

Course Materials for GEN-AI

Northeastern University

These materials have been prepared and sourced for the course **GEN-AI** at Northeastern University. Every effort has been made to provide proper citations and credit for all referenced works.

If you believe any material has been inadequately cited or requires correction, please contact me at:

Instructor: Ramin Mohammadi
r.mohammadi@northeastern.edu

Thank you for your understanding and collaboration.

Homework - 1

Problem 1 (5 points)

Find the derivative of the softmax function.

Problem 2 (10 points)

Consider a model where the prior distribution over the parameters is a normal distribution with mean zero and variance σ^2 , so that

$$Pr(\phi) = \prod_{j=1}^J \text{Norm}(0, \sigma_j^2; \phi_j),$$

where j indexes the model parameters. We now maximize $\prod_{i=1}^n Pr(y_i | x_i, \phi) Pr(\phi)$. Show that the associated loss function of this model is equivalent to L2 regularization.

Problem 3 (10 points)

Show that the weight decay parameter update with decay rate λ :

$$\phi \leftarrow (1 - \lambda)\phi - \alpha \frac{\partial L}{\partial \phi},$$

on the original loss function $L[\phi]$ is equivalent to a standard gradient update using L2 regularization so that the modified loss function $\tilde{L}[\phi]$ is:

$$\tilde{L}[\phi] = L[\phi] + \frac{\lambda}{2\alpha} \sum_k \phi_k^2,$$

where ϕ are the parameters, and α is the learning rate.

Problem 4 (8 points)

A network consists of three 1D convolutional layers. At each layer, a zero-padded convolution with kernel size three, stride one, and dilation one is applied. What size is the receptive field of the hidden units in the third layer?

Problem 5 (8 points)

A network consists of three 1D convolutional layers. At each layer, a zero-padded convolution with kernel size seven, stride one, and dilation one is applied. What size is the receptive field of hidden units in the third layer?

Problem 6 (8 points)

Consider a convolutional network with 1D input x . The first hidden layer H_1 is computed using a convolution with kernel size five, stride two, and a dilation rate of one. The second hidden layer H_2 is computed using a convolution with kernel size three, stride one, and a dilation rate of one. The third hidden layer H_3 is computed using a convolution with kernel size five, stride one, and a dilation rate of two. What are the receptive field sizes at each hidden layer?

Problem 7 (12 points)

A surface is guaranteed to be convex if the eigenvalues of the Hessian $H[\phi]$ are positive everywhere. In this case, the surface has a unique minimum, and optimization is easy. Find an algebraic expression for the Hessian matrix,

$$H[\phi] = \begin{bmatrix} \frac{\partial^2 L}{\partial \phi_0^2} & \frac{\partial^2 L}{\partial \phi_0 \partial \phi_1} \\ \frac{\partial^2 L}{\partial \phi_1 \partial \phi_0} & \frac{\partial^2 L}{\partial \phi_1^2} \end{bmatrix},$$

for the linear regression model. Prove that this function is convex by showing that the eigenvalues are always positive. This can be done by showing that both the trace and the determinant of the matrix are positive.

Linear Regression:

Consider applying gradient descent to the 1D linear regression model. The model $f(x_i, \phi)$ maps a scalar input x to a scalar output y and has parameters $\phi = (\phi_0, \phi_1)^T$, which represent the y-intercept and the slope:

$$y = f(x_i, \phi) = \phi_0 + \phi_1 x_i.$$

Given a dataset $\{(x_i, y_i)\}$ containing I input/output pairs, we choose the least squares loss function:

$$L[\phi] = \sum_{i=1}^I (f(x_i, \phi) - y_i)^2 = \sum_{i=1}^I (\phi_0 + \phi_1 x_i - y_i)^2,$$

Problem 8 (12 points)

The logistic regression model uses a linear function to assign an input to one of two classes $y \in \{0, 1\}$. For a 1D input and a 1D output, it has two parameters, ϕ_0 and ϕ_1 , and is defined by:

$$Pr(y = 1|x) = \text{sig}(\phi_0 + \phi_1 x),$$

where sig is the logistic sigmoid function:

$$\text{sig}(z) = \frac{1}{1 + \exp(-z)}.$$

- (i) Plot y against x for this model for different values of ϕ_0 and ϕ_1 and explain the qualitative meaning of each parameter.
- (ii) What is a suitable loss function for this model?
- (iii) Compute the derivatives of this loss function with respect to the parameters.
- (iv) Generate ten data points from a normal distribution with mean -1 and standard deviation 1 and assign them the label $y = 0$. Generate another ten data points from a normal distribution with mean 1 and standard deviation 1 and assign these the label $y = 1$. Plot the loss as a heatmap in terms of the two parameters ϕ_0 and ϕ_1 .
- (v) Is this loss function convex? How could you prove this?

