# Solution 1

## Solution 1 (a)

### Step 1: Identify characteristics of the dataspace

We are given 10 dimensional vectors where each element can be any real number ($x_i \in \mathbb{R}$):

$\therefore$ we can express the dataspace $\chi$ as: $\chi = \mathbb{R}^{10}$

## Solution 1 (b)

### Step 1: Identify characteristics of the dataspace

We are given 3 dimensional vectors where each element is zero or one ($x_i \in [0, 1]$):

$\therefore$ we can express the dataspace $\chi$ as: $\chi = [0, 1]^3$

## Solution 2

### Solution 2 (a)

**Step 1: Define Euclidean distance ($\ell_2$)**

$$\ell_2 = \|p - q\|_2 = \sqrt{\sum_{i=1}^{n}(p_i - q_i)^2}$$

**Step 2: Compute $\ell_2$**

Let $p = 1$ and $q = 10$

$$\ell_2 = \sqrt{\sum_{i=1}^{n}(p_i - q_i)^2}$$

$$\ell_2 = \sqrt{\sum_{i=1}^{1}(1 - 10)^2}$$

$$\ell_2 = \sqrt{(-9)^2}$$

$$\ell_2 = 9$$

$\therefore \ell_2 = 9$

### Solution 2 (b)

**Step 1: Define Euclidean distance ($\ell_2$)**

$$\ell_2 = \|p - q\|_2 = \sqrt{\sum_{i=1}^{n}(p_i - q_i)^2}$$

**Step 2: Compute $\ell_2$**

Let $p = \begin{bmatrix} -1 \\ 12 \end{bmatrix}, q = \begin{bmatrix} 6 \\ -12 \end{bmatrix}$

$$\ell_2 = \sqrt{\sum_{i=1}^{n}(p_i - q_i)^2}$$

$$\ell_2 = \sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2}$$

$$\ell_2 = \sqrt{(-1 - 6)^2 + (12 - (-12))^2}$$

$$\ell_2 = \sqrt{(-7)^2 + (24)^2}$$

$$\ell_2 = \sqrt{625}$$

$$\ell_2 = 25$$

$\therefore \ell_2 = 25$

## Solution 2

### Solution 2 (c)

**Step 1: Define Euclidean distance $(\ell_2)$**

$$\ell_2 = \|p - q\|_2 = \sqrt{\sum_{i=1}^{n}(p_i - q_i)^2}$$

**Step 2: Compute $\ell_2$**

Let $p = \begin{bmatrix} 1 \\ 5 \\ -1 \end{bmatrix}, q = \begin{bmatrix} 5 \\ 2 \\ 11 \end{bmatrix}$

$$\ell_2 = \sqrt{\sum_{i=1}^{n}(p_i - q_i)^2}$$
$$\ell_2 = \sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2 + (p_3 - q_3)^2}$$
$$\ell_2 = \sqrt{(1 - 5)^2 + (5 - 2)^2 + (-1 - 11)^2}$$
$$\ell_2 = \sqrt{(-4)^2 + (3)^2 + (-12)^2}$$
$$\ell_2 = \sqrt{169}$$
$$\ell_2 = 13$$

$\therefore \ell_2 = 13$

## Solution 3

## Solution 3 (a)

**Step 1: Normalize the vector $x$**

Let $x = \begin{bmatrix} 10 \\ 15 \\ 25 \end{bmatrix}$

$$\sum_{i=1}^{3} x_i = x_1 + x_2 + x_3 = 10 + 15 + 25 = 50$$

Now, divide each entry by the total sum:

$$p = \frac{1}{50} \cdot x = \frac{1}{50} \begin{bmatrix} 10 \\ 15 \\ 25 \end{bmatrix} = \begin{bmatrix} 10/50 \\ 15/50 \\ 25/50 \end{bmatrix} = \begin{bmatrix} 0.2 \\ 0.3 \\ 0.5 \end{bmatrix}$$

$\therefore$ the result $(p)$ of scaling vertor $x$ is the following:

$$p = \begin{bmatrix} 0.2 \\ 0.3 \\ 0.5 \end{bmatrix}$$

## Solution 3 (b)

**Step 1: Define dimension of the probability simplex**

The dimension of vector $p$ is 3 and $k = n - 1$ where $k$ is the dimension of the probability simplex

