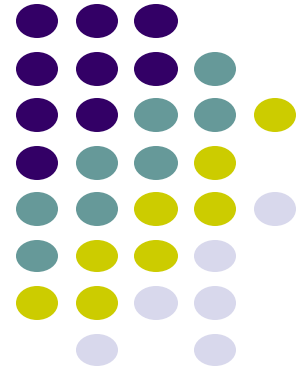


ErSE222: Machine learning in Geoscience

Feb. 16th, 2025

Tariq Alkhalifah and Omar Saad





Backpropagation (BCE)

General algorithm, called **backpropagation** → used to compute gradients of any complex function that be composed of simple, differentiable functions

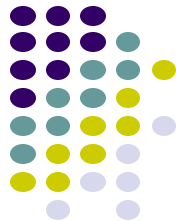
$$z = \mathbf{x}^T \boldsymbol{\theta}, \quad a = \sigma(z), \quad \mathcal{L} = -(y \log(a) + (1 - y) \log(1 - a))$$



Chain-
rule

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}} = \frac{\partial \mathcal{L}}{\partial a} \frac{\partial a}{\partial z} \frac{\partial z}{\partial \boldsymbol{\theta}}$$

Backpropagation (BCE)



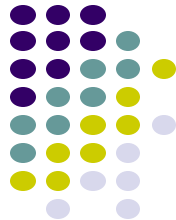
Let's now (for once) do the math:

$$\mathcal{L} = -(y \log(a) + (1 - y) \log(1 - a))$$

$$\frac{\partial \mathcal{L}}{\partial a} = -\frac{y}{a} + \frac{1 - y}{1 - a} = \frac{-y(1 - a) + (1 - y)a}{a(1 - a)}$$

$$\frac{\partial a}{\partial z} = a(1 - a) \quad \boxed{a = \sigma(z)}$$

Backpropagation (BCE)



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

$$u = 1 + e^{-x}$$

$$y = \frac{1}{u}$$

$$\frac{dy}{du} = -\frac{1}{u^2}$$

$$\frac{du}{dx} = -e^{-x}$$

Backpropagation (BCE)

$$\frac{dy}{dx} = \frac{dy}{du} \cdot \frac{du}{dx}$$

$$\frac{dy}{du} = -\frac{1}{u^2} \quad \frac{du}{dx} = -e^{-x}$$

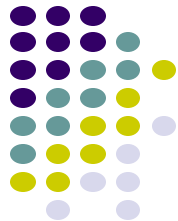
$$\frac{dy}{dx} = \frac{e^{-x}}{u^2}$$

$$\frac{dy}{dx} = \frac{e^{-x}}{(1+e^{-x})^2}$$

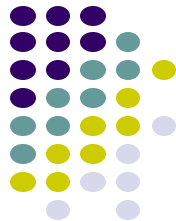
$$\sigma(x) = \frac{1}{1+e^{-x}}$$

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

$$\sigma'(x) = \sigma(x) \cdot (1-\sigma(x))$$



Backpropagation (BCE)



Let's now (for once) do the math:

$$\mathcal{L} = -(y \log(a) + (1 - y) \log(1 - a))$$

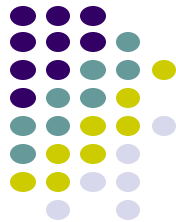
$$\frac{\partial \mathcal{L}}{\partial a} = -\frac{y}{a} + \frac{1 - y}{1 - a} = \frac{-y(1 - a) + (1 - y)a}{a(1 - a)}$$

$$\frac{\partial a}{\partial z} = a(1 - a) \quad \boxed{a = \sigma(z)}$$



$$\frac{\partial \mathcal{L}}{\partial z} = \frac{\partial \mathcal{L}}{\partial a} \frac{\partial a}{\partial z} = -y(1 - a) + (1 - y)a = a - y = dz$$

Backpropagation (BCE)



$$z = \mathbf{x}^T \boldsymbol{\theta}, \quad a = \sigma(z), \quad \mathcal{L} = -(y \log(a) + (1 - y) \log(1 - a))$$

Let's now (for once) do the math:

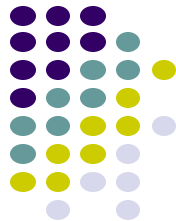
$$\frac{\partial \mathcal{L}}{\partial z} = \frac{\partial \mathcal{L}}{\partial a} \frac{\partial a}{\partial z} = -y(1 - a) + (1 - y)a = a - y = dz$$

$$\frac{\partial z}{\partial w_i} = x_i, \quad \frac{\partial z}{\partial b} = 1$$



$$\frac{\partial \mathcal{L}}{\partial w_i} = dz \cdot x_i = dw_i, \quad \frac{\partial \mathcal{L}}{\partial b} = dz = db$$

Backpropagation (BCE)



$$\mathbf{z} = \mathbf{X}_{train}^T \boldsymbol{\theta}$$

$$\mathbf{a} = \sigma(\mathbf{z})$$

} Forward

$$\mathbf{dz} = \mathbf{a} - \mathbf{y}$$

$$\mathbf{dw} = \frac{1}{N_s} \mathbf{X}_{train} \mathbf{dz}$$

$$db = \frac{1}{N_s} \mathbf{1}^T \mathbf{dz}$$

} Backward

$$\mathbf{w} \leftarrow \mathbf{w} - \alpha \mathbf{dw}$$

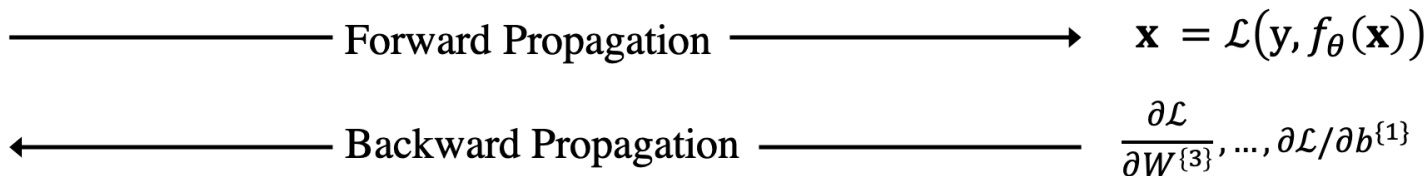
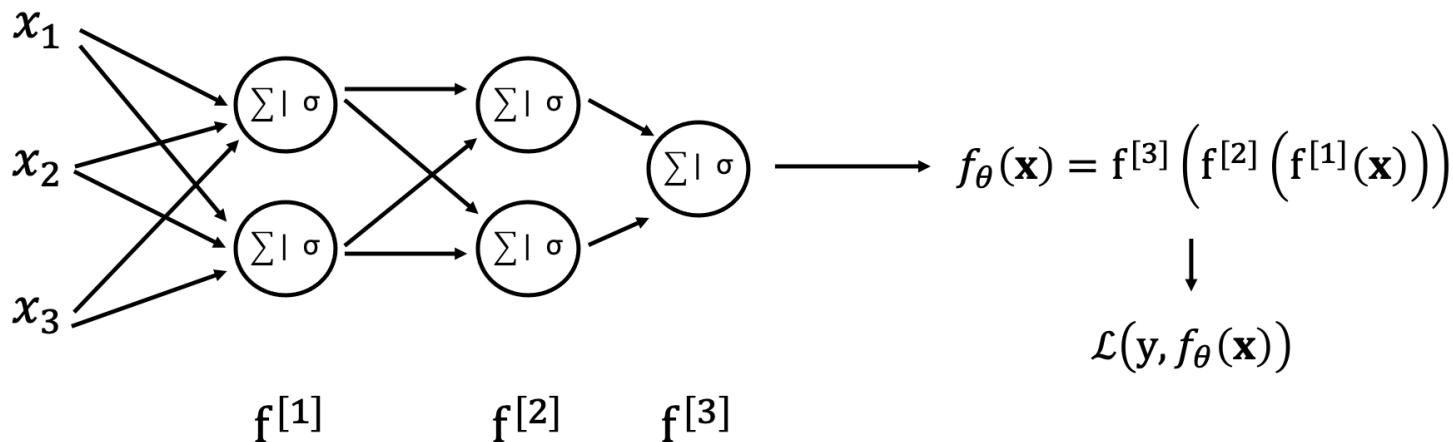
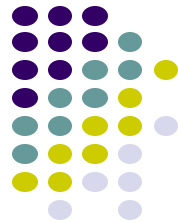
$$b \leftarrow b - \alpha db$$

