# Uncovering the Internal Representations of Neural Networks with Regression Analysis

Hiroshi Nonaka

May 12, 2025

#### Abstact

With the rise of AI technologies in different social sectors, understanding the internal states of deep neural networks is becoming a critical research field that tackles the black box problem. One possible approach is probing, a regression analysis that takes internal activations as independent variables and regresses target information. This method quantifies how informative the representations are in relation to the dependent variables. I implemented probing for a simple three-layer neural network for binary image classifications with three regression models: linear probability models, logistic regression, and neural networks. The results revealed that all three layers possess equally rich information for the binary image labels. Through further statistical analysis, I discovered that the representation vectors from deeper layers are more multicollinear. By applying Principal Component Analysis (PCA) to mitigate multicollinearity, I also found out that the vectors essentially acquired an equally useful representation. Heteroskedasticity tests showed that the data was heteroskedastic at a high significance level, and the training steps did not affect it.

#### 1 Introduction

#### 1.1 Overview

Neural networks are essential building components of recent AI models such as ChatGPT [1]. Although neural network predictions often outperform linear regression models because of their flexible representation capacity, their internal states are hard to interpret compared to linear models. For instance, the common multivariate linear regression can be written in the following form:

$$\hat{Y} = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n$$

where Y is an estimated value of a dependent variable Y, and  $\beta_i$  are the coefficients of independent random variables  $X_i$ . This is easy to interpret. In contrast, neural networks with m layers are often represented as:

$$\hat{Y} = (T_l \circ f_{l-1} \circ T_{l-1} \circ \dots f_1 \circ T_1)(\vec{X}) \tag{1}$$

where  $\circ$  is the function composition operation,  $\vec{X} = [X_1, X_2, \dots, X_n]$ ,  $T_i$  is a transformation such that  $T_i(\vec{V}) = A_i \vec{V} + b_i$  for some input vector  $\vec{V}$ , matrix  $A_i$  and bias vector  $b_i$ , and f is some function (mostly non-linear) that slightly modify outputs from  $T_i$ . Repetitively applying  $f \circ T$  gives rise to the non-linearity of neural network functions. We train neural networks so they learn patterns in a dataset through multiple training steps. During the training, the error between the output  $\hat{Y}$  and the true value Y is calculated. This error function is called a loss function, and neural networks update the values of their matrices  $A_i$  and bias vectors  $b_i$  by taking the gradient of the loss function. This is how training optimizes neural networks to fit to the given dataset.

#### 1.2 Motivation

Due to the multiple transformations and non-linear functions within neural networks, it becomes hard to interpret what the matrix and bias vector in each transformation represent compared to simple linear regression models, in which the coefficients denote how much each variable contributes to the estimated point  $\hat{Y}$ . The uninterpretability of neural networks is named the black box problem and has stimulated various research endeavors, including explainable AI (XAI) [6]. Given that AI systems are implemented in many different sectors today, understanding how AI thinks is one of the critical themes of ongoing research. I aim to tackle this problem through a popular approach called probing [5]. My research questions are mentioned in Section 4.

# 1.3 Approach

Probing applies regression analysis to measure how rich information an internal vector contains in relation to target data. Let an activation  $\vec{Y}_i = [Y_{i,1}, Y_{i,2}, \dots, Y_{i,m}]$  represent the output vector of the *i*-th transformation  $T_i$  in Equation 1:  $\vec{Y}_i = T_i(f_{i-1}(\vec{Y}_{i-1}))$ . Note that the first activation is produced from the observation vector:  $\vec{Y}_1 = T_1(\vec{X})$ . In prob-

ing, regression models P, namely probes, regress some target information using the vector  $\vec{Y}_i = [Y_{i,1}, Y_{i,2}, \dots, Y_{i,m}]$  as m independent variables. If P produces good predictions, this indicates that  $\vec{Y}_i$  are good independent variables, or more specifically, contain rich information about the target.

#### 2 Notations

Here, I outline notations that frequently appear in this paper.  $T_i$  indicates the *i*-th transformation layer in a neural network.  $f_i$  represents a function applied to the output from  $T_i$ . Most models uses a single non-linear function throughout their pipeline, so I occasionally write f to denote the unique function.  $\vec{X} = [X_1, X_2, ..., X_n]$  is an n-dimensional observation and an input to a neural network. An activation  $\vec{Y}_i = [Y_{i,1}, Y_{i,2}, ..., Y_{i,m}]$  is an m-dimensional vector from the i-th layer in the neural network. Therefore, each observation  $\vec{X}$  has corresponding activations  $\vec{Y}_1$ ,  $\vec{Y}_2$  through  $\vec{Y}_l$  for each transformation in Equation 1. In my paper, P is a probe/function/regression model that takes activations and maps them to a probability p:  $p = P(\vec{Y}_i)$ . One example of P is logistic regression.

