Scientific Computing April 28, 2025 Announcements > Homework 6 due Friday, May 2, 11:59pm -> Final exam is take-have only assigned Fri, May 2 due Fri, May 9 Office Hours: Today Mont Fri -> Loss Functions -> Backpropagation 9:30am - 10:30am Cudahy 307

How do we measure how good or bad a NN is m terms of matching known data? Linear Regression: Mean Squared Ervor, E. (actual-predicted)² 0.4 32 knots in this very small example

(Loss) 5 The "score" of a NN relative to particular training data. Change weights or biases: loss goes up (bad) or down (good)

7 55 7	wo types of problems we'll use NNs fer:
· · · · · · · · · · · · · · · · · · · ·	(1) Regression - predict output values based on in put values * Predict home price based on zip code, sq.ft, # bedrooms, # bathrooms, crime rate, school quali
· · · ·	* Predict # of bike rentals based on day, weather, holiday, etc.
· · · ·	(2) Classification - classify input into categories * MNIST climates
· ·	* Predict whether a patient has diabetes predicibet

We use diffe	rent loss functions for each type of
problem.	Scoring function for a NN, smaller is better
- Actually, fiv	st some notation.
- The loss fun pair at a input/outp	ction is defined not for one input/output time, but for a whole batch of rut pairs. [Remember "batching" from our last lecture.]
botch of mputs	NN (borteh) (055 for of outputs) (055 -> that botch (one #)

botch of mputs NN (of of function) (oss for outputs (one #) Remember each input is a whole vector (one # per mput neuron) and some for each output. y, y2, ..., yn. all vertors

The goal of a loss function is to measure how far apart the actual output \hat{y} is from the desired output y, and "training" = "make the loss get smaller"

* Regression. Loss function #1: Mean Squared Error (MSE) For a NN whose output layer has 1 neuron and for a batch of n input/output pars: $loss = \frac{1}{n} \left(\sum_{i=1}^{n} \left(y_i - \hat{y}_i \right)^2 \right)$ nean of that, over the whole For a NN whose output layer has k neurons, and for a botch of n input/output pars: $loss = \frac{1}{n \cdot k} \left(\sum_{i=1}^{n} \sum_{j=1}^{k} (y_{ij} - \bar{y}_{ij})^2 \right) \quad \begin{array}{c} y_{ij} \text{ is the } j^{\text{th}} \\ \text{Component of the} \\ \text{vector } y_i \end{array}$

Example:

>>> y = np.round(np.random.randn(3,5),2); y array([[1.9 , 0.24, -0.38, -0.86, 2.49], [-0.22, 0.17, -0.74, 1.12, 1.06],[-2.39, 0.98, -1.87, -1.62, 0.23]])>>> yhat = np.round(np.random.randn(3,5),2); yhat array([[0.03, 0.43, -0.25, 0.61, -0.92],[-1.84, -0.35, 2.32, 0.82, 0.51],[-1.11, -0.19, 0.48, 2.1, 0.6]])>>> (y - yhat)**2 array([[3.4969, 0.0361, 0.0169, 2.1609, 11.6281], [2.6244, 0.2704, 9.3636, 0.09, 0.3025],[1.6384, 1.3689, 5.5225, 13.8384, 0.1369]])>>> np.sum((y-yhat)**2) / 15 np.float64(3.4996599999999995) >>> 3 output neurons, batch of 5 mput/output pairs

* Regression. Loss function #2: Mean Absolute Error (MAE) For a NN whose output layer has 1 neuron and for a batch of n input/output pars: $loss = \frac{1}{n} \left(\sum_{i=1}^{n} |y_i - \hat{y}_i| \right)$ For a NN whose output layer has k neurons, and for a botch of n input/output pars: $loss = \frac{1}{n \cdot k} \left(\sum_{i=1}^{n} \sum_{j=1}^{k} |y_{ij} - \widehat{y}_{ij}| \right) \quad \begin{array}{c} y_{ij} \text{ is the } j^{\text{th}} \\ \text{Component of the } \\ \text{vector } y_i \end{array} \right)$

