilon_x



Local linear approximation of f, with slope a

y = f(x)

y

x

# Gradient = Derivative in Multiple Dimensions

# Gradient Descent

# Derivatives of Loss w.r.t All Parameters

Recall that Each Node's Output Can be Expressed as a Function of the Prior Nodes' Outputs

$$y_1 = \varphi \left( x_1 \cdot w_{1,1} + x_2 \cdot w_{1,2} + \cdots + b_1 \right)$$

$$y_2 = \varphi \left( x_1 \cdot w_{2,1} + x_2 \cdot w_{2,2} + \cdots + b_2 \right)$$

...



Start at the final nodes in the network and work backwards
- We calculate partial derivatives w.r.t. their inputs / weights.
- Then, use those partial derivatives and work backward into earlier layers to get partial derivatives w.r.t. *their* inputs / weights, and so on.

# Simplifying Gradients:
# The Computation Graph

$$J = 3(a \cdot b + c)$$

$a = 6$

$b = 4$

$$u = a \cdot b$$

$$z = u + c$$

$$J = 3(z)$$

$c = 2$

# Backpropagation = Working Backwards

$$J = 3(a \cdot b + c)$$

$a = 6$

$b = 4$

| $u = a \cdot b$ | $z = u + c$ | $J = 3(z)$ |

$c = 2$

$$\frac{dJ}{dz} = 3$$

# Backpropagation = Work Backwards

$$\frac{dJ}{dz} = 3$$

$$J = 3(a \cdot b + c)$$

$a = 6$

$b = 4$

| $u = a \cdot b$ |

| $z = u + c$ |

| $J = 3(z)$ |

$c = 2$

$$\frac{dJ}{du} = \frac{dJ}{dz} \cdot \frac{dz}{du}$$

$$\frac{dJ}{du} = 3 \cdot \frac{dz}{du}$$

$$\frac{dJ}{du} = 3 \cdot 1 = 3$$

# Backpropagation = Work Backwards

$$\frac{dJ}{dz} = 3$$

$$\frac{dJ}{du} = 3$$

$$J = 3(a \cdot b + c)$$

$a = 6$

$b = 4$

| $u = a \cdot b$ | $z = u + c$ | $J = 3(z)$ |

$c = 2$

$$\frac{dJ}{dc} = \frac{dJ}{dz} \cdot \frac{dz}{dc}$$

$$\frac{dJ}{dc} = 3 \cdot \frac{dz}{dc}$$

$$\frac{dJ}{dc} = 3 \cdot 1 = 3$$

# Backpropagation = Work Backwards

$$\frac{dJ}{dz} = 3$$

$$J = 3(a \cdot b + c)$$

$$\frac{dJ}{du} = 3$$

$$\frac{dJ}{dc} = 3$$



$a = 6$

$b = 4$

$u = a \cdot b$

$z = u + c$

$J = 3(z)$

$c = 2$

$$\frac{dJ}{da} = \frac{dJ}{du} \cdot \frac{du}{da}$$

$$\frac{dJ}{da} = 3 \cdot b$$

$$\frac{dJ}{da} = 3 \cdot 4 = 12$$

# Backpropagation = Work Backwards

$$\frac{dJ}{dz} = 3$$

$$J = 3(a \cdot b + c)$$

$$\frac{dJ}{du} = 3$$

$$\frac{dJ}{da} = 12$$

$$\frac{dJ}{dc} = 3$$



$a = 6$

$b = 4$

$u = a \cdot b$

$z = u + c$

$J = 3(z)$

$$\frac{dJ}{db} = \frac{dJ}{du} \cdot \frac{du}{db}$$

$$\frac{dJ}{db} = 3 \cdot a$$

$$\frac{dJ}{da} = 3 \cdot 6 = 18$$

$c = 2$

We thus update our parameters, a, b, and c, subtracting each's gradients*epsilon from its current value. Epsilon is the learning rate.

# Single Node with Sigmoid & Cross-Entropy Loss (i.e., Logistic Regression)



$x_1 = 2$

$w_1$

$b$

$\varphi\,()$

$\mathcal{L}(y, \hat{y})$

$w_2$

$x_2 = 3$

Remember that $\varphi$ here is just a placeholder for the argument to the loss function. It happens to be a sigmoid transformation of 'something', i.e., $\varphi(wx+b)$, but it doesn't really matter. We just represent it with some variable name and calculate an expression for the derivative.