## 3 Literature Review

In this section, I aim to introduce some prior works on probing and simultaneously deepen the audience's understanding of its method through concrete examples.

## 3.1 Machine Learning in Regression Analysis

Before discussing probing techniques, it is worth noting that machine learning algorithms are closely tied to regression analysis in many fields. Decision trees [10] and their families, such as random forest and gradient boosting, are popular regression algorithms for their accuracy and interpretability. Ridge and Lasso regressions [13], regularized linear regression models, are also important linear models in machine learning. These machine learning models are applied in many different areas, such as political sciences [7], chemistry [4], and economics [15].

## 3.2 Probing Neural Networks

Gurnee and Tegmark [8] applied Ridge regression and neural nets to measure the spatiotemporal representations within large language models (LLMs). LLMs are neural-network-based AI models that deal with language, such as ChatGPT. For spatial information, they gave LLMs prompts asking for the latitude and longitude of a specific place in the world, and for temporal information, the LLMs were asked about historic dates related to well-known books, figures, and et cetera. Activations produced by those stimuli were collected and used as independent variables to regress the latitudes and longitudes of the places or the associated years to the historic events. The researchers discovered that the linear model can predict the dependent variables with high accuracy, concluding that the spatio-temporal representations within LLMs are linearly related to the target values.

Multiple research works have applied probing to investigate representations of neural networks in various modalities, like text and image. McCoy and Leshinskaya [12] applied regression analysis to investigate how an LLM understands subject-object relations. Specifically, they focused on the has-color relation; for instance, the subject is bananas, and the relation between the subject and object is have the color of, then the expected object is yellow. They gave prompt queries asking the color of a particular subject, collected activations from the fifth layer of the LLM  $(\vec{Y_5})$ , and regressed the final output  $\vec{Y_l}$ .

Alain and Bengio [3] introduced linear classifier probes as a tool to understand the information captured in intermediate layers of neural networks. They collected activations from each layer of their neural network for image classification. Results revealed that the error of probes decreased as they used activations from deeper (later) layers as independent variables. Rashtchian et al. [17] probed image encoders of foundation models (AI models processing multimodal data) to quantify how the encoders generalize image representations. They collected activation vectors by giving the same images, but with different visual effects, and discovered that the representations that the encoders obtain from these differentially edited images converge well. Ilharco et al. [9] focused on the multimodal relations between contextual text representations with corresponding visual features by regressing image features based on activations from text-only LLMs. Results indicated that language representations are strong signals for predicting visual features.

Akhondzadeh et al. [2] leveraged linear, Bayesian, and pairwise probing to investigate the relations between the representations of graph-based models with important chemical properties like functional groups and odors. Ngo and Kim [14] trained probes to align language representations in LLMs with sound representations from audio models.

# 4 Research Directions and Questions

Despite the considerable number of prior studies, few have focused on the statistical analysis, such as heteroskedasticity test and multicollinearity, of representations along with probing. In this work, I aim to (1) mildly replicate existing research on probing techniques and (2) apply statistical analysis to internal representations for further understanding than probing. The research questions in this work are (1) which activation outputs  $\vec{Y}_i$  contain rich information, and (2) how the amount of training affects representation learning.

## 5 Data Collection

## 5.1 Training and Architecture

To collect activation vectors, I first trained my neural network with 2,000 images of cats and dogs from the CIFAR10 dataset [11]. The model classifies whether a given image is likely to be a cat or a dog. The training pipeline is as follows: (1) give the model an image and have

it predict probability p that the image is a cat; (2) calculate binary cross-entropy between p and its true label as an error score; (3) repeat step (1) and (2) for 1,600 training images; (4) calculate a validation score with 400 unseen data; and (5) repeat step (1) through (4) for 100 times with the same training and validation data. Each repetition of steps (1) through (4) is called an *epoch*.

The model consists of a *convolution* part and a neural net part. Convolution layers compress a  $64 \times 64 \times 3$  (height, width, and channels) image into  $13 \times 13 \times 16$ . Note that the convolution layers are also trainable.

