

EngSci 721

Inverse Problems and Learning From Data

Oliver Maclaren (oliver.maclaren@auckland.ac.nz)

1. Basic concepts [5 lectures + 1 Tutorial]

Forward vs inverse problems. Well-posed vs ill-posed problems. Algebra and calculus of inverse problems (left and right inverses, generalised and pseudo inverses, resolution operators, matrix calculus). Representing higher dimensional problems (image data etc).

2. Instability and regularisation in linear and nonlinear problems [6 lectures + 1 Tutorial]

Instability and related issues for generalised inverses. Introduction to regularisation and trade-offs. Tikhonov regularisation. Higher-order Tikhonov regularisation. Sparsity and regularisation using different norms. Truncated singular value decomposition. Iterative regularisation, including stochastic/mini-batch gradient descent.

3. Further topics [3 lectures + 1 Tutorial]

Regularisation parameter choice, including statistical and machine learning views of regularisation. Confidence sets for linear and nonlinear models. Physics-informed machine learning and neural networks.

Module overview

Inverse Problems and Learning From Data (*Oliver Maclaren*)

[~14 lectures/3 tutorials]

Lecture 14: Neural models (Nonlinear part II)

Topics:

- An attempt to learn about neural networks via matrix calculus!

EngSci 721 : Lecture 14: Neural models

i.e. Nonlinear II

- An attempt to apply matrix calculus & understand neural nets ^

See also (loosely based on)

Linear models (again)

Consider

$$y_i = w_0 + w_1 x_i \quad \text{for observations } i=1, \dots, m$$

i.e:

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_m \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$$

The Matrix Calculus You Need For Deep Learning

[Terence Parr](#) and [Jeremy Howard](#)

(Terence is a tech lead at Google and ex-Professor of computer/data science in [University of San Francisco's MS in Data Science program](#). You might know Terence as the creator of the ANTLR parser generator. For more material, see Jeremy's [fast.ai courses](#) and University of San Francisco's Data Institute [in-person version of the deep learning course](#).)

Please send comments, suggestions, or fixes to [Terence](#).

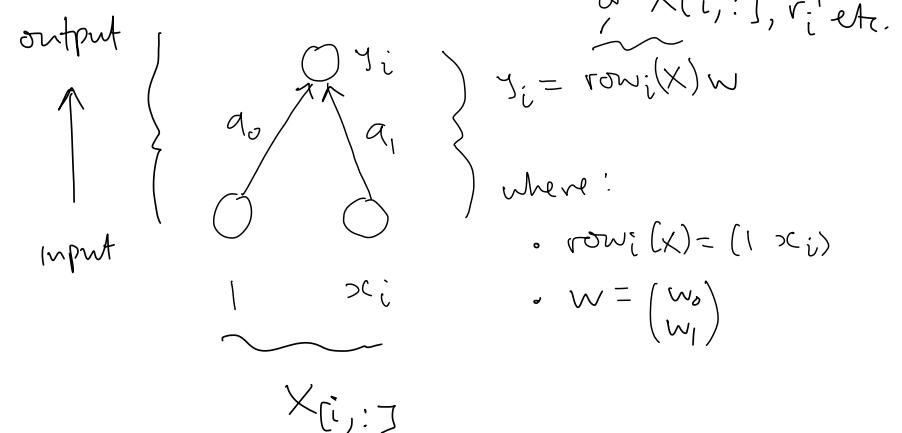
[Printable version](#) (This HTML was generated from markup using [bookish](#)). A [Chinese version](#) is also available (content not verified by us).

Abstract

This paper is an attempt to explain all the matrix calculus you need in order to understand the training of deep neural networks. We assume no math knowledge beyond what you learned in calculus 1, and provide links to help you refresh the necessary math where needed. Note that you do **not** need to understand this material before you start learning to train and use deep learning in practice; rather, this material is for those who are already familiar with the basics of neural networks, and wish to deepen their understanding of the underlying math. Don't worry if you get stuck at some point along the way—just go back and reread the previous section, and try writing down and working through some examples. And if you're still stuck, we're happy to answer your questions in the [Theory category at forums.fast.ai](#). Note: There is a [reference section](#) at the end of the paper summarizing all the key matrix calculus rules and terminology discussed here.

