

CAREFUL OF:
dimensions



[SUPERVISED LEARNING:]

REGRESSION: predicting continuous #

training data: $D_{train} = \{(x^1, y^1), \dots, (x^d, y^d)\}$

feature vector $x = [x^1, \dots, x^d]^T \in \mathbb{R}^d$

label $y \in \mathbb{R} \leftarrow$ scalar

n dpts, each w/ d dim feats & scalar label

LINEAR REG: $h(x, \theta, \theta_0) = \theta^T x + \theta_0$

SQUARED LOSS: $L(g, a) = (g - a)^2$

MSE: $J(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n (\theta^T x_i + \theta_0 - y_i)^2$

Jacobian: $\nabla J(\theta) = \frac{2}{n} \sum_{i=1}^n (x_i - y_i) x_i^T$

Gradient Descent: $\theta^{t+1} = \theta^t - \eta \nabla J(\theta^t)$

Objective Func. J is convex

Gradient Descent: gradually update soln to improve when can't solve OLS or RL analytically

Repeat until: $|J(\theta^t) - J(\theta^{t-1})| < \epsilon$

Global Minimizer: $\theta^* = \arg \min_{\theta} J(\theta)$

Objective Func. J is convex

Objective Func. J is convex

Objective Func. J is convex

Objective Func. J is convex

Objective Func. J is convex

Objective Func. J is convex

Objective Func. J is convex

Objective Func. J is convex

Objective Func. J is convex

Objective Func. J is convex

Objective Func. J is convex

Objective Func. J is convex

Objective Func. J is convex

Objective Func. J is convex

Objective Func. J is convex

Objective Func. J is convex

Objective Func. J is convex

Objective Func. J is convex

Objective Func. J is convex

Objective Func. J is convex

Objective Func. J is convex

Objective Func. J is convex

Objective Func. J is convex

Objective Func. J is convex

Objective Func. J is convex

Objective Func. J is convex

Objective Func. J is convex

Objective Func. J is convex

Objective Func. J is convex

Objective Func. J is convex

Objective Func. J is convex

Objective Func. J is convex

Objective Func. J is convex

REGULARIZATION:

Objective Func J(theta): trying to minimize

J(theta, theta_0, theta_1): $\frac{1}{n} \sum_{i=1}^n (\theta^T x_i + \theta_0 - y_i)^2 + \lambda ||\theta||^2$

Stochastic GD: C each step, randomly choose pt to calculate gradient

Linear Classifier:

h(x, theta, theta_0) = sign(theta^T x + theta_0)

0-1 loss: $L_{0-1}(g, a) = \begin{cases} 0 & \text{if } g=a \\ 1 & \text{otherwise} \end{cases}$

Negative Log Likelihood (NLL): $J(\theta) = -\sum_{i=1}^n \log(\sigma(\theta^T x_i + \theta_0))$

Objective Func of LLC: $J(\theta, \theta_0, D) = \frac{1}{n} \sum_{i=1}^n \text{NLL}(\sigma(\theta^T x_i + \theta_0), y_i) + \lambda ||\theta||^2$

Linear Classifier:

h(x, theta, theta_0) = sign(theta^T x + theta_0)

0-1 loss: $L_{0-1}(g, a) = \begin{cases} 0 & \text{if } g=a \\ 1 & \text{otherwise} \end{cases}$

Negative Log Likelihood (NLL): $J(\theta) = -\sum_{i=1}^n \log(\sigma(\theta^T x_i + \theta_0))$

Objective Func of LLC: $J(\theta, \theta_0, D) = \frac{1}{n} \sum_{i=1}^n \text{NLL}(\sigma(\theta^T x_i + \theta_0), y_i) + \lambda ||\theta||^2$

Linear Classifier:

h(x, theta, theta_0) = sign(theta^T x + theta_0)

0-1 loss: $L_{0-1}(g, a) = \begin{cases} 0 & \text{if } g=a \\ 1 & \text{otherwise} \end{cases}$

Negative Log Likelihood (NLL): $J(\theta) = -\sum_{i=1}^n \log(\sigma(\theta^T x_i + \theta_0))$

Objective Func of LLC: $J(\theta, \theta_0, D) = \frac{1}{n} \sum_{i=1}^n \text{NLL}(\sigma(\theta^T x_i + \theta_0), y_i) + \lambda ||\theta||^2$

Linear Classifier:

h(x, theta, theta_0) = sign(theta^T x + theta_0)

0-1 loss: $L_{0-1}(g, a) = \begin{cases} 0 & \text{if } g=a \\ 1 & \text{otherwise} \end{cases}$

Negative Log Likelihood (NLL): $J(\theta) = -\sum_{i=1}^n \log(\sigma(\theta^T x_i + \theta_0))$

Objective Func of LLC: $J(\theta, \theta_0, D) = \frac{1}{n} \sum_{i=1}^n \text{NLL}(\sigma(\theta^T x_i + \theta_0), y_i) + \lambda ||\theta||^2$

Linear Classifier:

h(x, theta, theta_0) = sign(theta^T x + theta_0)

0-1 loss: $L_{0-1}(g, a) = \begin{cases} 0 & \text{if } g=a \\ 1 & \text{otherwise} \end{cases}$

Negative Log Likelihood (NLL): $J(\theta) = -\sum_{i=1}^n \log(\sigma(\theta^T x_i + \theta_0))$

Objective Func of LLC: $J(\theta, \theta_0, D) = \frac{1}{n} \sum_{i=1}^n \text{NLL}(\sigma(\theta^T x_i + \theta_0), y_i) + \lambda ||\theta||^2$

Linear Classifier:

h(x, theta, theta_0) = sign(theta^T x + theta_0)

0-1 loss: $L_{0-1}(g, a) = \begin{cases} 0 & \text{if } g=a \\ 1 & \text{otherwise} \end{cases}$

Negative Log Likelihood (NLL): $J(\theta) = -\sum_{i=1}^n \log(\sigma(\theta^T x_i + \theta_0))$

Objective Func of LLC: $J(\theta, \theta_0, D) = \frac{1}{n} \sum_{i=1}^n \text{NLL}(\sigma(\theta^T x_i + \theta_0), y_i) + \lambda ||\theta||^2$

Linear Classifier:

h(x, theta, theta_0) = sign(theta^T x + theta_0)

0-1 loss: $L_{0-1}(g, a) = \begin{cases} 0 & \text{if } g=a \\ 1 & \text{otherwise} \end{cases}$