The neural network has three layers:  $T_1, T_2$  and  $T_3$ . Each  $T_i$  maps  $f(\vec{Y}_{i-1})$  or  $\vec{X}$  to a vector of 60 dimensions. Therefore,  $T_1: \mathbb{R}^{2704=13\cdot13\cdot16} \to \mathbb{R}^{60}, \ T_2: \mathbb{R}^{60} \to \mathbb{R}^{60}$ , and  $T_3: \mathbb{R}^{60} \to \mathbb{R}^{60}$ . A non-linear function f called Rectified Linear Unit follows  $T_1, T_2$ , and  $T_3$ . One final layer  $T_4$  transfoms the 60-dimensional vector from the previous transformation  $T_3$  to a scalar. The logistic function  $\sigma$  turns the scalar into probability p.

The overall model architecture is outlined below:

$$p = (\sigma \circ T_4 \circ f \circ T_3 \circ f \circ T_2 \circ f \circ T_1 \circ conv)(\vec{X})$$

where conv is the convolution part, and  $\vec{X}$  is an image vector.

#### 5.2 Data Requirement

After each *epoch* during training, I gave the same 2,000 images to the network and collected  $\vec{Y_1}$ ,  $\vec{Y_2}$ , and  $\vec{Y_3}$  for *each image*. Therefore, I collected 600,000(= 2,000 images × 3 layers × 100 epochs) activations in total. I restricted the output dimension of  $T_1$ ,  $T_2$ , and  $T_3$  to 60 in order for  $\vec{Y_i}$  to have the same dimension size, which, later in probing, prevents the variability of regression performance due to the different number of independent variables. In short, the data I collected is cross-sectional data consisting of 300 .csv files (3 files for each of the three layers, multiplied by 100 epochs) with a table of 60 independent variables ( $\vec{Y_i}$ ) and 2,000 observations for each file (refer to Figure 1 for details). I did not implement any control or instrument variables because of the nature of my research. The details of data preprocessing are explained in Section 6.2.

## 5.3 Assumptions and Diagnostics

For statistical analysis, I focused on two metrics, heteroskedasticity and multicollinearity. For heteroskedasticity, I ran the Breusch-Pagan test to calculate p-values for the null hypothesis of homoskedasticity. The left plot in Figure 2 shows that the entire data is heteroskedastic at a very high significance level (p-value <<<0.001). Only the first layer experiences rapid increases of p-value, but this trend is inconsistent.

VIF of the second and third layers are also extremely high. VIFs of the third layer activations converges around 3 as the epoch approaches the end of training. Therefore, the activations from the first layer are mildy multicollinear, and those of the second and third layers have high multicollinearity. However, I can observe the downward trend of VIF as the amount

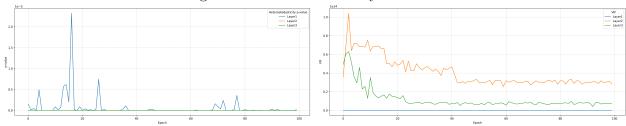
Figure 1: .csv file of the first layer at epoch 1