$\therefore$ vector $p$ lies in the probability simplex($\Delta_2$) for $k = 2$

## Solution 4

### Step 1: Understand what scaling means and the constraints of $\Delta_2$

Scaling a vector means multiplying all entries by the same positive constant $c > 0$.
The probability simplex $\Delta_2$ is defined as:

$$\Delta_2 = \left\{ x \in \mathbb{R}^3 : x_i \geq 0 \text{ for all } i = 1, 2, 3, \text{ and } \sum_{i=1}^{3} x_i = 1 \right\}$$

For a point to be scalable to $\Delta_2$, after scaling it must satisfy:

- All components must be non-negative

- The sum of components must equal 1

### Step 2: Find a point that violates the constraints

Since we need a 2-dimensional point, let's consider $x = \begin{bmatrix} a \\ b \end{bmatrix}$ where we interpret this as a 3-dimensional

vector $\begin{bmatrix} a \\ b \\ 0 \end{bmatrix}$.

For any point with a negative component, scaling cannot make it non-negative.

Example: Let $x = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$

If we scale by any $c > 0$: $cx = \begin{bmatrix} c \\ -2c \end{bmatrix}$

The second component $-2c < 0$ for any $c > 0$, so this scaled vector cannot satisfy the non-negativity constraint of $\Delta_2$.

$\therefore$ Final Answer Example: $x = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$ (or any 2D point with at least one negative component)

Graduate Level Explanation The key insight is that scaling preserves the sign of each component - it cannot transform negative values to positive ones. Since the probability simplex requires all components to be non-negative, any vector containing negative components cannot be scaled to lie within it. This geometric constraint is fundamental in applications like topic modeling or mixture models, where negative probabilities are meaningless. The orthant restriction (non-negativity) combined with the sum constraint defines the simplex as a bounded convex polytope, and scaling operations represent rays from the origin that either intersect this polytope or miss it entirely.

Explanation for 5 year old Imagine you have a recipe with ingredients, but one of the "amounts" is negative - like "negative 2 cups of flour." No matter how much you shrink or grow the recipe (scaling), you'll still need a negative amount of flour, which doesn't make sense! The probability simplex is like a rule that says "all ingredients must be positive amounts," so any recipe with negative ingredients can never follow this rule, no matter how you scale it.

## Solution 5

**Step 1: Define $\Delta_3$ and identify key properties**

The probability simplex $\Delta_3$ is defined as:

$$\Delta_3 = \left\{ x \in \mathbb{R}^4 : x_i \geq 0 \text{ for all } i = 1, 2, 3, 4, \text{ and } \sum_{i=1}^{4} x_i = 1 \right\}$$

This is a 3-dimensional simplex (tetrahedron) embedded in 4-dimensional space. However, due to the constraint $\sum x_i = 1$, we can visualize it in 3D by using three coordinates and letting the fourth be determined by the constraint.

The given points are the vertices of the simplex:

- $\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$ corresponds to $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ in 4D

- $\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$ corresponds to $\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$ in 4D

- $\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$ corresponds to $\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$ in 4D

**Step 2: Determine coordinates and sketch description**

The fourth vertex (not shown) would be $\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$, which corresponds to $\begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$ in our 3D representation (since the fourth coordinate is $1 - 0 - 0 - 0 = 1$).

The most central point (centroid) of the simplex is:

$$\text{centroid} = \frac{1}{4} \left( \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix} \right) = \begin{bmatrix} 1/4 \\ 1/4 \\ 1/4 \\ 1/4 \end{bmatrix}$$

In 3D coordinates, this is $\begin{bmatrix} 1/4 \\ 1/4 \\ 1/4 \end{bmatrix}$.
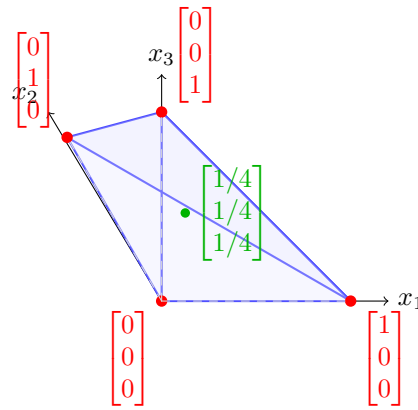
**Sketch Description:**