$$\frac{d\mathcal{L}}{d\varphi} = -\frac{y}{\varphi} + \frac{1-y}{1-\varphi}$$

$$\frac{d\mathcal{L}}{d\varphi} = \frac{\varphi(1-y) - y(1-\varphi)}{\varphi(1-\varphi)}$$

$$\frac{d\mathcal{L}}{d\varphi} = \frac{\varphi - \varphi y - y + \varphi y}{\varphi(1-\varphi)}$$

$$\frac{d\mathcal{L}}{d\varphi} = \frac{\varphi - y}{\varphi(1-\varphi)}$$

# Single Node with Sigmoid & Cross-Entropy Loss (i.e., Logistic Regression)

$x_1 = 2$

$w_1$

$b$

$\varphi\,()$

$\mathcal{L}(y, \hat{y})$

$w_2$

$x_2 = 3$

**Now we calculate derivative of the sigmoid with respect to its argument, z.**

$\varphi(z) = (1 + e^{-z})^{-1}$

$\varphi'(z) = -1 \cdot (1 + e^{-z})^{-2} \cdot (0 + e^{-z} \cdot -1)$

$\varphi'(z) = (1 + e^{-z})^{-2} \cdot e^{-z}$

$\varphi'(z) = \varphi(z) \cdot (1 - \varphi(z))$

$$\frac{d\mathcal{L}}{dz} = \frac{d\mathcal{L}}{d\varphi} \cdot \frac{d\varphi}{dz}$$

$$\frac{d\mathcal{L}}{dz} = \frac{\varphi - y}{\varphi(1 - y)} \cdot \frac{d\varphi}{dz}$$

$$\frac{d\mathcal{L}}{dz} = \frac{\varphi - y}{\varphi(1 - y)} \cdot \varphi(1 - \varphi)$$

$$\frac{d\mathcal{L}}{dz} = \varphi - y$$

# Single Node with Sigmoid & Cross-Entropy Loss (i.e., Logistic Regression)

$x_1 = 2$

$w_1 = 5$

$b=1$

$x_2 = 3$

$w_2 = 6$

$\varphi\,()$

$\mathcal{L}(y = 1, \hat{y})$

$$\frac{d\mathcal{L}}{dw_1} = \frac{d\mathcal{L}}{dz} \cdot \frac{dz}{dw_1}$$

$$\frac{d\mathcal{L}}{dw_1} = (\varphi - y) \cdot x_1$$

$$w_{1,new} = w_{1,old} - (\frac{d\mathcal{L}}{dw_{1,old}} \cdot \varepsilon)$$

**And, finally, we can get gradient of loss with respect to weights and bias. For example, for the first weight…**

**Evaluate $\varphi$ based on current values of parameters and the data.**

**Finally, update the weights…**

# Keras and Tensorflow

## 1. Tensorflow

- A Python platform for working with tensors, implementing automatic differentiation, providing access to repositories of (well-known) pre-trained models.

## 2. Keras

- A higher-level API that wraps common usage patterns with Tensorflow functions, pre-defined loss functions, optimization algorithms, etc.
- Keras simplifies data scientists' interaction with Tensorflow.

| Keras | Deep learning development: layers, models, optimizers, losses, metrics... |
| TensorFlow | Tensor manipulation infrastructure: tensors, variables, automatic differentiation, distribution... |
| CPU | GPU | TPU | Hardware: execution |

# Tensorflow GradientTape: AutoDiff

## 1. Gradient Tape

- A Tensorflow function that automates the calculation of derivatives.
- It constructs a computation graph in the background and implements codified rules for calculating derivatives of functions.
- You could technically use gradient tape to implement a gradient descent algorithm for many optimization problems.

# Keras Layers

Layers are the Key Building Block of NNs in Keras

- There are a few subclasses of the Layers class: e.g., Dense is the one we have seen so far, i.e., layers.Dense(), but we also have convolutional layers, max-pooling layers, recurrent layers, and so on. There are many pre-defined layers in Keras. See: https://keras.io/api/layers/.

- These are different architectural components that can be mixed and matched in different ways to create different network topologies.

- It is also possible to construct custom layers.

# Sequential vs. Functional API

## We Have Only Used Sequential API So Far

- Sequential is easy to work with but is also very inflexible. Can only really handle basic feed-forward networks. It automatically figures out the shape of each layer's output tensor and specifies the next layer's input shape accordingly.

## Functional API Let's You Construct Any Topology You Want

- But – we will look at the difference in how each API is used, syntactically.

*Syntactic Simplicity*

*Architectural Flexibility*

Sequential API                                    Functional API

# Quiz Prep: Review of Concepts!

# Optimizers

## Keras Supports 8 Optimizers

- SGD = Stochastic Gradient Descent
- Momentum
- Ftrl (2010) = Follow the Regularized Leader
- Adagrad and Adadelta (2012) = Adaptive Gradient Des
- RMSprop (~2012) = Root Mean Squared propagation
- Adam (2015) = Adadelta / RMSProp with Momentum.
    - Adamax, Nadam are extensions to Adam.