} Update

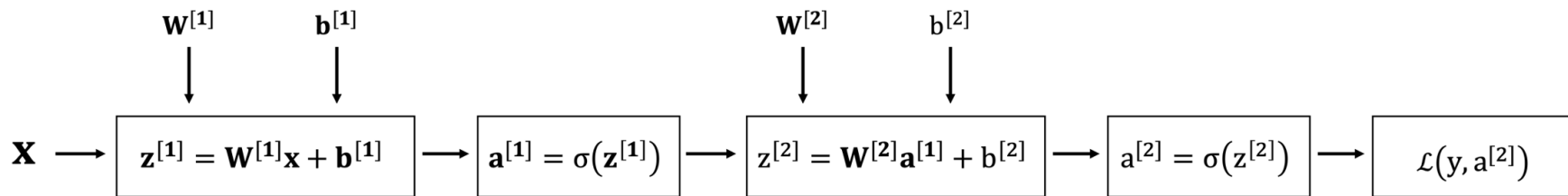
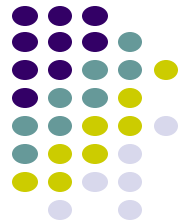


MLP

Backpropagation (BCE)



Backpropagation (BCE)



And write in reversed order (evaluated right to left):

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{[1]}} = \frac{\partial \mathbf{z}^{[1]}}{\partial \mathbf{W}^{[1]}} \frac{\partial \mathbf{a}^{[1]}}{\partial \mathbf{z}^{[1]}} \frac{\partial \mathbf{z}^{[2]}}{\partial \mathbf{a}^{[1]}} \frac{\partial \mathbf{a}^{[2]}}{\partial \mathbf{z}^{[2]}} \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[2]}}$$



Effective order of operations

Backpropagation (BCE)

$$\mathbf{z} = \mathbf{X}_{train}^T \boldsymbol{\theta}$$

$$\mathbf{a} = \sigma(\mathbf{z})$$

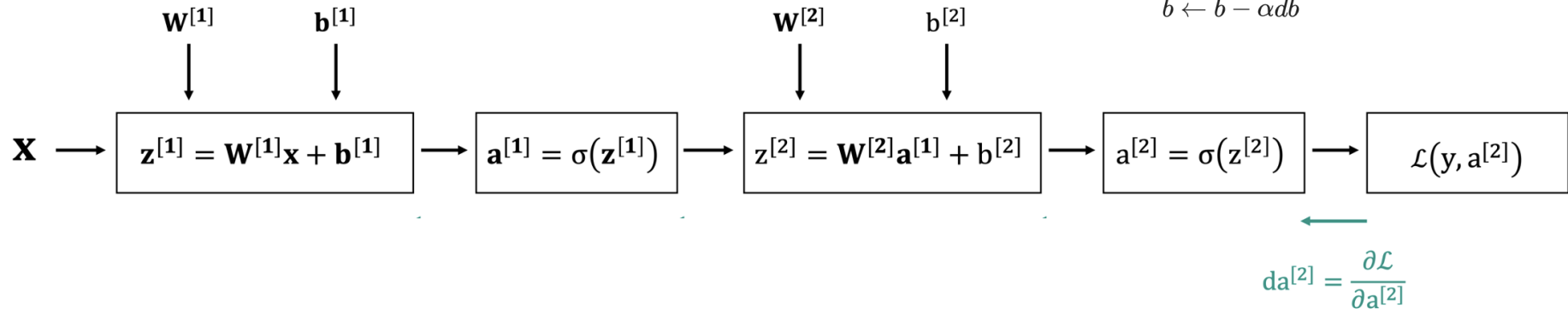
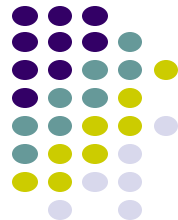
$$\mathbf{dz} = \mathbf{a} - \mathbf{y}$$

$$d\mathbf{w} = \frac{1}{N_s} \mathbf{X}_{train} d\mathbf{z}$$

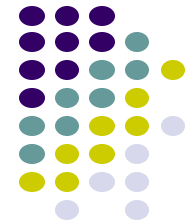
$$db = \frac{1}{N_s} \mathbf{1}^T d\mathbf{z}$$

$$\mathbf{w} \leftarrow \mathbf{w} - \alpha d\mathbf{w}$$

$$b \leftarrow b - \alpha db$$



Backpropagation (BCE)



$$\mathbf{z} = \mathbf{X}_{train}^T \boldsymbol{\theta}$$

$$\mathbf{a} = \sigma(\mathbf{z})$$

$$d\mathbf{z} = \mathbf{a} - \mathbf{y}$$

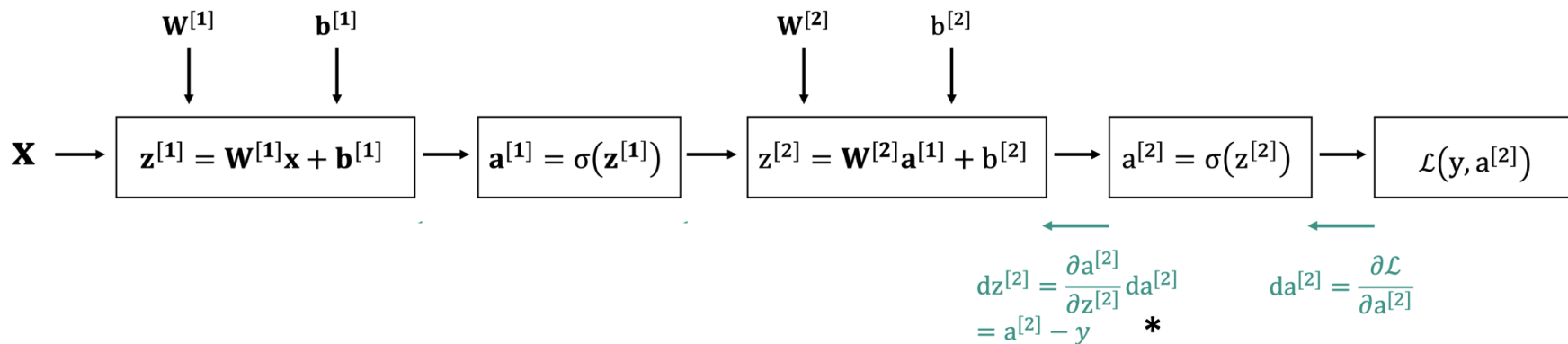
$$d\mathbf{w} = \frac{1}{N_s} \mathbf{X}_{train} d\mathbf{z}$$

$$db = \frac{1}{N_s} \mathbf{1}^T d\mathbf{z}$$

$$\mathbf{w} \leftarrow \mathbf{w} - \alpha d\mathbf{w}$$

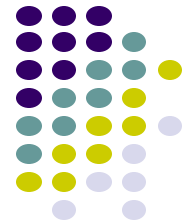
$$b \leftarrow b - \alpha db$$

$$\frac{\partial \mathcal{L}}{\partial z} = \frac{\partial \mathcal{L}}{\partial a} \frac{\partial a}{\partial z} = -y(1-a) + (1-y)a = a - y = dz$$



* Assuming \mathcal{L} =BCE

Backpropagation (BCE)



