

Training neural networks

Vikas Chandrakant Raykar

Microsoft

Collateral

notes - <https://vikasraykar.github.io/deeplearning>

code - <https://github.com/vikasraykar/deeplearning-doj>

IOAI 2025 Syllabus

Section 2: Neural Networks & Deep Learning

Topic	Subtopic	Category
Neural Networks	Perceptron Basics	Both
	Gradient Descent	Both
	Backpropagation	Both
	Activation Functions (ReLU, Sigmoid, Tanh)	Both
	Cost Functions (MSE, MAE, Cross Entropy, etc.)	Both
Deep Learning	Multi-Layer Perceptrons (MLP)	Both
	Stochastic Gradient Descent (SGD), Mini-Batch Gradient Descent	Both
	Momentum Methods (Adam, AdamW)	Practice
	Adaptive Learning Rates	Practice
	Convergence and Learning Rates	Both
	Weight Regularization	Practice
	Early Stopping	Practice
	Dropout, Gaussian Noise	Practice
	Weight Initialization	Practice
	Batch Normalization	Practice
	Autoencoders and Sparse Encoders	Practice

Training neural networks

The goal of training is to find the value of the **parameters** of a neural network **model** to make **effective predictions**.

Training neural networks

We choose the **model parameters** by **optimizing a loss function**.

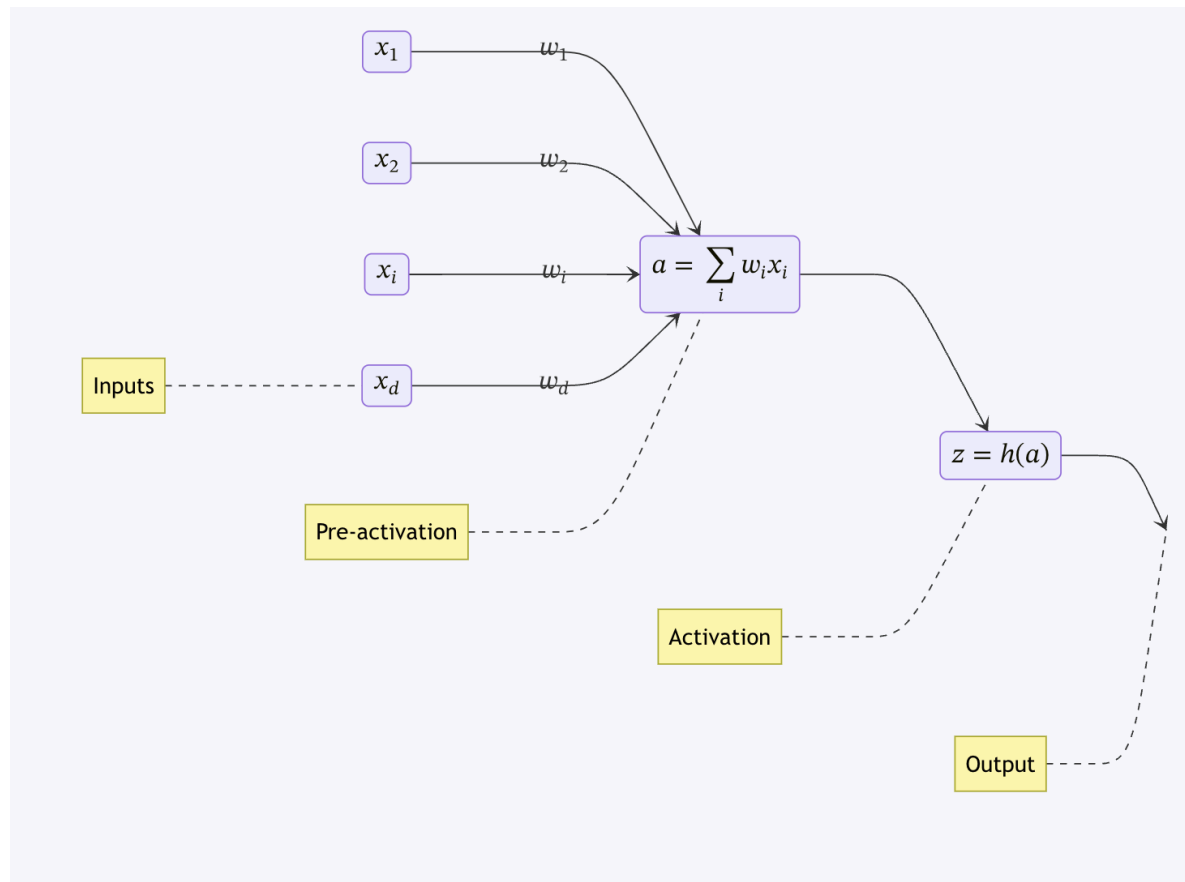
- Model and parameters
- Loss functions
- Gradient Descent
- Optimizers
- Normalization
- Regularization
- Training loop
- Backpropagation and Automatic differentiation

Training neural networks

- Model and parameters
- Loss function
- Gradient Descent
- Optimizers
- Normalization
- Regularization
- Training loop
- Backpropagation and Automatic differentiation

Single Layer Networks

For simplicity we will discuss single layer networks for regression and classification.



Linear Regression

Linear Regression

Linear Regression assumes a **linear relationship** between the target $y \in \mathbb{R}$ and the features $\mathbf{x} \in \mathbb{R}^d$.

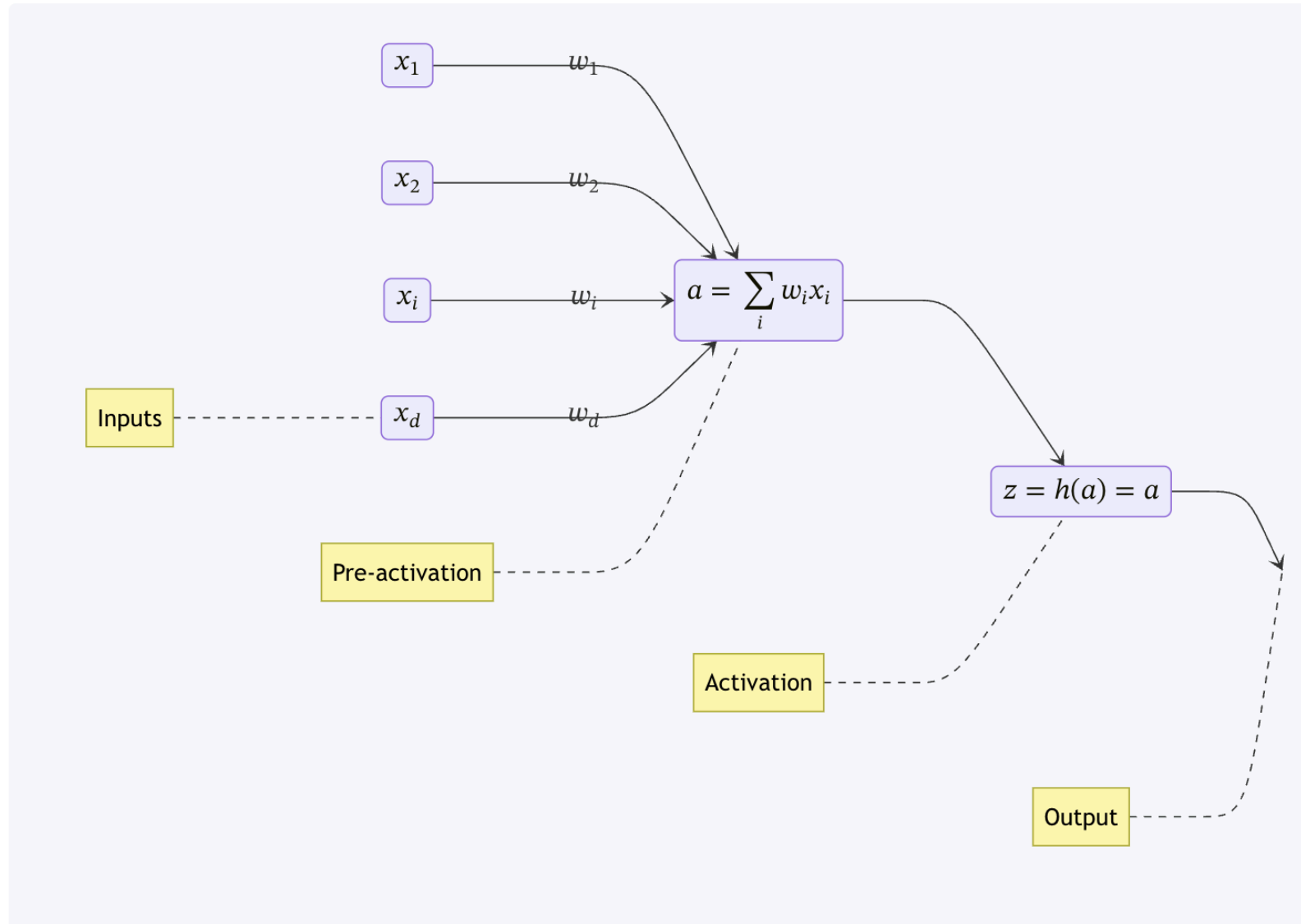
$$y = f(\mathbf{x}) = w_1x_1 + w_2x_2 + \dots + w_dx_d + b = \mathbf{w}^T \mathbf{x} + b,$$

where $\mathbf{w} \in \mathbb{R}^d$ is the d -dimensional *weight vector* and $b \in \mathbb{R}$ is the *bias term*.

Without loss of generalization we can ignore the bias term as it can be subsumed into the feature matrix by appending a column of ones.

$$y = f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$$

Linear Regression is a single layer neural network for regression.



Probabilistic model

The probability of y for a given feature vector ($\mathbf{x} \in \mathbb{R}^d$) is modelled as

$$\Pr[y|\mathbf{x}, \mathbf{w}] = \mathcal{N}(y|\mathbf{w}^T \mathbf{x}, \sigma^2)$$

where $\mathbf{w} \in \mathbb{R}^d$ are the weights/**parameters** of the model and \mathcal{N} is the **normal** distribution with mean $\mathbf{w}^T \mathbf{x}$ and variance σ^2 .

The prediction is given by

$$\mathbb{E}[y|\mathbf{x}, \mathbf{w}] = \mathbf{w}^T \mathbf{x}$$

Negative log likelihood

Given a dataset $\mathcal{D} = \{\mathbf{x}_i \in \mathbb{R}^d, \mathbf{y}_i \in \mathbb{R}\}_{i=1}^N$ containing n examples we need to estimate the parameter vector \mathbf{w} by maximizing the likelihood of data.