index label	scalar_0	scalar_1	scalar_2	scalar_3	scalar_4	scalar_5	scalar_6	scalar_7	scalar_8	scalar_52	scalar_53	scalar_54	scalar_55	scalar_56	scalar_57	scalar_58	scalar_50
0 1	-0.1806216670036320	0.238051952958107	-0.23040471971035000	0.2444995939731600	-0.224036306142807	-0.23286768794059800	0.40625232458114600	-0.45800262689590500	0.24776443839073	-0.3539388179779050	-0.3796703815460210	-0.2730218172073360	0.6700180768966680	0.24066977202892300	0.04509800872531130	-0.2485959678888320	0.11127142608165700
1 1	0.25909456610679600	0.025823909789323800	-0.2945656180381780	0.716274619102478	-0.310443639755249		0.16475407778216800	-0.32931967241859400	0.37751957774162	-0.22143204510211900	-0.2026701122522350	-0.24072836339473700	-0.540337065723877	-0.3383471966789900	0.6426206231117250	-0.38953566551208500	-0.16938401758670800
2 0	0.07127939164638520	0.15038348734378800	-0.1940902958202360	0.09208225458860400	-0.23451805114746100	-0.14938011765480000	0.3071085810961320	-0.30271488428115800	0.066796572638511	-0.24862990178642300	-0.33354413509368900	-0.2728663980960850	0.4940780669253080	0.04606788605451580	0.13055066764354700	-0.38924288749994800	0.1278744339942930
3 0	-0.07842961115932470	-0.04762215539813040	-0.7541811496217040	-0.4809391498565670	-0.515975832939148	-1,175797938300540	-0.44280582966397100	-0.7047756314277650	0.479679256677626	-0.3705447316169740	-0.5571001768112180	-0.2814175486564640	1.1031765907805200	0.2157389988401410	-0.8824149370193480	-1.1720149517059300	0.880963454532623
4 0	-0.061542364060878800	-0.16863088309764900	-0.3474876284596900	0.0032825618962315100	-0.34706827998161300	-0.001375948078930380	-0.2217320054769520	-0.3081977367401120	0.44666700197219	-0.3171845078468320	-0.33373409509658800	-0.0846555083990097	0.5567300319671630	-0.051794134090410000	-0.12725801765918700	-0.3900648627090450	0.2972768843173980
5 0	0.41749128899302700	-0.1074407547712330	-0.24033965170383600	0.2001992201805120	-0.28399658203125	0.042694512754678700	0.17685754597187000	-0.296296089887619	0.380857464266823	-0.22650395333766900	-0.33410676969364600	-0.18828488886356400	0.015215903520584100	-0.055375076830287100	0.07547566294670110	-0.04442659765481950	0.1997092515230180
6 0	0.6730656027793880	0.990655243396759	-0.45839348435401900	0.9568737745285030	-0.4012591540813450	0.1084270030260090	-0.10609669942951200	-0.14305035769939400	0.51419019899096	-0.18039974570274400	-0.2763514518737790	-0.17388403415679900	-0.5259617838668820	-0.40429896116256700	0.6202533841133120	-0.17390725016593900	-0.4173938930034640
7 0	0.2758157551288610	0.20073196291923500	-0.30904469776153900	-0.0007705451250076300	-0.13389921188354500	-0.4602925181388860	0.23953910171985600	-0.1180797996311260	0.40896729871490	-0.2789744734764100	-0.33012664318084700	-0.05675767734646800	0.26791778206825300	-0.002238065004348760	0.12983858585357700	-0.4240480066345760	0.0151141192764044
8 0	0.773219645023346	0.01706048659980300	-0.036968085914850200	0.8076162904303280	-0.34229186177253700	0.2895367741584780	0.7080490589141850	-0.229110449552536	0.128244680146477	-0.3432069420814510	-0.4426083564758300	-0.2785899043083190	-0.27088332176206500	-0.20661267638206500	0.5814420580963950	0.09412574768066410	-0.509296165237427
9 1	0.20727244018508400	0.08729682862758640	-0.5235839674758910	0.09087687730789190	-0.404738187788917	-0.37272348999997710	0.019256948941905600	-0.3468703627586070	0.33456200361251	-0.4703217148780820	-0.3891850709915160	-0.26104068756103500	0.28491583466529900	-0.4239004850387570	0.020990837602615400	-0.928119421005249	0.056147811710834
10 1	0.19982080161571500	0.027240820229053500	-0.5126203894615170	0.6585033535967340	-0.3920000196500240	0.17882852448854200	0.640214741230011	-0.38505464792251600	0.131356745958326	-0.30347615480423000	-0.5655096769332990	-0.17140811681747400		-0.4196825163268040	0.6556157469749450		-0.07924704253673559
11 1	0.38982293009758000	0.027972225099802000	-0.2735860047747800	0.06440538913011550	-0.23562952879091000	-0.19977356490473100	0.002979222988099880	-0.30567702651023900	0.25306424752725	-0.2430526614189150	-0.33455249667167700	-0.21270961629685000	0.4820329248906180	-0.15256863832473800	0.12415005266666400	-0.4248061180114750	0.2421190291643140
12 1	0.31986838579177900	-0.3176800310611730	-0.577718198299408	0.5730097889900210	-0.4688889682292940	1.0196634531021100	0.8371493816375730	-0.8747255802154540	0.9466013908386			-0.39032283425331100	0.6216617822647100	-0.6591458916984120		0.0034043043851852400	
13 1	0.22530657052963800	0.09646274149417880	-0.4689100980758670	-0.38229262828829900	-0.5227727293968200	-0.6466355821609500	0.13448722960541500	-0.5687738253044130	0.4954983890056		-0.40217864513397200		0.9091739078262330	-0.2371070663002470	-0.3357539176940920		0.1812322288751600
14 1	0.2266487181186680	-0.08985801041126250	-0.1823870688676830	-0.02595733106136320	-0.4284425377845760	-0.4984428882596880	0.3085671067237850	-0.447387158870697	-0.180618822574615		-0.5071409940719600	-0.27314165234565700	0.15022769570350600	-0.22291846573352800	0.2907170355319980		-0.0659877300262456
15 1	0.35751116275787400	0.08254578709802360	-0.31903275847435	0.49954652786254900	-0.4432121813297270	0.07680703899588780	0.19916433095932000	-0.42428305745124800	0.317549377679825		-0.2921089231967930	-0.2946103513240610	-0.19711783528327900	-0.30964179134368900	0.5197453496840330		-0.2442376911640170
16 0	1.386624813079830	0.3525449335575100	-0.21645677089691200	1.6877480745315600	-0.3242713212966920	1.093017339706420	1,709019860949710	-0.410988450950354		-0.4351404308272770	-0.7601648569107090	-0.3565865457057950	-1.08914053440094	-0.41500188920166000	1.5126643180847200	0.396284262207489	