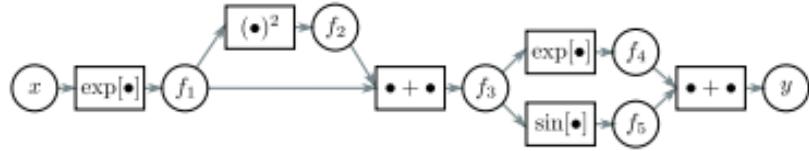


Figure 1: Image for question 9 and 10

Problem 9 (15 points)

This problem explores computing derivatives on general acyclic computational graphs. Consider the function:

$$y = \exp(\exp[x] + \exp[x]^2) + \sin(\exp[x] + \exp[x]^2).$$

We can break this down into a series of intermediate computations so that:

$$\begin{aligned} f_1 &= \exp[x], \\ f_2 &= f_1^2, \\ f_3 &= f_1 + f_2, \\ f_4 &= \exp[f_3], \\ f_5 &= \sin[f_3], \\ y &= f_4 + f_5. \end{aligned}$$

The associated computational graph is depicted in figure 1. Compute the derivative $\frac{\partial y}{\partial x}$ by reverse-mode differentiation. In other words, compute in order:

$$\frac{\partial y}{\partial f_5}, \frac{\partial y}{\partial f_4}, \frac{\partial y}{\partial f_3}, \frac{\partial y}{\partial f_2}, \frac{\partial y}{\partial f_1}, \text{ and } \frac{\partial y}{\partial x},$$

using the chain rule in each case to make use of the derivatives already computed.

Problem 10 (15 points)

For the same function as in problem 9, compute the derivative $\frac{\partial y}{\partial x}$ by forward-mode differentiation. In other words, compute in order:

$$\frac{\partial f_1}{\partial x}, \frac{\partial f_2}{\partial x}, \frac{\partial f_3}{\partial x}, \frac{\partial f_4}{\partial x}, \frac{\partial f_5}{\partial x}, \text{ and } \frac{\partial y}{\partial x},$$

using the chain rule in each case to make use of the derivatives already computed. Why do we not use forward-mode differentiation when we calculate the parameter gradients for deep networks?

Problem ① - Find the derivative of softmax

For a vector $\mathbf{z} = (z_1, z_2, \dots, z_n)$ the softmax function can be defined as:

$$\text{softmax}(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{j=1}^n e^{z_j}}$$

For each element $s_i = \text{softmax}(\mathbf{z})_i$ represents the probability assigned to class "i".

Each output in a neural network using softmax as a classification function will depend on all the inputs, that's why softmax's derivative is not a simple vector but a Jacobian Matrix

$$J_{ij} = \frac{\partial s_i}{\partial z_j}$$

Probabilidad del elemento pertenecer a la clase I

} ~~p/ la probabilidad de los logits.~~

This derivative has 2 cases

↳ ① when $i=j$ // Diagonal

$$\frac{\partial s_i}{\partial z_i} = s_i \cdot (1 - s_i)$$

↳ ② when $i \neq j$ // outside diagonal

$$\frac{\partial s_i}{\partial z_j} = -s_i \cdot s_j$$

Writing both cases:

$$\frac{\partial s_i}{\partial z_j} = s_i \cdot (\delta_{ij} - s_j) \quad \text{where } \delta_{ij} = 1 \text{ if } i=j \\ \delta_{ij} = 0 \text{ if } i \neq j$$

so,

$$J = \text{diag}(s) - s \cdot s^T$$

$\text{diag}(s)$ = diagonal matrix with $[s_1, s_2 \dots s_n]$ in the diagonal

and $s \cdot s^T$ is the outer product.

Problem 2 (10 points)

Consider a model where the prior distribution over the parameters is a normal distribution with mean zero and variance σ^2 , so that

$$Pr(\phi) = \prod_{j=1}^J \text{Norm}(0, \sigma_j^2; \phi_j),$$

where j indexes the model parameters. We now maximize $\prod_{i=1}^n Pr(y_i | x_i, \phi) Pr(\phi)$. Show that the associated loss function of this model is equivalent to L2 regularization.

To avoid overfitting, we penalized big weights. In order to do that, two approaches can be used

↳ Regularization approach (L2) : a penalty term is added to the loss function.

$$\text{Loss} = \text{Loss}_0 + \underbrace{\gamma (\|\phi\|)^2}_{\text{this term penalizes greater weights.}}$$

↳ Bayesian Approach : Assuming that the ϕ comes from a normal distribution centered at zero.

$$\text{Prior} : Pr(\phi) = \text{Norm}(0, \sigma^2; \phi)$$

With this formula

$$Pr(\phi) = \prod_{j=1}^J \text{Norm}(0; \sigma_j^2; \phi_j)$$

and also considering that we want to maximize

$$\prod_{i=1}^n Pr(y_i | x_i, \phi) \cdot Pr(\phi)$$

Taking logarithm (MAP)

$$\log \prod_{i=1}^n \Pr(y_i | x_i, \phi) + \log \Pr(\phi)$$



$$= \sum_{i=1}^n \log \Pr(y_i | x_i, \phi) + \log \prod_{j=1}^J \text{Norm}(0, \sigma_j^2; \phi_j)$$

Expanding the prior gaussian:

$$\log \text{Norm}(0, \sigma_j^2; \phi_j) = \log \left(\frac{1}{\sqrt{2\pi\sigma_j^2}} \cdot e^{-\frac{\phi_j^2}{2\sigma_j^2}} \right)$$

$$= -\frac{\phi_j^2}{2\sigma_j^2} - \frac{1}{2} \log(2\pi\sigma_j^2)$$

Summing up the parameters

$$\log \Pr(\phi) = \sum_{j=1}^J \left(-\frac{\phi_j^2}{2\sigma_j^2} \right) + \text{constant}$$