<https://explained.ai/matrix-calculus/>

We can write as a 'network diagram':



Constant or 'bias' terms

Given an $m \times n$ matrix X , let's introduce the notation

$$\bar{X} = [I_m \ X] = \text{hcat}((I_m, X))$$

$$= \begin{bmatrix} & \\ & \vdots & X \\ & & \end{bmatrix} \quad (\text{Think 'closure' or 'completion' of } X \text{ by } 1)$$

where I_m is an $m \times 1$ column

vector of ones. i.e we have added

an extra column of all ones to X

giving an $m \times (n+1)$ matrix, or

an $m \times \bar{n}$ matrix, where $\bar{n} = n+1$

Also, given a coefficient or weight vector

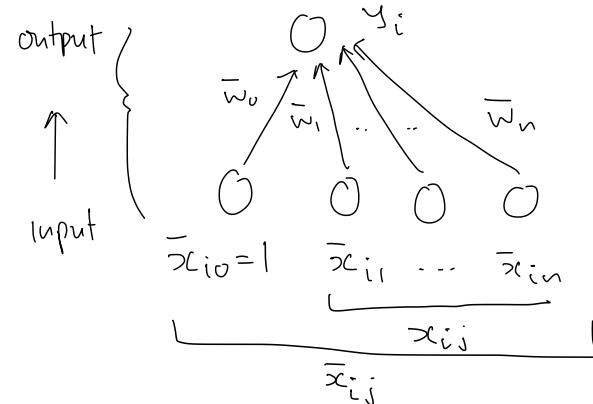
w of length n , define

$$\bar{w} = \begin{bmatrix} w_0 \\ w \end{bmatrix} \left\{ \begin{array}{l} 1 \\ \vdots \\ n \end{array} \right\} \quad n+1 = \bar{n}$$

Thus, $y_i = w_0 + w_i x_i, i=1,\dots,m$ becomes

$$m \times 1 \left\{ \overbrace{\begin{bmatrix} y = [1 \ X] \begin{bmatrix} w_0 \\ w \end{bmatrix} = \bar{X} \bar{w} \end{bmatrix}}^1 \right\}$$

More generally, write



for $x_{ij} = i^{\text{th}} \text{ obs. of } j^{\text{th}} \text{ feature}, j=1, \dots, n$

$\bar{x}_{ij} = \text{" } \text{" } \text{" } , j=0, \dots, n$

i.e. $\bar{x}_{ij} = x_{ij}, j=1, \dots, n$

& $\bar{X} = [\bar{x}_{ij}]$ matrix

$$\Rightarrow \overbrace{\begin{bmatrix} y = \bar{X} \bar{w} \end{bmatrix}}^1 \quad \underbrace{m \times (\bar{m} \times \bar{n}) \cdot (\bar{n})}$$

→ so far just a linear model

Nonlinearity & element-wise functions

Next consider a nonlinear transformation of the output of the linear part instead:

$$y_i = f(w_0 + w_1 x_i), \text{ for nonlinear } f,$$

$$\text{or } y_i = f(\text{row}_i(\bar{x})\bar{w}).$$

For vector inputs, $z = \begin{pmatrix} z_1 \\ z_2 \\ \vdots \\ z_m \end{pmatrix}$ say, (or z_0, z_m etc)

define f to act [element-wise]

$$f \text{ for vector input} \quad \left\{ \begin{array}{l} f(z) := \begin{pmatrix} f(z_1) \\ f(z_2) \\ \vdots \\ f(z_m) \end{pmatrix} \quad (\text{similarly for row vectors}) \\ \text{for matrices} \end{array} \right.$$

Then

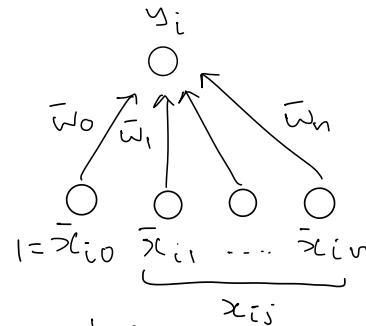
$$\boxed{y = f(\bar{x}\bar{w})} \quad \begin{matrix} \uparrow \\ f \text{ for scalar input} \end{matrix}$$

Defines a model which is part linear with a subsequent elementwise nonlinear transformation of the output

→ see later for more on element-wise f .