17 1	0.12201926112175000	0.1767570823431020	-0.39058204164505000	-0.0864125382900238	-0.3446931838989260	0.11605633957481400	-0.019767310470342600	-0.34622955322265600	0.410319656133652		-0.15077152848243700 -0.31022942066192600		0.3806799650192260 -0.10116390863922900	-0.1368672251701360 -0.10153555870056200	0.12411914765834800	-0.35342609882354700 -0.200678533434868000	
18 0	0.183485746383667	0.2510315775871280	-0.40054364671707200	0.2427864968776700	-0.2916232943534850	-0.172395259141922	0.0009698085486888890	-0.24458740651607500	0.381693929433823		-0.27103888008262900	-0.18450775742530800	-0.10116390883922600	-0.10153555870056200	0.17501966026763600		-0.1122037904303280
19 1	0.30902315068244900	0.2925279140472410	-0.357525110244751	0.7634008526802060	-0.26611432433128400	0.3345284163951870	0.07663780450820920	-0.27419421076774900	0.47773060202596	-1.1081730127334600		-0.18450775742530800	1.642326831817630	0.47127866744995100	0.8847113847732540		-0.296265308141708
20 1	0.28486692905426000	0.9624137282371520	-0.5235768556594850	0.8805583119382400	-0.5868474245071410	-0.7637255191802980	1.9514774084091200	-1,0468136072158800	0.060258477926254		-1.39249849319438 -0.6540349721908570	-1.1938680907592800 -0.5702263712883000	.0 94103592831817630	-0.907583832740385100	13832881777191200	1.061432794370900	
21 1	1,0997468519210800	-0.22630716860294300	-0.4516868691308690	1.7310093641281100	-0.5945833921432500	1.782225489616390	1,5429710149765000	-0.5582174062728880			-0.16099096722183200			-0.06995539421537400	0.01855584130063080		-0.1117647898278250
22 1	0.13756310939788800	0.1292366236448290	-0.13999612791633600	0.02229611575603490	-0.2233031690120700	-0.249485682113266	0.2193225622177120	-0.17764116823673200	0.152093946933746		-0.7314572334289550	-0.5875479002990180	0.9042062473297120	-0.22767136991024000	-0.2433871030807500	-1.9421751499178000	
23 1	-0.020161796568885600	0.1094130426645280	-0.6033146977424620	-0.5587277412414550	-0.6816045045852660	-0.8245342373847960	0.09062930941581730	-0.6765416264534000	-0.132378876209256		-0.20308298580104800		-0.16263073682785000	-0.11419809729219000	0.31597812304999900	-0.08663086563151090	
24 0	0.44274425506591800	0.34589120745658900	-0.2855651242828370	0.4958667756126960	-0.19244349002838100	0.011160935275256600	0.10797879099845900	-0.2441529631614690			-0.40584099292755100		0.47134116291999900	-0.09999784319400790	0.31644508242907100	-0.48017311096191400	
25 0	0.5142162442207340	0.01765696331858840	-0.19129341840744000	-0.10929680296857100	-0.32696817947387700	-0.10343368351459600	-0.08060242235660550	-0.6093207597732540			-0.4403841496513920		0.38809665713119600	-0.1372072845697400	0.17197073996067000		0.00845113582909107
26 1	0.004775624722242360	-0.22621287405490900	-0.5996531844139100	0.3877355456352230	-0.6064339876174930	0.4526548385620120	0.12400354444980600	-0.5647415518760680			-0.4925215542316440	-0.3839457929134370	0.247898131606963	-0.29395759106682400	-0.1631847321987150	-0.8453869819641110	
27 0	0.2478729784488680	-0.034320179373025900	-0.5150821805000310	-0.3110767602920530	-0.4615599513053890	-0.5440896153450010	-0.2930900643157960	-0.48962876200676000			-0.5027818063763120	-0.1711081862449650	0.558029294013977	-0.05808150768280030	-0.2511945962905880	-0.6447329521179000	0.3052671849727630
28 0	0.11289404332637800	-0.22849927842617000	-0.23577523231506300	-0.006443098187446590	-0.406571626663208	-0.28491079807281500	-0.1280091553926470	-0.41924622654914900	0.155230373144150	-0.14149418473243700	-0.18658339997032200	-0.15294326841831200	-0.20629757642746000	-0.2871888279914890	0.21802721917629200	-0.20742063224315600	-0.1067532490491871
29 0	0.16638296842575100	0.112979456782341	-0.2496940642595290	0.3563551604747770	-0.25542157888412500	-0.03948894515633580	0.3152730464935300	-0.16512851417064700	0.130261689424515	-0.36512213945388800	-0.6330290634422300	-0.1882133036851880	-0.835198163986206	-0.4176303744316100	1.119140386581420	-0.08972938358783720	0.0825790315866470
30 0	0.9128257036209110	0.4222291111946110	-0.3815732002258300	1.2524895668029800	-0.23628352582454700	0.756817638874054	1.0525587797164900	-0.3361946940422060	0.61265778541564	-0.23408790992031100	-0.32058462500572200	-0.23116451501846300	-0.5164337158203130	-0.2539595365524290	0.6219285726547240	-0.14462578127384200	-0.1331446468830110
31 1	0.3824794590473180	0.2634543478488920	-0.2900915241241460	0.8067054491043090	-0.2569068372249600	0.46744024753570600	0.44067293406532800	-0.2996390347480770	0.5034121274948	-0.21832291781902300	-0.2877565622329710	-0.2065703570842740	-0.07424116134643560	-0.09419800341129300	0.16209056973457300	-0.24114640057086900	-0.06166597455739980
32 0	0.5680714845657350	0.1357696291540150	-0.2625271677970890	0.5131095051765440	-0.331266313791275	0.15244382619857800	0.23171024024486500	-0.26764315366745000	0.383948564529419	-0.37436962127685500	-0.4437200129032140	-0.2146567553281780	0.3511119484901430	-0.20664279162883800	0.042507022619247400	-0.18075691163539900	-0.3458870947360990
33 0	0.2827852666378020	0.05722084897418020	-0.13761602342128800	0.06408745050430300	-0.17366260290145900	-0.04491389915347100	0.39127373696373500	-0.36058443784713700	0.394024461507797	-0.42121440172195400	-0.6728751063346860	-0.32759854197502100	-0.3371215760707860	-0.21343019604682900	0.9773392081260680	0.15996858477592500	-0.4914964735507971
34 1	1.0980379581451400	0.24320665001869200	-0.20708459615707400	1.09586763381958	-0.19992893934249900	0.7316392084094540	1.2555079460144000	-0.6515689492225650	0.467025369405746	-0.25224143269677900	-0.49110883474350000	-0.12585359811782900	-0.32929004187583900	-0.13533231616020200	0.5881228108406070	-0.14764265716075900	0.1517637073993680