* Regression. Example: $y = [i], \quad \hat{y} = [s] \in Off by 4$ $MSE = \frac{1}{2} \left((5-1)^2 + (1-1)^2 \right) = 8$ $MAE = \frac{1}{2}(15 - 11 + 11 - 11) = 2$ $y = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \quad \begin{cases} y = \begin{bmatrix} 3 \\ 3 \end{bmatrix} = \begin{bmatrix} 3 \\ 3 \end{bmatrix} = \begin{bmatrix} 2 \\ 3 \end{bmatrix} = \begin{bmatrix} 3 \\ 2 \end{bmatrix} = \begin{bmatrix} 2 \\ 3 \end{bmatrix} = \begin{bmatrix} 2 \\$ MAE = = = (13-11+13-11) = 2 < same

* *	Cla Re	255 2Ca	fia -	that that		if w Le.q.	Je . 1	ar O	e diqi	so	647 4	ng	ì	npu IN	its IS		nt)		Ł	· · ·		CX	°5	· · ·	· · ·	· · ·	· · · · · · · · · · · · · · · · · · ·
• •	H	her	v †	here	e	an	2	k	Ov	dp	nt	• •	nei	A re	MS		•	• •	•	•	• •	•	• •	•	•	•	• •
· · ·	W	le ac Dro	pas tru bal	s H Vertica Xility	Ven N	n Fuu dist	col ncti ribe	lecton on tia	lve tc	ly	H tui	hvc rn	ug t	h he	H	12	nte	sof	tn c			•	· ·		•	•	· · ·
• •	•			J			•	•••		л ́	• •	• •	•	•	•	• •	•	• •	•	•	• •	•	• •	•	•	•	• •
• •	•	•	L	kn	Du	itput	- '•\				•	• •	•	•	•	• •	•	• •	•	•	• •	•		•	•	•	• •
• •	•		() U	fory	S	Gui	M	4			•		•	•	•	• •	•	• •	•	•		•	• •	•	•	•	
• •	•	• •			<i>–</i>					•	•	• •	•	•	•	• •	•	• •	•	•	• •	•	• •	•	•	•	• •
• •	•			• •	•	• • •	•	• •		•	•	• •	•	•	•		•	• •	•	•		•		• •	•	•	• •
• •			• •	• •	• •		٠			•	•	• •	•	•	•	• •			•			•	•	· •		٠	
		• •			•		٠	• •		•	•	• •	•	•	•	• •	٠	• •	٠	٠	• •	٠	• •	•	•	•	• •
	•	• •	• •	• •	•		٠	• •	• •	•	•	• •	٠	•	•	• •	٠	• •	٠	٠	• •	٠	• •	•	•	•	• •
• •	٠	• •	• •	• •	• •	• • •	•	• •	• •	٠	•	• •	٠	•	٠	• •	•	• •	٠	•	•••	•	• •	٠	•	٠	•••
• •	•	• •	• •	• •	• •		•	• •	• •	•	•	• •	•	•	•	• •	•	• •	•	•	•••	•	• •	٠	•	٠	•••
• •	٠	• •	• •	• •	• •	• • •	٠	• •	• •	•	•	• •	•	•	•	• •	•	• •	•	•	• •	•		•	•	•	• •
• •	٠	• •	• •	• •	• •	• • •	•	• •	• •	•	•	• •	•	•	•	• •	•	• •	•	•	• •	•	• •	•	٠	•	• •

Previously:

Unlike our other activation functions, Softmax works on the whole vector at once, not one value at a time individually. $\begin{cases} x_{i} \\ x_{z} \\ x_{s} \\ \vdots \\ x_{k} \\ x_{k}$ e×e/s Obvious these add up to 1 because the denominator is their sum. Obviously >0 because et >0.