Negative Log Likelihood (NLL): $J(\theta) = -\sum_{i=1}^n \log(\sigma(\theta^T x_i + \theta_0))$

Objective Func of LLC: $J(\theta, \theta_0, D) = \frac{1}{n} \sum_{i=1}^n \text{NLL}(\sigma(\theta^T x_i + \theta_0), y_i) + \lambda ||\theta||^2$

Linear Classifier:

h(x, theta, theta_0) = sign(theta^T x + theta_0)

0-1 loss: $L_{0-1}(g, a) = \begin{cases} 0 & \text{if } g=a \\ 1 & \text{otherwise} \end{cases}$

Negative Log Likelihood (NLL): $J(\theta) = -\sum_{i=1}^n \log(\sigma(\theta^T x_i + \theta_0))$

Objective Func of LLC: $J(\theta, \theta_0, D) = \frac{1}{n} \sum_{i=1}^n \text{NLL}(\sigma(\theta^T x_i + \theta_0), y_i) + \lambda ||\theta||^2$

Linear Classifier:

h(x, theta, theta_0) = sign(theta^T x + theta_0)

0-1 loss: $L_{0-1}(g, a) = \begin{cases} 0 & \text{if } g=a \\ 1 & \text{otherwise} \end{cases}$

Negative Log Likelihood (NLL): $J(\theta) = -\sum_{i=1}^n \log(\sigma(\theta^T x_i + \theta_0))$

Objective Func of LLC: $J(\theta, \theta_0, D) = \frac{1}{n} \sum_{i=1}^n \text{NLL}(\sigma(\theta^T x_i + \theta_0), y_i) + \lambda ||\theta||^2$

Linear Classifier:

MATRIX MATH:

Scalar x vector: $\frac{\partial y}{\partial x} = \left[\frac{\partial y_1}{\partial x_1}, \dots, \frac{\partial y_m}{\partial x_n} \right]$

Vector x vector: $\frac{\partial y}{\partial x} = \left[\frac{\partial y_1}{\partial x_1}, \dots, \frac{\partial y_m}{\partial x_n} \right]^{m \times n}$

denominator layout: $\left(\frac{\partial b}{\partial a} \right)_{ij} = \frac{\partial b_j}{\partial a_i}$

multivar. chain rule: $\nabla_{\theta} J(\theta) = \frac{\partial J}{\partial \theta} = \frac{\partial J}{\partial z} \frac{\partial z}{\partial \theta}$

EVALUATING:

LEARNING ALG: procedure that takes data on \mathcal{D} & returns hypothesis h from hypothesis class H .

VALIDATION: train on new training set then evaluate h on validation set w/ no overfitting w/ training set.

CROSS-VAL: re-use train data which can be hard to get

Cross-Validate (D, k): Find λ w/ lowest avg. error

1. divide D into k chunks D_1, D_2, \dots, D_k (of roughly equal size)

2. for $i = 1$ to k (for each candidate of λ):

3. train h_i on $D \setminus D_i$ (withholding chunk D_i as the validation set)

4. compute "test" error $E_i(h_i)$ on withheld data D_i

5. return $\frac{1}{k} \sum_{i=1}^k E_i(h_i)$ validation return $\lambda^* = \arg \min_{\lambda} E_{val}(\lambda)$

6. return final model h^* w/ D_{train} using reg. (bic solve for λ^*)

7. evaluate h^* on D_{test}

LINEAR LOGISTIC CLASSIFIER:

SIGMOID $\sigma(\cdot)$: confidence/estimated likelihood

- that x belongs to + class

- monotonic, elegant gradient (smooth)

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

approaches 1 as $z \rightarrow +\infty$

approaches 0 as $z \rightarrow -\infty$

higher z mag = more confident = push closer to 0 & 1

$\sigma'(z) = \sigma(z)(1 - \sigma(z))$

FOR $y=1$: $L = -\log(\sigma)$

$\sigma=0$ correct & confident

larger z = smaller loss

smaller z = larger loss

FOR $y=0$: $L = -\log(1-\sigma)$

$\sigma=0$ correct & confident

larger z = highest loss

smaller z = lower loss

smaller z = lower loss

smaller z = lower loss

smaller z = lower loss

smaller z = lower loss

smaller z = lower loss

smaller z = lower loss

smaller z = lower loss

smaller z = lower loss

smaller z = lower loss

smaller z = lower loss

smaller z = lower loss

smaller z = lower loss

smaller z = lower loss

smaller z = lower loss

smaller z = lower loss

smaller z = lower loss

wk 1-4

Cat tu: 6, 3, 90

when $X^T X$ nearly not invertible:

there are many solns that fit the data almost equally well.

OLS chooses a high sloping one.

Ex: $\begin{bmatrix} 1 & 1 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$

$\hat{\theta} = (d+1)x + 1$

$X = nx(d+1)$

$Y = nx$

when well defined, $\hat{\theta}$ is unique minimizer of $J(\theta)$

Vectorized MSE:

$\hat{\theta} = (X^T X)^{-1} X^T Y$

$\hat{\theta} = \frac{1}{n} X^T X^{-1} X^T Y$

$\hat{\theta} = \frac{1}{n} X^T X^{-1} X^T Y$

$\hat{\theta} = \frac{1}{n} X^T X^{-1} X^T Y$

$\hat{\theta} = \frac{1}{n} X^T X^{-1} X^T Y$

$\hat{\theta} = \frac{1}{n} X^T X^{-1} X^T Y$

$\hat{\theta} = \frac{1}{n} X^T X^{-1} X^T Y$

$\hat{\theta} = \frac{1}{n} X^T X^{-1} X^T Y$

$\hat{\theta} = \frac{1}{n} X^T X^{-1} X^T Y$

$\hat{\theta} = \frac{1}{n} X^T X^{-1} X^T Y$

$\hat{\theta} = \frac{1}{n} X^T X^{-1} X^T Y$

$\hat{\theta} = \frac{1}{n} X^T X^{-1} X^T Y$

$\hat{\theta} = \frac{1}{n} X^T X^{-1} X^T Y$

$\hat{\theta} = \frac{1}{n} X^T X^{-1} X^T Y$

$\hat{\theta} = \frac{1}{n} X^T X^{-1} X^T Y$

$\hat{\theta} = \frac{1}{n} X^T X^{-1} X^T Y$

$\hat{\theta} = \frac{1}{n} X^T X^{-1} X^T Y$

$\hat{\theta} = \frac{1}{n} X^T X^{-1} X^T Y$

$\hat{\theta} = \frac{1}{n} X^T X^{-1} X^T Y$

$\hat{\theta} = \frac{1}{n} X^T X^{-1} X^T Y$

$\hat{\theta} = \frac{1}{n} X^T X^{-1} X^T Y$

$\hat{\theta} = \frac{1}{n} X^T X^{-1} X^T Y$

$\hat{\theta} = \frac{1}{n} X^T X^{-1} X^T Y$

$\hat{\theta} = \frac{1}{n} X^T X^{-1} X^T Y$

$\hat{\theta} = \frac{1}{n} X^T X^{-1} X^T Y$

$\hat{\theta} = \frac{1}{n} X^T X^{-1} X^T Y$

$\hat{\theta} = \frac{1}{n} X^T X^{-1} X^T Y$

$\hat{\theta} = \frac{1}{n} X^T X^{-1} X^T Y$

$\hat{\theta} = \frac{1}{n} X^T X^{-1} X^T Y$

$\hat{\theta} = \frac{1}{n} X^T X^{-1} X^T Y$

$\hat{\theta} = \frac{1}{n} X^T X^{-1} X^T Y$

$\hat{\theta} = \frac{1}{n} X^T X^{-1} X^T Y$

$\hat{\theta} = \frac{1}{n} X^T X^{-1} X^T Y$

$\hat{\theta} = \frac{1}{n} X^T X^{-1} X^T Y$

$\hat{\theta} = \frac{1}{n} X^T X^{-1} X^T Y$

$\hat{\theta} = \frac{1}{n} X^T X^{-1} X^T Y$

SOFTMAX:

$P_k = \frac{e^{z_k}}{\sum_{j=1}^K e^{z_j}}$

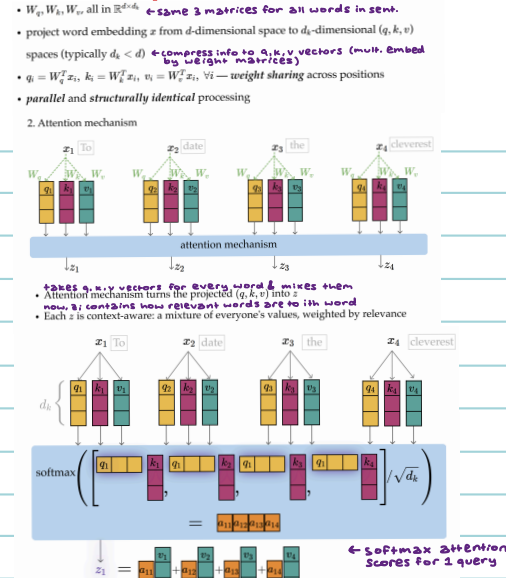
higher score = e^z

training: compared to true label w/ cross entropy loss

$L = -\sum_{k=1}^K y_k \log(P_k)$

</

TRANSFORMER ARCHITECTURE



EX: LINEAR LOGISTIC

(c) Again, assume $d = 3$. Consider a data point with $x = [10, -1, 2]^T$ and $y = +1$. A linear logistic classifier with $\theta = [1, -1, 2]^T$ and $\theta_0 = 0$ correctly predicts the label for this point. Recall that the negative log-likelihood (NLL) loss is:

$$L(y, y) = -[y \log g + (1 - y) \log (1 - g)],$$

where $g = \sigma(\theta^T x + \theta_0)$. $\rightarrow g = \sigma(2) \approx 0.88$ is close to 1: small loss
 $\rightarrow g = \sigma(-2) \approx 0.12$ is close to 0: large loss

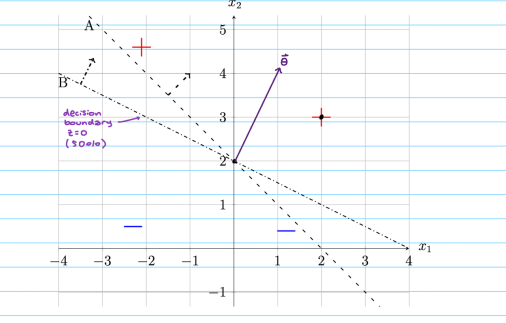
Consider logistic classifiers defined by the following parameters. Which classifier achieves the lowest NLL loss on this data point?

- i. Which classifier achieves the lowest NLL loss on this data point?
- Solution:**
- $\theta = [1, -1, 2]^T, \theta_0 = 0 \rightarrow z = 15$
 - $\theta = 10 * [1, -1, 2]^T, \theta_0 = 10 \rightarrow z = 10 - 15 + 10 = 5$
 - $\theta = 20 * [1, -1, 2]^T, \theta_0 = 20 \rightarrow z = 30$
 - $\theta = -30 * [1, -1, 2]^T, \theta_0 = -30 \rightarrow z = -90$

ii. Which classifier achieves the highest NLL loss on this data point?

- Solution:**
- $\theta = [1, -1, 2]^T, \theta_0 = 0$
 - $\theta = 10 * [1, -1, 2]^T, \theta_0 = 10$
 - $\theta = 20 * [1, -1, 2]^T, \theta_0 = 20$
 - $\theta = -30 * [1, -1, 2]^T, \theta_0 = -30$

EX: LINEAR LOGISTIC CLASSIFIER



Which of the following could be possible θ, θ_0 values for the hyperplane described by model B?

- A. $\theta = \begin{bmatrix} -2 \\ -4 \end{bmatrix}, \theta_0 = 8$
- B. $\theta = \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \theta_0 = -8$
- C. $\theta = \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \theta_0 = -4$
- D. $\theta = \begin{bmatrix} 4 \\ 2 \end{bmatrix}, \theta_0 = -8$

What value of z does model B assign to the point (2, 3)? What is the numerical probability outputted by the classifier?

$z = \theta^T x + \theta_0 = 1*2 + 2*3 - 8 = 4 - 8 = -4$

$\sigma(-4) = \frac{e^{-4}}{1 + e^{-4}} \approx 0.08$

EX: SOFTMAX

Suppose that we have a model defined by the following matrix:

$$\theta = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix}$$

Consider the data point $x = [1, -1, 1]^T$. Compute $z = \theta^T x$ and determine which class will be assigned to x .

$z = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}$

\rightarrow class 1 will be assigned to class 1.

Does the provided model defined by θ perfectly separate the data as desired in the graph above? Explain.