The table includes expected labels and vector values as columns.  $scalar_{-j}$  denotes the j-th value  $(Y_{i,j})$  in the vector  $(\vec{Y}_i = [Y_{i,1}, \dots, Y_{i,j}, \dots, Y_{i,m}])$ . Indexes on the leftmost column number different 2,000 activations.

of training increases and the mean VIF converges to certain value points (this trend is also observable in the third layer).

Figure 2: Heteroskedasticity and VIF



The left line plot is the p-values for the heteroskedasticity test, and the right graph shows the mean VIF within data from each layer and epoch. Note that the y-axis of the left graph is scaled by 1e-5 (1.0 in the y-axis is equal to 0.000001), and VIF is scaled by 1e14 = 100,000,000,000,000.

# 6 Methodological Framework

#### 6.1 Probes

I prepared three regression models as probes. The first is linear probability models:

$$p = \beta_0 + \beta_1 Y_{i,1} + \beta_2 Y_{i,2} + \dots + \beta_{60} Y_{i,60}$$

Although this model cannot capture the trend of cumulative Gaussian distribution functions, it serves as a good measure of the linearity of representations because it directly applies linear regression for probability p. In other words, if a linear probability model predicts probabilities with high accuracy, the independent variables  $Y_{i,1}, Y_{i,2}, \ldots, Y_{i,60}$  have high linearity with

target labels. We clipped any p > 1 as 1 and p < 0 as 0.