In practice we minimize the **negative log likelihood**.

Let $\mu_i = \mathbf{w}^T \mathbf{x}_i$ be the model prediction for each example in the training dataset.

The negative log likelihood (NLL) is given by

$$L(\mathbf{w}) = - \sum_{i=1}^N \log [\Pr[y_i | \mathbf{x}_i, \mathbf{w}]] = \frac{N}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \mu_i)^2$$

Mean Squared Error loss

This is equivalent to minimizing the **Mean Squared Error** (MSE) loss.

$$L(\mathbf{w}) = \frac{1}{2N} \sum_{i=1}^N (y_i - \mu_i)^2 = \frac{1}{2N} \sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{x}_i)^2$$

We need to choose the model parameters that optimizes (minimizes) the loss function.

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} L(\mathbf{w})$$

```
torch.nn.MSELoss
```

Vectorization

We often stack all the n examples into a *design matrix* $\mathbf{X} \in \mathbb{R}^{n \times d}$, where each row is one instance. The predictions for all the n instances $\mathbf{y} \in \mathbb{R}^n$ can be written conveniently as a matrix-vector product.

$$\mathbf{y} = \mathbf{X}\mathbf{w}$$

The loss function using matrix notation can be written as follows.

$$L(\mathbf{w}) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2$$

Quiz - Compute the gradient of the loss function.

Linear Regression Vectorization

We often stack all the n examples into a *design matrix* $\mathbf{X} \in \mathbb{R}^{n \times d}$, where each row is one instance. The predictions for all the n instances $\mathbf{y} \in \mathbb{R}^n$ can be written conveniently as a matrix-vector product.

$$\mathbf{y} = \mathbf{X}\mathbf{w}$$

The loss function using matrix notation can be written as follows.

$$L(\mathbf{w}) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2$$

The gradient of the loss function is given by

$$\nabla_{\mathbf{w}} L(\mathbf{w}) = \mathbf{X}^T (\mathbf{X}\mathbf{w} - \mathbf{y})$$

Quiz - Compute the analytic solution.

Logistic Regression

Logistic Regression

The probability of the positive class ($y = 1$) for a given feature vector ($\mathbf{x} \in \mathbb{R}^d$) is given by

$$\Pr[y = 1 | \mathbf{x}, \mathbf{w}] = \sigma(\mathbf{w}^T \mathbf{x})$$

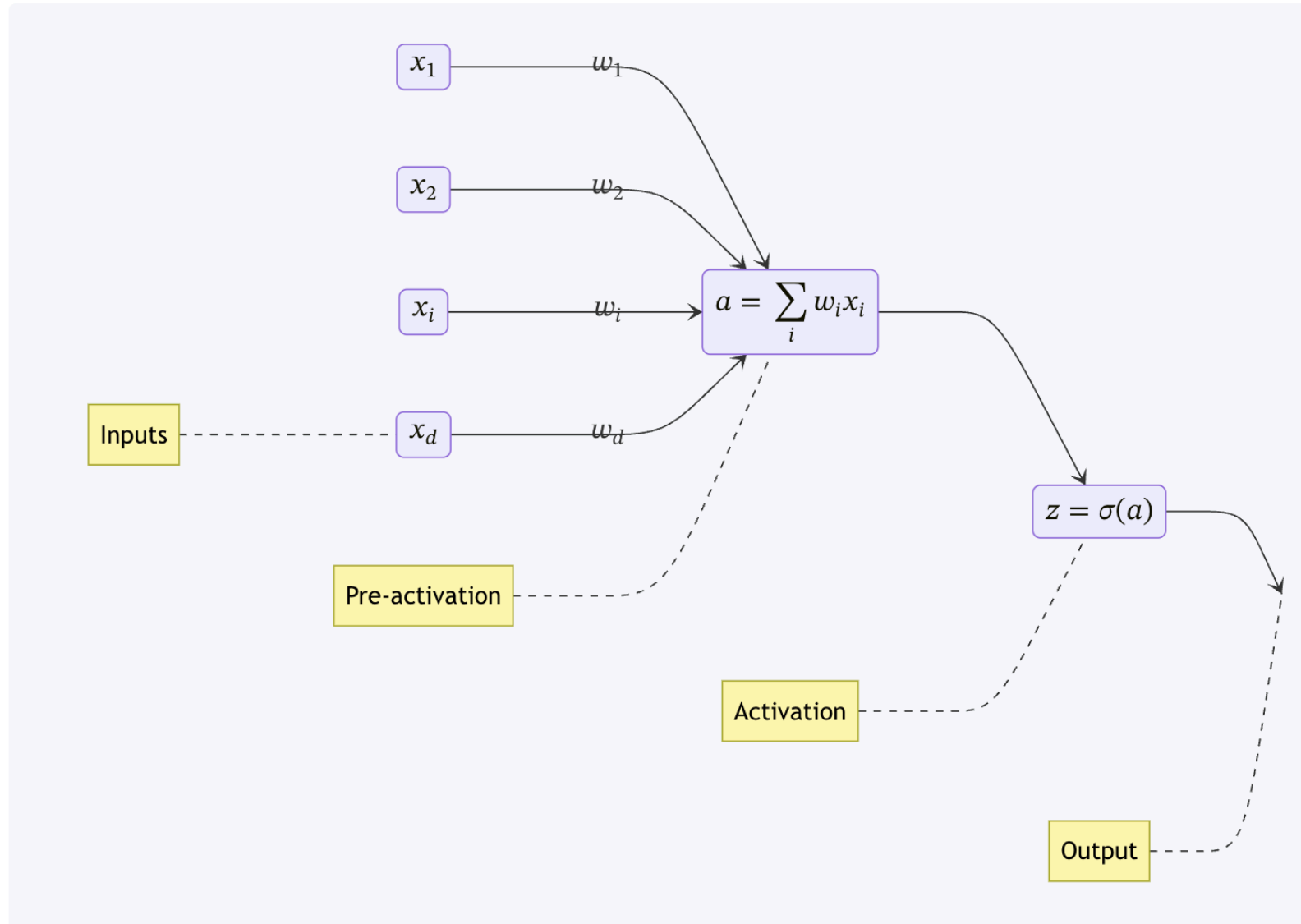
where $\mathbf{w} \in \mathbb{R}^d$ are the weights/**parameters** of the model and σ is the **sigmoid** activation function defined as

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

Without loss of generalization we ignore the bias term as it can be incorporated into the feature vector.

Quiz - *Plot the sigmoid function.*

Logistic Regression is a single layer neural network for binary classification.



Negative log likelihood

Given a dataset $\mathcal{D} = \{\mathbf{x}_i \in \mathbb{R}^d, \mathbf{y}_i \in [0, 1]\}_{i=1}^N$ containing n examples we need to estimate the parameter vector \mathbf{w} by maximizing the likelihood of data.

In practice we minimize the **negative log likelihood**.

Let $\mu_i = \Pr[y_i = 1 | \mathbf{x}_i, \mathbf{w}] = \sigma(\mathbf{w}^T \mathbf{x}_i)$ be the model prediction for each example in the training dataset.

The the negative log likelihood (NLL) is given by

$$L(\mathbf{w}) = - \sum_{i=1}^N \log [\mu_i^{y_i} (1 - \mu_i)^{1-y_i}] = - \sum_{i=1}^N [y_i \log(\mu_i) + (1 - y_i) \log(1 - \mu_i)]$$

Binary cross entropy loss

This is referred to as the **Binary Cross Entropy** (BCE) loss.

$$L(\mathbf{w}) = - \sum_{i=1}^N [y_i \log(\mu_i) + (1 - y_i) \log(1 - \mu_i)]$$

We need to choose the model parameters that optimizes (minimizes) the loss function.

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} L(\mathbf{w})$$

`torch.nn.BCELoss`

`torch.nn.BCEWithLogitsLoss`

Quiz - *Why is this called cross entropy ?*

Quiz - *Compute the gradient of the loss function.*

Gradient

The gradient of the loss function is given by

$$\nabla_{\mathbf{w}} L(\mathbf{w}) = \mathbf{X}^T (\sigma(\mathbf{X}\mathbf{w}) - \mathbf{y})$$

Compare with the gradient of linear regression.

$$\nabla_{\mathbf{w}} L(\mathbf{w}) = \mathbf{X}^T (\mathbf{X}\mathbf{w} - \mathbf{y})$$

Softmax regression

Multi-class logistic regression

Given k classes the probability of class i for a given feature vector ($\mathbf{x} \in \mathbb{R}^d$) is given by

$$\Pr[y = i | \mathbf{x}, (\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_k)] = \frac{\exp(\mathbf{w}_i^T \mathbf{x})}{\sum_{j=1}^k \exp(\mathbf{w}_j^T \mathbf{x})}$$

where $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_k \in \mathbb{R}^d$ are the weight vector or **parameters** of the model for each class.