$$\mathbf{z} = \mathbf{X}_{train}^T \boldsymbol{\theta}$$

$$\mathbf{a} = \sigma(\mathbf{z})$$

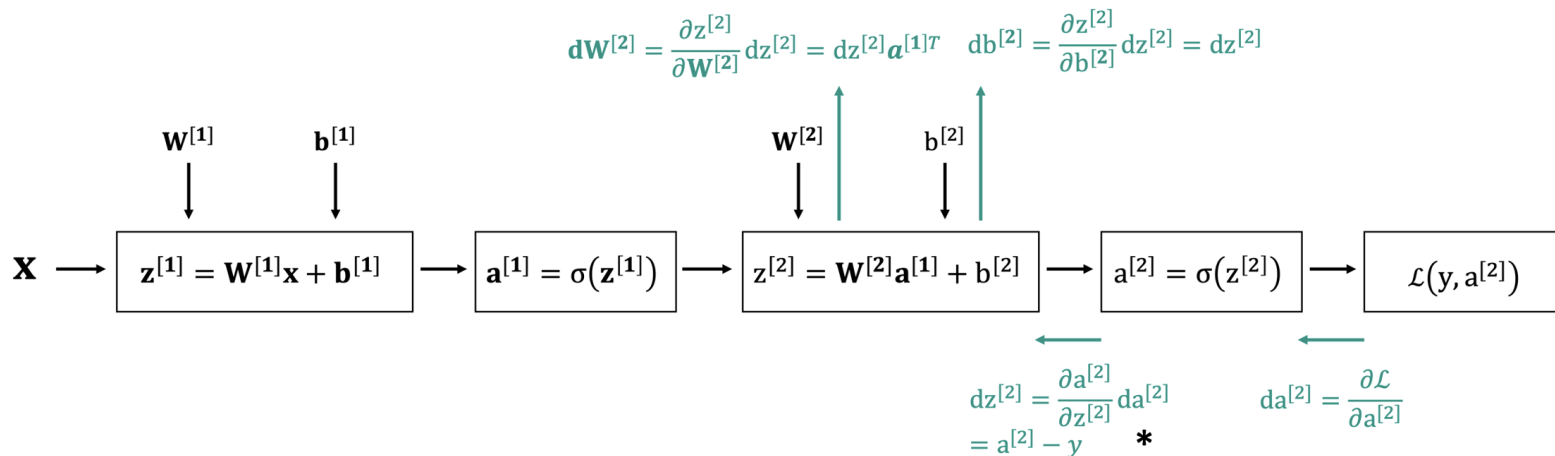
$$d\mathbf{z} = \mathbf{a} - \mathbf{y}$$

$$d\mathbf{w} = \frac{1}{N_s} \mathbf{X}_{train} d\mathbf{z}$$

$$db = \frac{1}{N_s} \mathbf{1}^T d\mathbf{z}$$

$$\mathbf{w} \leftarrow \mathbf{w} - \alpha d\mathbf{w}$$

$$\frac{\partial \mathcal{L}}{\partial w_i} = dz \cdot x_i = dw_i, \quad \frac{\partial \mathcal{L}}{\partial b} = dz = db$$



* Assuming $\mathcal{L}=\text{BCE}$

Backpropagation (BCE)

$$\frac{\partial \mathcal{L}}{\partial w_i} = dz \cdot x_i = dw_i, \quad \frac{\partial \mathcal{L}}{\partial b} = dz = db$$

$$\mathbf{z} = \mathbf{X}_{train}^T \boldsymbol{\theta}$$

$$\mathbf{a} = \sigma(\mathbf{z})$$

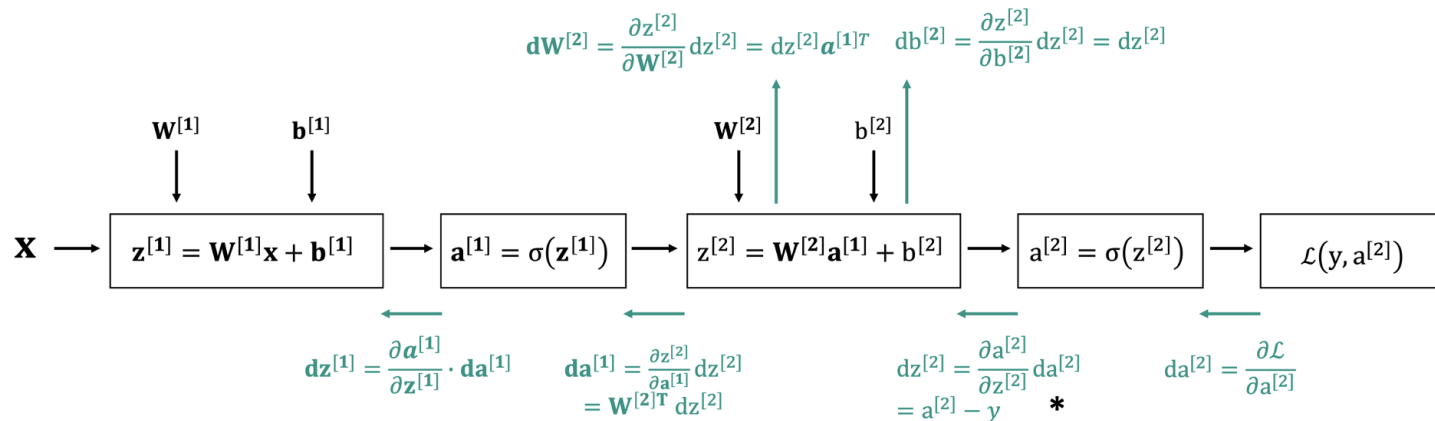
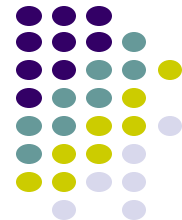
$$\mathbf{dz} = \mathbf{a} - \mathbf{y}$$

$$d\mathbf{w} = \frac{1}{N_s} \mathbf{X}_{train} d\mathbf{z}$$

$$db = \frac{1}{N_s} \mathbf{1}^T d\mathbf{z}$$

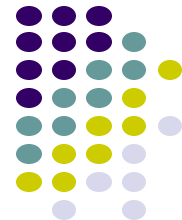
$$\mathbf{w} \leftarrow \mathbf{w} - \alpha d\mathbf{w}$$

$$b \leftarrow b - \alpha db$$



* Assuming \mathcal{L} =BCE

Backpropagation (BCE)



$$\mathbf{z} = \mathbf{X}_{train}^T \boldsymbol{\theta}$$

$$\mathbf{a} = \sigma(\mathbf{z})$$

$$d\mathbf{z} = \mathbf{a} - \mathbf{y}$$

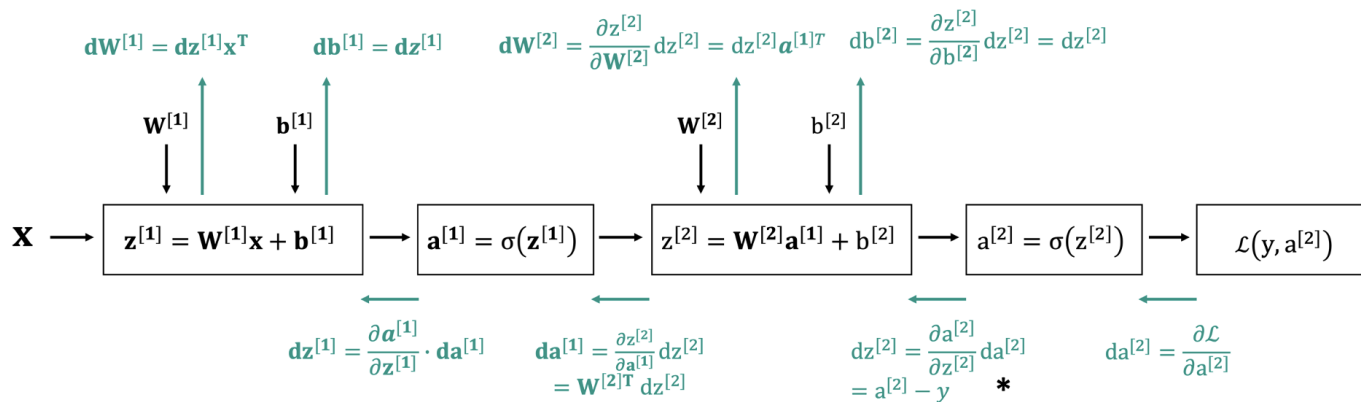
$$d\mathbf{w} = \frac{1}{N_s} \mathbf{X}_{train} d\mathbf{z}$$

$$db = \frac{1}{N_s} \mathbf{1}^T d\mathbf{z}$$

$$\mathbf{w} \leftarrow \mathbf{w} - \alpha d\mathbf{w}$$

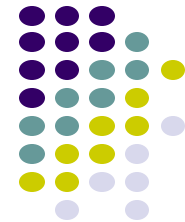
$$b \leftarrow b - \alpha db$$

$$\frac{\partial \mathcal{L}}{\partial w_i} = dz \cdot x_i = dw_i, \quad \frac{\partial \mathcal{L}}{\partial b} = dz = db$$

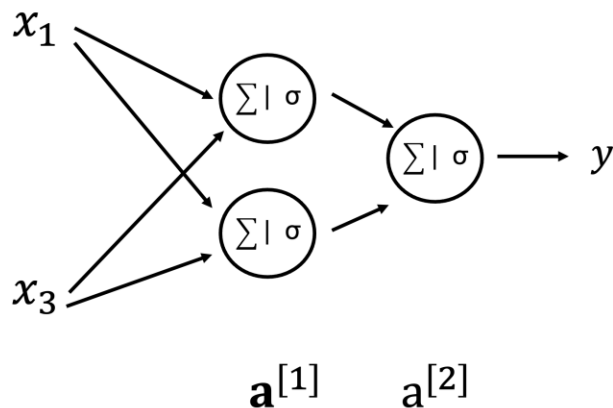


* Assuming $\mathcal{L}=\text{BCE}$

Initialization



- **Zero (or constant) initialization:** symmetry problem (aka symmetric gradients)