The second model is the logistic regression model. Logistic regression captures the natural shape of the cumulative Gaussian distribution functions and is a popular model for binary classification.

$$p = \frac{1}{1 + exp(-(\beta_0 + \beta_1 Y_{i,1} + \beta_2 Y_{i,2} + \dots + \beta_{60} Y_{i,60}))}$$

Finally, I introduce simple three-layer neural networks. Neural networks serve as a good indicator of non-linearity in activations. If the accuracy of linear probability models and logistic regression is low while the neural networks regress well, there is a non-linear relationship between the activations and true labels. I applied the three models to each of the 300 .csv files and recorded the accuracy. I split 2,000 observations into training and test datasets.

#### 6.2 Data Preprocessing

Most values of my data range from 0.3 and -0.3, which resulted in the poor convergence of the regression models. Therefore, I standardized the values. Standardization is given by  $x_{std} = \frac{x-\mu}{\sigma}$  where  $\mu$  is the mean of the independent variable X, which sample x belongs to, and  $\sigma$  is the standard deviation of the variable. In the follow-up research mentioned later, I also performed principal component analysis (PCA) to reduce the dimensions of the observations from 60 to 10. PCA essentially removes unnecessary, sparse representations and compresses data to lower dimensions while retaining most of the information contained in the original data.

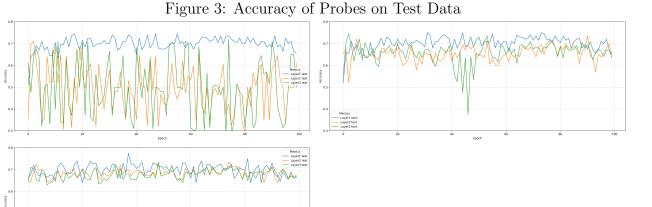
#### 7 Results

# 7.1 Probing Results

Figure 3 shows the results of probing. The linear probability models suffer from fluctuating accuracy. Logistic regression is more robust than linear probability models, but there still exists a large decline in accuracy. In contrast, neural networks regressed the true labels more stably. Moreover, the overall results show that the first layer contains more stable representations, while the activations from the second and third layers resulted in slightly lower or significantly lower accuracy in the logistic models and linear probability models. These graphs present counterintuitive results, as I initially speculated that representations would grow richer as the layer goes deeper, closer to the output layer. However, in Section 4, I identified the multicollinearity in the activation data and hypothesized that it degraded the accuracy of the probes.

# 7.2 Probing Results with PCA

Speculating that the high multicollinearity contributes to the low accuracy with the second and third layers, I performed PCA to reduce dimensionality. PCA mitigates multicollinearity by creating new variables that are linearly unrelated to each other [16]. Probing results



The upper-left graph is the linear probability models, the upper-right is the logistic models, and the bottom-left is the neural networks. The blue lines represent the scores of probes on the activations from the first layer, the orange lines correspond to the second layer, and the green lines track that of the third layer.

after PCA are outlined in Figure 4. I can observe clear improvements in the accuracy scores across the three probes, especially in the linear probability models. Furthermore, the accuracy scores across different layers generally stay at almost the same level, indicating that activations from the first, second, and third layers possess equally rich information. Moreover, the stable accuracy scores of the linear probability models indicate that there is a solid linear relationship between the activations and true labels.

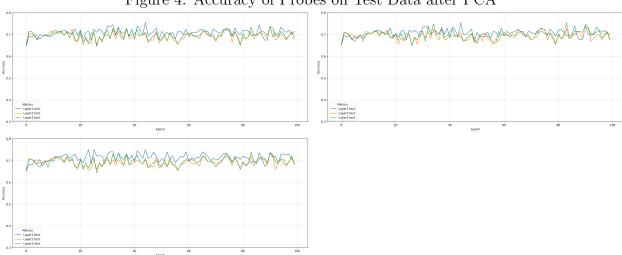


Figure 4: Accuracy of Probes on Test Data after PCA

The upper-left graph is the linear probability models, the upper-right is the logistic models, and the bottom-left is the neural networks.

#### 8 Conclusion

In this study, I investigated how the internal representations of neural networks in binary image classification change depending on the depth of layers and training amount. I discovered that the layers later in the network pipeline possess representations with higher multicollinearity. However, the multicollinearity in every layer declines and converges to certain points of value as the amount of training increases. Further regression analysis on the data after PCA revealed that the representations from different layers possess the equivalent level of information about the true binary labels.