The complete MAP objective is:

$$\sum_{i=1}^n \log \Pr(y_i | x_i, \phi) - \sum_{j=1}^J \frac{\phi_j^2}{2\sigma_j^2}$$

↓
considering a minimization
problem

$$\mathcal{L} = -\sum_{i=1}^n \log \Pr(y_i | x_i, \phi) + \sum_{j=1}^J \phi_j^2 / 2\sigma_j^2$$

$$= -\sum_{i=1}^n \log \Pr(y_i | x_i, \phi) + \frac{1}{2\sigma^2} \|\phi\|^2$$

this is a negative log-likelihood + L₂ regularization with $\lambda = \frac{1}{2\sigma^2}$

Problem 3 (10 points)

Show that the weight decay parameter update with decay rate λ :

$$\phi \leftarrow (1 - \lambda)\phi - \alpha \frac{\partial L}{\partial \phi},$$

on the original loss function $L[\phi]$ is equivalent to a standard gradient update using L2 regularization so that the modified loss function $\tilde{L}[\phi]$ is:

$$\tilde{L}[\phi] = L[\phi] + \frac{\lambda}{2\alpha} \sum_k \phi_k^2,$$

where ϕ are the parameters, and α is the learning rate.

we have two ways of regularizing

① weight decay

$$\phi \leftarrow (1 - \underbrace{\lambda}_{\text{decay rate}}) \cdot \phi - \alpha \cdot \frac{\partial L}{\partial \phi}$$

At each iteration, we multiply the weights * $(1 - \lambda)$ before subtracting the gradient.

② L₂ regularization .

$$\tilde{L}[\phi] = L[\phi] + \frac{\lambda}{2\alpha} \sum_K \phi_K^2$$

we add a penalty term to the loss and gradient descent is used on this new loss.

so, considering the standard gradient with L_2

$$\hookrightarrow \varphi \leftarrow \varphi - \alpha \cdot \frac{\partial \tilde{L}}{\partial \varphi}$$

we compute the \tilde{L} gradient

$$\hookrightarrow \frac{\partial \tilde{L}}{\partial \varphi} = \frac{\partial L}{\partial \varphi} + \left(\frac{\lambda}{2\alpha}\right) \cdot 2\varphi = \frac{\partial L}{\partial \varphi} + (\lambda/2)\varphi$$

Substituting in the update

$$\hookrightarrow \varphi \leftarrow \varphi - \alpha \cdot [\frac{\partial L}{\partial \varphi} + (\lambda/\alpha)\varphi]$$

$$\varphi \leftarrow \varphi - \alpha \cdot \frac{\partial L}{\partial \varphi} - \lambda\varphi$$

$$\boxed{\varphi \leftarrow (1-\lambda)\varphi - \alpha \cdot \frac{\partial L}{\partial \varphi}} // \text{ weight decay.}$$

Problem 4 (8 points)

A network consists of three 1D convolutional layers. At each layer, a zero-padded convolution with kernel size three, stride one, and dilation one is applied. What size is the receptive field of the hidden units in the third layer?

Receptive field \rightarrow we want to compute how many pixels from the original image a neuron in the 3rd convolutional layer "sees".



Kernel size = 3

stride "s" = 1

dilation (spacing between kernel elements) = 1

Padding "p" = 0

The formula for calculating RF for a single convolutional layer \rightarrow $RF_{\text{layer}} = \underbrace{K}_{\text{Kernel}} + \underbrace{(K-1) \times (d-1)}_{\text{dilation effect}}$

Replacing with dilation = 1 :

$$RF_{\text{layer}} = K + (K-1) \cdot (0) = K = 3$$

Layer 1

$$RF_{\text{layer} 1} = 3$$

Layer 2

$$RF_2 = RF_1 + (K-1) \cdot \text{stride}_1$$

$$RF_2 = 3 + (3-1) \cdot 1$$

$$RF_2 = 5$$

Layer 3

$$RF_3 = RF_2 + (K-1) \cdot \text{stride}_2$$

$$RF_3 = 5 + (3-1) \cdot 1$$

RF₃ = 7 // The receptive field of units in the 3rd layer is 7 (each neuron in layer 3 sees and is influenced by 7 consecutive pixels from the original input).

Problem 5 (8 points)

A network consists of three 1D convolutional layers. At each layer, a zero-padded convolution with kernel size seven, stride one, and dilation one is applied. What size is the receptive field of hidden units in the third layer?

Kernel size = 7

stride = 1

dilation = 1

padding = zero-padded

Computing RF

• $RF_1 = 7$ // each neuron in layer 1 sees 7 pixels.

$$\cdot RF_2 = RF_1 + (K-1) \cdot \text{stride}_1$$

$$RF_2 = 7 + (7-1) \cdot 1$$

$$RF_2 = 13$$

$$\cdot RF_3 = RF_2 + (K-1) \cdot \text{stride}_2$$

$$RF_3 = 13 + (7-1) \cdot 1$$

$$\boxed{RF_3 = 19} \quad // \text{The receptive field of units in the 3rd layer is 19.}$$

Problem 6 (8 points)