Initialization:

$$\{\mathbf{W}^{[1]} = c^{[1]}\mathbf{1}, \mathbf{b}^{[1]} = c^{[1]}\mathbf{1}, \mathbf{W}^{[2]} = c^{[2]}\mathbf{1}, \mathbf{b}^{[2]} = c^{[2]}\mathbf{1}\}$$

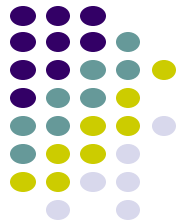


$$a_1^{[1]} = a_2^{[1]} \longrightarrow dz_1^{[1]} = dz_2^{[1]}$$

$$\longrightarrow dW_{11}^{[1]} = dW_{21}^{[1]}, dW_{12}^{[1]} = dW_{22}^{[1]}$$

$$\longrightarrow W_{11}^{[1]} = W_{21}^{[1]}, W_{12}^{[1]} = W_{22}^{[1]}$$

Initialization



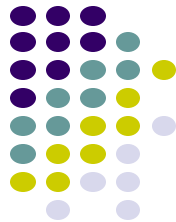
- **Random initialization:** general idea

$$w_{i,j}^{[.]} \sim \mathcal{N}(0, 0.01)$$

Not too small: to avoid slow training
Not too high: to avoid unstable training

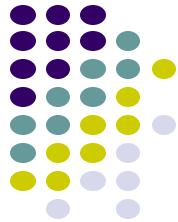
$$b_i^{[.]} = 0 \quad (\text{or constant - e.g. for ReLU})$$

Initialization



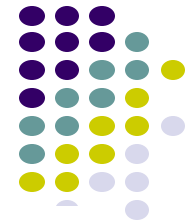
- **Random initialization:** more appropriate choices (depend on #units in layer)

- Uniform distributior $w_{ij}^{[k]} \sim \mathcal{U}(-1/\sqrt{N^{[k]}}, 1/\sqrt{N^{[k]}})$
- Xavier initialization $w_{ij}^{[k]} \sim \mathcal{N}(0, 1/N^{[k]})$
- He initialization $w_{ij}^{[k]} \sim \mathcal{N}(0, 2/N^{[k]})$

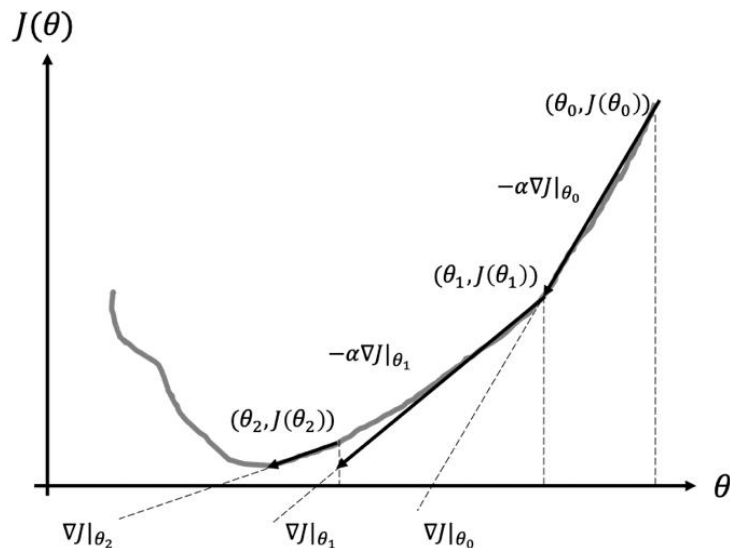


Optimization

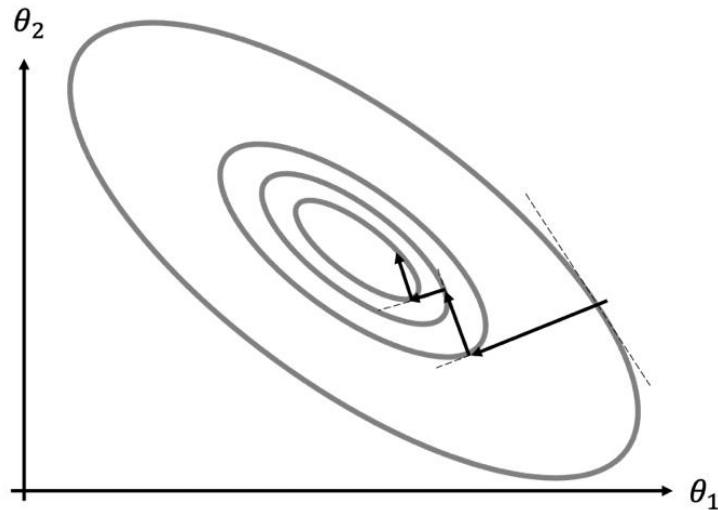
Gradient-based optimization



1d function



2d function



Use local information of $f(x)$ to find its minimum/maximum

Gradient-based optimization



Gradient descent algorithm:

Initialization: choose $\theta \in \mathbb{R}$

For $i = 0, \dots, N - 1$;

1. Compute update direction $d_i = -\nabla J|_{\theta_i}$
2. Estimate step-length α_i
3. Update $\theta_{i+1} = \theta_i + \alpha_i d_i$

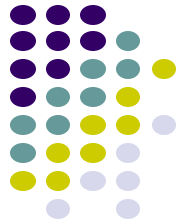
Choice of step-length:

- Constant or with some schedule (in DL)
- Exact
- Backtracking

Extension to n-dimensional problems:

In DL applications, $\theta = [\theta_1, \theta_2, \dots, \theta_N]$

$$\nabla J = [\delta J / \delta \theta_1, \delta J / \delta \theta_2, \dots, \delta J / \delta \theta_N].$$



Stochastic gradient descent (SGD)

Batched gradient descent:

$$\boldsymbol{\theta}_{i+1} = \boldsymbol{\theta}_i - \alpha_i \nabla J = \boldsymbol{\theta}_i - \frac{\alpha_i}{N_s} \sum_{j=1}^{N_s} \nabla \mathcal{L}_j$$

Stochastic gradient descent:

$$\boldsymbol{\theta}_{i+1} = \boldsymbol{\theta}_i - \alpha_i \nabla \mathcal{L}_j$$

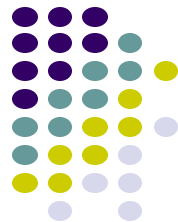
Mini-batch gradient descent:

$$\boldsymbol{\theta}_{i+1} = \boldsymbol{\theta}_i - \frac{\alpha_i}{N_b} \sum_{j=1}^{N_b} \nabla \mathcal{L}_j$$

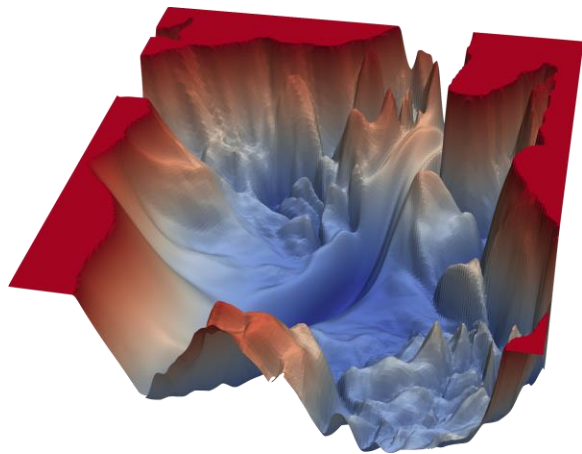
Mini-batch size



Limitations of SGD



- **Ill-conditioning:** In an ill-conditioned problem, the gradients may change dramatically with small changes in the parameter values. This results in steep gradients in certain directions and shallow gradients in others.
- **Local minima:** the landscape of NN functionals is generally non-convex and populated with a multitude of local minima