My research indicates that probing techniques are a suitable approach to understand *how* AI thinks by quantifying what kind of representations it acquires. The results also present answers to the research questions: (1) every layer contains equally rich information, and (2) the number of training steps does not affect accuracy or heteroskedasticity but makes multicollinearity converge.

This research contributes to the further understanding of deep neural network models not only through probing but also via statistical analysis. Limitations posed to this work are that I solely focused on neural networks for binary image classification. I also failed to replicate similar results with Alain and Bengio [3], who argued that deeper layers contain richer representations in binary image classification. Thus, further studies should be done to apply my analysis methods to neural networks with broader tasks and perform stricter replication studies.

#### References

- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2023.
- [2] Mohammad Sadegh Akhondzadeh, Vijay Lingam, and Aleksandar Bojchevski. Probing graph representations. In Francisco Ruiz, Jennifer Dy, and Jan-Willem van de Meent, editors, *Proceedings of The 26th International Conference on Artificial Intelligence and Statistics*, volume 206 of *Proceedings of Machine Learning Research*, pages 11630–11649. PMLR, 25–27 Apr 2023.
- [3] Guillaume Alain and Yoshua Bengio. Understanding intermediate layers using linear classifier probes, 2018.
- [4] Vy Anh Tran, Le Thi Nhu Quynh, Thu-Thao Thi Vo, Phuc An Nguyen, Ta Ngoc Don, Yasser Vasseghian, Hung Phan, and Sang-Wha Lee. Experimental and computational investigation of a green knoevenagel condensation catalyzed by zeolitic imidazolate framework-8. *Environmental Research*, 204:112364, 2022.
- [5] Yonatan Belinkov. Probing classifiers: Promises, shortcomings, and advances. *Computational Linguistics*, 48(1):207–219, 04 2022.
- [6] Prashant Gohel, Priyanka Singh, and Manoranjan Mohanty. Explainable ai: current status and future directions, 2021.
- [7] Justin Grimmer, Margaret E. Roberts, and Brandon M. Stewart. Machine learning for social science: An agnostic approach. *Annual Review of Political Science*, 24, 2021.
- [8] Wes Gurnee and Max Tegmark. Language models represent space and time, 2024.
- [9] Gabriel Ilharco, Rowan Zellers, Ali Farhadi, and Hannaneh Hajishirzi. Probing contextual language models for common ground with visual representations. In Kristina Toutanova, Anna Rumshisky, Luke Zettlemoyer, Dilek Hakkani-Tur, Iz Beltagy, Steven Bethard, Ryan Cotterell, Tanmoy Chakraborty, and Yichao Zhou, editors, *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5367–5377, Online, June 2021. Association for Computational Linguistics.
- [10] Yacine Izza, Alexey Ignatiev, and Joao Marques-Silva. On explaining decision trees, 2020.
- [11] Alex Krizhevsky. pages 32–33.
- [12] Michael B. McCoy and Anna Leshinskaya. Relational composition during attribute retrieval in GPT is not purely linear. In NeurIPS 2024 Workshop on Compositional Learning: Perspectives, Methods, and Paths Forward, 2024.
- [13] L.E. Melkumova and S.Ya. Shatskikh. Comparing ridge and lasso estimators for data analysis. *Procedia Engineering*, 201:746–755, 2017. 3rd International Conference "In-

- formation Technology and Nanotechnology", ITNT-2017, 25-27 April 2017, Samara, Russia.
- [14] Jerry Ngo and Yoon Kim. What do language models hear? probing for auditory representations in language models, 2024.
- [15] Saeed Nosratabadi, Amirhosein Mosavi, Puhong Duan, Pedram Ghamisi, Ferdinand Filip, Shahab S. Band, Uwe Reuter, Joao Gama, and Amir H. Gandomi. Data science in economics: Comprehensive review of advanced machine learning and deep learning methods. *Mathematics*, 8(10), 2020.
- [16] Lexi V Perez. Principal component analysis to address multicollinearity. Whitman College: Walla Walla, WA, USA, 99362, 2017.
- [17] Cyrus Rashtchian, Charles Herrmann, Chun-Sung Ferng, Ayan Chakrabarti, Dilip Krishnan, Deqing Sun, Da-Cheng Juan, and Andrew Tomkins. Probing the equivariance of image embeddings. In NeurIPS 2023 Workshop on Distribution Shifts: New Frontiers with Foundation Models, 2024.