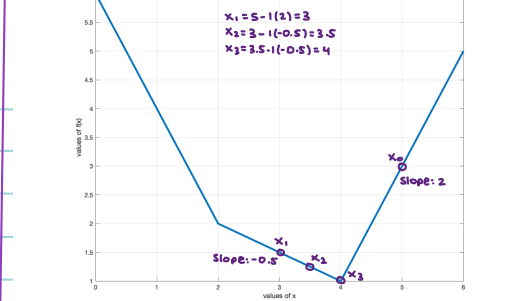
$z = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} \rightarrow \logit = 1$

\rightarrow class 1 is largest when $-x_1 - 3x_2 + x_3 > 0$

class 1: $x_1 > 0$, class 2: $-3 \leq x_2 \leq 0$

EX: 6D CALCULATION

First, John applies gradient descent to a piecewise-linear function $f: \mathbb{R} \rightarrow \mathbb{R}$ with the (partial) graph shown on the figure below:



At points $x = 2$ and $x = 4$, John uses $\nabla f(2) = -0.5$ and $\nabla f(4) = 0$.

i. Starting from the initial guess $x^{(0)} = 5$, and using step size $\eta = 1$, what will be the values of $x^{(1)}, x^{(2)}$, and $x^{(3)}$?

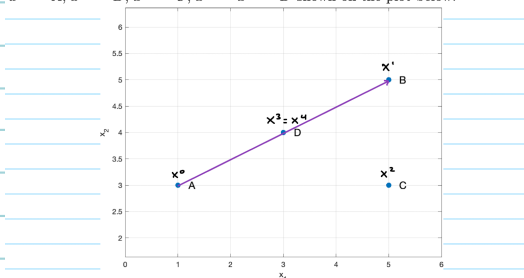
Solution: Answer: $x^{(1)} = 3, x^{(2)} = 3.5, x^{(3)} = 4$. Reasoning: by inspection of the graph, $\nabla f(x^{(0)}) = 2$, hence $x^{(1)} = 5 - 1 \cdot 2 = 3$. Next, $\nabla f(x^{(1)}) = -0.5$, hence $x^{(2)} = 3 - 1 \cdot (-0.5) = 3.5$. Finally, $\nabla f(x^{(2)}) = -0.5$, hence $x^{(3)} = 3.5 - 1 \cdot (-0.5) = 4$.

6D OSCILLATION

John discovers that, starting with $x^{(0)} = 1$, there are many values of $\eta > 0$ for which the gradient descent iterations produce oscillations of period 2 within the range (0, 6) (i.e., $x^{(k+2)} = x^{(k)} \in (0, 6)$ for all $k = 0, 1, 2, \dots$). Find all such values of η .

Solution: Answer: $\eta \in (1.5, 2.5)$. Reasoning: equality $x^{(k+2)} = x^{(k)}$ requires $\nabla f(x^{(k)}) = -\nabla f(x^{(k+1)})$. Since $\nabla f(x^{(0)}) = 2$, we need $\nabla f(x^{(1)}) = -2$, which means $x^{(1)} = x^{(0)} - \eta \nabla f(x^{(0)}) = 1 + 2\eta \in (4, 6)$. Hence $\eta \in (1.5, 2.5)$.

After mastering one-dimensional optimization, Anh applies gradient descent to a smooth function $f: \mathbb{R}^2 \rightarrow \mathbb{R}$, with $\eta = 0.1$, resulting in the sequence of points $x^{(0)} = A, x^{(1)} = B, x^{(2)} = C, x^{(3)} = D$ shown on the plot below:



i. Find $\nabla f(x^{(0)}), \nabla f(x^{(1)}), \nabla f(x^{(2)}),$ and $\nabla f(x^{(3)})$.

$\nabla f(x^{(0)}) = \begin{bmatrix} 2 \\ 0 \end{bmatrix}$

$\nabla f(x^{(1)}) = \begin{bmatrix} -2 \\ 0 \end{bmatrix}$

$\nabla f(x^{(2)}) = \begin{bmatrix} 2 \\ 0 \end{bmatrix}$

$\nabla f(x^{(3)}) = \begin{bmatrix} -2 \\ 0 \end{bmatrix}$

Consider running gradient descent (GD) to learn θ . We initialize GD with $\theta_{\text{initial}} = 4$ and use a fixed learning rate of $\eta = 0.02$. After just one iteration of the GD update, which of the following is a possible value for the resulting updated parameter θ_{new} ? Briefly justify your answer.

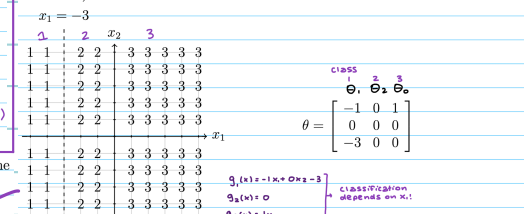
- $\theta_{\text{new}} = 4$
- $\theta_{\text{new}} = 3.76$
- $\theta_{\text{new}} = 3.44$
- $\theta_{\text{new}} = 2.56$
- $\theta_{\text{new}} = 0.88$
- Not enough information to determine any possible value of θ_{new} .

Consider running stochastic gradient descent (SGD) to learn θ . We initialize SGD with $\theta_{\text{initial}} = 4$ and use a fixed learning rate of $\eta = 0.02$. After just one iteration of the GD update, which of the following is a possible value for the resulting updated parameter θ_{new} ? Briefly justify your answer.

- $\theta_{\text{new}} = 4$
- $\theta_{\text{new}} = 3.76$
- $\theta_{\text{new}} = 3.44$
- $\theta_{\text{new}} = 2.56$
- $\theta_{\text{new}} = 0.88$
- Not enough information to determine any possible value of θ_{new} .

EX: SOFTMAX

Examine the classification problem represented by the graph below, where points are labeled with their class: 1, 2, or 3. Also given below is a model represented by the matrix θ (note that this is θ and so the first column represents θ_1, θ_2 , and θ_3 for class 1).



Write out the definition of the hypothesis visualized in the graph given above. Then, determine activation functions f^1, f^2 and weights such that $h_2(x) = \sigma(\theta_2^T(x) + \theta_0)$, where $\phi_2(x)$ is the same transformation from part (e).

$g = f^2(w^1 x_1 + w_2 x_2 + w_3 x_3 + w_0)$

$\theta_2 = \begin{bmatrix} -1 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix}$

$w^1 = \sigma, w^2 = \sigma, w^3 = \sigma, w_0 = 0$

BINARY & MULTI-CLASS CLASSIFICATION

This problem explores the relationship between binary and multi-class classification, and shows how the multi-class formulation reduces to the binary one when $K = 2$.