Neurons!

(analogy to brain blah blah)



represents:

$$y_i = f(\text{row}_i(\bar{x})\bar{w})$$

i.e. $y = f(\bar{x}\bar{w})$ for vector $y: \frac{\text{dim}}{m}$ matrix $\bar{x}: m \times n$ vector $\bar{w}: n$

\bar{w} : weights

\bar{x} : inputs/features

f : nonlinear activation function

Note : If f is Heaviside, then 'perceptron' (1958)

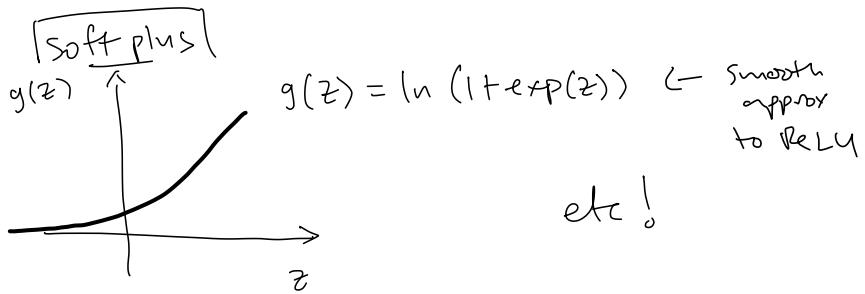
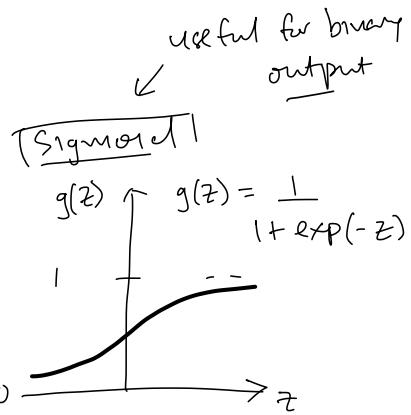
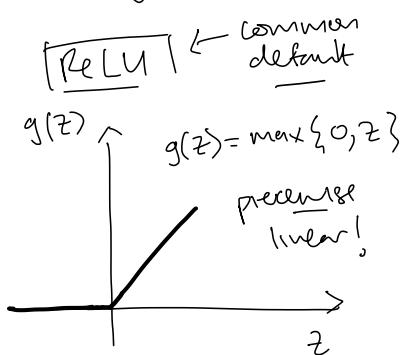
If f sigmoid, is 'logistic regression'

Nonlinearity

If f is linear then its just a linear model. This is also true for more complex models.

→ f should be nonlinear to be more exciting!

Typical:



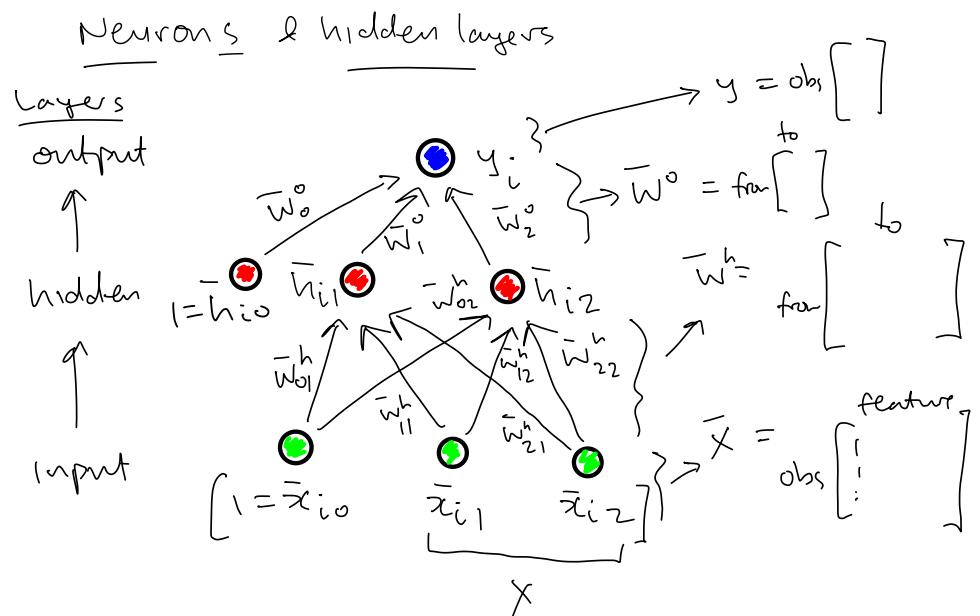
Neurons & hidden layers

A neural network typically contains many, many neurons ('regressions'!)

→ a feedforward neural network
(also called multi-layer perceptrons)

consists of a series of layers;
each with multiple neurons:

- each 'layer' has multiple nonlinear (linear) regressions
- the 'next' layer uses the outputs of the previous layer as the inputs/features for its collection of regressions
- layers that aren't observable inputs or outputs are called hidden
- 'deep': many hidden layers



First hidden layer : or 'bias' term when weighted

$$\begin{cases} \bar{h}_{i0} = 1 = x_{i0} \text{ (constant feature)} \\ \bar{h}_{i1} = f(\text{row}_i(\bar{x}) \bar{w}^h[:, 1]) = h_{i1} \\ \bar{h}_{i2} = f(\text{row}_i(\bar{x}) \bar{w}^h[:, 2]) = h_{i2} \end{cases}$$

$$\Rightarrow \boxed{\boxed{\boxed{\quad}}} = f\left(\boxed{1} \quad \boxed{\boxed{\boxed{\quad}}}\right)$$

$$\boxed{\boxed{\boxed{\quad}}} = \boxed{\boxed{\quad}} f\left(\boxed{1} \quad \boxed{\boxed{\boxed{\quad}}}\right)$$

dim: $\begin{cases} \bar{x} : m \times \bar{n} \\ \bar{w}^h : \bar{n} \times p \end{cases}$

So: $\bar{H}[i, 0] = \bar{h}_{i0} = 1$

$\bar{H}[i, 1:2] = [\bar{h}_{i1} \ \bar{h}_{i2}] = f(\text{row}_i(\bar{x}) \bar{w}^h)$

$$\Leftrightarrow \bar{H} = \boxed{i \downarrow \left[\begin{array}{c|c} \boxed{1} & f(\bar{x} \bar{w}^h) \\ \vdots & \vdots \\ 1 & \end{array} \right]}$$

$$= \boxed{\left[\begin{array}{c} \boxed{1} \\ \vdots \\ 1 \end{array} \right] f\left(\begin{array}{c|c} \bar{x} & \bar{w}^h \\ \hline \boxed{1} & \boxed{\boxed{\boxed{\quad}}} \end{array} \right)}$$

$$= \boxed{\left[\begin{array}{c} \boxed{1} \\ \vdots \\ 1 \end{array} \right] f\left(\begin{array}{c} \boxed{\quad} \\ \vdots \\ \boxed{\quad} \end{array} \right) f\left(\begin{array}{c} \boxed{\quad} \\ \vdots \\ \boxed{\quad} \end{array} \right)}$$

$\bar{x} \bar{w}^h[:, 1]$ $\bar{x} \bar{w}^h[:, 2]$

since f element-wise
(& applies to col element-wise too)

Also $H = f(\bar{x} \bar{w}^h)$

↑ completed
not completed
(need extra 1 vector for \bar{H})

Hidden layer cont'd

We've arranged as collection of individual regressions, one per hidden neuron.