Stacking the weights vectors $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_k \in \mathbb{R}^d$ into a **weight matrix** $\mathbf{W} \in \mathbb{R}^{d \times k}$ we can write the **pre-activation** vector $\mathbf{a} \in \mathbb{R}^k$ as follows.

$$\mathbf{a} = \mathbf{W}^T \mathbf{x}$$

The **activation** vector $\mathbf{z} \in \mathbb{R}^k$ is given by

$$\mathbf{z} = \text{softmax}(\mathbf{a})$$

and the **softmax** activation function is defined as

$$\text{softmax}(\mathbf{a})_i = \frac{\exp(\mathbf{a}_i)}{\sum_{j=1}^k \exp(\mathbf{a}_j)}$$

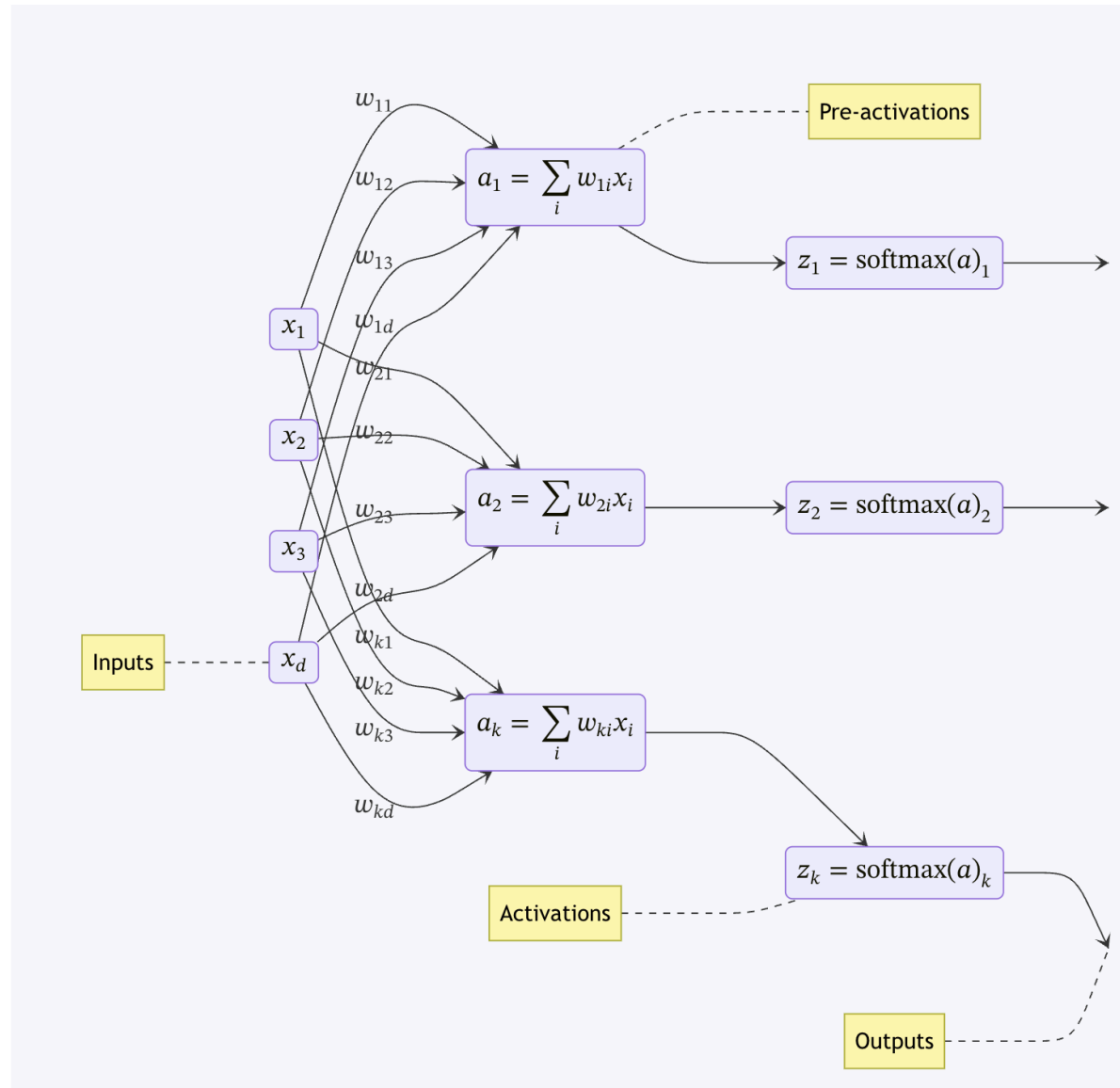
Hence

$$\Pr[y = i | \mathbf{x}, \mathbf{W}] = \text{softmax}(\mathbf{W}^T \mathbf{x})_i$$

We often stack all the n examples into a *design matrix* $\mathbf{X} \in \mathbb{R}^{n \times d}$, where each row is one instance. The predictions for all the n instances $\mathbf{y} \in \mathbb{R}^{n \times k}$ can be written conveniently as a matrix-vector product.

$$\mathbf{y} = \text{softmax}(\mathbf{X}\mathbf{W})$$

Softmax Regression is a single layer neural network for multi-class classification.



Likelihood

Given a dataset $\mathcal{D} = \{\mathbf{x}_i \in \mathbb{R}^d, \mathbf{y}_i \in [1, 2, \dots, k]\}_{i=1}^n$ containing n examples we need to estimate the parameter vector \mathbf{W} by maximizing the likelihood of data.

Let μ_i^j be the model prediction for class j for each example i in the training dataset.

$$\mu_i^j = \Pr[y_i = j | \mathbf{x}_i, \mathbf{W}] = \text{softmax}(\mathbf{W}^T \mathbf{x}_i)_j$$

Let y_i^j be the corresponding true label.

The negative log likelihood (NLL) is given by

$$L(\mathbf{W}) = - \sum_{i=1}^n \sum_{j=1}^k y_i^j \log \mu_i^j$$

Loss function

Cross Entropy loss

$$L(\mathbf{W}) = - \sum_{i=1}^n \sum_{j=1}^k y_i^j \log \mu_i^j$$

Compare to the earlier **Binary Cross Entropy** loss

$$L(\mathbf{w}) = \sum_{i=1}^n [y_i \log(\mu_i) + (1 - y_i) \log(1 - \mu_i)]$$

`torch.nn.CrossEntropyLoss`

Summary

	Model	Parameters	Loss function
Linear Regression	$\mu = \mathbf{w}^T \mathbf{x}$	$\mathbf{w} \in \mathbb{R}^{d+1}$	MSE $\frac{1}{2n} \sum_{i=1}^N (y_i - \mu_i)^2$
Logistic Regression	$\mu = \sigma(\mathbf{w}^T \mathbf{x})$	$\mathbf{w} \in \mathbb{R}^{d+1}$	BCE $-\sum_{i=1}^n [y_i \log(\mu_i) + (1 - y_i) \log(1 - \mu_i)]$
Softmax regression	$\mu = \text{softmax}(\mathbf{W}^T \mathbf{x})$	$\mathbf{W} \in \mathbb{R}^{d+1 \times k}$	CE $-\sum_{i=1}^n \sum_{j=1}^k y_i^j \log \mu_i^j$

Pytorch loss functions

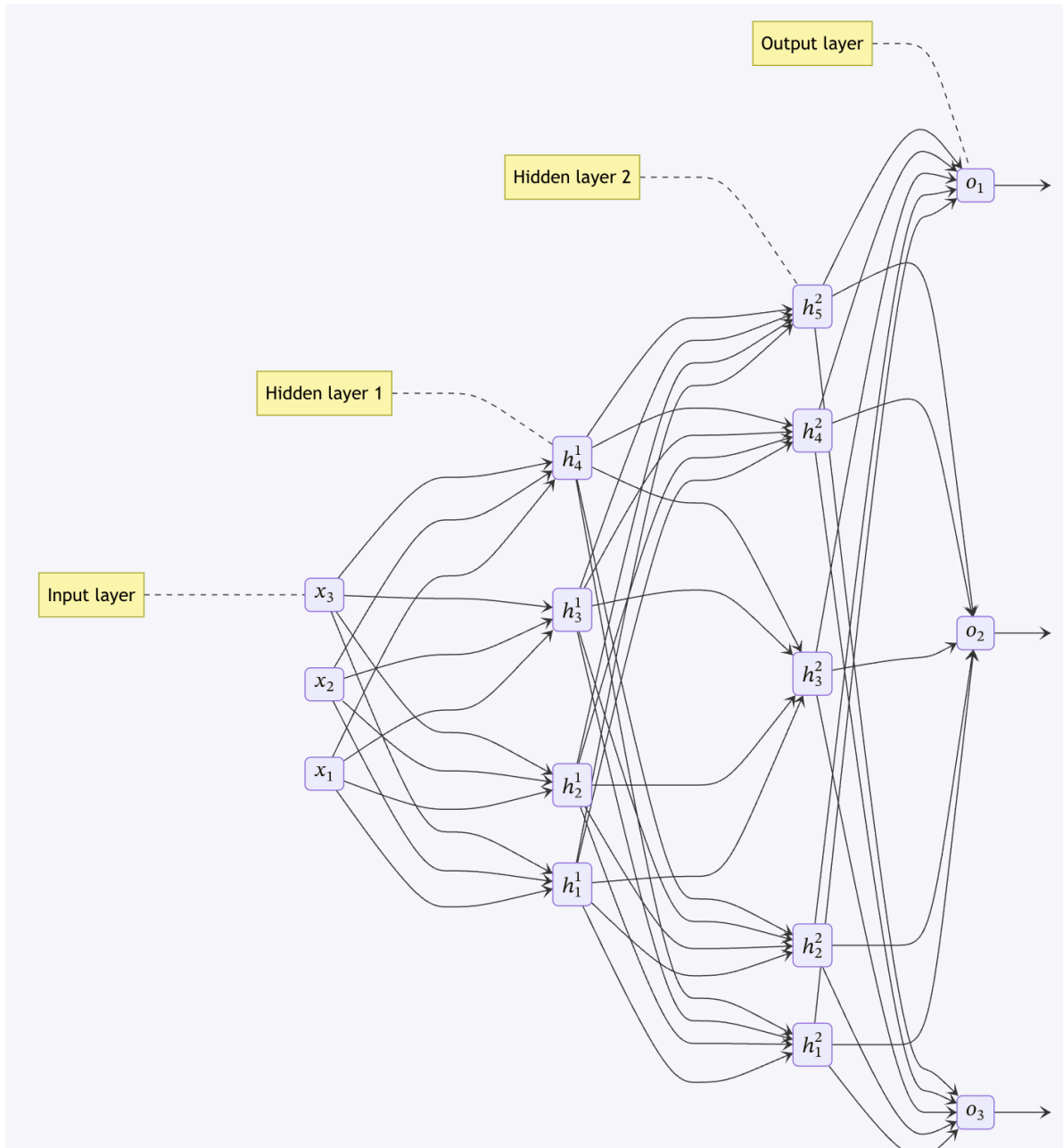
Quiz - Can we use MSE for Logistic Regression ?