Previously:

$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	0.041 0.111 0.091 0.184 $*$ 18.4% chance 0.136 He digit is 0.101 0.061 α 4
--	--

* Classification Goal of a loss func	toon: measur the the the	e The div actual cla reducted pr	stance ssification obability	ætween and distr.
Ex: Actual: [0 0 Predicted: [0.01 0.03 0 Another predicted: [0 0 -0	0 0 1 0.1 0.07 0.33 01 0.03 0.97	0.02 0.0		00] 0.030] 0.010]
Second one is much the loss should	closer to to the lawer.	the actual	avsuer,	50
. .	. .		· · · · · ·	· · · · · · · · ·

* Classification Ex: Achal: [0 0 0 0 1 0 0 0 0] Predicted: [0.01 0.03 0.1 0.07 0.33 0.11 0.01 0.21 0.03 al Another predicted: [0 0 -01 0.03 0.92 0.02 0.01 0 0.01 9] Lots of ways to devise a loss function to capture this closeness. The most popular is "Categorical Cross-Entropy". For I sample with k outputs: k=10 for digits loss = - Éyi logyi Oor $= -(y_1 \cdot \log(y_1) + y_2 \cdot \log(y_2) + \dots + y_k \cdot \log(y_k))$ = - log(ýy) (in this sample)

* Classification		
Ex: Actual: [0 0 0 0 1 0 0	0	0 0]
Predicted: [0.01 0.03 0.1 0.07 0.33 0.11 0.01	0.21	0.03 0.1]
Arother predicted: [0 0 -01 0.03 0.97 0.02 0.01	O	<i>v.01</i> 0]
$logs = -5' \text{Gi} \cdot log(\hat{y}_i)$	· · ·	· · · · · · · · ·
	· · ·	
O for the "wrong classes" I for the "correct class"		
So, this simplifies to $loss = -log(\hat{y}_c)$ where \hat{y}_c is the confidence level of the correct class	· · · · · · · · · · · ·	. .
	· · · ·	· · · · · · · ·

* Classification	· · · · · · · ·	· · · · ·	· · · · · · ·
Ex: Actual: [0 0 0 0 1		 . 0	0 0]
Predicted: [0.01 0.03 0.1 0.07 0.33 0.	-11 0.01	0.21	0.03 01]
Another predicted: [0 0 -01 0.03 0.92 0).02 0.01		0.01 0
$loss = -log(\hat{u}_c)$	· · · · · · · ·		· · · · · · ·
$\sum loc = -loc (0.22) \approx 0.48$	· · · · · · · ·	· · · ·	· · · · · · ·
· · · · · · · · · · · · · · · · · · ·	smaller	· · · · ·	
→ loss = - log(0.92) ≈ 0.03 €		· · · ·	· · · · · · ·
J			
· · · · · · · · · · · · · · · · · · ·			

***	F	Cla			fi	(<u></u>	le f	a he	~ 2 ~	• • • • • •	0 10 11	f No No	5	0		ne L	2	hi e	ole an di	e ch d			ata	cl	n // /		f Or gr	n vt	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	ur H	ようう	/6	wr Pa	j. iv		ļç I		j		st Sd	L	<i>[</i>],	y.	· · · · · · · · · ·
•	h		2	+	1n C	is , 01	<u>``</u>	f.	en	- M		b a 5	?	D	•	S	.	q	t de	51	s S C S C S	s		d	۲۶ Pi	ta	vic	e bi		2 2 2 1	d	ee st		bu		Jr.	>0		e	Le	•	•	•	•
•	• •	•	•	-	h)0	rk	S		Ŵ	ė		•	Ņ		Ę		à	5	ce		•	•	•		•	•	•			•	•	•	•	•	•	•	•	•	•	•	•	•	•
•		•	٠		•				,	٠	•	٠	٠			ľ	•	٠	•	٠	•	•	٠	•	•	•	•	٠		•	•	•	٠	•	٠	•	٠	٠	•	٠	•	•	•	٠
٠	• •	٠	٠	•	٠	•	•	•	•	•	•	٠	•	٠			•	•	٠	•	٠	•	٠	•	•	•	•	•	•	٠	•	•	•	•	•	•	•	٠	•	•	٠	•	•	٠
•		•	•	•	•	•	•	•	•	•	•	•	•	•			•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•
		•	•	•	•	•			•	•	•	•	•				•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•		•	•	•	•	•	•	•	•	•	•	•
•		•	٠	•		٠	•		,	•				•			•	•	•	٠	•	٠	٠	٠		٠	٠	٠		٠	•		٠		٠		٠		•	•	•	•	٠	•
•		٠	•		٠	•	•	•	,	•	•	•	•	•			•	•	•	•	•	•	•	•	•	•	•	•	•	•		٠	•	•	•	•	•	٠	•		•	•	•	•
•		•	٠	٠	•	٠	•		,	•	•	٠	•	•			•	•	•	٠	•	٠	٠	٠	•	٠	٠	٠	•	٠	٠	•	٠	•	٠	•	٠	•	•	٠	•	٠	•	•
•	• •	٠	•	•	٠	٠	•	•	,	•	•	٠	•	•			•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	٠	•	٠	•	•	•	٠	•	•	٠	•	•	•
•	• •	•	•	•	•	•	•	•	•	•	•	•	٠	•			•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	٠
•	• •	•	•	•	•	•	•	•	•	•	•	•	•	•			•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•