Note though:

$$\bar{H} = [I f(\bar{x}\bar{w}^h)] = [I f(\bar{x}\bar{w}^h I)]$$

where I = vector of ones

I = identity matrix

$$\text{vec}(\bar{H}) = \begin{bmatrix} I \\ \text{vec}(f(\bar{x}\bar{w}^h I)) \end{bmatrix}$$

$$\text{vec}(\bar{x}\bar{w}^h I) = (I \otimes \bar{x}) \text{vec}(\bar{w}^h)$$

& $\text{vec}(f(\cdot)) = f(\text{vec}(\cdot))$ if f element-wise

$$(I \otimes \bar{x}) \text{vec}(\bar{w}^h) = \begin{bmatrix} \bar{x} & 0 \\ 0 & \bar{x} \end{bmatrix} \begin{bmatrix} \bar{w}_{01}^h \\ \bar{w}_{11}^h \\ \bar{w}_{21}^h \\ \bar{w}_{02}^h \\ \bar{w}_{12}^h \\ \bar{w}_{22}^h \end{bmatrix}$$

\Rightarrow

So:

$$\text{vec}(\bar{H}) = \begin{bmatrix} I \\ \bar{H}[:, 1] \\ \bar{H}[:, 2] \end{bmatrix} = \begin{bmatrix} I \\ \bar{h}_1 \\ \bar{h}_2 \end{bmatrix}$$

constant feature for use
as input to next
layer

element-wise nonlinear

where:

$$\begin{bmatrix} \bar{H}[:, 1] \\ \bar{H}[:, 2] \end{bmatrix} = \tilde{f} \left(\begin{bmatrix} \bar{x} & 0 \\ 0 & \bar{x} \end{bmatrix} \begin{bmatrix} \bar{w}_{01}^h \\ \bar{w}_{11}^h \\ \bar{w}_{21}^h \\ \bar{w}_{02}^h \\ \bar{w}_{12}^h \\ \bar{w}_{22}^h \end{bmatrix} \right)$$

linear

$$(\Leftrightarrow H = f((I \otimes \bar{x}) \text{vec}(\bar{w}^h)) \Rightarrow$$

i.e. one big 'linear' regression for weights of prev. layer
... with a nonlinear transformation of output ... & ...
an extra constant vector for later

(conceptually, if not computationally)

Output layer

For each observation i , we use the outputs of the hidden layer (here) as the inputs of the output layer

→ do same again! single output here though

→ possibly different activation function, g

→ include constant feature ('bias') from prev. layer x weight.

$$\boxed{y_i = g(\bar{H}_{[i,:]} \bar{w}^o)}$$

w^o : w^{output}
: $\bar{P} \times 1$

$$\boxed{y = g(\bar{H} \bar{w}^o)}$$

$$i \downarrow \begin{bmatrix} \square \\ \square \\ \vdots \\ \square \end{bmatrix} = g \left(\begin{bmatrix} 1 \\ \bar{x} \\ \vdots \\ 1 \end{bmatrix}, \begin{bmatrix} \square \end{bmatrix} \right)$$

} single
'regression'
for single
output node

Binary classification: $y \in \{0, 1\}$, class probability,

→ use e.g. logistic, range = $(0, 1)$

Composite

$$\begin{aligned} y &= g(\bar{H} \bar{w}^o) \\ H &= f(\bar{x} \bar{w}^h) \end{aligned} \quad \left. \begin{array}{l} \text{Note:} \\ y = g([I + H] \bar{w}^o) \\ H = f(\bar{x} \bar{w}^h) \end{array} \right\}$$

$$\text{so } \boxed{y = g([I f(\bar{x} \bar{w}^h)] \bar{w}^o)} \quad | \quad \left. \begin{array}{l} \text{model} \\ \text{prediction} \end{array} \right\}$$

unknowns \bar{w}^h, \bar{w}^o } $\bar{n}p + \bar{p} = npt + pti$
given \bar{x} unknowns
outputs y

want $\boxed{D_{\bar{w}^h} y}$ & $\boxed{D_{\bar{w}^o} y}$

i.e. $D_{\bar{w}} y = [D_{\bar{w}^h} y \ D_{\bar{w}^o} y]$

→ derivative of prediction wrt param.

→ combine with ...

'Loss function' $\ell(y; d)$ (minimise)

compare $y(x)$ output to data d associated
with x , i.e. (x, d) is full data.

