# Exam study-guide: Deep learning
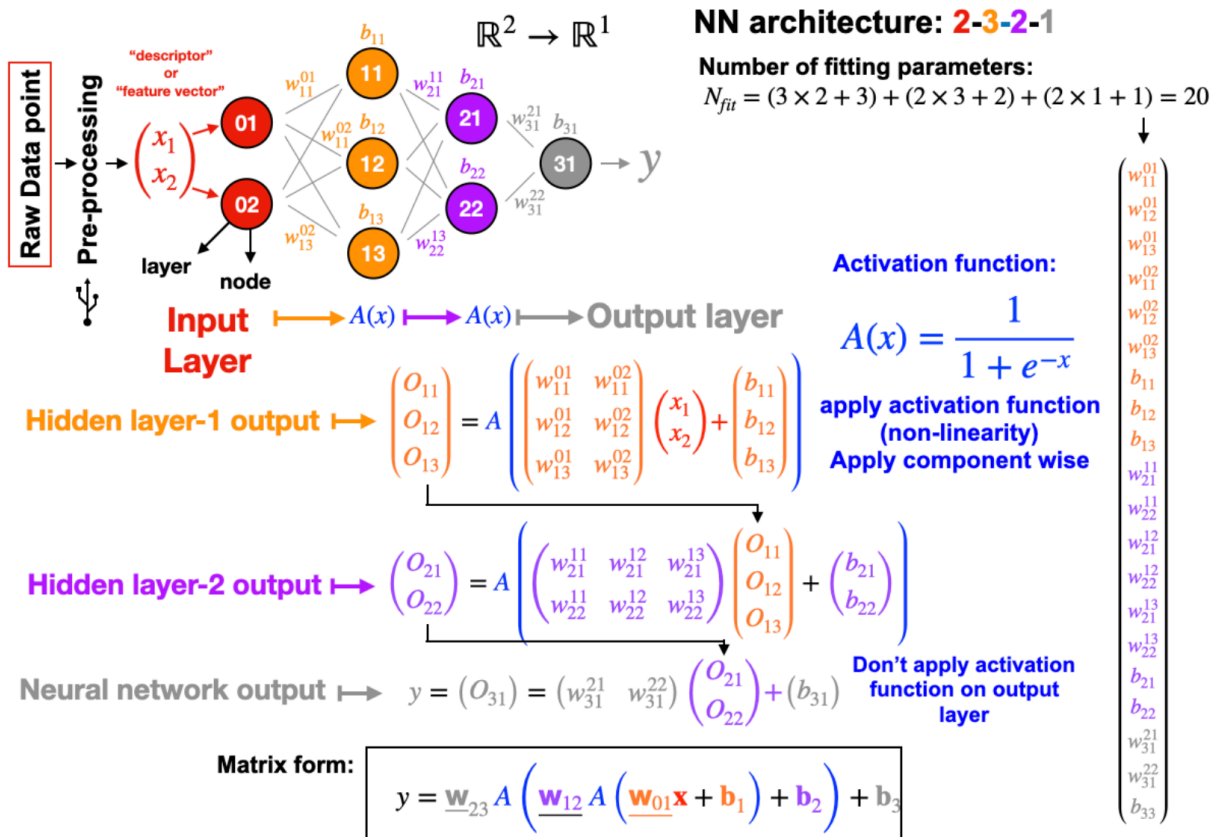
## Shallow Neural Networks 🔗

- Be able to draw shallow models (e.g. linear and logistic regression as simple networks)
  - Study the following notes section: click here

## Dense Feed-Forward Neural Networks

- Understand the architecture of dense feed-forward neural networks.
  - Study the following notes section: click here
- **For example**:
- **Simple Neural network for scalar regression: 2 feature input and 1 target output**
  - Note a **better** choice for hidden layer activate would be RELU here (to help combat the vanishing gradient effect)