Consider a one-hot- K classifier where the input is $x \in \mathbb{R}^d$ and the parameters are $\theta \in \mathbb{R}^{d \times K}$. Let θ_k denote the k th column. The logits are computed as

$$z_k = \theta_k^T x,$$

and the softmax probabilities are

$$p_k = \frac{e^{z_k}}{\sum_{j=1}^K e^{z_j}}$$

(a) How many independent values are needed to specify a probability distribution over K classes? How many outputs does softmax actually produce? How do we reconcile this difference? Does binary classification encounter the same issue?

Solution: A probability distribution over K classes must sum to 1, so only $K - 1$ values are required to fully specify it. Softmax produces K outputs, which introduces one redundant degree of freedom.

In binary classification with $K = 2$, this issue does not typically arise because we use a single logit instead of two. However, we could formulate binary classification with two logits, in which case the same redundancy would appear.

We use K outputs because: (1) it provides a symmetric representation across all classes, and (2) it matches the K -dimensional one-hot encoding of class labels, making the loss computation straightforward. We could eliminate the redundancy by fixing one logit (e.g., $z_K = 0$), but the symmetric formulation is simpler in practice.

(b) Now consider $K = 2$. Suppose you train both a binary classifier (parameter $\theta \in \mathbb{R}^d$) and a 2-class softmax classifier (parameters $\theta_1, \theta_2 \in \mathbb{R}^d$) on the same dataset.

After training, the two classifiers produced the same decision boundary. Derive how θ relates to θ_1 and θ_2 .

Note: Assume class 1 in the multi-class setting corresponds to class 1 (positive) in binary.

Hint: Write p_1 in the multi-class setting and rewrite it to match $p = \sigma(z)$.

Solution: In the multi-class setting, the probability of class 1 is

$$p_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2}}$$

Divide both the numerator and denominator by e^{z_2} :

$$p_1 = \frac{1}{1 + e^{z_1 - z_2}}$$

Define $z = z_1 - z_2$. Then

$$p_1 = \frac{1}{1 + e^{-z}} = \sigma(z),$$

which is exactly the binary logistic function. This shows that both classifiers produce the same decision boundary: we predict class 1 when $p_1 > 0.5$, i.e., when $z > 0$, which is equivalent to $z_1 > z_2$. The classification depends only on the relative score between classes.

Therefore, we have

$$z = z_1 - z_2 = \theta_1^T x - \theta_2^T x = (\theta_1 - \theta_2)^T x.$$

Comparing with the binary form $z = \theta^T x$, we get

$$\theta = \theta_1 - \theta_2.$$

CNN DIMS:

EX: 3x3x3 RGB input, Conv layer: 8 filter 5x5, 10 filters, stride 1 pad 0.

Total PARAMS: each filter 5x5x3 = 75 + 1 (bias) weights per filter. 10 filters = 750 + 10 = 760 PARAMS

FULLY CONV: paramS = (input dim * output dim) + output dim

CNN DIMS:

Layer	Input Dimensions	Output Dimension
convolutional layer with 9 filters, stride 2, with zero padding	$n^k = \left\lfloor \frac{n^i + 2p - k}{s} \right\rfloor$ $32 \times 32 \times 3$ $= \left\lfloor \frac{32 + 2(1 - 1) - 1}{2} \right\rfloor = 16$	$16 \times 16 \times 9$
ReLU activation	$16 \times 16 \times 9$	$16 \times 16 \times 9$
4 x 4 max pooling, stride 4	$16 \times 16 \times 9$	$4 \times 4 \times 9$
Flattening layer	$4 \times 4 \times 9$	144×1
Fully connected layer with 8 neurons	144×1	8×1
Softmax activation	8×1	8×1

How many weight parameters are learned in the first convolutional layer?

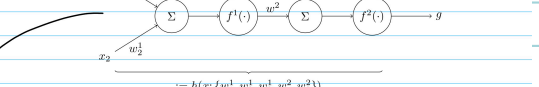
Solution: 9 filters of size $3 \times 3 \times 3$, each with a bias, leads to 252 weight parameters.

Given the size of the network output, how many image categories is this network designed to classify?

Solution: 8 neurons \rightarrow 8 classes.

EX: HYPOTHESIS (JAN)

(f) Consider the following computation graph:



Write out the definition of the hypothesis visualized in the graph given above. Then, determine activation functions f^1, f^2 and weights such that $h_2(x) = \sigma(\theta_2^T(x) + \theta_0)$, where $\phi_2(x)$ is the same transformation from part (e).

$g = f^2(w^1 x_1 + w_2 x_2 + w_3 x_3 + w_0)$

$\theta_2 = \begin{bmatrix} -1 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix}$

$w^1 = \sigma, w^2 = \sigma, w^3 = \sigma, w_0 = 0$

Does the provided model defined by θ perfectly separate the data as desired in the graph above? Explain.

$z = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} \rightarrow \logit = 1$

\rightarrow class 1 is largest when $-x_1 - 3x_2 + x_3 > 0$

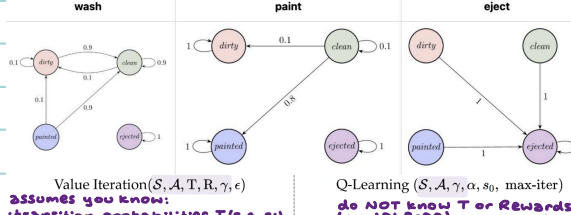
class 1: $x_1 > 0$, class 2: $-3 \leq x_2 \leq 0$

12-12 REINFORC. LEARN:
 • AGENT INTERACTS w/ ENVIRONMENT
 • TAKES ACTIONS, GET REWARDS
 • MAXIMIZE LONG-TERM REWARDS
 • TRANSITIONS & REW. ARE UNKNOWN!