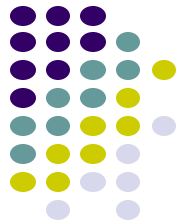


Model identifiability

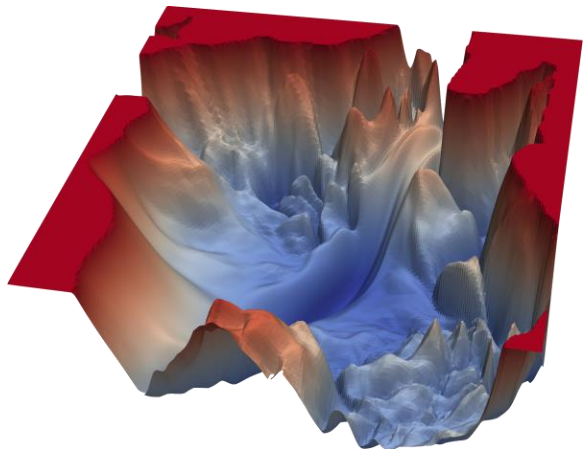
θ_{gm} Single set of params for optimal performance

$\theta_{gm}, \theta_{lm,1}, \dots$ Multiple models with similar performance

Limitations of SGD



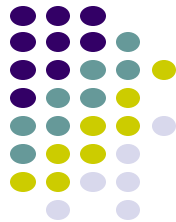
- **Saddle points (and other flat regions):** landscapes associated to the training of deep neural networks may have much fewer local minima than we tend to believe... the main hinder to slow training is represented by saddle points



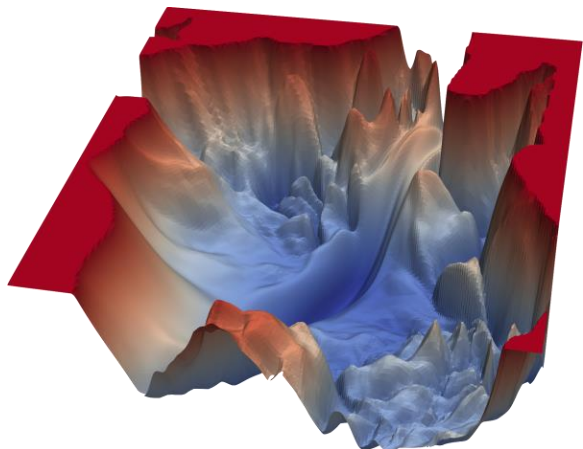
$$J(\theta) \rightarrow 0$$

Vanishing gradients

Limitations of SGD



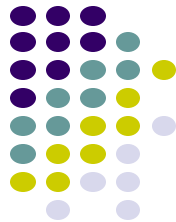
- **Cliffs:** steep regions in the NN landscape with high gradient



$J(\theta) \rightarrow \infty$ Exploding gradients

$\nabla J(\theta_i) = \min(\nabla J(\theta_i), th)$ Gradient clipping

Cooling



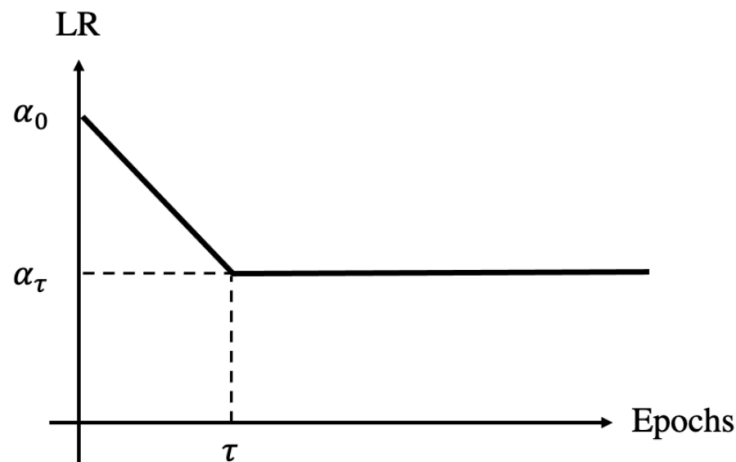
- Linearly decaying learning rate

$$\alpha_i = (1 - \beta)\alpha_0 + \beta\alpha_\tau \quad i < \tau$$

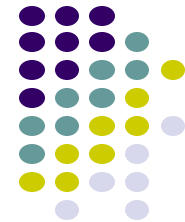
$$\alpha_i = \alpha_\tau \quad i \geq \tau$$

$$\beta = i/\tau$$

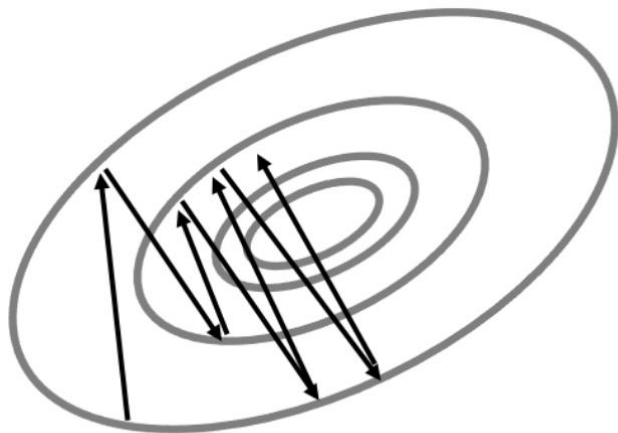
$$\tau \approx 100N_{epochs}, \alpha_\tau = \alpha_0/100$$



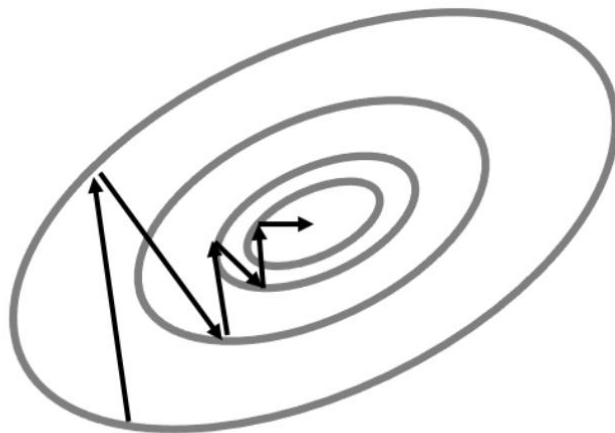
Cooling



Fixed learning rate



Cooling strategy



SGD with momentum



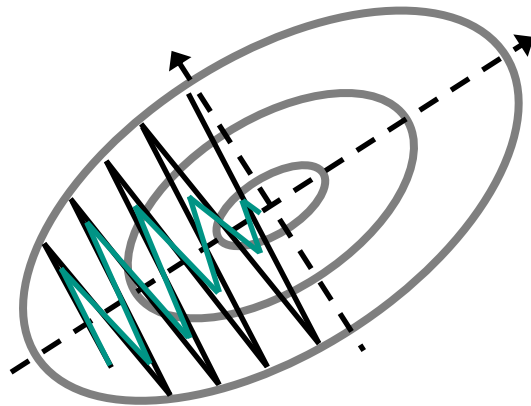
Idea of momentum (from Polyak and Nesterov in 60'): use update that is an exponentially decaying moving average of the past gradients.

Velocity Momentum term $[0,1)$

↓ ↙

$$\mathbf{v}_{i+1} = \gamma \mathbf{v}_i - \mathbf{g}_{i+1} = \gamma \mathbf{v}_i - \frac{\alpha}{N_b} \sum_{j=1}^{N_b} \nabla \mathcal{L}_j$$

$$\boldsymbol{\theta}_{i+1} = \boldsymbol{\theta}_i + \mathbf{v}_{i+1}$$



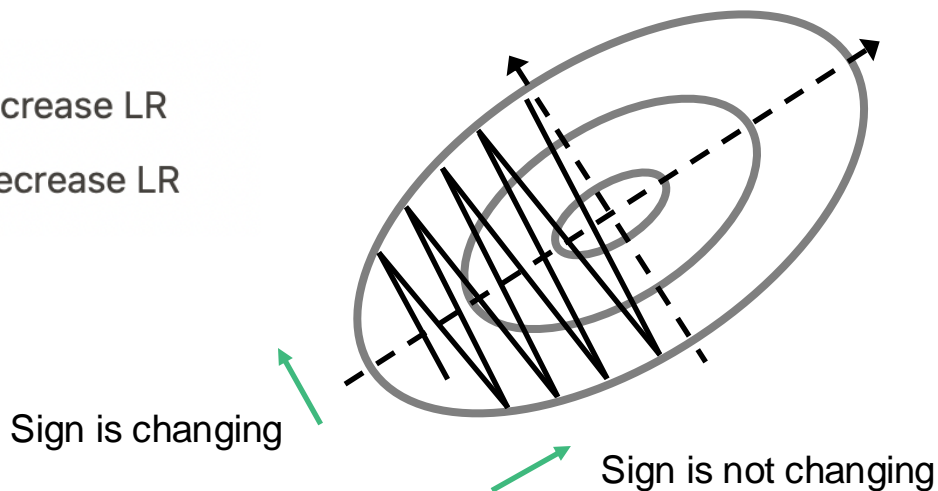
Adaptive Learning rates

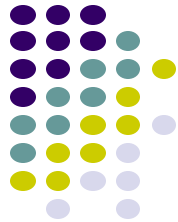


Idea of adaptive LR : instead of scaling the entire gradient, different learning rates are adaptively applied to different directions