Multilayer perceptrons

A 3-layer multilayer perceptron.

$$\begin{aligned}\mathbf{X} &= \mathbf{X} \\ \mathbf{H}^{(1)} &= g_1 \left(\mathbf{X} \mathbf{W}^{(1)} + \mathbf{b}^{(1)} \right) \\ \mathbf{H}^{(2)} &= g_2 \left(\mathbf{H}^{(1)} \mathbf{W}^{(2)} + \mathbf{b}^{(2)} \right) \\ \mathbf{O} &= \mathbf{H}^{(2)} \mathbf{W}^{(3)} + \mathbf{b}^{(3)}\end{aligned}$$

g is a nonlinear **activation function**



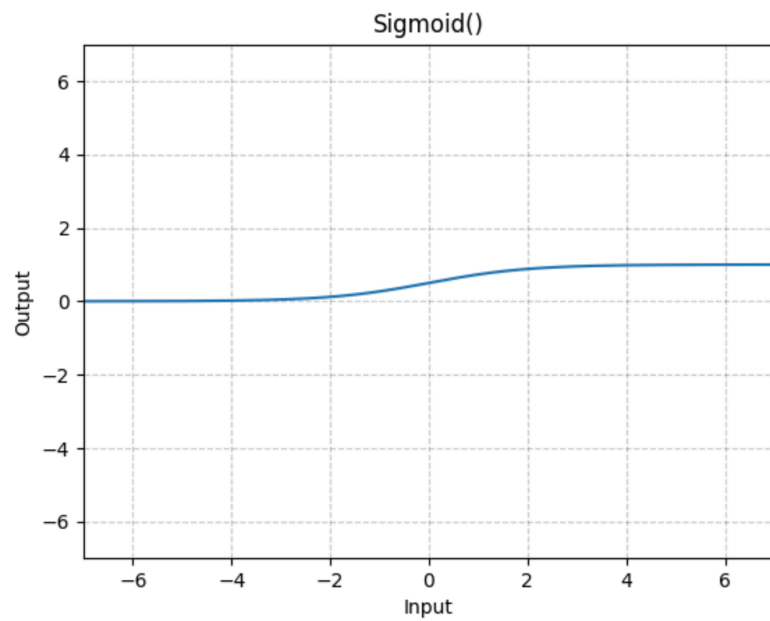
Activation functions

`torch.nn`

Sigmoid

Sigmoid/Logistic

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

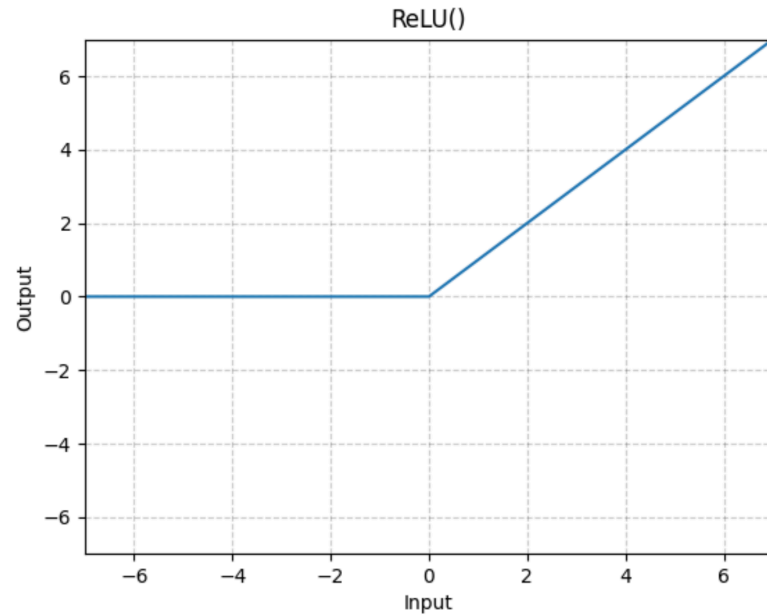


The derivative is given by $\sigma'(z) = \sigma(z)(1 - \sigma(z))$.

ReLU

Rectified Linear Unit (ReLU)

$$\text{ReLU}(z) = \max(z, 0)$$

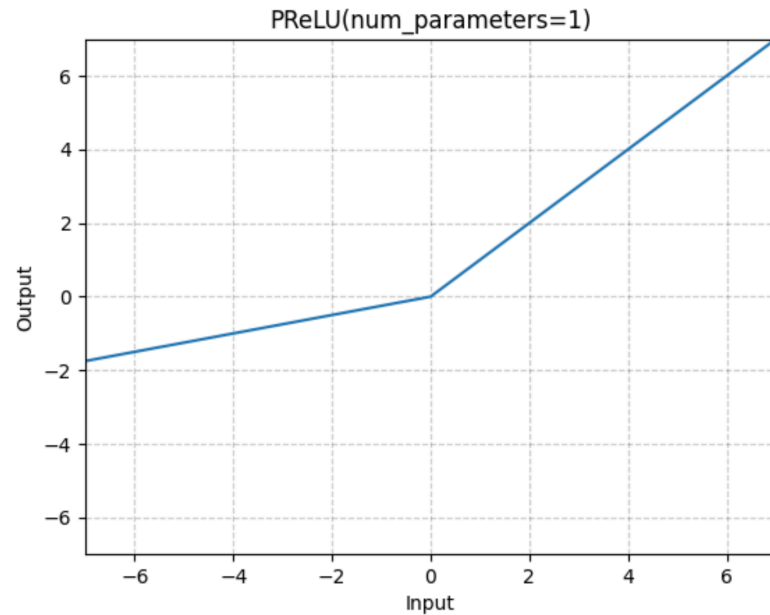


Nair and Hinton, 2010

pReLU

parameterized Rectified Linear Unit (pReLU)

$$\text{pReLU}_{\alpha}(z) = \max(z, 0) + \alpha \min(z, 0)$$

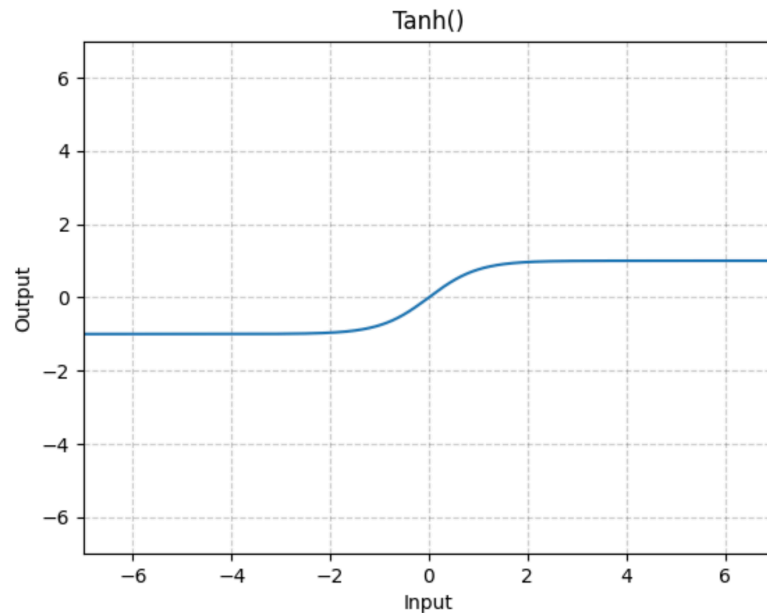


He et al., 2015

Tanh

Hyperbolic tangent.

$$\tanh(z) = \frac{1 - \exp(-2z)}{1 + \exp(-2z)}$$



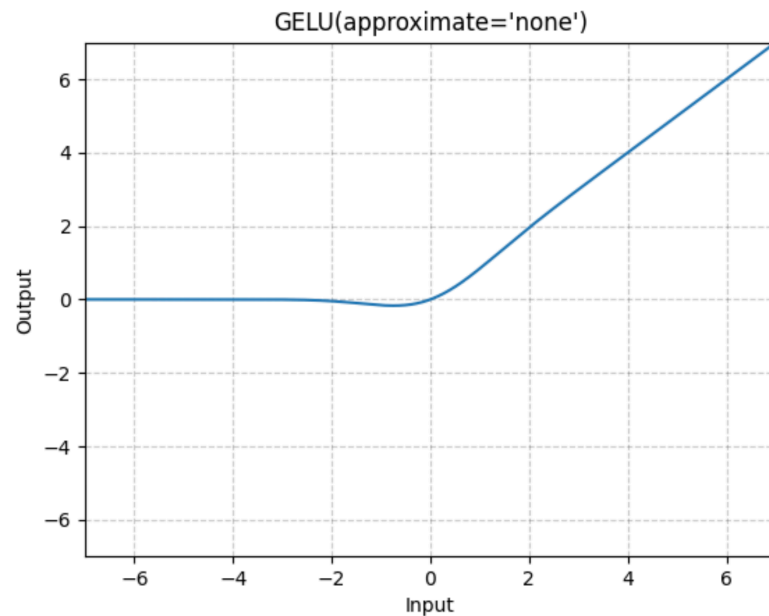
The derivative is given by $\tanh'(z) = 1 - \tanh^2(z)$

GeLU

Gaussian error Linear Unit/Smooth ReLU

$$\text{GeLU}(z) = z\Phi(z)$$

where $\Phi(z)$ is the standard Gaussian cumulative distribution.



A brief primer on entropy, cross-entropy and perplexity.

Entropy

The **entropy** of a discrete random variable X with K states/categories with distribution $p_k = \Pr(X = k)$ for $k = 1, \dots, K$ is a measure of uncertainty and is defined as follows.

$$H(X) = \sum_{k=1}^K p_k \log_2 \frac{1}{p_k} = - \sum_{k=1}^K p_k \log_2 p_k$$

The term $\log_2 \frac{1}{p}$ quantifies the notion of surprise or uncertainty and hence entropy is the average uncertainty.

The unit is bits ($\in [0, \log_2 K]$) (or nats in case of natural log).

Quiz : *What is the discrete distribution with maximum entropy?*

Quiz : *What is the discrete distribution with minimum entropy?*

Entropy

The **entropy** of a discrete random variable X with K states/categories with distribution $p_k = \Pr(X = k)$ for $k = 1, \dots, K$ is a measure of uncertainty and is defined as follows.

$$H(X) = \sum_{k=1}^K p_k \log_2 \frac{1}{p_k} = - \sum_{k=1}^K p_k \log_2 p_k$$

The term $\log_2 \frac{1}{p}$ quantifies the notion of surprise or uncertainty and hence entropy is the average uncertainty. The unit is bits ($\in [0, \log_2 K]$) (or nats in case of natural log).

The discrete distribution with maximum entropy ($\log_2 K$) is uniform.

The discrete distribution with minimum entropy (0) is any delta function which puts all mass on one state/category.

Binary entropy

For a binary random variable $X \in \{0, 1\}$ with $\Pr(X = 1) = \theta$ and $\Pr(X = 0) = 1 - \theta$ the entropy is as follows.

$$H(\theta) = -[\theta \log_2 \theta + (1 - \theta) \log_2(1 - \theta)]$$

The range is $H(\theta) \in [0, 1]$ and is maximum when $\theta = 0.5$.