So, whether we're doing regression or classification, we have a "loss function" that measures: for a botch of inputs, how far is the actual from what the training data says it should be Lover is better. Gool: pick weights and biases that make the loss as low as possible on your traming data, and hope that the resulting NN performs well on whatever new data you have for it. How do we find good weights and biases?

Bad idea: Pick randon #s for them over and over again until some combination is good. Interesting Idea: Do Hill Climbing or Simulated Annealing. Change some /all of the weights and brases by a little, see if the loss of the new NN is better or worse. This is not ever done in practice, and I've never seen it discussed much. Why? My guesses: * too slow * so many parameters that it's hand to go downhill * the method people do use is good and fast Next topic: Backpropagation

 Have which a join which random weights and biases Feed all of the known inputs (training data) through in one batch and get all their outputs Compute the loss between the expected artputs (y) and the actual outputs (g) * Use CALCULUS to Figure out how to tweak the weights + biases a little to make the loss go 	Big	picture: A Stack will will will made
* Use CALCULUS to figure out how to tweak The weights + biases a little to make the loss go	 . .<	 Freed all of the known inputs (training data) through in one batch and get all their outputs Compute the loss between the expected artputs (y) and the actual outputs (j)
	· · · · · ·	* Use CALCULUS to figure out how to tweak The weights + biases a little to make the loss go

CALCULUS Let's recall the idea of <u>Gradient Descent</u>. The vector $\begin{bmatrix} \partial f \\ \partial x \\ \frac{\partial f}{\partial y} \end{bmatrix}$ tells you the livection n^{3} Suppose you have a function f(x,y)=z. direction of steepest ascent from any point (x,y). $\begin{bmatrix} -\delta f \\ \overline{\delta x} \end{bmatrix}$ is the direction of steepest $\begin{bmatrix} -\delta f \\ \overline{\delta y} \end{bmatrix}$ descent. The negative version -Vf =

CALCULUS $-\frac{\delta f}{\delta x}$ is the direction of steepest $-\frac{\delta f}{\delta y}$ descent. The negative version This means if you're enly going to walk I step in the mountains, then are step in that divection will get you as low as possible among all possible one steps. E_{x} : $f(x_{iy}) = x_{cos}(x)(s_{in}(y))^{2}$. $\int \frac{\delta}{\delta x} (x \cos(x) (\sin(y))^2)$ (-1:5,25) $-\frac{\delta}{\delta y} (x\cos(x)(\sin(y))^2)$

CALCULUS E_{x} : $f(x_{iy}) = x_{cos}(x)(s_{mly}))^{2}$. (-1.5, 2.5 0.6- $\frac{\partial f}{\partial x} = (smly)^2 \cdot (cos(x) - xsm(x))$ -0.2--0.4 $\frac{\delta f}{\delta y} = \chi(os(x), 2 \leq m(y) \cos(y))$ $(-\nabla f)(-1.5, 2.5) = \int \sin(2.5)^2 \cdot (\cos(-1.5) + 1.5 \sin(1.5))$ $(-1.5)\cos(-1.5)\cdot 2\cdot \sin(2.5)\cos(2.5)$ 0.5105. -0.1017.