Consider a convolutional network with 1D input x . The first hidden layer H_1 is computed using a convolution with kernel size five, stride two, and a dilation rate of one. The second hidden layer H_2 is computed using a convolution with kernel size three, stride one, and a dilation rate of one. The third hidden layer H_3 is computed using a convolution with kernel size five, stride one, and a dilation rate of two. What are the receptive field sizes at each hidden layer?

$$\begin{array}{lll} H_1: \text{Kernel} = 5 & \text{stride} = 2 & \text{dilation} = 1 \\ H_2: \text{Kernel} = 3 & \text{stride} = 1 & \text{dilation} = 1 \\ H_3: \text{Kernel} = 5 & \text{stride} = 1 & \text{dilation} = 2 \end{array}$$

$$\underbrace{\text{input}}_{\text{RF}_0 = 1} \rightarrow \text{stride}_0 = 1$$

$$\begin{array}{ll} \underbrace{H_1}_{\text{effective kernel size}} \rightarrow & K_{\text{effective},1} = K + (K-1) \cdot (d-1) \\ & = 5 + (5-1) \cdot (1-1) \\ & = 5 \end{array}$$

$$\begin{array}{l} \frac{RF_1}{RF_0} \\ RF_1 = RF_0 + (K_{\text{effective},1} - 1) \cdot \text{stride}_0 \\ = 1 + (5-1) \cdot 1 \\ = 5 \end{array}$$

$$\underbrace{RF_1 \rightarrow 5}_{,}$$

Layer 2 - H₂

Effective Kernel size 2

$$\begin{aligned} K_{\text{effective}}_2 &= K + (K-1) \cdot (d-1) \\ &= 3 + (3-1) \cdot (1-1) \\ &= 3 \end{aligned}$$

RF₂

$$\begin{aligned} RF_2 &= RF_1 + (K_{\text{effective}}_2 - 1) \cdot \text{stride}_1 \\ &= 5 + (3-1) \cdot 2 \\ &= 9 \end{aligned}$$

RF₂ = 9

Layer 3 → H₃

Effective Kernel size H₃

$$\begin{aligned} K_{\text{effective}}_3 &= K + (K-1) \cdot (d-1) \\ &= 5 + (5-1) \cdot (2-1) \\ &= 9 \end{aligned}$$

RF₃

$$RF_3 = RF_2 + (K_{\text{effective}}_3 - 1) \cdot \text{stride}_2$$

$$RF_3 = 9 + (9-1) \cdot 1$$

RF₃ = 17

$$RF_1 = 5$$

$$RF_2 = 9$$

RF₃ = 17

Problem 7 (12 points)

A surface is guaranteed to be convex if the eigenvalues of the Hessian $H[\phi]$ are positive everywhere. In this case, the surface has a unique minimum, and optimization is easy. Find an algebraic expression for the Hessian matrix,

$$H[\phi] = \begin{bmatrix} \frac{\partial^2 L}{\partial \phi_0^2} & \frac{\partial^2 L}{\partial \phi_0 \partial \phi_1} \\ \frac{\partial^2 L}{\partial \phi_1 \partial \phi_0} & \frac{\partial^2 L}{\partial \phi_1^2} \end{bmatrix},$$

linear regression

for the linear regression model. Prove that this function is convex by showing that the eigenvalues are always positive. This can be done by showing that both the trace and the determinant of the matrix are positive.

Linear Regression:

Consider applying gradient descent to the 1D linear regression model. The model $f(x_i, \phi)$ maps a scalar input x to a scalar output y and has parameters $\phi = (\phi_0, \phi_1)^T$, which represent the y-intercept and the slope:

$$y = f(x_i, \phi) = \phi_0 + \phi_1 x_i.$$

Given a dataset $\{(x_i, y_i)\}$ containing I input/output pairs, we choose the least squares loss function:

$$L[\phi] = \sum_{i=1}^I (f(x_i, \phi) - y_i)^2 = \sum_{i=1}^I (\phi_0 + \phi_1 x_i - y_i)^2,$$

for this problem I have to proof that the Loss function of linear regression is convex, that means that the surface is convex and positive in both ways so it is easily optimizable and has a unique solution.

Model $y = \underbrace{\phi_0}_{\text{interc.}} + \underbrace{\phi_1 x_i}_{\text{slope * input}}$

Parameters $\phi = (\phi_0, \phi_1)^T$

Loss Function $L(\phi) = \sum_{i=1}^I (\phi_0 + \phi_1 x_i - y_i)^2$ // sum of squares errors
of linear regression

First derivatives

* with respect to φ_0

$$\frac{\partial L}{\partial \varphi_0} = \frac{\partial}{\partial \varphi_0} \cdot [\sum_i (\varphi_0 + \varphi_1 x_i - y_i)^2]$$

Applying chain rule

$$= \sum_i 2(\varphi_0 + \varphi_1 x_i - y_i) \cdot \underbrace{\frac{\partial}{\partial \varphi_0} (\varphi_0 + \varphi_1 x_i - y_i)}_{= 1}$$

$$\frac{\partial L}{\partial \varphi_0} = 2 \sum_i (\varphi_0 + \varphi_1 x_i - y_i) // \text{residual error}$$