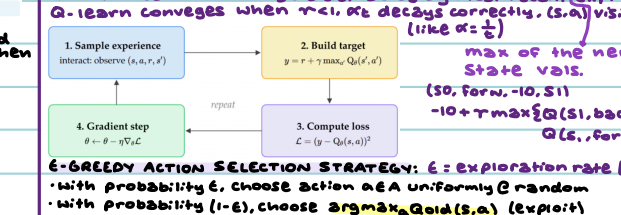
MARKOV DECISION PROCESS: (S, A, T, R, γ, ε)
 • S: state space (current situation)
 • A: action space (possible actions)
 • T(S, a, S'): probability of S to S' taking a
 • R(S, a): reward for S taking a
 • γ ∈ [0, 1]: discount factor, how much you care for future rewards
 • π(S): policy (strategy), deterministic π(a) or probabilistic π(s|a)

VALUE FUNC: how good is a state?
 → long term sum calc. expectation
 $V_{\pi}^{\pi}(s) = R(s, \pi(s)) + \gamma \sum_{s'} T(s, \pi(s), s') V_{\pi}^{\pi}(s')$
 current reward + discounted future reward
 sum of future values weighted by the probability of reaching that next state s'

GEOMETRIC SERIES: $\sum_{i=0}^{\infty} a(r)^{i-1} = \frac{a}{1-r}$
 α = 0: Q-table never updates
 TABULAR Q-LEARNING: store values in table
 Q-learning approx. Q(S, a) w/ func. Q(S, a) w/ func. Q(S, a) w/ func.
 Q-new(S, a) = (1-α)Q-old(S, a) + α(r + γ max_{a'} Q-old(S', a'))



Q-VALUES: Expected discounted reward
 Starting from state s with action a then following policy π Bellman Recursion:
 $V_{\pi}^{\pi}(s) = R(s, \pi_h(s)) + \gamma \sum_{s'} T(s, \pi_h(s), s') V_{\pi}^{\pi}(s')$
 do what a says, then back to π for future
 $Q_{\pi}^{\pi}(s, a) = R(s, a) + \gamma \sum_{s'} T(s, a, s') V_{\pi}^{\pi}(s')$



Value Iteration (S, A, T, R, γ, ε)
 ASSUMES you know:
 • transition probabilities T(S, a, S')
 • reward func. R(S, a)

Bellman:
 for s ∈ S, a ∈ A:
 1. for a ∈ S, a ∈ A:
 2. Q(s, a) ← 0
 3. while True:
 4. for s' ∈ S, a' ∈ A:
 5. Q(s', a') ← R(s', a') + γ max_a Q(s', a')
 6. if max_{s, a} |Q(s, a) - Q(s, a)| < ε:
 7. return Q(s, a)
 8. Q(s, a) ← Q(s, a)

OPTIMAL POLICY π*: policy that yields highest value V_{π*} from every state (may not be unique)
 $V_{\pi^*}(s) = \max_a [R(s, a) + \gamma \sum_{s'} T(s, a, s') V_{\pi^*}(s')]$
 $Q_{\pi^*}(s, a) = \max_{a'} [R(s, a) + \gamma \sum_{s'} T(s, a, s') V_{\pi^*}(s')]$

ε-GREEDY ACTION SELECTION STRATEGY: ε = exploration rate [0, 1]
 • with probability ε, choose action a ∈ A uniformly @ random
 • with probability (1-ε), choose argmax_a Q(s, a) (exploit)
 • start: ↑ ε (more exploratory), but later ↓ ε (more greedy)
 • trades off exploration vs. exploitation
 POLICY GRADIENTS: param. π with θ, run π_θ, update θ s.t. → no Q-table needed
 → "good" trajectories are more likely to happen
 → learns policy directly

Q-Learning Algorithm for Finding Optimal Policy
 Q-learning: Begin in state s and then repeat:

- Take action a, selected random with probability ε and greedy with probability 1-ε.
- Observe the reward r and next state s'.
- Update Q(s, a) with learning rate α to merge old belief and new target.
 $Q(s, a) \leftarrow (1-\alpha)Q(s, a) + \alpha(r + \gamma \max_{a'} Q(s', a'))$
 ← this one is ASSUMING you know dist. calc. → larger dominates
 $= Q(s, a) + \alpha((r + \gamma \max_{a'} Q(s', a')) - Q(s, a))$
- Build the policy
 $\pi_h^*(s) = \arg \max_a Q_h^*(s, a)$

Define the optimal state-action value functions Q_{π*}^*(s, a):
 the expected sum of discounted rewards, obtained by
 • starting in state s
 • take action a, for one step
 • act optimally thereafter for the remaining (h - 1) steps

Q* satisfies:
 $Q_{\pi^*}^*(s, a) = R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q_{\pi^*}^*(s', a')$
 $V_{\pi^*}^{\pi^*}(s) = \max_a Q_{\pi^*}^*(s, a)$ $V_{\pi^*}^{\pi^*}(s) = \max_a Q_{\pi^*}^*(s, a)$
 $Q_{\pi^*}^*(s, a) = R(s, a) + \gamma \sum_{s'} T(s, a, s') V_{\pi^*}^{\pi^*}(s')$
 $\pi_{\pi^*}^*(s) = \arg \max_a Q_{\pi^*}^*(s, a)$
 $\sum_s \sum_a T(s, a, s') = \sum_s T(s, s' | S) = 1$ ← trans. prob. → stochastic

K-MEANS CLUSTERING: (unsupervised)
 Structural discovery w/o labels
 1. choose k
 2. initialize k cluster centers
 3. assign each pt to nearest center
 4. recompute centers ← based on pts in its cluster
 5. repeat
 minimizes $\sum_i \|x_i - \mu_c\|^2$ (total sq. dist. bet pts & clusters)
 → sensitive to initialization and k larger vals. dominate
 → relative scale of features matters b/c use distance
 → return groupings (neither regression nor classification)

NON-PARAMETRIC MODELS: no fixed functional form for h.
 Complexity grows with data
 → interpretability, insight, speed, adaptivity

k-Means Clustering
 Iterate until convergence or some max iteration
 Assign each data point to the cluster of its closest centroid
 Update each centroid to be the average of the data points assigned to it
 Stop when converged or we've reached max iteration

CNN EX:

Layer	Input Dimensions	Output Dimension
convolutional layer with 9 filters, stride 2, with zero padding	32x32x3	16x16x9 Z ¹
ReLU activation	16x16x9	16x16x9 A ¹
4 x 4 max pooling, stride 4	16x16x9	4 x 4 x 9 A ²
Flattening layer	4 x 4 x 9	144 x 1 A ³
Fully connected layer with 8 neurons	144 x 1	8 x 1 Z ⁴
Softmax activation	8 x 1	8 x 1 9

K-NEAREST NEIGHBOR: (supervised)
 Training: none. Just memorize the training set.
 Small k = var (overfit)
 large k = bias (underfit)
 Predicting: for a new x_{new}, take the majority (class) or mean (regression) label of its k nearest neighbors.
 → sensitive to feat. scale b/c uses dist. calc. → larger dominates
 Parameters learned: the entire training set (its features and labels).
 becomes param. (caution when:
 • large n and dimensionality
 • irrelevant/noisy features
 • features on diff. scales
 Hyperparameters:
 • k: number of neighbors considered.
 • distance metric (typically Euclidean or Manhattan).
 • tie-breaking scheme (typically at random).
 Consider a one-dimensional dataset with only one feature x. The points are in two classes given by the value of y^{(j)}.