Quiz : *Plot the binary entropy..*

Cross entropy

Cross entropy is the average number of bits needed to encode the data from a source p when we model it using q .

$$H(p, q) = - \sum_{k=1}^K p_k \log_2 q_k$$

Perplexity

$$\text{PPL}(p, q) = 2^{H(p, q)}$$

$$\text{PPL}(p, q) = e^{H(p, q)}$$

KL Divergence

The **Kullback-Leibler** (KL) divergence or **relative entropy** measures the dissimilarity between two probability distributions p and q .

$$KL(p, q) = \sum_{k=1}^K p_k \log_2 \frac{p_k}{q_k}$$

$$KL(p, q) = H(p, q) - H(p, p) \geq 0$$

Mutual Information

$$I(X, Y) = KL(P(X, Y) \| P(X)P(Y))$$

Training neural networks

- **Model and parameters**
 - Single layer neural networks (Linear Regression, Logistic Regression, Softmax Regression)
 - Multilayer neural networks
 - Activation functions (ReLU, pReLU, Sigmoid, Tanh, GeLU)
- **Loss functions**
 - Mean Squared Error loss, Binary Cross Entropy loss, Cross Entropy Loss
 - Concept of entropy and cross entropy
- **Gradient Descent**
- Optimizers
- Normalization
- Regularization
- Training loop
- Backpropagation and Automatic differentiation

Parameter estimation

Let \mathbf{w} be a vector of all the parameters for a model.

Let $L(\mathbf{w})$ be the loss function (or error function).

We need to choose the model parameters that optimizes (minimizes) the loss function.

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} L(\mathbf{w})$$

We will do this via **gradient descent**.

Gradient Descent

Steepest descent.

Let $\nabla L(\mathbf{w})$ be the **gradient vector**, where each element is the partial derivative of the loss function with respect to each parameter.

The gradient vector points in the **direction of the greatest rate of increase of the loss function**.

So to minimize the loss function we take small steps in the direction of $-\nabla L(\mathbf{w})$.

At the minimum $\nabla L(\mathbf{w}) = 0$.

Stationary points

$\nabla L(\mathbf{w}) = 0$ are known as stationary points, which can be either be

- minima (local or global)
- maxima (local or global)
- saddle point

The necessary and sufficient condition for a local minima is

1. The gradient of the loss function should be zero.
2. The Hessian matrix should be positive definite.

Gradient descent

- Batch Gradient Descent
- Stochastic Gradient Descent
- Min-batch Stochastic Gradient Descent

For now we will assume the gradient is given. For deep neural networks the gradient can be computed efficiently via **backpropagation**(which we will revisit later).

Batch Gradient Descent

We take a small step in the direction of the **negative gradient**.

$$\mathbf{w}^t \leftarrow \mathbf{w}^{t-1} - \eta \nabla L(\mathbf{w}^{t-1})$$

The parameter $\eta > 0$ is called the **learning rate** and determines the step size at each iteration. This update is repeated multiple times (till convergence).

```
for epoch in range(n_epochs):  
    dw = gradient(loss, data, w)  
    w = w - lr * dw
```

Each step requires that the **entire training data** be processed to compute the gradient $\nabla L(\mathbf{w}^{t-1})$. For large datasets this is not computationally efficient.

Stochastic Gradient Descent

In general most loss functions can be written as sum over each training instance.

$$L(\mathbf{w}) = \sum_{i=1}^N L_i(\mathbf{w})$$

In Stochastic Gradient Descent (SGD) we update the parameters **one data point at a time**.

$$\mathbf{w}^t \leftarrow \mathbf{w}^{t-1} - \eta \nabla L_i(\mathbf{w}^{t-1})$$

A complete passthrough of the whole dataset is called an **epoch**.

```
for epoch in range(n_epochs):  
    for i in range(n_data):  
        dw = gradient(loss, data[i], w)  
        w = w - lr * dw
```

Stochastic Gradient Descent

- SGD is much faster and more computationally efficient, but it has noise in the estimation of the gradient.
- Since it updates the weight frequently, it can lead to big oscillations and that makes the training process highly unstable.

Bottou, L. (2010). [Large-Scale Machine Learning with Stochastic Gradient Descent](#).
In: Lechevallier, Y., Saporta, G. (eds) Proceedings of COMPSTAT'2010. Physica-Verlag HD.

Mini-batch Stochastic Gradient Descent

Using a single example results in a very noisy estimate of the gradient.

So we use a small random subset of data called **mini-batch** of size B (**batch size**) to compute the gradient.

$$\mathbf{w}^t \leftarrow \mathbf{w}^{t-1} - \eta \nabla L_{batch}(\mathbf{w}^{t-1})$$

```
for epoch in range(n_epochs):  
    for mini_batch in get_batches(data, batch_size):  
        dw = gradient(loss, mini_batch, w)  
        w = w - lr * dw
```

Mini-batch Stochastic Gradient Descent

Mini-batch SGD is the most commonly used method and is sometimes referred to as just SGD.

- Typical choices of the batch size are $B=32, 64, 128, 256, \dots$ (power of 2)
- In practice we do a random shuffle of the data per epoch.

In practice, mini-batch SGD is the most frequently used variation because it is both computationally cheap and results in more robust convergence.

```
torch.optim.SGD
```

```
optimizer = optim.SGD(model.parameters(), lr=1e-3)
```

Training neural networks

- **Model and parameters**
 - Single layer neural networks (Linear Regression, Logistic Regression, Softmax Regression)
 - Multilayer neural networks
 - Activation functions (ReLU, pReLU, Sigmoid, Tanh, GeLU)
- **Loss functions**
 - Mean Squared Error loss, Binary Cross Entropy loss, Cross Entropy Loss
 - Concept of entropy and cross entropy
- **Gradient Descent**
 - Batch Gradient Descent
 - Stochastic Gradient Descent
 - Mini-batch Stochastic Gradient Descent
- **Optimizers**
- Backpropagation and Automatic differentiation
- Normalization
- Regularization
- Training loop

Optimizers

Improvements over min-batch stochastic gradient descent for **faster convergence** and **stability**.

- Adding momentum.
- Adaptive learning rates.
 - Adagrad
 - RMSProp
 - Adam

Adding momentum

One of the basic improvements over SGD comes from adding a **momentum** term.

At every time step, we update **velocity** by decaying the previous velocity by a factor of $0 \leq \mu \leq 1$ (called the **momentum** parameter) and adding the current gradient update.

$$\mathbf{v}^{t-1} \leftarrow \mu \mathbf{v}^{t-2} - \eta \nabla L(\mathbf{w}^{t-1})$$

Then, we update our weights in the direction of the velocity vector.

$$\mathbf{w}^t \leftarrow \mathbf{w}^{t-1} + \mathbf{v}^{t-1}$$

```
for epoch in range(n_epochs):  
    for mini_batch in get_batches(data, batch_size):  
        dw = gradient(loss, mini_batch, w) # gradient  
        v = momentum * v - lr * dw # velocity  
        w = w + v
```

Adding momentum

We now have two hyper-parameters **learning rate** and **momentum**.

Typically we set the momentum parameter to 0.9.

One interpretation of momentum to increase the effective learning rate from η to $\frac{\eta}{(1-\mu)}$.

```
torch.optim.SGD
```

```
optimizer = optim.SGD(model.parameters(), lr=0.01, momentum=0.9)
```

Adding momentum

- We can now escape local minima or saddle points because we keep moving downwards even though the gradient of the mini-batch might be zero.
- Momentum can also help us reduce the oscillation of the gradients because the velocity vectors can smooth out these highly changing landscapes.
- It reduces the noise of the gradients and follows a more direct walk down the landscape.

Why Momentum Really Works

Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. 2013. [On the importance of initialization and momentum in deep learning](#). In Proceedings of the 30th International Conference on International Conference on Machine Learning - Volume 28 (ICML'13). JMLR.org, III-1139-III-1147.

Adding momentum

Adaptive Learning Rates

Different learning rate for each parameter.

- Adagrad
- RMSProp
- Adam
- AdamW

Adagrad

Adaptive gradient.

AdaGrad reduces each learning rate parameter over time by using the accumulated sum of squares of all the derivatives calculated for that parameter.

$$\mathbf{w}^t \leftarrow \mathbf{w}^{t-1} - \frac{\eta}{\sqrt{\mathbf{r}^t} + \delta} \odot \nabla L(\mathbf{w}^{t-1})$$

where \mathbf{r}^t is the running sum of the squares of the gradients and δ is a small constant to ensure numerical stability.

$$\mathbf{r}^t = \mathbf{r}^{t-1} + (\nabla L(\mathbf{w}^t))^2$$

Adagrad

```
for epoch in range(n_epochs):  
    for mini_batch in get_batches(data, batch_size):  
        dw = gradient(loss, mini_batch, w) # gradient  
        r += dw*dw # Accumulated squared gradients  
        w = w - lr * dw / (r.sqrt() + delta)
```

`torch.optim.Adagrad`

```
optimizer = torch.optim.Adagrad(model.parameters(), lr=0.01, eps=1e-10)
```


Adagrad

We can see that when the gradient is changing very fast, the learning rate will be smaller. When the gradient is changing slowly, the learning rate will be bigger.

A drawback of Adagrad is that as time goes by, the learning rate becomes smaller and smaller due to the monotonic increment of the running squared sum.

John Duchi, Elad Hazan, and Yoram Singer. 2011. [Adaptive Subgradient Methods for Online Learning and Stochastic Optimization](#). J. Mach. Learn. Res. 12, null (2/1/2011), 2121–2159.

RMSProp

Root Mean Square Propagation, Leaky AdaGrad

Since AdaGrad accumulates the squared gradients from the beginning, the associated weight updates can become very small as training progresses.

RMSProp essentially replaces it with an **exponentially weighted average**.

$$\mathbf{r}^t = \alpha \mathbf{r}^{t-1} + (1 - \alpha) (\nabla L(\mathbf{w}^t))^2$$

where $0 < \alpha < 1$.

$$\mathbf{w}^t \leftarrow \mathbf{w}^{t-1} - \frac{\eta}{\sqrt{\mathbf{r}^t} + \delta} \odot \nabla L(\mathbf{w}^{t-1})$$

RMSProp

```
for epoch in range(n_epochs):  
    for mini_batch in get_batches(data, batch_size):  
        dw = gradient(loss, mini_batch, w) # gradient  
        r += alpha * r + (1-alpha) * dw*dw # Accumulated squared gradients  
        w = w - lr * dw / (r.sqrt() + delta)
```

`torch.optim.RMSprop`

```
optimizer = torch.optim.RMSProp(model.parameters(), lr=0.01, alpha=0.99, eps=1e-8)
```

Typically we set the $\alpha = 0.9$.