- Draw 3D coordinate axes labeled $x_1$, $x_2$, $x_3$

- The simplex is a triangular region (actually a tetrahedron) with vertices at:

  - $(1, 0, 0)$ on the $x_1$-axis
  - $(0, 1, 0)$ on the $x_2$-axis
  - $(0, 0, 1)$ on the $x_3$-axis
  - $(0, 0, 0)$ at the origin (representing the fourth vertex)

- Connect these four points to form a tetrahedron

- Mark the centroid at $(1/4, 1/4, 1/4)$

$\therefore$ Final Answer The most central point in $\Delta_3$ has coordinates: $\begin{bmatrix} 1/4 \\ 1/4 \\ 1/4 \end{bmatrix}$ (or $\begin{bmatrix} 0.25 \\ 0.25 \\ 0.25 \end{bmatrix}$)

Note: The fourth coordinate is $1 - 1/4 - 1/4 - 1/4 = 1/4$

Graduate Level Explanation The 3-simplex $\Delta_3$ is a regular tetrahedron when embedded properly in 3D space. Each vertex represents a pure probability distribution (all probability mass on one outcome), while interior points represent mixed distributions. The centroid represents the uniform distribution over 4 outcomes. This geometric structure is fundamental in Bayesian statistics, where it represents the space of all possible probability distributions over 4 categories. The simplex's convex hull property ensures that convex combinations of probability distributions remain valid probability distributions, making it a natural space for optimization in machine learning algorithms like EM or variational inference.

Explanation for 5 year old Imagine a pyramid made of triangles (like the ones in Egypt, but pointier)! Each corner of the pyramid represents putting all your marbles in one basket - like having 100% chocolate ice cream, or 100% vanilla, or 100% strawberry, or 100% mint. The very center of the pyramid is where you have equal amounts of all four flavors - 25% of each! Any point inside the pyramid represents some mix of the four flavors that adds up to 100%.

Figure 1: The probability simplex $\Delta_3$ showing vertices and centroid

## Solution 6

**Step 1: Recall the formulas for $\ell_1$ distance and KL divergence**

IDK.. LOG IS BASE e NOT 2 The $\ell_1$ distance (Manhattan distance) between two vectors $u$ and $v$ is:

$$\|u - v\|_1 = \sum_{i=1}^{n} |u_i - v_i|$$

The KL divergence between two probability distributions $P$ and $Q$ is:

$$KL(P\|Q) = \sum_{i=1}^{n} P_i \log\left(\frac{P_i}{Q_i}\right)$$

where we use the convention that $0\log(0/Q_i) = 0$ and $P_i\log(P_i/0) = \infty$ if $P_i > 0$.

**Step 2: Calculate each distance and divergence**

**Part (i):** $\|p - q\|_1$ where $p = \begin{bmatrix} 1/2 \\ 1/4 \\ 1/8 \\ 1/8 \end{bmatrix}$ and $q = \begin{bmatrix} 1/4 \\ 1/4 \\ 1/4 \\ 1/4 \end{bmatrix}$

$$p - q = \begin{bmatrix} 1/2 - 1/4 \\ 1/4 - 1/4 \\ 1/8 - 1/4 \\ 1/8 - 1/4 \end{bmatrix} = \begin{bmatrix} 1/4 \\ 0 \\ -1/8 \\ -1/8 \end{bmatrix}$$

$$\|p - q\|_1 = |1/4| + |0| + |-1/8| + |-1/8| = 1/4 + 0 + 1/8 + 1/8 = 1/4 + 1/4 = 1/2$$