ENTROPY:
 Weighted Average Entropy (WAE)
 Set of points in leaf m:
 Fraction of points in leaf m with label k
 Entropy of leaf m:
 WAE of the split:
 entropy H := -∑_{class} P_c (log_2 P_c)
 K-CALC. L & R H, then: c iterates over 3 classes

#weight params:
 Solution: 9 filters of size 3 x 3 x 3, each with a bias, leads to 252 weight parameters.
 During training, you are running stochastic gradient descent to update the weights of the network using backpropagation. Suppose that you are using NLL Multiclass Loss L. Derive the backpropagation rule for the weight parameters W^1 of the convolutional layer. Explicitly write the partial derivative factors according to the chain rule using variable names from problem (a).
 Solution:
 $\frac{\partial \mathcal{L}}{\partial W^1} = \frac{\partial \mathcal{L}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^1} \frac{\partial Z^1}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^1} \frac{\partial A^1}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^2} \frac{\partial Z^2}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^2} \frac{\partial A^2}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^3} \frac{\partial Z^3}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^3} \frac{\partial A^3}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^4} \frac{\partial Z^4}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^4} \frac{\partial A^4}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^5} \frac{\partial Z^5}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^5} \frac{\partial A^5}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^6} \frac{\partial Z^6}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^6} \frac{\partial A^6}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^7} \frac{\partial Z^7}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^7} \frac{\partial A^7}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^8} \frac{\partial Z^8}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^8} \frac{\partial A^8}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^9} \frac{\partial Z^9}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^9} \frac{\partial A^9}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{10}} \frac{\partial Z^{10}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{10}} \frac{\partial A^{10}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{11}} \frac{\partial Z^{11}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{11}} \frac{\partial A^{11}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{12}} \frac{\partial Z^{12}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{12}} \frac{\partial A^{12}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{13}} \frac{\partial Z^{13}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{13}} \frac{\partial A^{13}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{14}} \frac{\partial Z^{14}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{14}} \frac{\partial A^{14}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{15}} \frac{\partial Z^{15}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{15}} \frac{\partial A^{15}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{16}} \frac{\partial Z^{16}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{16}} \frac{\partial A^{16}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{17}} \frac{\partial Z^{17}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{17}} \frac{\partial A^{17}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{18}} \frac{\partial Z^{18}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{18}} \frac{\partial A^{18}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{19}} \frac{\partial Z^{19}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{19}} \frac{\partial A^{19}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{20}} \frac{\partial Z^{20}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{20}} \frac{\partial A^{20}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{21}} \frac{\partial Z^{21}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{21}} \frac{\partial A^{21}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{22}} \frac{\partial Z^{22}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{22}} \frac{\partial A^{22}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{23}} \frac{\partial Z^{23}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{23}} \frac{\partial A^{23}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{24}} \frac{\partial Z^{24}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{24}} \frac{\partial A^{24}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{25}} \frac{\partial Z^{25}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{25}} \frac{\partial A^{25}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{26}} \frac{\partial Z^{26}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{26}} \frac{\partial A^{26}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{27}} \frac{\partial Z^{27}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{27}} \frac{\partial A^{27}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{28}} \frac{\partial Z^{28}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{28}} \frac{\partial A^{28}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{29}} \frac{\partial Z^{29}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{29}} \frac{\partial A^{29}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{30}} \frac{\partial Z^{30}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{30}} \frac{\partial A^{30}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{31}} \frac{\partial Z^{31}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{31}} \frac{\partial A^{31}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{32}} \frac{\partial Z^{32}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{32}} \frac{\partial A^{32}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{33}} \frac{\partial Z^{33}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{33}} \frac{\partial A^{33}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{34}} \frac{\partial Z^{34}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{34}} \frac{\partial A^{34}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{35}} \frac{\partial Z^{35}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{35}} \frac{\partial A^{35}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{36}} \frac{\partial Z^{36}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{36}} \frac{\partial A^{36}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{37}} \frac{\partial Z^{37}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{37}} \frac{\partial A^{37}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{38}} \frac{\partial Z^{38}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{38}} \frac{\partial A^{38}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{39}} \frac{\partial Z^{39}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{39}} \frac{\partial A^{39}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{40}} \frac{\partial Z^{40}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{40}} \frac{\partial A^{40}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{41}} \frac{\partial Z^{41}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{41}} \frac{\partial A^{41}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{42}} \frac{\partial Z^{42}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{42}} \frac{\partial A^{42}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{43}} \frac{\partial Z^{43}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{43}} \frac{\partial A^{43}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{44}} \frac{\partial Z^{44}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{44}} \frac{\partial A^{44}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{45}} \frac{\partial Z^{45}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{45}} \frac{\partial A^{45}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{46}} \frac{\partial Z^{46}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{46}} \frac{\partial A^{46}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{47}} \frac{\partial Z^{47}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{47}} \frac{\partial A^{47}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{48}} \frac{\partial Z^{48}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{48}} \frac{\partial A^{48}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{49}} \frac{\partial Z^{49}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{49}} \frac{\partial A^{49}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{50}} \frac{\partial Z^{50}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{50}} \frac{\partial A^{50}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{51}} \frac{\partial Z^{51}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{51}} \frac{\partial A^{51}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{52}} \frac{\partial Z^{52}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{52}} \frac{\partial A^{52}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{53}} \frac{\partial Z^{53}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{53}} \frac{\partial A^{53}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{54}} \frac{\partial Z^{54}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{54}} \frac{\partial A^{54}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{55}} \frac{\partial Z^{55}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{55}} \frac{\partial A^{55}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{56}} \frac{\partial Z^{56}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{56}} \frac{\partial A^{56}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{57}} \frac{\partial Z^{57}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{57}} \frac{\partial A^{57}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{58}} \frac{\partial Z^{58}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{58}} \frac{\partial A^{58}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{59}} \frac{\partial Z^{59}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{59}} \frac{\partial A^{59}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{60}} \frac{\partial Z^{60}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{60}} \frac{\partial A^{60}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{61}} \frac{\partial Z^{61}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{61}} \frac{\partial A^{61}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{62}} \frac{\partial Z^{62}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{62}} \frac{\partial A^{62}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{63}} \frac{\partial Z^{63}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{63}} \frac{\partial A^{63}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{64}} \frac{\partial Z^{64}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{64}} \frac{\partial A^{64}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{65}} \frac{\partial Z^{65}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{65}} \frac{\partial A^{65}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{66}} \frac{\partial