Hinton, 2012. Neural Networks for Machine Learning. [Lecture 6a](#).

Adam

Adaptive moments.

If we combine RMSProp with momentum we obtain the most popular Adam optimization method.

Kingma, D.P. and Ba, J., 2014. [Adam: A method for stochastic optimization](#). arXiv preprint arXiv:1412.6980.

Adam

Adam maintains an exponentially weighted average of the first and the second moments.

$$\mathbf{s}^t = \beta_1 \mathbf{s}^{t-1} + (1 - \beta_1) (\nabla L(\mathbf{w}^t))$$

$$\mathbf{r}^t = \beta_2 \mathbf{r}^{t-1} + (1 - \beta_2) (\nabla L(\mathbf{w}^t))^2$$

We correct for the bias introduced by initializing \mathbf{s}^0 and \mathbf{r}^0 to zero.

$$\hat{\mathbf{s}}^t = \frac{\mathbf{s}^t}{1 - \beta_1^t} \quad \hat{\mathbf{r}}^t = \frac{\mathbf{r}^t}{1 - \beta_2^t}$$

The updates are given as follows.

$$\mathbf{w}^t \leftarrow \mathbf{w}^{t-1} - \frac{\eta}{\sqrt{\hat{\mathbf{r}}^t} + \delta} \odot \hat{\mathbf{s}}^t$$

Adam

```
for epoch in range(n_epochs):  
    for mini_batch in get_batches(data, batch_size):  
        dw = gradient(loss, mini_batch, w) # gradient  
        s += beta1 * s + (1-beta1) * dw # Accumulated gradients  
        r += beta2 * r + (1-beta2) * dw*dw # Accumulated squared gradients  
        s_hat = s / (1-beta1**t)  
        r_hat = r / (1-beta2**t)  
        w = w - lr * s_hat / (r_hat.sqrt() + delta)
```

`torch.optim.Adam`

```
optimizer = optim.Adam(model.parameters(), lr=0.001, betas=(0.9,0.99), eps=1e-08)
```

Typically we set the $\beta_1 = 0.9$ and $\beta_2 = 0.99$.

AdamW

AdamW proposes a modification to Adam that improves regularization by adding **weight decay**. At each iteration we pull the parameters towards zero.

$$\mathbf{w}^t \leftarrow \mathbf{w}^t - \eta \lambda \mathbf{w}^t$$

```
torch.optim.AdamW
```

```
optimizer = optim.AdamW(model.parameters(), lr=0.001, betas=(0.9,0.99), eps=1e-08, weight_decay=0.01)
```

Ilya Loshchilov, Frank Hutter, [Decoupled Weight Decay Regularization](#), ICLR 2019.

Adam and AdamW are the most widely used optimizers.

Learning rate schedule

A small learning rate leads to slow convergence while a large learning rate leads to instability (due to divergent oscillations).

In practice we start with a large learning rate and then reduce it over time.

$$\mathbf{w}^t \leftarrow \mathbf{w}^{t-1} - \eta^{t-1} \nabla L(\mathbf{w}^{t-1})$$

Learning rate schedule

	Learning rate schedule
Linear	$\eta^t = \left(1 - \frac{t}{K}\right)\eta^0 + \left(\frac{t}{K}\right)\eta^K$
Power	$\eta^t = \eta^0 \left(1 + \frac{t}{s}\right)^c$
Exponential	$\eta^t = \eta^0 c^{\frac{t}{s}}$

```
from torch.optim import SGD
from torch.optim.lr_scheduler import ExponentialLR

optimizer = SGD(model.parameters(), lr=0.01, momentum=0.9)
scheduler = ExponentialLR(optimizer, gamma=0.9)
```

Parameter initialization

Initialization before starting the gradient descent.

Avoid all parameters set to same value. (**symmetry breaking**)

Uniform distribution in the range $[-\epsilon, \epsilon]$

Zero-mean Gaussian $\mathcal{N}(0, \epsilon^2)$

```
nn.init
```

Training neural networks

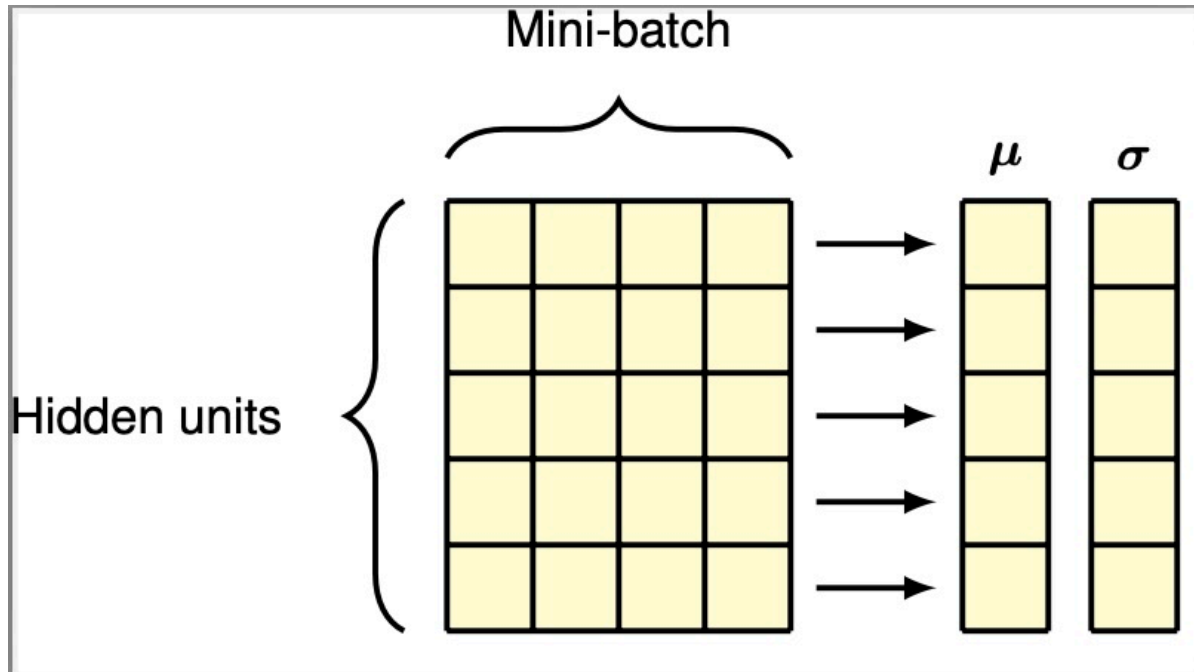
- **Model and parameters**
 - Single layer neural networks (Linear Regression, Logistic Regression, Softmax Regression)
 - Multilayer neural networks, Activation functions (ReLU, pReLU, Sigmoid, Tanh, GeLU)
- **Loss functions**
 - Mean Squared Error loss, Binary Cross Entropy loss, Cross Entropy Loss
 - Concept of entropy and cross entropy
- **Gradient Descent**
 - Batch, Stochastic, Mini-batch Stochastic Gradient Descent
- **Optimizers**
 - Momentum and Adaptive learning rates (Adagrad, RMSProp, Adam)
 - Learning rate schedule
 - Parameter initialization
- **Normalization**
- **Regularization**
- **Training loop**
- **Backpropagation and Automatic differentiation**

Normalization

Normalization is sometimes important for effective training and mitigating vanishing and exploding gradients.

- Batch Normalization
- Layer Normalization

Batch normalization



In batch normalization the mean and variance are computed across the mini-batch separately for each feature/hidden unit.

For a mini-batch of size B

$$\mu_i = \frac{1}{B} \sum_{n=1}^B a_{ni} \quad \sigma_i^2 = \frac{1}{B} \sum_{n=1}^B (a_{ni} - \mu_i)^2$$

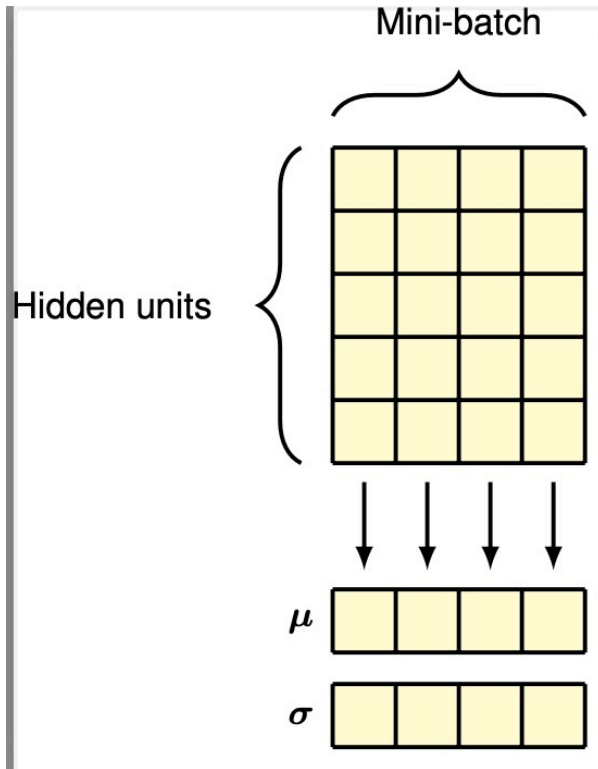
We normalize the pre-activations as follows.

$$\hat{a}_{ni} = \frac{a_{ni} - \mu_i}{\sqrt{\sigma_i^2 + \delta}}$$
$$\tilde{a}_{ni} = \gamma_i \hat{a}_{ni} + \beta_i$$

```
torch.nn.BatchNorm1d
```

```
m = nn.BatchNorm1d(num_features)
```

Layer normalization



In layer normalization the mean and variance are computed across the feature/hidden unit for each example separately.

$$\mu_n = \frac{1}{M} \sum_{i=1}^M a_{ni}$$
$$\sigma_n^2 = \frac{1}{M} \sum_{i=1}^M (a_{ni} - \mu_n)^2$$

We normalize the pre-activations as follows.

$$\hat{a}_{ni} = \frac{a_{ni} - \mu_n}{\sqrt{\sigma_n^2 + \delta}}$$
$$\tilde{a}_{ni} = \gamma_n \hat{a}_{ni} + \beta_n$$

```
torch.nn.LayerNorm
```

```
layer_norm = nn.LayerNorm(normalized_shape)
```

Regularization

Dropout

Dropout is one of the most effective form of regularization that is widely used.