Delta-Bar-Delta algorithm

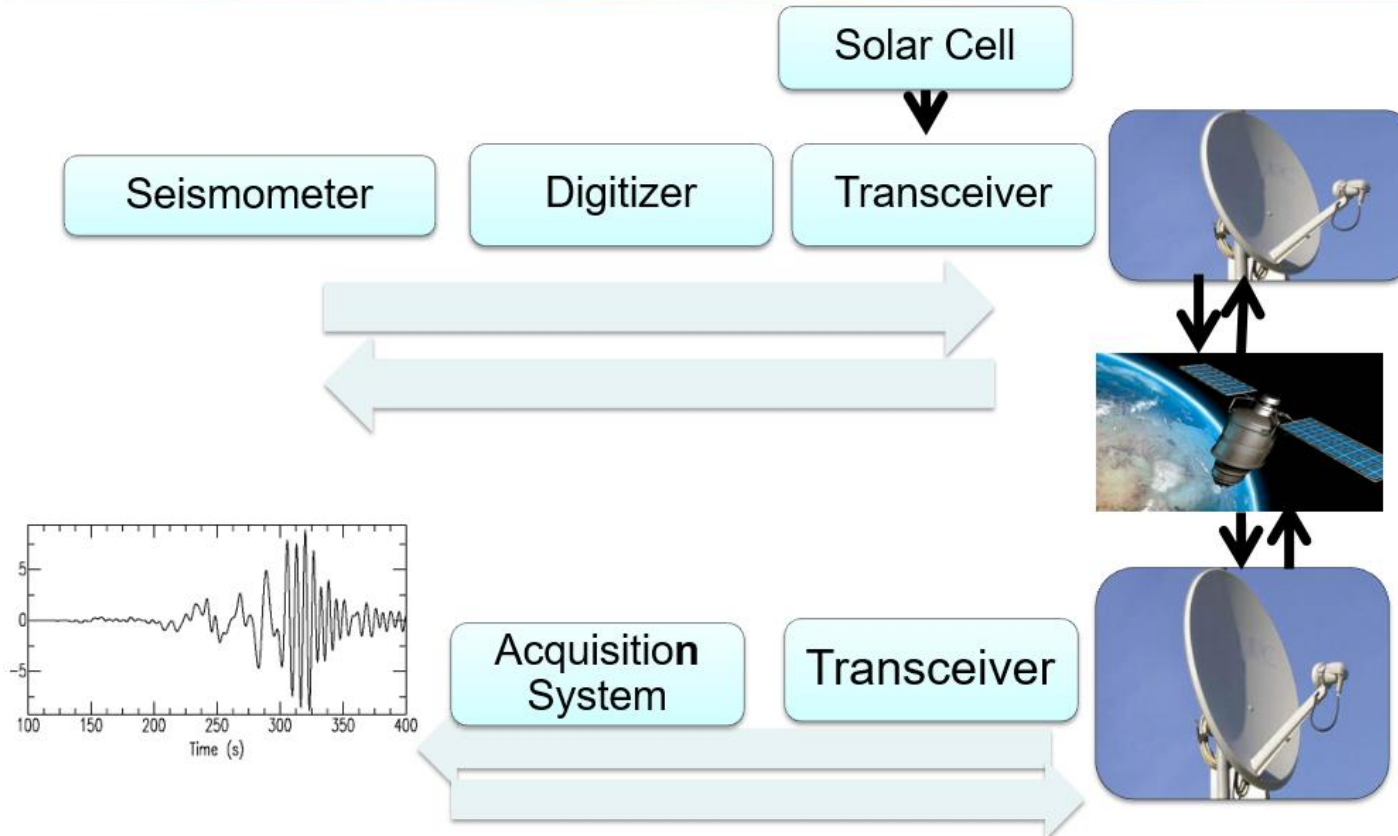
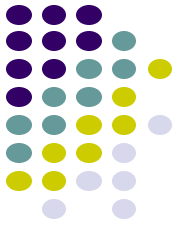
- if $\text{sign}\{g_{i+1}^j\} = \text{sign}\{g_i^j\}$, increase LR
- if $\text{sign}\{g_{i+1}^j\} \neq \text{sign}\{g_i^j\}$, decrease LR



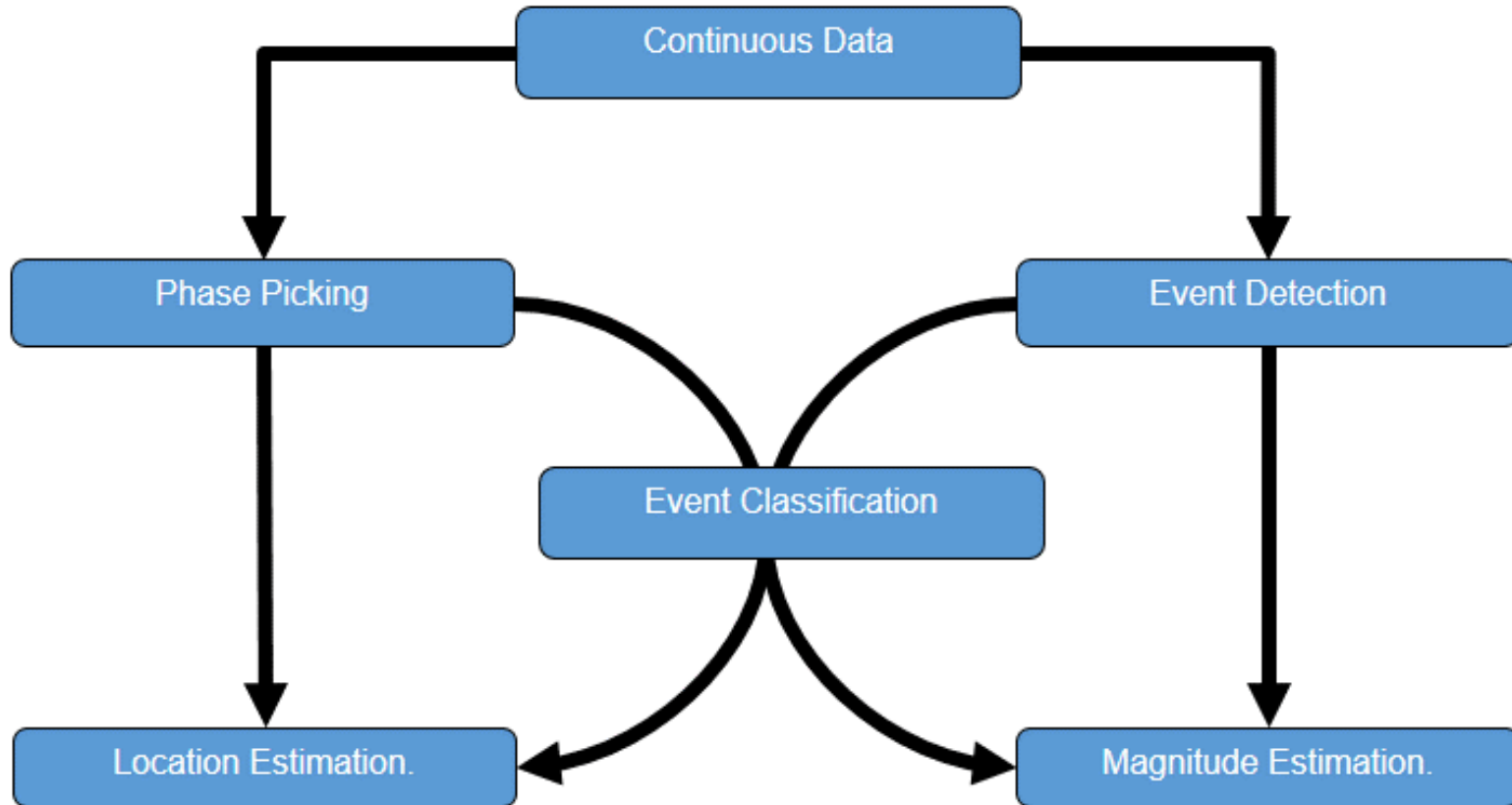
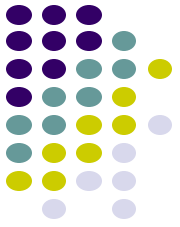