* with respect to φ_1

$$\frac{\partial L}{\partial \varphi_1} = \frac{\partial}{\partial \varphi_1} [\sum_i (\varphi_0 + \varphi_1 x_i - y_i)^2]$$

chain rule:

$$\frac{\partial L}{\partial \varphi_1} = \sum_i 2(\varphi_0 + \varphi_1 x_i - y_i) \cdot (\underbrace{\frac{\partial}{\partial \varphi_1} (\varphi_0 + \varphi_1 x_i - y_i)}_{x_i})$$

$$\frac{\partial L}{\partial \varphi_1} = \sum_i 2(\varphi_0 + \varphi_1 x_i - y_i) \cdot x_i$$

residual * input

Now, I have to compute the 2nd derivatives since these are the Hessian Elements that \ni to the Hessian Matrix.

$$\bullet H[1,1] : \frac{\partial^2 \ell}{\partial \varphi_0^2}$$

$$\frac{\partial^2 L}{\partial \varphi_0^2} = \frac{\partial}{\partial \varphi_0} \left[2 \sum_i (\varphi_0 + \varphi_1 x_i - y_i) \right]$$

$$= 2 \sum_i \frac{\partial}{\partial \varphi_0} (\varphi_0 + \varphi_1 x_i - y_i)$$

$$= 2 \sum_i 1$$

$$= 2I \quad // \text{ 2 times the number of data}$$

since $I > 0 \rightarrow \text{Positive} \checkmark$

$$\bullet H[2,2] : \frac{\partial^2 \ell}{\partial^2 \varphi_1}$$

$$\frac{\partial^2 L}{\partial \varphi_1^2} = \frac{\partial}{\partial \varphi_1} \left[2 \sum_i (\varphi_0 + \varphi_1 x_i - y_i) \cdot x_i \right]$$

$$= 2 \sum_i x_i \cdot \frac{\partial}{\partial \varphi_1} (\varphi_0 + \varphi_1 x_i - y_i)$$

$$= 2 \sum_i x_i \cdot x_i$$

$$= 2 \sum_i x_i^2 \quad // \text{ 2 times the sum of inputs}^2$$

since $x_i^2 \geq 0$ always $\rightarrow \text{Positive} \checkmark$

$$\bullet H[1,2] = H[2,1] = \frac{\partial^2 \ell}{\partial \varphi_0 \partial \varphi_1}$$

$$\begin{aligned}
 \frac{\partial^2 L}{\partial \theta_0 \partial \theta_1} &= \frac{\partial}{\partial \theta_0} \left[2 \sum_i (\theta_0 + \theta_1 x_i - y_i) \cdot x_i \right] \\
 &= 2 \sum_i x_i \cdot \frac{\partial}{\partial \theta_0} (\theta_0 + \theta_1 x_i - y_i) \\
 &= 2 \sum_i x_i \cdot 1 \\
 &= 2 \sum_i x_i \quad // \quad 2 \times \text{sum of the inputs}
 \end{aligned}$$

Hessian Matrix

$$H = \begin{bmatrix} 2I & 2 \sum_i x_i \\ 2 \sum_i x_i & 2 \sum_i x_i^2 \end{bmatrix} = 2 \begin{bmatrix} I & \sum_i x_i \\ \sum_i x_i & \sum_i x_i^2 \end{bmatrix}$$

Now, I have to prove positive eigenvalues (Trace method).
(sum of diagonal elements)

$$\text{Trace}(H) = 2(I + \sum_i x_i^2) > 0 \quad \checkmark$$

$I > 0$ (we have datapoints) \rightarrow trace is positive

$$\text{Trace} = \lambda_1 + \lambda_2 > 0 \quad \checkmark$$

For the 2nd condition : prove positive eigenvalues (Determinant)

$\text{Determinant} = (\text{diagonal product}) - (\text{OFF diagonal product})^2$

$$\begin{aligned}
 \text{Det}(H) &= (2I)(2\sum_i x_i^2) - (2\sum_i x_i)^2 \\
 &= 4I \cdot \sum_i x_i^2 - 4(\sum_i x_i)^2 \\
 &= 4 \underbrace{[I \cdot \sum_i x_i^2 - (\sum_i x_i)^2]}_{\substack{\text{I have to show that this is } > 0 \\ (\text{Cauchy - Inequality})}
 \end{aligned}$$

IF I consider real data :

$$I \cdot \sum_i x_i^2 - (\sum_i x_i)^2 > 0$$

Therefore :

$$\begin{aligned}
 \text{Det}(H) &= 4 [I \cdot \sum_i x_i^2 - (\sum_i x_i)^2] > 0 \quad \checkmark \\
 | \\
 \text{is positive}
 \end{aligned}$$

IF $\det > 0$ AND trace > 0 , both eigenvalues > 0

Therefore H is positive definite $\therefore L(\phi)$ is convex \therefore

Linear regression has a unique global minimum.