Z^{66}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{66}} \frac{\partial A^{66}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{67}} \frac{\partial Z^{67}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{67}} \frac{\partial A^{67}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{68}} \frac{\partial Z^{68}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{68}} \frac{\partial A^{68}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{69}} \frac{\partial Z^{69}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{69}} \frac{\partial A^{69}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{70}} \frac{\partial Z^{70}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{70}} \frac{\partial A^{70}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{71}} \frac{\partial Z^{71}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{71}} \frac{\partial A^{71}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{72}} \frac{\partial Z^{72}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{72}} \frac{\partial A^{72}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{73}} \frac{\partial Z^{73}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{73}} \frac{\partial A^{73}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{74}} \frac{\partial Z^{74}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{74}} \frac{\partial A^{74}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{75}} \frac{\partial Z^{75}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{75}} \frac{\partial A^{75}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{76}} \frac{\partial Z^{76}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{76}} \frac{\partial A^{76}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{77}} \frac{\partial Z^{77}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{77}} \frac{\partial A^{77}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{78}} \frac{\partial Z^{78}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{78}} \frac{\partial A^{78}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{79}} \frac{\partial Z^{79}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{79}} \frac{\partial A^{79}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{80}} \frac{\partial Z^{80}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{80}} \frac{\partial A^{80}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{81}} \frac{\partial Z^{81}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{81}} \frac{\partial A^{81}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{82}} \frac{\partial Z^{82}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{82}} \frac{\partial A^{82}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{83}} \frac{\partial Z^{83}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{83}} \frac{\partial A^{83}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{84}} \frac{\partial Z^{84}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{84}} \frac{\partial A^{84}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{85}} \frac{\partial Z^{85}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{85}} \frac{\partial A^{85}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{86}} \frac{\partial Z^{86}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{86}} \frac{\partial A^{86}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{87}} \frac{\partial Z^{87}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{87}} \frac{\partial A^{87}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{88}} \frac{\partial Z^{88}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{88}} \frac{\partial A^{88}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{89}} \frac{\partial Z^{89}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{89}} \frac{\partial A^{89}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{90}} \frac{\partial Z^{90}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{90}} \frac{\partial A^{90}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{91}} \frac{\partial Z^{91}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{91}} \frac{\partial A^{91}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{92}} \frac{\partial Z^{92}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{92}} \frac{\partial A^{92}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{93}} \frac{\partial Z^{93}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{93}} \frac{\partial A^{93}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{94}} \frac{\partial Z^{94}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{94}} \frac{\partial A^{94}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{95}} \frac{\partial Z^{95}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{95}} \frac{\partial A^{95}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{96}} \frac{\partial Z^{96}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{96}} \frac{\partial A^{96}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{97}} \frac{\partial Z^{97}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{97}} \frac{\partial A^{97}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{98}} \frac{\partial Z^{98}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{98}} \frac{\partial A^{98}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{99}} \frac{\partial Z^{99}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{99}} \frac{\partial A^{99}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{100}} \frac{\partial Z^{100}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{100}} \frac{\partial A^{100}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{101}} \frac{\partial Z^{101}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{101}} \frac{\partial A^{101}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{102}} \frac{\partial Z^{102}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{102}} \frac{\partial A^{102}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{103}} \frac{\partial Z^{103}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{103}} \frac{\partial A^{103}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{104}} \frac{\partial Z^{104}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{104}} \frac{\partial A^{104}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{105}} \frac{\partial Z^{105}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{105}} \frac{\partial A^{105}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{106}} \frac{\partial Z^{106}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{106}} \frac{\partial A^{106}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{107}} \frac{\partial Z^{107}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{107}} \frac{\partial A^{107}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{108}} \frac{\partial Z^{108}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{108}} \frac{\partial A^{108}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{109}} \frac{\partial Z^{109}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{109}} \frac{\partial A^{109}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{110}} \frac{\partial Z^{110}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{110}} \frac{\partial A^{110}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{111}} \frac{\partial Z^{111}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{111}} \frac{\partial A^{111}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{112}} \frac{\partial Z^{112}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{112}} \frac{\partial A^{112}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{113}} \frac{\partial Z^{113}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{113}} \frac{\partial A^{113}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{114}} \frac{\partial Z^{114}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{114}} \frac{\partial A^{114}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{115}} \frac{\partial Z^{115}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{115}} \frac{\partial A^{115}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{116}} \frac{\partial Z^{116}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{116}} \frac{\partial A^{116}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{117}} \frac{\partial Z^{117}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{117}} \frac{\partial A^{117}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{118}} \frac{\partial Z^{118}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{118}} \frac{\partial A^{118}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{119}} \frac{\partial Z^{119}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{119}} \frac{\partial A^{119}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{120}} \frac{\partial Z^{120}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{120}} \frac{\partial A^{120}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{121}} \frac{\partial Z^{121}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{121}} \frac{\partial A^{121}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{122}} \frac{\partial Z^{122}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{122}} \frac{\partial A^{122}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{123}} \frac{\partial Z^{123}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{123}} \frac{\partial A^{123}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{124}} \frac{\partial Z^{124}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{124}} \frac{\partial A^{124}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{125}} \frac{\partial Z^{125}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{125}} \frac{\partial A^{125}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{126}} \frac{\partial Z^{126}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{126}} \frac{\partial A^{126}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{127}} \frac{\partial Z^{127}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{127}} \frac{\partial A^{127}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{128}} \frac{\partial Z^{128}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{128}} \frac{\partial A^{128}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{129}} \frac{\partial Z^{129}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{129}} \frac{\partial A^{129}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{130}} \frac{\partial Z^{130}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{130}} \frac{\partial A^{130}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{131}} \frac{\partial Z^{131}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{131}} \frac{\partial A^{131}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{132}} \frac{\partial Z^{132}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{132}} \frac{\partial A^{132}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{133}} \frac{\partial Z^{133}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{133}} \frac{\partial A^{133}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial Z^{134}} \frac{\partial Z^{134}}{\partial W^1} \frac{\partial \mathcal{L}}{\partial A^{134}} \frac{\partial A^{13$



A sheet of white paper with a vertical red margin line on the left side and horizontal blue lines for writing. There are three circular punch holes along the left edge.