**Part (ii):** $\|q - r\|_1$ where $q = \begin{bmatrix} 1/4 \\ 1/4 \\ 1/4 \\ 1/4 \end{bmatrix}$ and $r = \begin{bmatrix} 1/2 \\ 0 \\ 1/4 \\ 1/4 \end{bmatrix}$

$$q - r = \begin{bmatrix} 1/4 - 1/2 \\ 1/4 - 0 \\ 1/4 - 1/4 \\ 1/4 - 1/4 \end{bmatrix} = \begin{bmatrix} -1/4 \\ 1/4 \\ 0 \\ 0 \end{bmatrix}$$

$$\|q - r\|_1 = |-1/4| + |1/4| + |0| + |0| = 1/4 + 1/4 + 0 + 0 = 1/2$$

**Part (iii):** $KL(p\|q)$

$$KL(p\|q) = \sum_{i=1}^{4} p_i \log\left(\frac{p_i}{q_i}\right)$$

$$= \frac{1}{2}\log\left(\frac{1/2}{1/4}\right) + \frac{1}{4}\log\left(\frac{1/4}{1/4}\right) + \frac{1}{8}\log\left(\frac{1/8}{1/4}\right) + \frac{1}{8}\log\left(\frac{1/8}{1/4}\right)$$

$$= \frac{1}{2}\log(2) + \frac{1}{4}\log(1) + \frac{1}{8}\log(1/2) + \frac{1}{8}\log(1/2)$$

$$= \frac{1}{2}\log(2) + 0 + \frac{1}{8}(-\log(2)) + \frac{1}{8}(-\log(2))$$

$$= \frac{1}{2}\log(2) - \frac{1}{4}\log(2) = \frac{1}{4}\log(2)$$

**Part (iv):** $KL(q\|r)$

$$KL(q\|r) = \sum_{i=1}^{4} q_i \log\left(\frac{q_i}{r_i}\right)$$

$$= \frac{1}{4}\log\left(\frac{1/4}{1/2}\right) + \frac{1}{4}\log\left(\frac{1/4}{0}\right) + \frac{1}{4}\log\left(\frac{1/4}{1/4}\right) + \frac{1}{4}\log\left(\frac{1/4}{1/4}\right)$$

Since $r_2 = 0$ and $q_2 = 1/4 > 0$, we have $\log(q_2/r_2) = \log(1/4/0) = +\infty$.
Therefore: $KL(q\|r) = +\infty$

$\therefore$ Final Answers

1. $\|p - q\|_1 = \frac{1}{2}$

2. $\|q - r\|_1 = \frac{1}{2}$

3. $KL(p\|q) = \frac{1}{4}\log(2) \approx 0.173$

4. $KL(q\|r) = +\infty$

Graduate Level Explanation The $\ell_1$ distance is symmetric and satisfies the triangle inequality, making it a proper metric on the probability simplex. Interestingly, both pairs have the same $\ell_1$ distance despite having different structures. The KL divergence, however, is asymmetric and not a true metric. $KL(p\|q)$ is finite because $q$ has full support (no zero entries), but $KL(q\|r) = \infty$ because $r$ has a zero entry where $q$ has positive probability. This illustrates a fundamental property: KL divergence from a distribution with full support to one with restricted support is infinite, making it useful for detecting when probability mass is assigned to impossible events in the reference distribution.

Explanation for 5 year old The $\ell_1$ distance is like counting how much you need to move marbles between jars to make them the same - it's always the same no matter which direction you go. But KL divergence is like asking "how surprised would I be?" If one jar is completely empty but you expected marbles there, you'd be infinitely surprised! That's why one answer is infinity - it's like expecting something that's impossible.

## Solution 7

**Part a: Dimensionality for each of the representations (raw pixel, HoG, VGG-last-fc, VGG-last-conv)**

| Feature Type | Dimensionality |
|---|---|
| Raw Pixel | 3072 |
| HoG | 512 |
| VGG-last-fc | 4096 |
| VGG-last-conv | 512 |

**Part b: Test accuracies for 1-nearest neighbor classification using the various representations (raw pixel, HoG, VGG-last-fc, VGG-last-conv, random-VGG-last-fc, random-VGG-last-conv).**

| Feature Type | 1-NN test accuracy (%) |
|---|---|
| Raw Pixel | 35.4 |
| HoG | 36.6 |
| VGG-last-fc | 92.1 |
| VGG-last-conv | 92.0 |
| random VGG-last-fc | 39.1 |
| random VGG-last-conv | 40.6 |