Problem 8 (12 points)

The logistic regression model uses a linear function to assign an input to one of two classes $y \in \{0, 1\}$. For a 1D input and a 1D output, it has two parameters, ϕ_0 and ϕ_1 , and is defined by:

$$Pr(y = 1|x) = \text{sig}(\phi_0 + \phi_1 x),$$

where sig is the logistic sigmoid function:

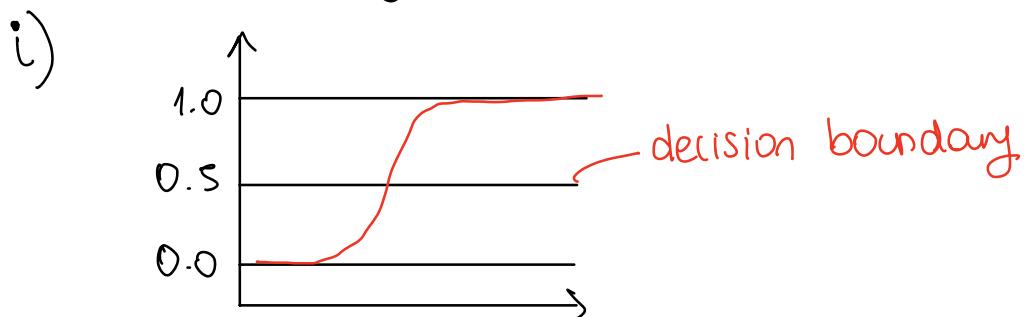
$$\text{sig}(z) = \frac{1}{1 + \exp(-z)}.$$

- (i) Plot y against x for this model for different values of ϕ_0 and ϕ_1 and explain the qualitative meaning of each parameter.
- (ii) What is a suitable loss function for this model?
- (iii) Compute the derivatives of this loss function with respect to the parameters.
- (iv) Generate ten data points from a normal distribution with mean -1 and standard deviation 1 and assign them the label $y = 0$. Generate another ten data points from a normal distribution with mean 1 and standard deviation 1 and assign these the label $y = 1$. Plot the loss as a heatmap in terms of the two parameters ϕ_0 and ϕ_1 .
- (v) Is this loss function convex? How could you prove this?

① Model $Pr(y = 1|x) = \text{sig}(\phi_0 + \phi_1 x)$

where $\text{sig}(z) = 1 / (1 + \exp(-z))$

Parameters $\left\{ \begin{array}{l} \phi_0 = \text{bias (shifts decision boundary)} \\ \phi_1 = \text{slope (transition)} \end{array} \right.$



Parameters $\left\{ \begin{array}{l} \phi_0 = 0 \text{ centered the curve at } x=0 \\ \phi_0 > 0 \text{ curve shifts left.} \\ \phi_0 < 0 \text{ curve shifts right} \end{array} \right.$

↳ $\begin{cases} \Theta_1 \text{ controls curve steepness} \\ \Theta_1 \text{ large } \rightarrow \text{steep transition, ergo, confident predictions} \\ \Theta_1 \text{ small } \rightarrow \text{gradual transition} \\ \Theta_1 < 0 \rightarrow \text{flipped curve} \end{cases}$

ii) suitable loss function

↳ Log-loss (binary cross entropy) since measures how well predicted probabilities match true labels, penalizes confident wrong predictions and is standard for binary classification as well as convex for logistic regression.

iii) Derivative of log-loss

Considering the sigmoid property

$$\frac{\partial}{\partial z} [\text{sig}(z)] = \text{sig}(z) \cdot (1 - \text{sig}(z))$$

$$\text{Let } p_i = \text{sig}(\Theta_0 + \Theta_1 x_i)$$

Derivative with respect to Θ_0 :

$$L = -\sum_i [y_i \cdot \log(p_i) + (1-y_i) \cdot \log(1-p_i)]$$

chain rule to each term

$$\frac{\partial}{\partial \Theta_0} [y_i \cdot \log(p_i)] = y_i \cdot \left(\frac{1}{p_i}\right) \cdot \frac{\partial p_i}{\partial \Theta_0}$$

Find $\partial p_i / \partial \Theta_0$

$$\frac{\partial p_i}{\partial \Theta_0} = p_i(1-p_i) \cdot 1 = p_i \cdot (1-p_i)$$

$$1^{\text{st}} \text{ term} \cdot \left(\frac{1}{p_i} \right) \cdot p_i (1-p_i) = y_i \cdot (1-p_i)$$

$$2^{\text{nd}} \text{ term} \rightarrow - (1-y_i) p_i$$

Combining

$$\frac{\partial L}{\partial \theta_0} = - \sum_i [y_i (1-p_i) - (1-y_i) p_i]$$

$$\frac{\partial L}{\partial \theta_0} = - \sum_i [y_i - p_i]$$

$$\frac{\partial L}{\partial \theta_0} = \sum_i (p_i - y_i)$$

Derivative w.r.t θ_1 :

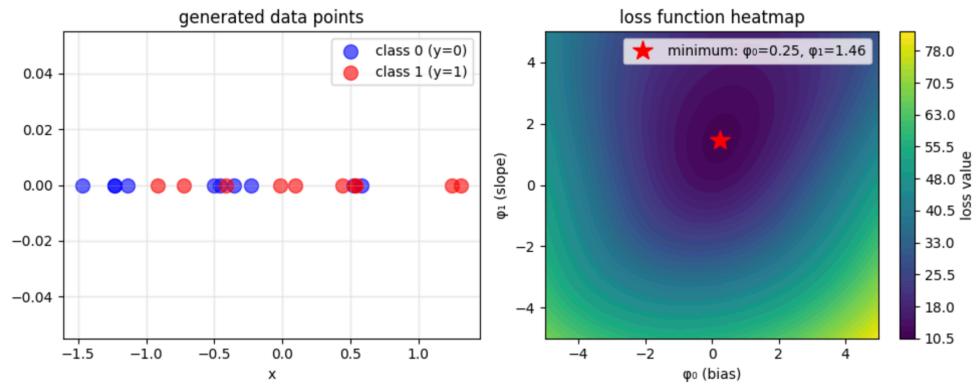
Same than before but:

$$\frac{\partial p_i}{\partial \theta_1} = p_i (1-p_i) \cdot x_i$$

Following same steps

$$\frac{\partial L}{\partial \theta_1} = \sum_i (p_i - y_i) \cdot x_i$$

IV



V As we can see from the heatmap, the Loss Function is convex ($\phi_0 = 0.25, \phi_1 = 1.46$; I did all this with code). This can be proven by showing that the Hessian matrix is positive semi-definite. The 2nd derivatives $\frac{\partial^2 L}{\partial \phi_0^2} = \sum_i p_i(1-p_i)$ and $\frac{\partial^2 L}{\partial \phi_1^2} =$

$\sum_i p_i(1-p_i) \cdot x_i^2$ are always positive since $0 < p_i < 1$ for sigmoid. Therefore the loss has a unique global minimum and gradient descent will always converge.

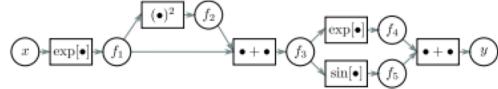


Figure 1: Image for question 9 and 10

Problem 9 (15 points)

This problem explores computing derivatives on general acyclic computational graphs. Consider the function:

$$y = \exp(\exp[x] + \exp[x]^2) + \sin(\exp[x] + \exp[x]^2).$$

We can break this down into a series of intermediate computations so that:

$$\begin{aligned} f_1 &= \exp[x], \\ f_2 &= f_1^2, \\ f_3 &= f_1 + f_2, \\ f_4 &= \exp[f_3], \\ f_5 &= \sin[f_3], \\ y &= f_4 + f_5. \end{aligned}$$

The associated computational graph is depicted in figure 1. Compute the derivative $\frac{\partial y}{\partial x}$ by reverse-mode differentiation. In other words, compute in order:

$$\frac{\partial y}{\partial f_5}, \frac{\partial y}{\partial f_4}, \frac{\partial y}{\partial f_3}, \frac{\partial y}{\partial f_2}, \frac{\partial y}{\partial f_1}, \text{ and } \frac{\partial y}{\partial x},$$

using the chain rule in each case to make use of the derivatives already computed.

I need to compute $\frac{\partial y}{\partial x}$ for the nested function using backpropagation. My approach is to break the complex function into simple intermediate steps.

$$y = \exp(\exp[x] + \exp[x]^2) + \sin(\exp[x] + \exp[x]^2)$$

Breaking down the computational graph :

$$f_1 = \exp[x]$$

$$f_2 = f_1^2$$

$$f_3 = f_1 + f_2$$

$$f_4 = \exp[f_3]$$

$$f_5 = \sin[f_3]$$

$$y = f_4 + f_5$$

$$\textcircled{a} \quad \frac{\partial y}{\partial f_5} \rightarrow y = f_4 + f_5 \quad \frac{\partial y}{\partial f_5} = 1$$

$$\textcircled{b} \quad \frac{\frac{\partial y}{\partial f_4}}{y = f_4 + f_5}$$

$$\frac{\partial y}{\partial f_4} = 1$$

$$\textcircled{c} \quad \frac{\frac{\partial y}{\partial f_3}}{f_3 \rightarrow 2 \text{ childrens}}$$

f_4 f_5
 Accumulate gradients from both parts.

$$\frac{\partial y}{\partial f_3} = \frac{\partial y}{\partial f_4} \cdot \frac{\partial f_4}{\partial f_3} + \frac{\partial y}{\partial f_5} \cdot \frac{\partial f_5}{\partial f_3}$$

$$\text{since } f_4 = \exp [f_3] \rightarrow \frac{\partial f_4}{\partial f_3} = \exp [f_3] = f_4$$

$$f_5 = \sin [f_3] \rightarrow \frac{\partial f_5}{\partial f_3} = \cos [f_3]$$

$$\begin{aligned}\frac{\partial y}{\partial f_3} &= 1 \cdot f_4 + 1 \cdot \cos [f_3] \\ &= f_4 + \cos [f_3]\end{aligned}$$

$$\textcircled{D} \quad \underbrace{\frac{dy}{df_2}} \quad | \quad f_2 \text{ only affects } f_3 ; \text{ so :}$$

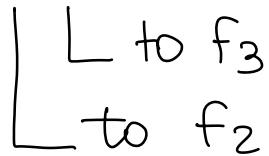
$$\frac{\partial y}{\partial f_2} = \frac{\partial y}{\partial f_3} \cdot \frac{\partial f_3}{\partial f_2}$$

$$\text{since } f_3 = f_1 + f_2 \rightarrow \frac{\partial f_3}{\partial f_2} = 1$$

$$\begin{aligned} \frac{\partial y}{\partial f_2} &= (f_4 + \cos[f_3]) \cdot 1 \\ &= f_4 + \cos[f_3] \end{aligned}$$