Application

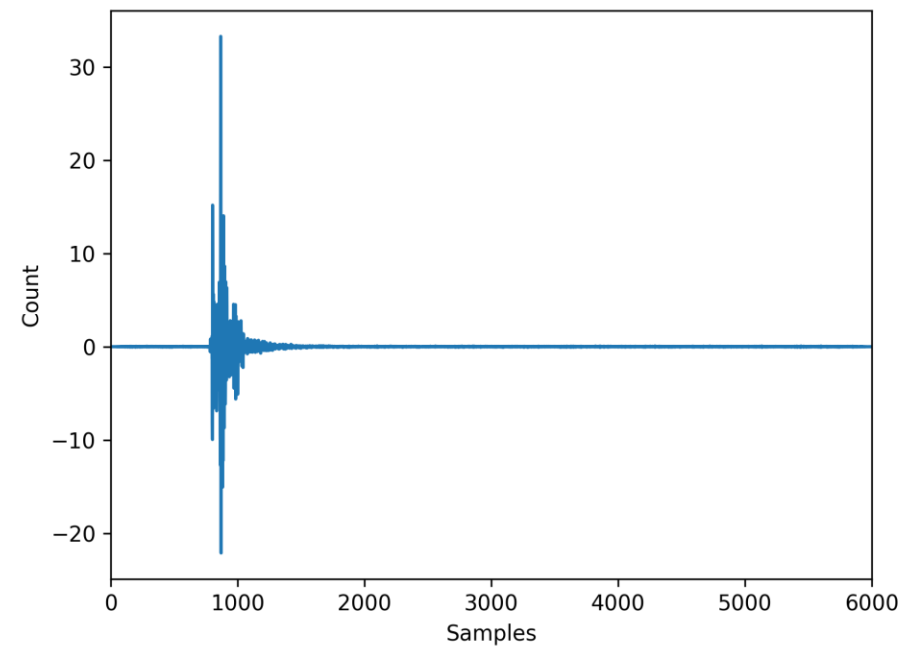
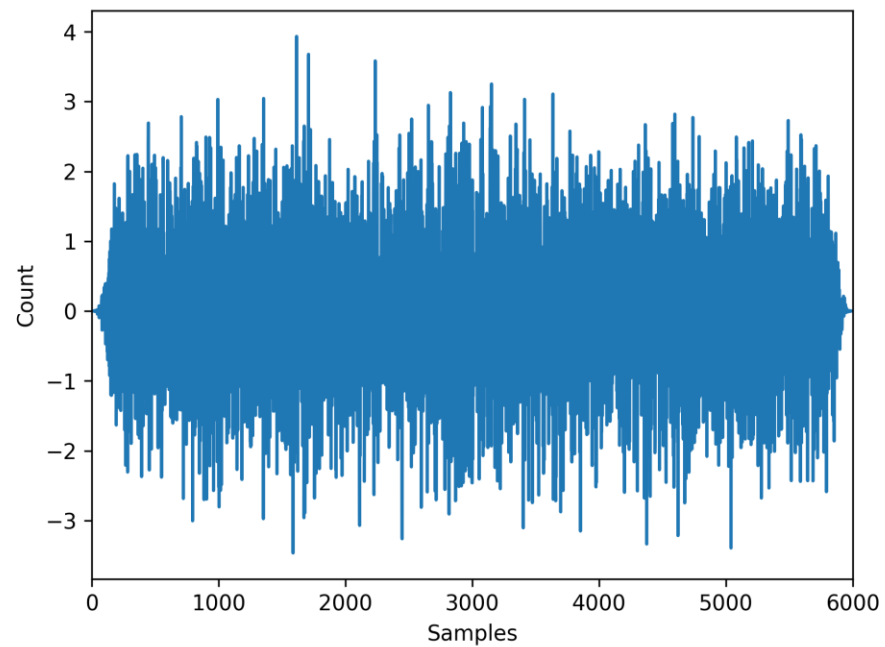
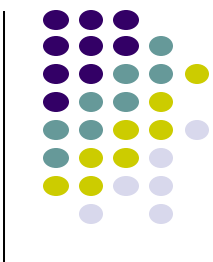
Earthquake Monitoring System



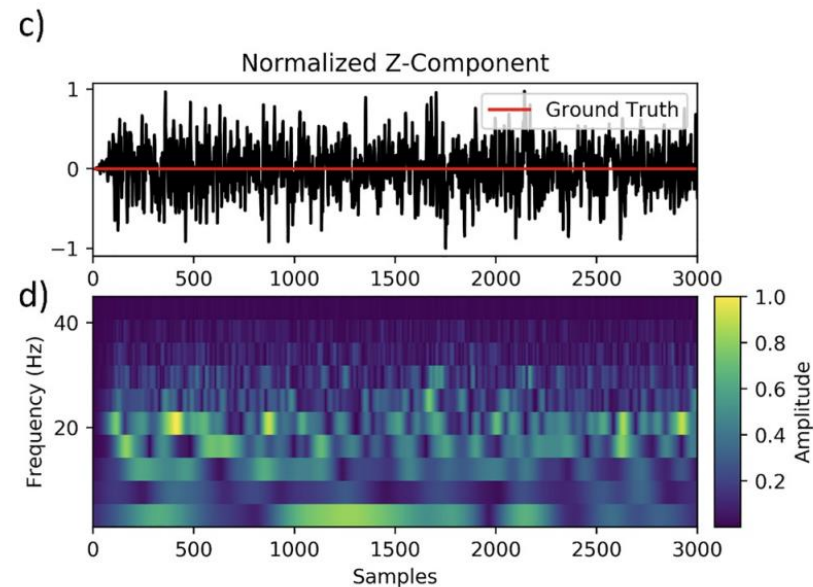
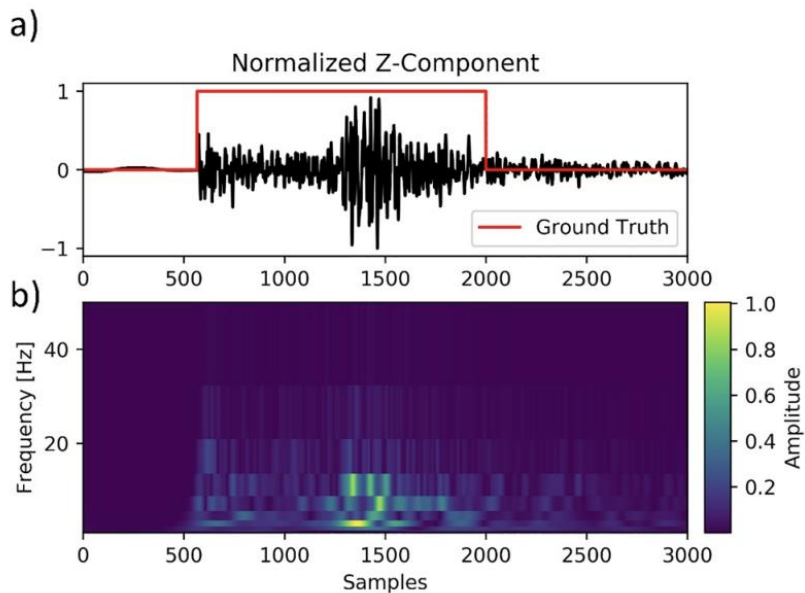
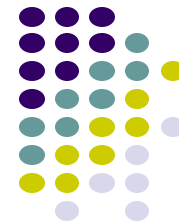
Earthquake Monitoring System