- **Problem Type:** Scalar regression (un-constrained output)
- **Mapping:** $^{\text{(input layer)}}\ \mathbb{R}^2 \to \mathbb{R}^1\ ^{\text{(output layer)}}$
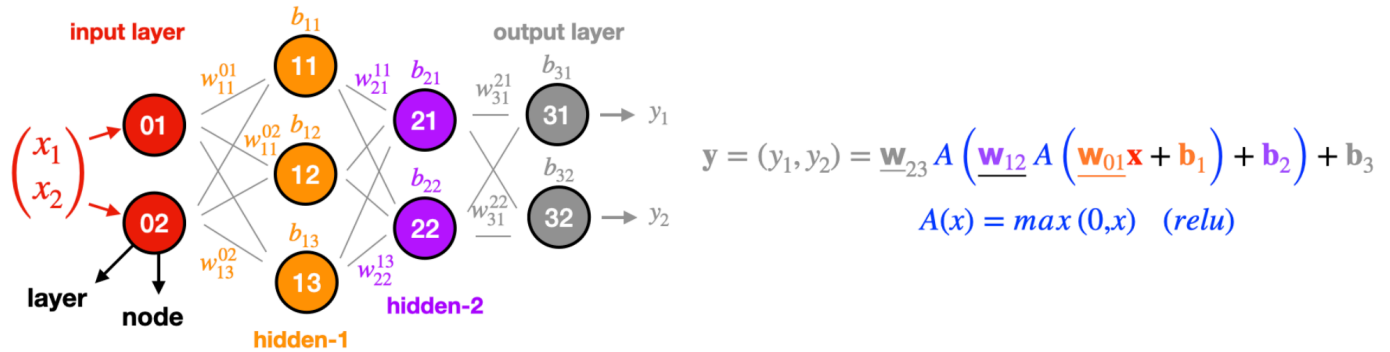- **Hidden layers**: $N_{hidden} = 2 \to [3,2]$

- **Hidden-layer activation**: Sigmoid
- **Output-layer activation**: Linear
- **Number of fitting parameters**: 20



- **Simple Neural network for vector regression: 2 feature input and 2 target output**

- **Problem Type:** Vector regression
- **Mapping:** $^{\text{(input layer)}} \mathbb{R}^2 \rightarrow \mathbb{R}^2 {}^{\text{(output layer)}}$
- **Hidden layers**: $N_{hidden} = 2 \rightarrow [3,2]$
- **Hidden-layer activation**: ReLu
- **Output-layer activation**: Linear
- **Number of fitting parameters**: 23

# Computational graph



$$\mathbf{y} = (y_1, y_2) = \underline{\mathbf{w}}_{23}\, A\left(\underline{\mathbf{w}}_{12}\, A\left(\underline{\mathbf{w}}_{01}\mathbf{x} + \mathbf{b}_1\right) + \mathbf{b}_2\right) + \mathbf{b}_3$$

$$A(x) = max\,(0,x) \quad (relu)$$

**Note**: We didn't have to use `relu` on the *hidden layers*, other activations would also work

## Be able to do the following:

- Determine the number of parameters in each layer, and in the total network, given node counts for various layers.
- Given node counts, activation functions, and weight matrices, be able to:
  - Draw a network accurately.
  - Perform step-by-step computations to evaluate a very small dense feed forward model:
    1. Matrix multiplication
    2. Vector addition (bias)
    3. Component-wise activation
- Understand and be able to draw common hidden-layer activation functions:
  - Sigmoid, Tanh, ReLU, Leaky ReLU
  - Clearly label critical points and characteristics of activation functions.
  - See the following section of the notes for examples: [click here](#)