CALCULUS (-1.5, 2.5 E_{X} : $f(x,y) = x\cos(x)(sinly))^{2}$ 0.2- $(-\nabla f)(-1.5, 2.5) = \begin{bmatrix} 0.5105 \\ -0.1017 \end{bmatrix}$ -0.2--0.4-This doesn't mean you move +0.51 on the x-axis an -0.10 on the y-axis. That would be a <u>huge</u> step Instead you move a small step, maybe too of that. $\begin{bmatrix} -1.5\\2.5\end{bmatrix} + \begin{bmatrix} 0.0051 \\ -0.0010 \\ -0.0000 \\ -0.0000 \\ -0.0000 \\ -0.0000 \\ -0.0000 \\ -0.0000 \\ -0.0000 \\$

•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	· · ·	•	•	•	•	•	•	•	•	•	•	•	•
•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	· · ·	•	•	•	•	•	•	•	•	•	•	•	•
•	•	4	0	6	Z . V	ng	d	ie	Å.	•	•	D	es		2N	·	•			21	Ņ	D	•	•	•	•	•	•	•	· ·	•	•	•	•	•	•	•	•	•	•	•	•
•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	· ·	•	•	•	•	•	•	•	•	•	•	•	•
•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	· ·	•	•	•	•	•	•	•	•	•	•	•	•
•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	• •	•	•	•	•	•	•	•	•	•	•	•	•
•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	· ·	•	•	•	•	•	•	•	•	•	•	•	•
•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	· ·	•	•	•	•	•	•	•	•	•	•	•	•
•	•	•		•	٠				•	•	•			•	•		•	•	•	٠	٠	•	٠	•	٠	•		•	٠		•	•	٠	•	٠	•	•			•	٠	•

We have been thinking about Neural Networks as big functions (input) NN > Coutput data data For now, instead of thinking of the input data as the Variables, we'll think of all the training data as a fixed constant and the weights and biases as the Variables. a single # input data > NN -> output > loss weights and brases

a single # input data > NN -> output > loss weights and brases UW5 How can we find the inputs that minimize the output? Gradient Descent. Suppose the neural network has weights w, w2,..., wm and biases b, bz,..., bn. Let l be the mean loss over the whole batched input. $-\nabla NN(w_{i},w_{2},...,w_{n},b_{i},b_{2},...,b_{n})?$ What is

$\ell = NNLw$	Wz(,	$w_{m,b_{1},b_{2},\dots,b_{n}}$	· · · · ·
How will w	Ne Con	mpute	
	<u>dr</u>	CALCULUS	· · · · · ·
$\nabla P N N = 0$	JL JWm	Specifically ? Hhe CHAIN RULE	· · · · ·
· · · · · · · · · · · ·	Sl Jb,		· · · · ·
· · · · · · · · · · · ·	<u>Sl</u>		· · · · ·
	Jbn -		· · · · ·

the CHAI	IN RULE	 	· · · · · · · ·	· · · · · · ·	· · · · · · · · · · ·
· · · · · · · · ·	· · · · · · · ·		t(+)	· · · · · · ·	· · · · · · · · · ·
<u>4</u>			· · · · · · · ·		· · · · · · · · · · ·
$\frac{\partial \Gamma}{\partial X}$ (3	$\mathcal{D} = \mathcal{D}$	neans if	you chang	e go fru	
	x=3 to	x = 3 + La	e little ma	e J	· · · · · · · · · · ·
	f(3) $f(3)$	f(3) + 5	• Ta little	MARO T	· · · · · · · · · ·
· · · · · · · · · ·					· · · · · · · · · · ·
· · · · · · · · · ·		· · · · · · · ·		· · · · · · ·	· · · · · · · · · · ·
	· · · · · · · ·	· · · · · · · ·	· · · · · · · ·	· · · · · · ·	