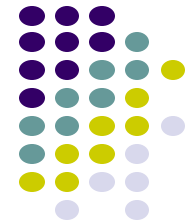
Seismic Event Classification



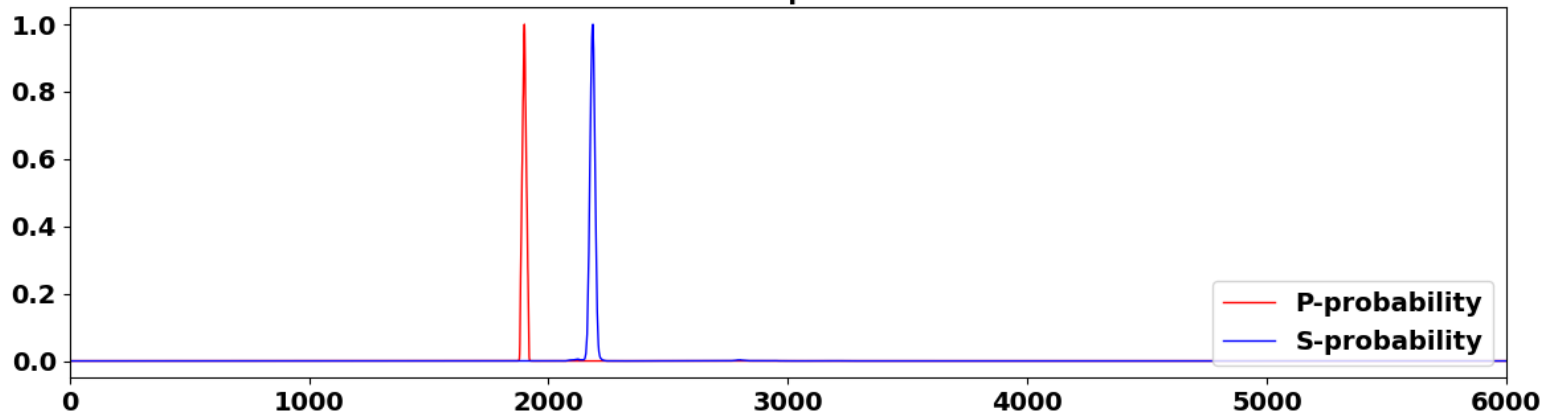
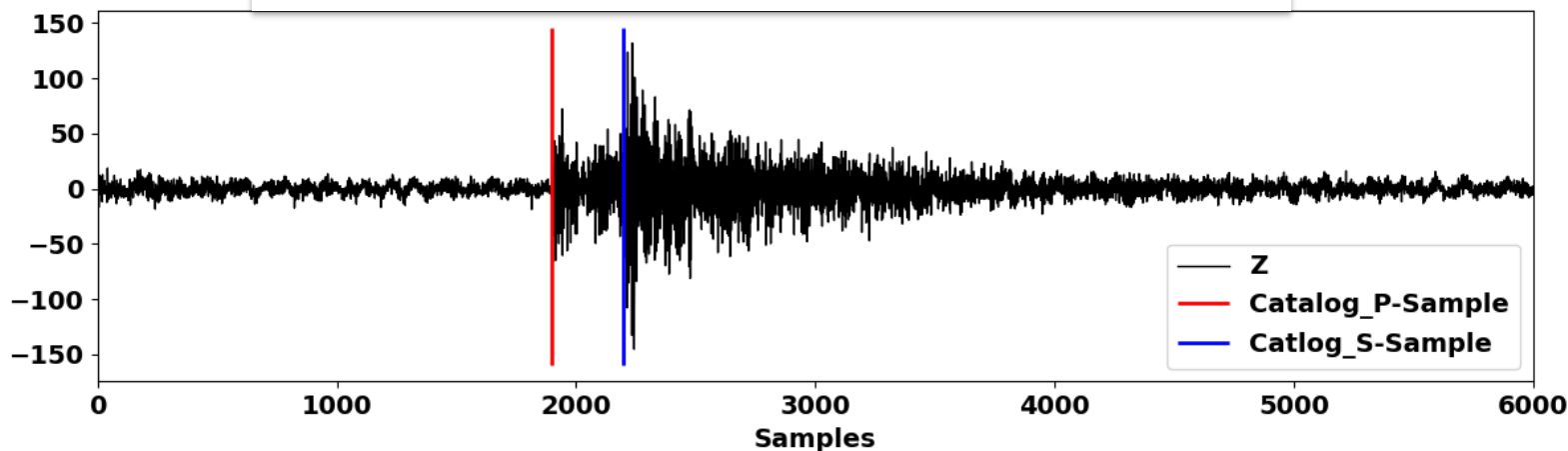
Sample-based Classification (Deteciton)



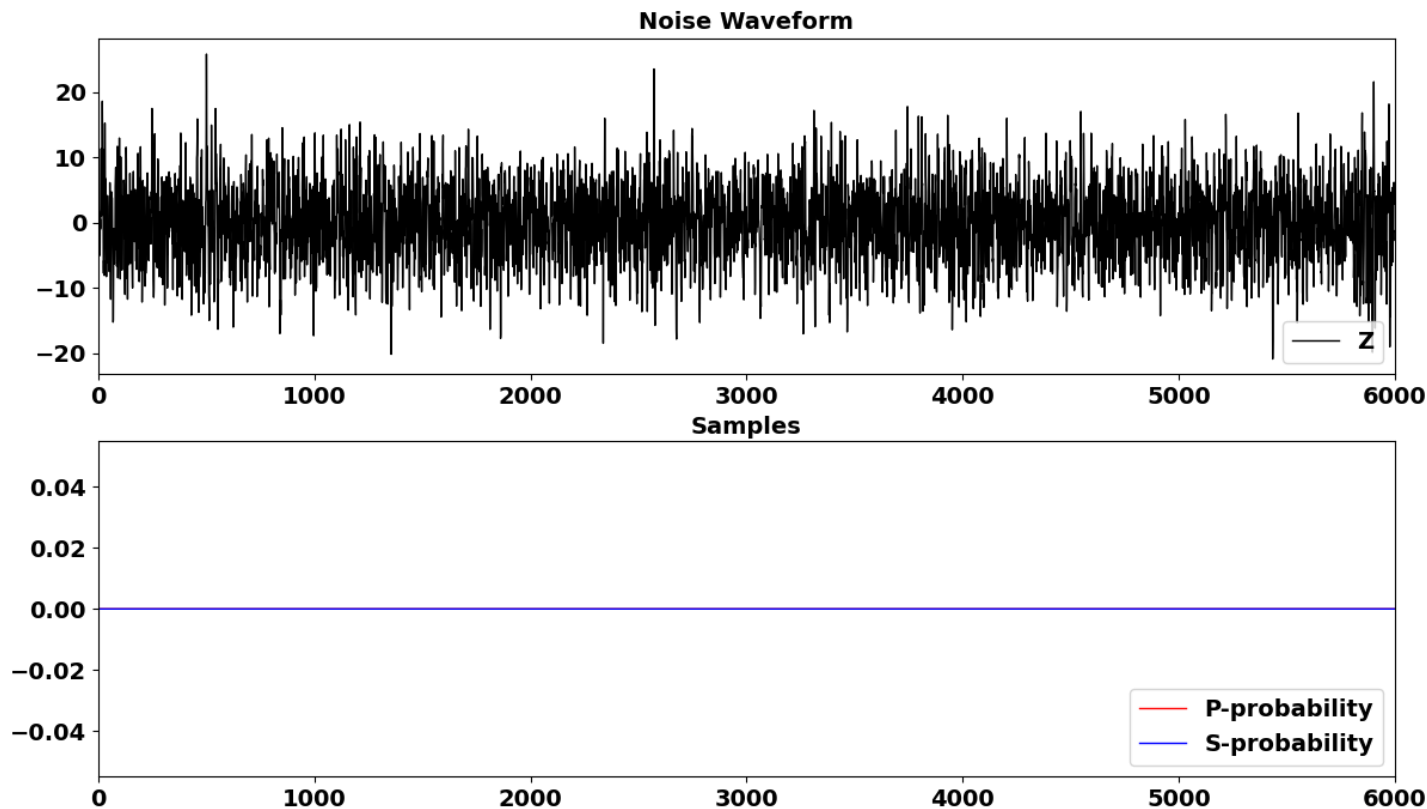
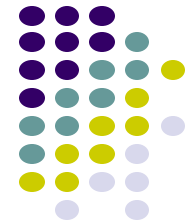
Sample-based Classification (Picking)



Recorded Earthquake



Sample-based Classification (Picking)



A wide-angle photograph of a desert landscape. In the foreground, there are prominent, layered rock formations, likely sandstone, with a warm, golden-brown hue. The rock layers are stacked horizontally, creating a textured appearance. Beyond the rocks, a vast, flat valley stretches out, with a winding road or path visible in the distance. The background features rolling hills and a hazy horizon under a clear sky. The overall scene is arid and expansive.

Thank You

30 17:50