## Solution 7

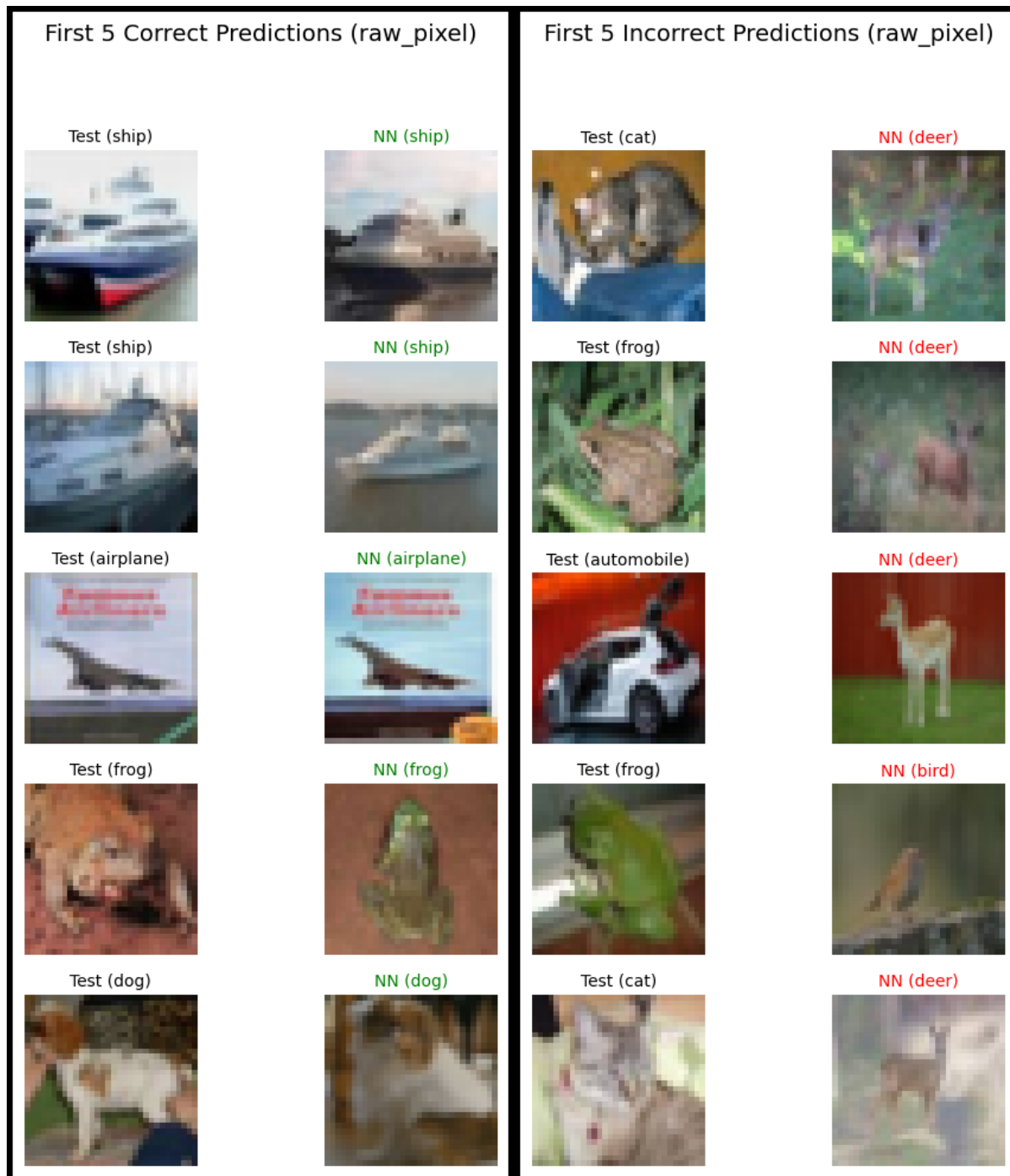**Part c: Raw Pixel correct/incorrect**



Figure 2: First five correct/incorrect images for Raw Pixel