$$\textcircled{e} \quad \underline{\frac{\partial y}{\partial f_1}}$$

$f_1 \rightarrow 2 \text{ childrens}$



$$\frac{\partial y}{\partial f_1} = \frac{\partial y}{\partial f_3} \cdot \frac{\partial f_3}{\partial f_1} + \frac{\partial y}{\partial f_2} \cdot \frac{\partial f_2}{\partial f_1}$$

$$\text{considering } f_3 = f_1 + f_2 \rightarrow \frac{\partial f_3}{\partial f_1} = 1$$

$$f_2 = f_1^2 \rightarrow \frac{\partial f_2}{\partial f_1} = 2f_1$$

$$\begin{aligned}\frac{\partial y}{\partial f_1} &= (f_4 + \cos[f_3] \cdot 1 + (f_4 + \cos[f_3]) \cdot 2f_1) \\ &= (f_4 + \cos[f_3])(1 + 2f_1)\end{aligned}$$

so, $\frac{\partial y}{\partial x} \rightarrow \frac{\partial y}{\partial f_1} \cdot \frac{\partial f_1}{\partial x}$

since $f_1 = \exp[x] \rightarrow \frac{\partial f_1}{\partial x} = \exp[x] = f_1$

$$\begin{aligned}\therefore \frac{\partial y}{\partial x} &= (f_4 + \cos[f_3])(1 + 2f_1) \cdot f_1 \\ &= f_1(f_4 + \cos[f_3])(1 + 2f_1)\end{aligned}$$

Problem 10 (15 points)

For the same function as in problem 9, compute the derivative $\frac{\partial y}{\partial x}$ by forward-mode differentiation. In other words, compute in order:

$$\frac{\partial f_1}{\partial x}, \frac{\partial f_2}{\partial x}, \frac{\partial f_3}{\partial x}, \frac{\partial f_4}{\partial x}, \frac{\partial f_5}{\partial x}, \text{ and } \frac{\partial y}{\partial x}$$

using the chain rule in each case to make use of the derivatives already computed. Why do we not use forward-mode differentiation when we calculate the parameter gradients for deep networks?

For this problem, I need to compute the same derivative $\frac{\partial y}{\partial x}$ than problem 9 but using forward mode differentiation. The difference is in the direction: I start at x and propagate derivatives forward to the graph until I reach y .

In Forward mode, I compute how changes in x affect each intermediate variable as I go forward. This contrasts with reverse mode where I computed how y depends on each variable working backwards.

computational graph

$$f_1 = \exp[x]$$

$$f_2 = f_1^2$$

$$f_3 = f_1 + f_2$$

$$f_4 = \exp[f_3]$$

$$f_5 = \sin[f_3]$$

$$y = f_4 + f_5$$

Applying the forward differentiation:

$$\left. \frac{\partial F_1}{\partial x} \right|_{f_1 = \exp[x]} \quad \frac{\partial F_1}{\partial x} = \exp[x] = f_1$$

$$\left. \frac{\partial f_2}{\partial x} \right|_{f_2 = f_1^2} \quad f_2 = f_1^2 \rightarrow \quad \begin{matrix} \text{using} \\ \text{chain} \\ \text{rule} \end{matrix}$$

$$\rightarrow \frac{\partial f_2}{\partial x} = \frac{\partial f_2}{\partial f_1} \cdot \frac{\partial f_1}{\partial x} = 2f_1 \cdot f_1 = 2f_1^2$$

$$\frac{\partial f_3}{\partial x} \quad | \quad f_3 = f_1 + f_2$$

$$\begin{aligned}\frac{\partial f_3}{\partial x} &= \frac{\partial f_1}{\partial x} + \frac{\partial f_2}{\partial x} \\ &= f_1 + 2f_1^2 \\ &= f_1(1 + 2f_1)\end{aligned}$$

$$\frac{\partial f_4}{\partial x} \quad |$$

$$f_4 = \exp[f_3]$$

$$\begin{aligned}\frac{\partial f_4}{\partial x} &= \frac{\partial f_4}{\partial f_3} \cdot \frac{\partial f_3}{\partial x} \\ &= \exp[f_3] \cdot f_1(1 + 2f_1) \\ &= f_4 \cdot f_1(1 + 2f_1)\end{aligned}$$

$$\frac{\partial f_5}{\partial x} \quad | \quad f_5 = \sin[f_3]$$

$$\frac{\partial f_5}{\partial x} = \frac{\partial f_5}{\partial f_3} \cdot \frac{\partial f_3}{\partial x}$$

$$\frac{\partial f_5}{\partial x} = \cos(f_3) \cdot f_1(1 + 2f_1)$$

Final answer

$$y = f_4 + f_5 \quad // \quad f_4 \text{ and } f_5 \text{ depends on } x$$

$$\frac{\partial y}{\partial x} = \frac{\partial f_4}{\partial x} + \frac{\partial f_5}{\partial x}$$

$$\begin{aligned} &= f_4 \cdot f_1(1 + 2f_1) + \cos[f_3] \cdot f_1(1+2f_1) \\ &= f_1(1 + 2f_1)(f_4 + \cos[f_3]). \end{aligned}$$

Forward is inefficient for deep neural networks since it implies one pass per variable, so you can just think that if we have 1.000.000 parameters, we will need 1000000 forward passes and, of course, the computational cost of this is really high; In backpropagation (reverse-mode) one backward pass all the parameters ($RT \rightarrow O(1)$).