## Learning Tasks and Output Layer Design

- Given specific learning tasks, know how to determine:
  - Number of output layer nodes.
  - Appropriate loss functions:
    - Mean Squared Error (MSE) for scalar/vector regression
    - Binary Cross-Entropy (BCE) for binary classification
    - Categorical Cross-Entropy (CCE) for multi-class classification
  - Choosing the correct activation functions for the output layer.
  - **For example:**
    - Vector regression mapping from $R^{50}$ to $R^{20}$ where $R^n$ is the continuous spaced defined by the intersection of $n$ mutually orthogonal coordinate axes (i.e. number lines)
      - Input layer size: 50
      - Output layer size: 20
      - Output activation: Linear
      - Loss: MSE
    - Sentiment analysis (binary classification) with a 10000 word vocabulary and OHE vectorization
      - Input layer size: 10000
      - Output layer size: 1 (probability of positive class)
      - Output activation: Sigmoid
      - Loss: BCE
    - MNIST, 10 class multi class classification problem (with 28x28 images (matrices), reshaped into (784,1) vectors )
      - Nodes in input layer: 784
      - Nodes in output layer: 10 node output layer
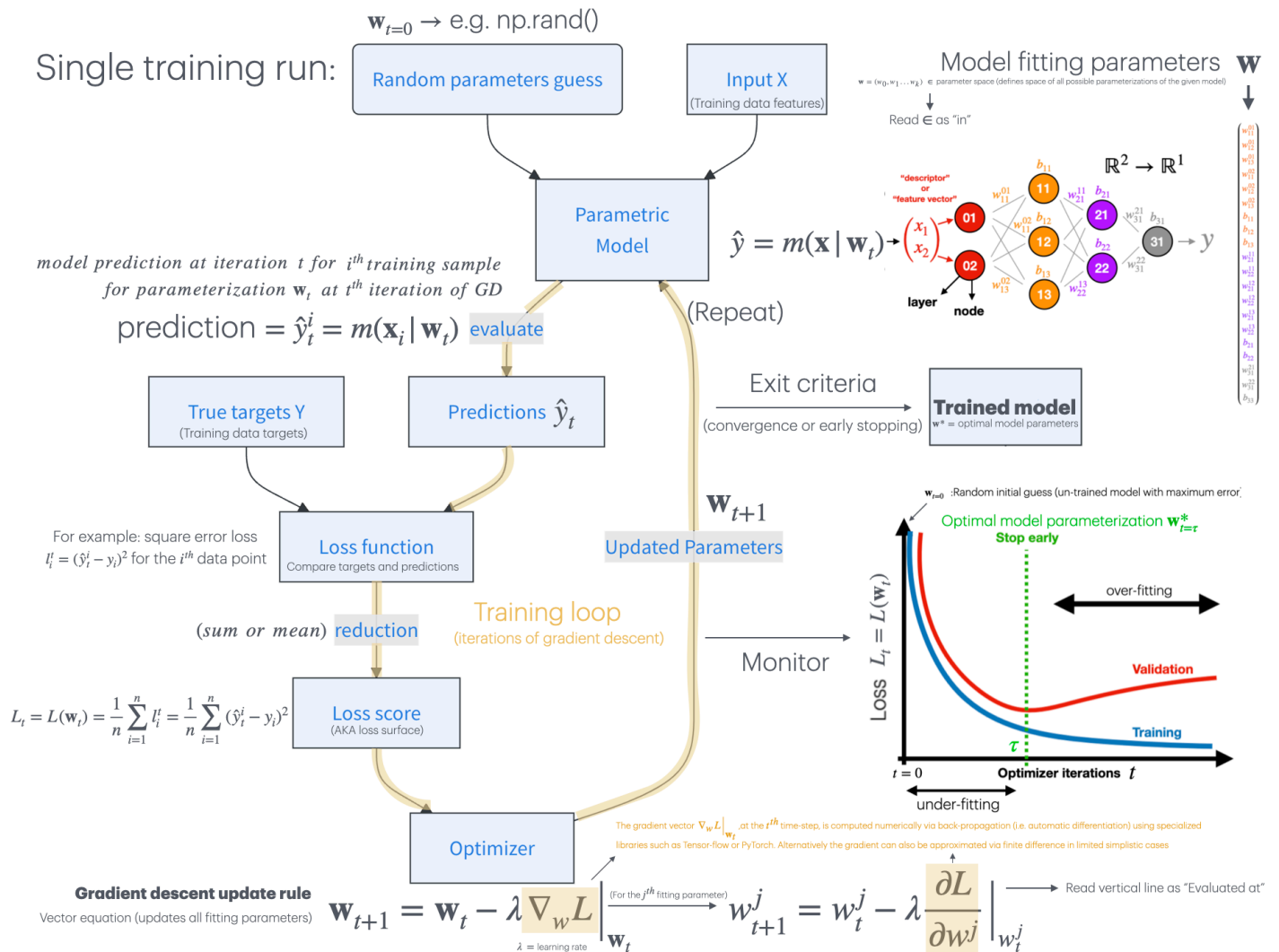      - Output activation: Soft-max activation