# SGD: Gradient Descent

Types of GD
- Batch GD = Use all the available training data in each pass.
  - Works well if the loss surface is smooth and lacks any saddle points / valleys.
- Stochastic GD = Mini-batch with batch size = 1.
  - If troughs / saddles exist, we move past them as our exploration of gradients for the model will vary withe a given observation that we are considering in an iteration.
  - Computationally quite burdensome but performs well on non-linear problems (eventually).
- Mini-batch GD = What we have been doing so far (randomly split the data in each epoch, into folds, and then cycle over the folds for training).
  - This is a happy-medium between batch and stochastic GD.

Role of Batch Size
- Empirically has been observed that smaller batches yield less overfitting (because of implicit noise in the training process – variance of the gradients obtained will go up).

# Batch (All) vs. Stochastic (1)

## Same Convergence
- If you have a convex surface, either approach will converge to the global optimum (no guarantee your problem is convex of course). Always converges at least to a local minimum.

## Tradeoffs
- Batch, each step is slower, more computationally burdensome, but convergence with fewer iterations; Need to be able to hold the entire dataset in memory.
- SGD makes noisier updates, and requires more iterations to converge, but a single iteration is quick. Only need one observation in memory at a time.

# Momentum

## Getting Past Local Minima

- SGD gets stuck in local minima; the idea of momentum is to make updates be a function of current gradient*learning rate, as well as some fraction (decay) of the update you made last iteration.
- This reduces updates to parameters where the gradients are flipping sign and amplifies updates to gradients that are going in a consistent direction (steeply descer
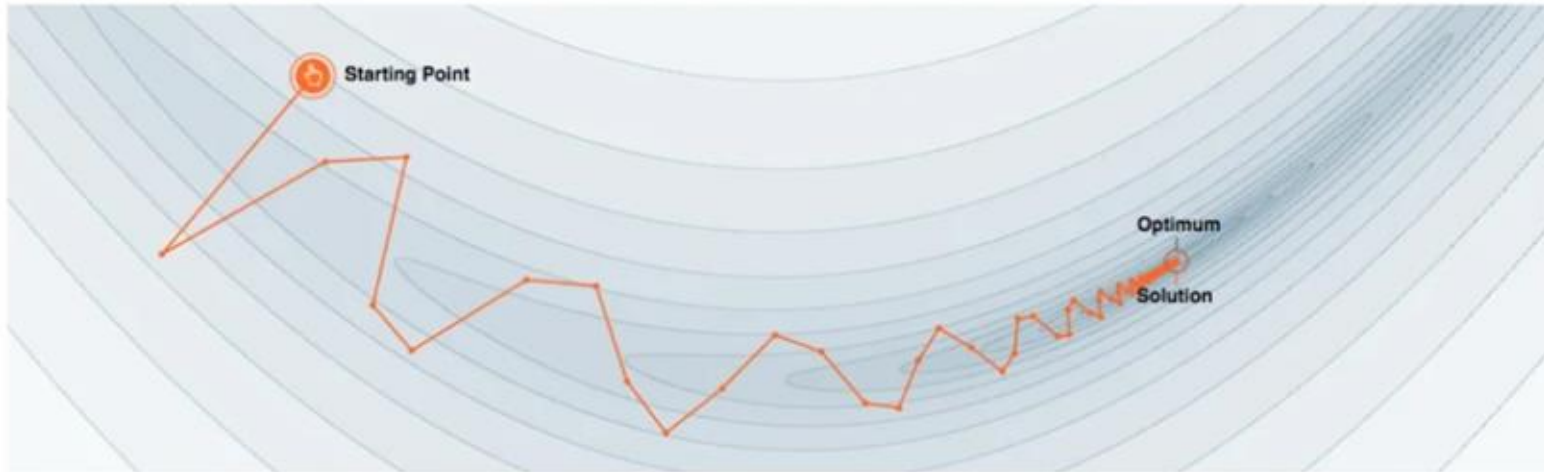


Figure: Optimization with momentum (Source: distill.pub)

# Adagrad & Adadelta (RMS Prop)

Adaptive Gradient Descent (Variable Learning Rate)
- We implicitly apply a high learning rate for features we have been updating very little so far (speed up movement through saddle points, for example).
- We implicitly apply a low learning rate for features we have been updating a lot so far.
- Technically learning rate is removed from the process, every update is a function of past updates.

Adadelta
- Same idea but we use a sliding window of previous updates to determine magnitude of current updates (rather than all prior updates).
- RMSProp is conceptually very similar but was independently developed (around the same time).

# Recap

## Building Blocks of NNs

- Tensors and Tensor Operations
- Activation Functions
- Loss Functions
- Backpropagation: Derivatives, Gradients & the Chain Rule

## Procedure of Minibatch Stochastic Gradient Descent

- Grab a batch of observations (samples)
- Predict their labels using current weights / bias terms.
- Calculate loss value.
- Calculate gradient of loss w.r.t. all weight / bias terms.
- Update each weight by subtracting its gradient*learning rate
- Cycle over the whole training dataset (each cycle is an epoch) repeatedly, until loss is small.