## Solution 7

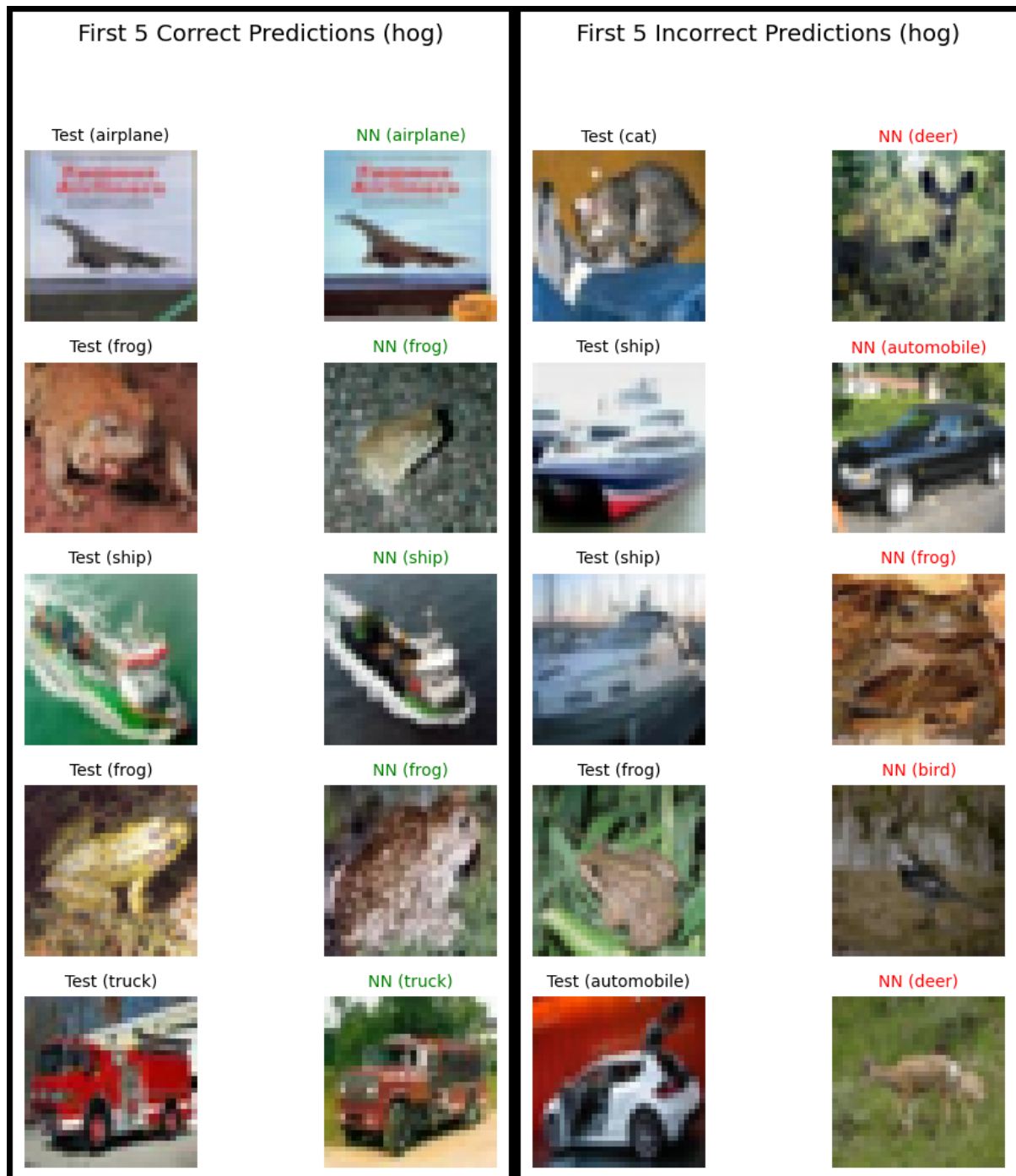**Part c: HoG correct/incorrect**



Figure 3: First five correct/incorrect images for HoG