The core idea is to delete nodes from the network, including their connections, at random during training.

Dropout is applied to both hidden and input nodes, but not outputs. It is equivalent to setting the output of a dropped node to zero.

[torch.nn.Dropout](#)

Early stopping

For good generalization training should be stopped at the point of smallest error with respect to the validation set.

Training neural networks

- **Model and parameters**
 - Single layer neural networks (Linear Regression, Logistic Regression, Softmax Regression)
 - Multilayer neural networks, Activation functions (ReLU, pReLU, Sigmoid, Tanh, GeLU)
- **Loss functions**
 - Mean Squared Error loss, Binary Cross Entropy loss, Cross Entropy Loss
 - Concept of entropy and cross entropy
- **Gradient Descent**
 - Batch, Stochastic, Mini-batch Stochastic Gradient Descent
- **Optimizers**
 - Momentum and Adaptive learning rates (Adagrad, RMSProp, Adam)
 - Learning rate schedule, Parameter initialization
- **Normalization**
 - Batch and layer normalization.
- **Regularization**
 - Dropout, Early stopping.
- **Training loop**
- **Backpropagation and Automatic differentiation**

Training loop

```
# Load the dataset.
train_dataset = SampleDataset(X_train, y_train)

# Preparing your data for training with DataLoaders.
batch_size = 64
train_dataloader = DataLoader(train_dataset, batch_size=batch_size, shuffle=True)

# Define the model class.
model = LogisticRegression(num_features=d)

# Loss function.
loss_fn = nn.BCELoss()

# Optimizer.
optimizer = SGD(model.parameters(), lr=0.01, momentum=0.9)

# Learning rate scheduler.
scheduler = ExponentialLR(optimizer, gamma=0.9)
```

```
# Run for a few epochs.
for epoch in range(n_epochs):
    # Iterate through the DataLoader to access mini-batches.
    for batch, (input, target) in enumerate(train_dataloader):
        # Prediction.
        output = model(input)

        # Compute loss.
        loss = loss_fn(output, target)

        # Compute gradient.
        loss.backward()

        # Gradient descent.
        optimizer.step()

        # Prevent gradient accumulation
        optimizer.zero_grad()

    # Adjust learning rate
    scheduler.step()
```

Training neural networks

- **Model and parameters**
 - Single layer neural networks (Linear Regression, Logistic Regression, Softmax Regression)
 - Multilayer neural networks, Activation functions (ReLU, pReLU, Sigmoid, Tanh, GeLU)
- **Loss functions**
 - Mean Squared Error loss, Binary Cross Entropy loss, Cross Entropy Loss
 - Concept of entropy and cross entropy
- **Gradient Descent**
 - Batch, Stochastic, Mini-batch Stochastic Gradient Descent
- **Optimizers**
 - Momentum and Adaptive learning rates (Adagrad,RMSProp,Adam)
 - Learning rate schedule, Parameter initialization
- **Normalization**
 - Batch and layer normalization.
- **Regularization**
 - Dropout, Early stopping.
- **Training loop**
- **Backpropagation and Automatic differentiation**

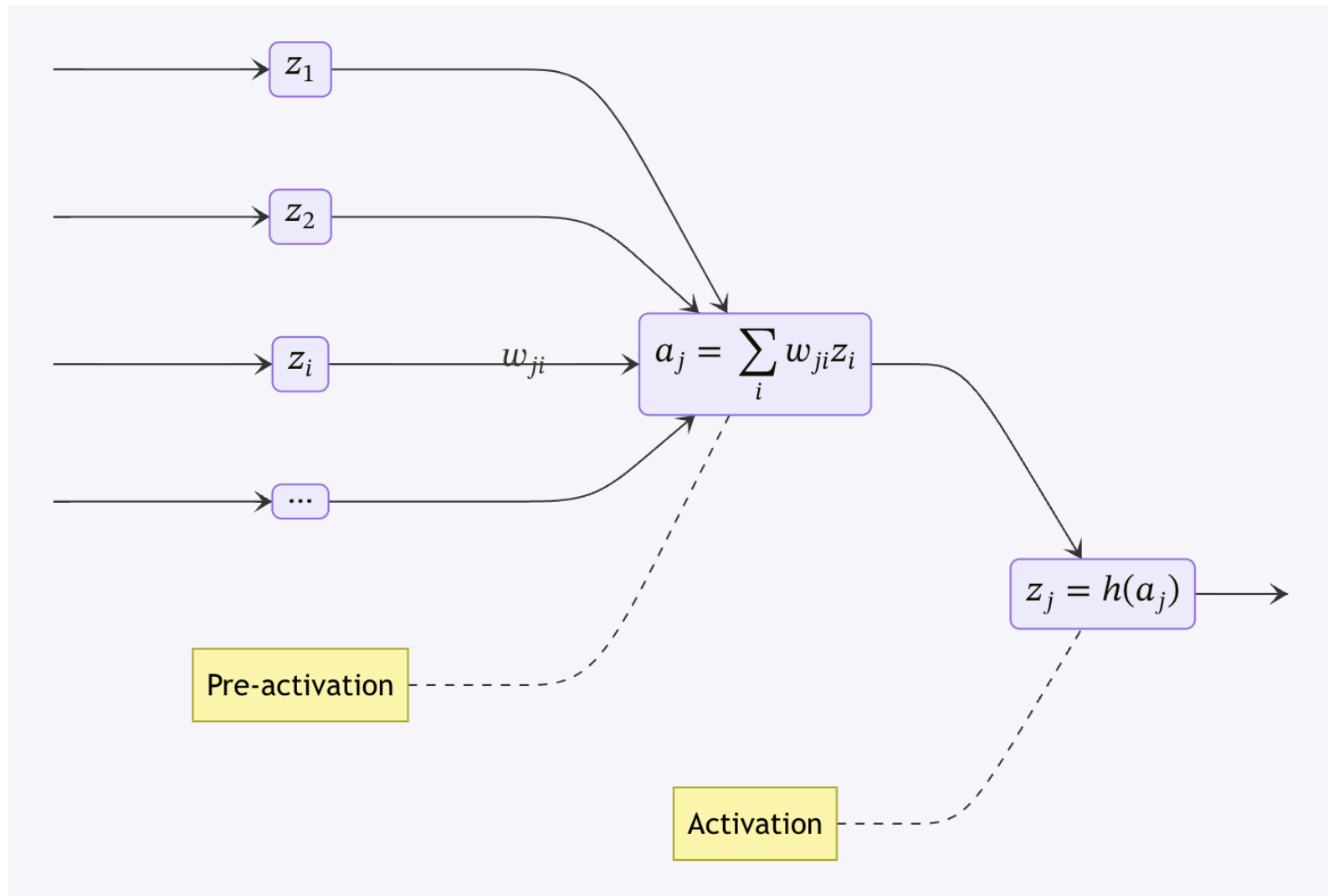
Backpropagation

Backprop, Error Backpropagation.

Backpropagation (or backprop) is an efficient technique to compute the gradient of the loss function.

It boils down to a **local message passing scheme** in which information is sent backwards through the network.

Forward propagation



Forward propagation

Let's consider a hidden unit in a general feed forward neural network.

Pre-activation

$$a_j = \sum_i w_{ji} z_i$$

Activation

$$z_j = h(a_j)$$

This process is called **forward propagation** since it is the forward flow of information through the network.

Gradients via chain rule

To compute the gradient of the loss function we use the chain rule.

$$\frac{\partial L_n}{\partial w_{ji}} = \frac{\partial L_n}{\partial a_j} \frac{\partial a_j}{\partial w_{ji}} = \delta_j z_i$$

where

$$\begin{aligned}\frac{\partial L_n}{\partial a_j} &= \delta_j \\ \frac{\partial a_j}{\partial w_{ji}} &= z_i\end{aligned}$$

δ_j are referred to as **errors**.

δ for the output units are based on the loss function.

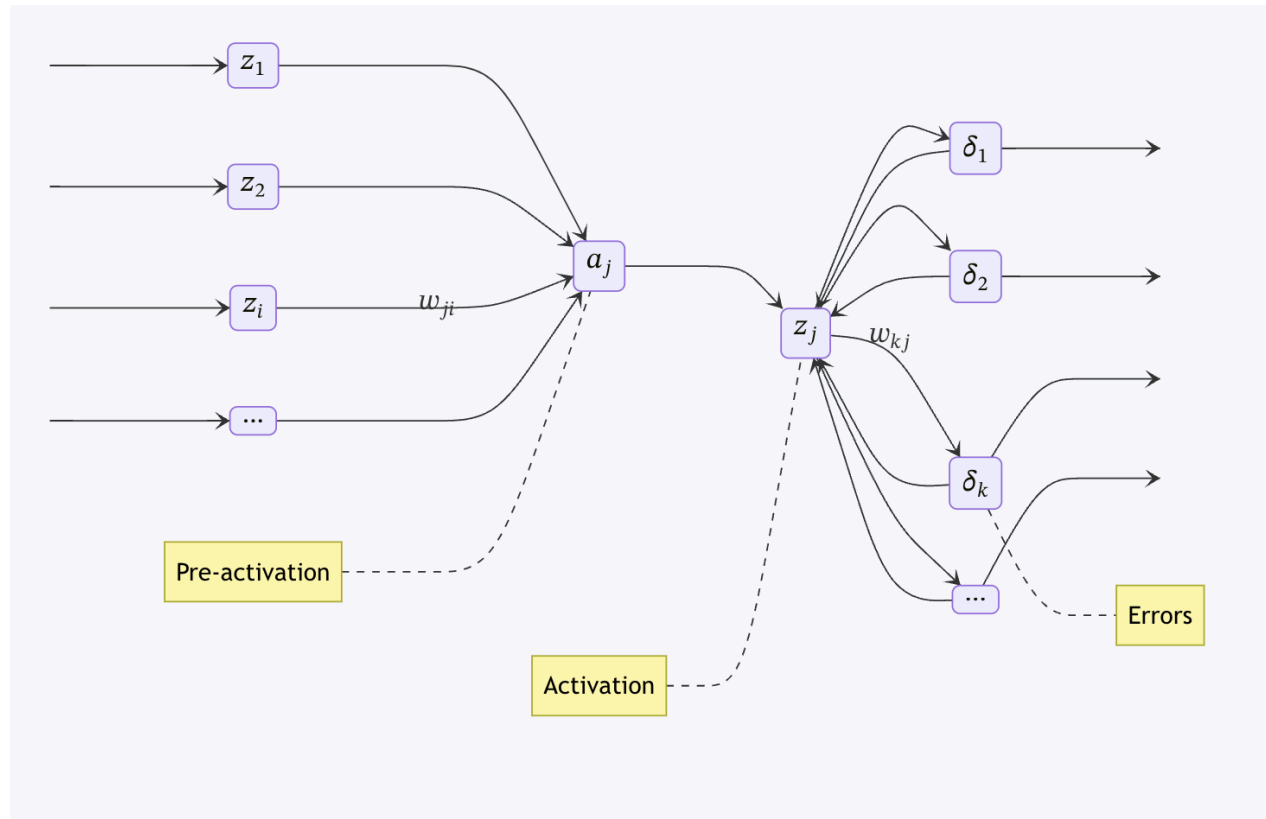
To evaluate the δ for the hidden units we again make use of the chain rule.

$$\delta_j := \frac{\partial L_n}{\partial a_j} = \sum_k \frac{\partial L_n}{\partial a_k} \frac{\partial a_k}{\partial a_j}$$

where the sum runs over all the units k to which j sends connections.

$$\delta_j = h'(a_j) \sum_k w_{kj} \delta_k$$

This tells us that the value of δ for a particular hidden unit can be obtained by propagating the δ backward from units higher up in the network.