the CHAIN Suppose	RULE $x > [-]$ h(x) = g(f(x))	$\frac{f}{f} \xrightarrow{f(r)} [$	g gl-		· · · · · · · · · · · · · · · · · · ·
<u>Chain</u>	ule: h'(x) = g'(x)	f(*)) · f'($(\mathbf{x}_{1}, \mathbf{x}_{2}, \mathbf{x}_{3})$	· · · · · ·	· ·
6^	dh_ Oh of			· · · · · ·	• •
		how much	ch g Cha Changes	inges a	· · ·
how much h when x che	ox Changes nges	how much when f liftle	ch g cha changes	inges a	· · · · · · · · · · · · · · · · · · ·
how much h when x che liftle	ox changes nges a how much f chan	how much when f liftle nges	ch g cha changes	inges a a	· · · · · · · · · · · · · · · · · · ·
how much h when x che liftle	ox changes nges a how much f chan when x changes	how much when f liftle nges c	ch g cha changes	inges .	· · · · · · · · · · · · · · · · · · ·
how much b when x che liftle	ox changes nges a how much f chan when x changes liffle	how much when f liftle nges c	cha changes	inges .	· · · · · · · · · · · · · · · · · · ·
how much h when x che liftle	ox changes nges a how much f chan when x changes liftle	how much when f liftle ngles G	ch g cha changes	inges .	

the CHAIN RULE Neural Networks are just really big multivariate compositions of smple functions (addition, multiplication, activation functions, loss functions) a single # input data > NN -> output > loss weights and brases We can apply the chain rule a lot to see how changing any weight or bias individually affects the loss

the CHAIN RULE a single # NN -> output > loss [input clater] -> [weights and brases We can apply the chain rule a lot to see how changing any weight or bias individually affects the loss <u>SR</u> That means we'll know the gradient and we can do gradient descent. <u>Jl</u> Jwm = und SL اطل ا <u>I</u> Jpu

No AF, no hidden layers, 3 mput, 1 output, only are point of frammy clata First example: Training data: (1,-1,2) -> 1.5 loss $l = (1.5 - 1.6)^2 = 0.01$ $\frac{1}{-1} \frac{w_{1}=-0.2}{10^{10}} \frac{16}{10^{10}} \frac{$ loss How do we adjust W, W, W, W, b to make l go down fer this training data?

Training data: (1,-1,2) -> 1.5 () 4, 2 2 $(-1)^{\frac{1}{2}-0.2}$ loss (2) 2 0 2 Want: <u>de</u> <u>de</u> <u>de</u> <u>de</u> <u>de</u> <u>du</u>, <u>du</u>, <u>du</u>, <u>du</u>, <u>de</u> $l = (\hat{y} - 1.5)^{\prime}.$ $\hat{y} = 1 \cdot \omega_1 - 1 \cdot \omega_2 + 2 \cdot \omega_3 + b$ $S_{0}, \frac{\partial l}{\partial \hat{y}} = 2[\hat{y}-1.5].$ = $2\cdot(0.1) = 0.2$ $\frac{\partial y}{\partial w} = 1$, $\frac{\partial y}{\partial w_2} = -1$, $\frac{\partial y}{\partial w_3} = 2$, $\frac{\partial y}{\partial b} = 1$ So, λ chain rule! $\partial \omega_{z}$ كطم

Training data: (1,-1,2) > 1.5 () 4, 50 ¥ 1.6 -1) w-=-02 loss 2 w2 w2 b= 0.5 Now adjust each weight a small amount in this direction. Let's do too. de der, DR - PNN = -0.2 0.2 Jun -0.4 $W_{c}: 0.1 \rightarrow 0.098$ dl -0-2 W_2 ; $-0.2 \rightarrow -0.198$ Jwz $\omega_3: 0.4 \rightarrow 0.396$) L 6:0.5 ->0.498 26

L loss 0.007396 (prev. 0.01) 41.50.048 ×1.586 $\frac{1}{2} \frac{\omega_{1} = -0.198}{\omega_{2} = -0.198} \frac{1.586}{1.586}$ b = 0.498Now repeat! Compute - VNW with these altered weights, adjust by too of it, and so on. * Python demo.