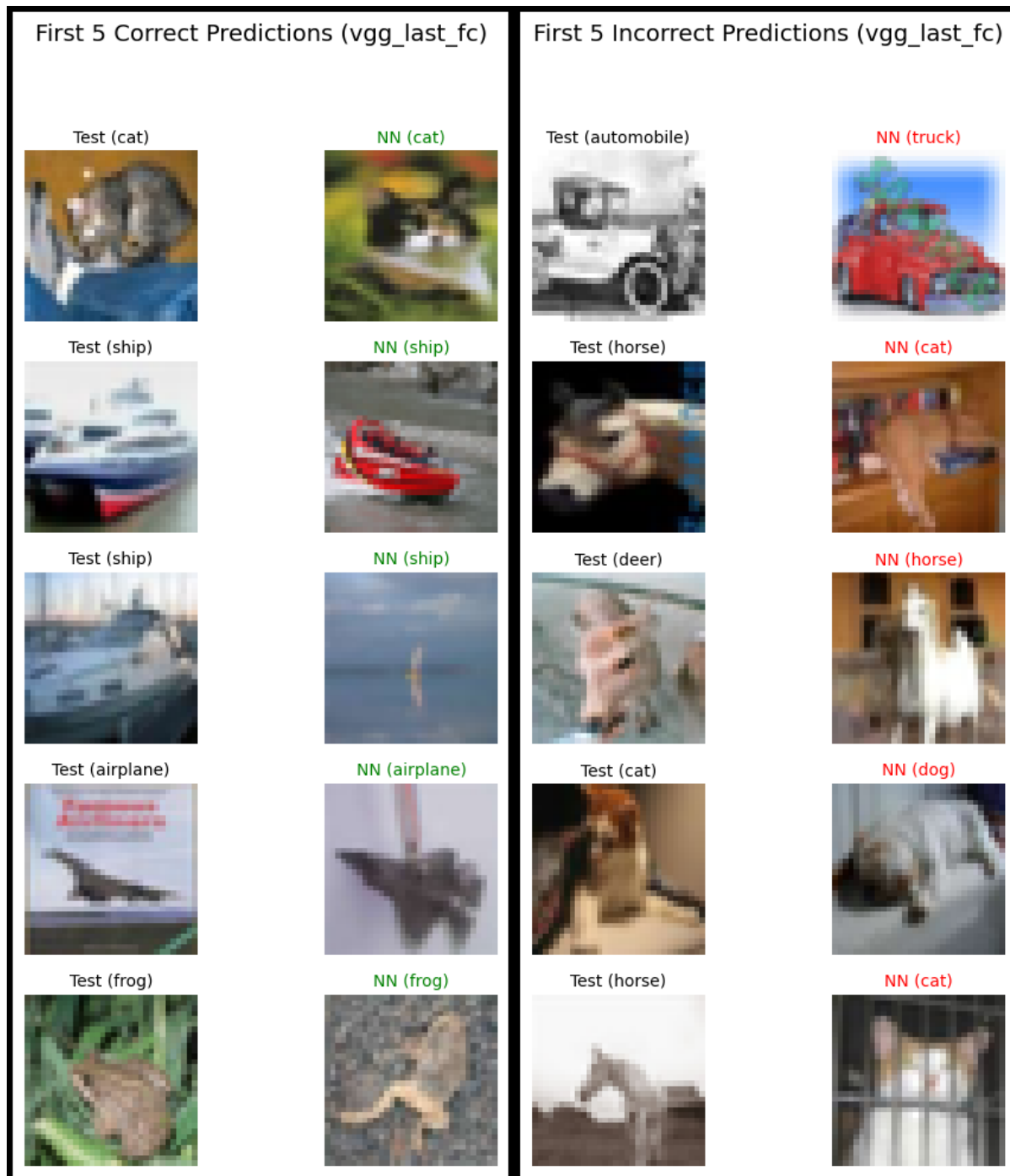## Solution 7

**Part c: VGG-last-fc correct/incorrect**



Figure 4: First five correct/incorrect images for VGG-last-fc