Forward propagation

For all hidden and output units
compute in **forward order**

$$a_j \leftarrow \sum_i w_{ji} z_i$$

$$z_j \leftarrow h(a_j)$$

Error evaluation

For all output units compute

$$\delta_k \leftarrow \frac{\partial L_n}{\partial a_k}$$

Backward propagation

For all hidden units compute in
reverse order

$$\delta_j \leftarrow h'(a_j) \sum_k w_{kj} \delta_k$$

$$\frac{\partial L_n}{\partial w_{ji}} \leftarrow \delta_j z_i$$

Automatic differentiation

Algorithmic differentiation, autodiff, autograd

There are broadly 4 approaches to compute derivatives.

Atılım Günes Baydin, Barak A. Pearlmutter, Alexey Andreyevich Radul, and Jeffrey Mark Siskind. 2017. [Automatic differentiation in machine learning: a survey](#). J. Mach. Learn. Res. 18, 1 (January 2017), 5595–5637.

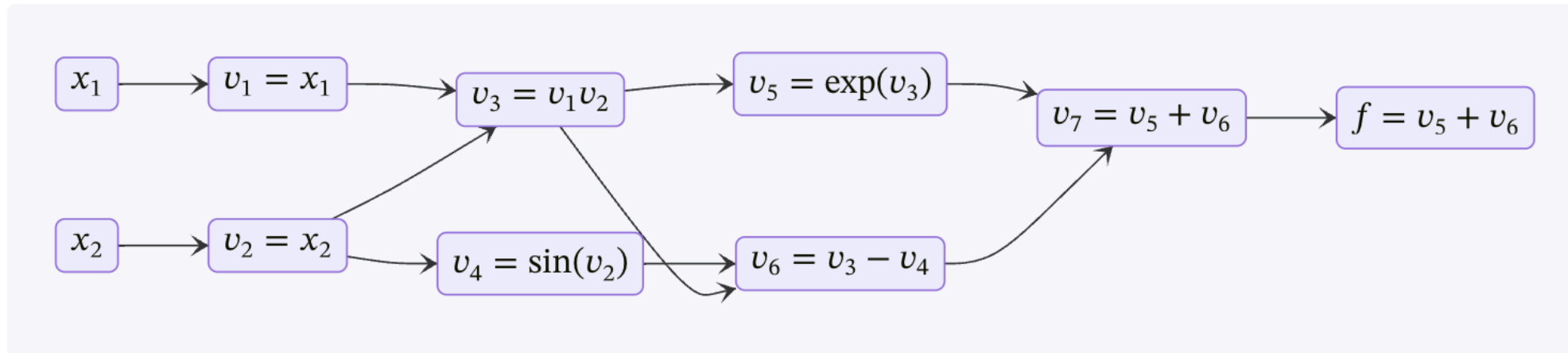
Approach	Pros	Cons
Manual derivation of backprop equations.	If done carefully can result in efficient code.	Manual process, prone to errors and not easy to iterate on models
Numerical evaluation of gradients via finite differences.	Sometimes used to check for correctness of other methods.	Limited by computational accuracy. Scales poorly with the size of the network.
Symbolic differentiation using packages like <code>sympy</code>		Closed form needed. Resulting expression can be very long (<i>expression swell</i>).
Automatic differentiation	Most preferred.	

Forward-mode automatic differentiation

Consider the following function.

$$f(x_1, x_2) = x_1x_2 + \exp(x_1x_2) - \sin(x_2)$$

When implemented in software the code consists of a sequence of operations than can be expressed as an **evaluation trace** of the underlying elementary operations. This trace can be visualized as a computation graph with respect to the following 7 **primal variables**.



We first write code to implement the evaluation of the primal variables.

$$v_1 = x_1$$

$$v_2 = x_2$$

$$v_3 = v_1 v_2$$

$$v_4 = \sin(v_2)$$

$$v_5 = \exp(v_3)$$

$$v_6 = v_3 - v_4$$

$$v_7 = v_5 + v_6$$

Primal and tangent variables

We augment each intermediate variable z_i (known as **primal** variable) with an additional variable representing the value of some derivative of that variable, which we denote as \dot{z}_i , known as **tangent** variable.

The tangent variables are generated automatically.

Not say we wish to evaluate the derivative $\partial f / \partial x_1$. First we define the tangent variables by

$$\dot{v}_i = \frac{\partial v_i}{\partial x_1}$$

Expressions for evaluating these can be constructed automatically using the chain rule of calculus.

$$\dot{v}_i = \frac{\partial v_i}{\partial x_1} = \sum_{j \in \text{parents}(i)} \frac{\partial v_i}{\partial v_j} \frac{\partial v_j}{\partial x_1} = \sum_{j \in \text{parents}(i)} \dot{v}_j \frac{\partial v_i}{\partial v_j}$$

where $\text{parents}(i)$ denotes the set of **parents** of node i in the evaluation trace diagram.

The associated equations and corresponding code for evaluating the tangent variables are generated automatically.

$$\dot{v}_1 = 1$$

$$\dot{v}_2 = 0$$

$$\dot{v}_3 = v_1 \dot{v}_2 + \dot{v}_1 v_2$$

$$\dot{v}_4 = \dot{v}_2 \cos(v_2)$$

$$\dot{v}_5 = \dot{v}_3 \exp(v_3)$$

$$\dot{v}_6 = \dot{v}_3 - \dot{v}_4$$

$$\dot{v}_7 = \dot{v}_5 + \dot{v}_6$$

To evaluate the derivative $\frac{\partial f}{\partial x_1}$ we input specific values of x_1 and x_2 and the code then executes the primal and tangent equations, numerically evaluating the tuples (v_i, \dot{v}_i) in **forward** order until we obtain the required derivative.

- The forward mode with slight modifications can handle multiple outputs in the same pass but the process has to be repeated for every parameter that we need the derivative.
- Since we are often in the regime of one output with millions of parameters this is not scalable for modern deep neural networks.
- We therefore turn to an alternative version based on the backwards flow of derivative data through the evaluation trace graph.

Reverse-mode automatic differentiation

Reverse-mode automatic differentiation is a generalization of the error backpropagation procedure we discussed earlier.

Primal and adjoint variables

As with forward mode, we augment each primal variable v_i with an additional variable called **adjoint** variable, denoted as \bar{v}_i .

$$\bar{v}_i = \frac{\partial f}{\partial v_i}$$

Expressions for evaluating these can be constructed automatically using the chain rule of calculus.

$$\bar{v}_i = \frac{\partial f}{\partial v_i} = \sum_{j \in \text{children}(i)} \frac{\partial f}{\partial v_j} \frac{\partial v_j}{\partial v_i} = \sum_{j \in \text{children}(i)} \bar{v}_j \frac{\partial v_j}{\partial v_i}$$

where $\text{children}(i)$ denotes the set of **children** of node i in the evaluation trace diagram.

The successive evaluation of the adjoint variables represents a flow of information backwards through the graph. For multiple parameters a single backward pass is enough.

$$\bar{v}_7 = 1$$

$$\bar{v}_6 = \bar{v}_7$$

$$\bar{v}_5 = \bar{v}_7$$

$$\bar{v}_4 = -\bar{v}_6$$

$$\bar{v}_3 = \bar{v}_5 v_5 + \bar{v}_6$$

$$\bar{v}_2 = \bar{v}_2 v_1 + \bar{v}_4 \cos(v_2)$$

$$\bar{v}_1 = \bar{v}_3 v_2$$

Autograd in pytorch

- A Gentle Introduction to `torch.autograd`
- The Fundamentals of Autograd

Training neural networks

- **Model and parameters**
 - Single layer neural networks (Linear Regression, Logistic Regression, Softmax Regression)
 - Multilayer neural networks, Activation functions (ReLU, pReLU, Sigmoid, Tanh, GeLU)
- **Loss functions**
 - Mean Squared Error loss, Binary Cross Entropy loss, Cross Entropy Loss
 - Concept of entropy and cross entropy
- **Gradient Descent**
 - Batch, Stochastic, Mini-batch Stochastic Gradient Descent
- **Optimizers**
 - Momentum and Adaptive learning rates (Adagrad,RMSProp,Adam)
 - Learning rate schedule, Parameter initialization
- **Normalization**
 - Batch and layer normalization.
- **Regularization**
 - Dropout, Early stopping.
- **Training loop**
- **Backpropagation and Automatic differentiation**

Short quiz

<https://vikasraykar.github.io/deeplearning/docs/training/quiz/>

Derive the gradient of the loss function for linear regression and logistic regression.

What is cross entropy ?

What is the most widely used optimizer ?

What are the typically used parameters of the optimizer ?

For SGD with momentum show that it increases the effective learning rate from η to $\frac{\eta}{(1-\mu)}$.

In **Attention Is All You Need** paper what is the optimizer and the learning rate scheduler used ?

What are the optimizers used in LLaMA and GPT-3 ?

What is the disadvantage of forward-mode automatic differentiation ?

What is the difference between batch and layer normalization ?

Why do we do `optimizer.zero_grad()` ?

Coding assignment

<https://vikasraykar.github.io/deeplearning/docs/training/coding/>

Thanks you and any questions ?

notes - <https://vikasraykar.github.io/deeplearning>

code - <https://github.com/vikasraykar/deeplearning-doj>