- Loss: CCE

## Metrics vs. Loss

- Understand the difference between:
  - Loss: directly optimized by the training process (seen by optimizer, used to find optimal model parameters).
  - Metrics: monitored but not directly optimized (not seen by optimizer).
- Be familiar with common metrics used for Regression (MSE, MAE, MAPE, etc) & Classification (accuracy, precision, recall, F1)

## Training Process Diagram

- Be able to describe, and annotate the following diagram which summarizes the training process, and possibly evaluate the progression for a very simple model, by taking a few steps of gradient descent manually (given the gradient vector)
- Understand training curves for training and validation loss, i.e. loss as a function of gradient decent iteration index $t$, understand what is happening in different regions of the curve (under-fitting, optimal fit, over-fitting), and the Concept of early stopping to avoid overfitting
- **Study the following workflow diagram carefully and ensure you fully understand it. If you do, you will have a solid grasp of how the vast majority of modern deep learning methods operate (e.g., MLP, CNN, RNN, GNN, LSTM, Auto-regression, ARIMA, Transformer models, etc.).** Remember, for parametric machine learning the goal is to learn optimal model parameters via optimization, i.e. **Training = Optimization**. Most models follow some variation of this workflow. In special cases (e.g., linear regression), the optimization problem can be solved analytically (by hand), and numerical optimization isn't necessary. However, this is an exception. Typically, we must solve the optimization numerically, such as through gradient descent or its variations, as outlined in the following diagram.



## Standard Deep Learning Workflow

- Understand the standard deep learning workflow, both **conceptually**, and as **implemented in Keras,** i.e. the following should be very familiar and intuitive to you at this point in the course (note how this is a very general process that applies to most parametric models in some form or another)
  - Implementation of the workflow in Keras can be found here: [click here](#)

1. **Raw Data Analysis**:

- Input data: from the input data X and target Y data be able to tell what type of design decisions can be made about the neural network modeling
    1. e.g. size of input layer, size of output layer, activation on output layer, loss function
2. **Data Preprocessing**: Normalization and One-hot encoding
3. **Data Splitting**: Training, Validation, Testing
4. **Parametric Model Definition**: $y_{pred} = \hat{y} = m(\mathbf{x}|\mathbf{w})$
    1. e.g. linear regression, logistic regression, or neural network
5. **Loss Function Definition**:
    - Loss function definition $L(\mathbf{w})$, for given learning task, to score the quality of a particular parameterization $\mathbf{w}$.
6. **Parameter Initialization**:
    - Random initial guess for model parameters $\mathbf{w}_{t=0}$ where $t$ is an iteration variable for gradient descent (this instantiates a particular version, i.e. parameterization, of the model). Since it is un-trained this random model parameterization will generally be bad (high loss/error)
7. **Training Loop (Gradient Descent)**: over iteration index $t$
    - Model is trained through the training data to optimize model parameters via gradient descent or one of it's "cousins" (Adam, RMSprop, etc)
    - Iteratively optimize parameters:
        - Full vector update: $\mathbf{w}_{t+1} = \mathbf{w}_t - \lambda \nabla_{\mathbf{w}} L(\mathbf{w}_t)$
            - This is the vector expression to update all fitting parameters at once
        - For a single fitting parameter, i.e. the $j^{th}$ fitting parameter, the formula is just $w_{t+1}^j = w_t^j - \lambda \frac{\partial L}{\partial w_t^j}$
    - Parameter stops updating when the gradient is zero.
    - Keep iterating until the convergence or the onset of overfitting (terminate run via early stopping if needed)
    - Monitor validation loss to avoid over fitting
8. **Final Model Evaluation**:
    - Final model transferability check: evaluate on the test data and compare metrics to validation and training set

## Regularization Techniques

- Understand the principles and effects of:
    - L1 regularization
    - L2 regularization
- Explain how these regularizations modify the loss function and influence model parameters.
- Note: Dropout regularization is not included on the exam.

## Back-propagation and Differentiation

- **Not required**: no back-propagation computations on the exam.