[call $d(x)$ & compare $y(x) \& d(x)$]

→ sum of squares (even for
classification!)

→ $-\log$ likelihood = $-\log$ prob. of
data given
param.

→ can use binary prob. dist.
with y = prob. parameter

$$\ell(\bar{w}; d) = \ell(y(\bar{w}); d)$$

want

$$\begin{aligned} D_{\bar{w}} \ell &= D_y \ell \cdot D_{\bar{w}} y \\ \text{vector} &\quad \text{vector} \quad \text{matrix (Jacobian)} \\ \bar{w} &= \text{array of } \underline{\text{all}} \text{ weights} \end{aligned}$$

Loss: simple case

$$\begin{aligned} \ell(y; d) &= \frac{1}{2} \|y - d\|^2 \\ &= \frac{1}{2} (y - d)^T (y - d) \\ &= \frac{1}{2} (y^T y - 2d^T y + d^T d) \end{aligned}$$

$$\begin{aligned} D_y \ell &= \frac{1}{2} (2y^T - 2d^T) \\ &= y^T - d^T \\ &= (y - d)^T \quad [\text{done}] \end{aligned}$$



Hard part: $D_{\bar{w}^h} y$!

Warm-up

$$\text{derivatives of } \left\{ \begin{array}{l} y = \underbrace{[1 \ H(\bar{w}^h)]}_{H} \bar{w}^o \\ H = \bar{x} \bar{w}^h \end{array} \right.$$

$$\begin{aligned} y &: m \times 1 \\ H &: m \times p \\ \bar{H} &: m \times \bar{p} \\ \bar{w}^h &: \bar{n} \times p \\ \bar{x} &: m \times \bar{n} \\ \bar{w}^o &: \bar{p} \end{aligned}$$

if no activation functions:



$$\begin{aligned} 1. \text{ via differentials:} \\ \bar{H} &= [1 \ \bar{x} \bar{w}^h] \\ \Rightarrow d\bar{H} &= [0 \ \bar{x} d\bar{w}^h] \\ &= [0 \ d\bar{H}] \end{aligned}$$

$$dy = \bar{H} d\bar{w}^o + d\bar{H} \bar{w}^o$$

$$= \bar{H} d\bar{w}^o + [0 \ \bar{x} d\bar{w}^h] \bar{w}^o$$

$$= \bar{H} d\bar{w}^o + \bar{x} d\bar{w}^h w^o, w^o = \begin{bmatrix} w_1^o \\ \vdots \\ w_p^o \end{bmatrix}$$

$$= \bar{H} \underbrace{d\bar{w}^o}_{\parallel} + (w^o \otimes \bar{x}) d\text{vec}(\bar{w}^h)$$

$$= \bar{H} \widetilde{d\text{vec}}(\bar{w}^o) + (w^o \otimes \bar{x}) \widetilde{d\text{vec}}(\bar{w}^h)$$

$$\Rightarrow D_{\bar{w}^o} y = \bar{H}, \quad D_{\bar{w}^h} y = w^o \otimes \bar{x}$$

$$\begin{array}{cccc} \sim & \sim & \sim & \sim \\ m \times \bar{p} & m \times \bar{p} & m \times \bar{n} p & m \times \bar{n} p \end{array} \checkmark$$

(no activation functions)



$$\text{unknowns: } \begin{cases} \bar{p} = p+1 \\ \bar{n} p = (n+1)p \end{cases}$$

2. via vectorised chain rule

$$y = [1 \ H(\bar{w}^h)] \bar{w}^o = \bar{H} \bar{w}^o$$

$$H = \bar{x} \bar{w}^h, \bar{H} = [1 \ \bar{x} \bar{w}^h]$$

$$D_{\bar{w}^h} y = \frac{\partial \text{vec}(y)}{\partial \text{vec}(\bar{w}^o)} = \bar{H} \quad (\bar{w}^o \text{ already vec})$$

$$D_{\bar{w}^h} y = \frac{\partial \text{vec}(y)}{\partial \text{vec}(\bar{w}^h)} = \frac{\partial \text{vec}(y)}{\partial \text{vec} H} \frac{\partial \text{vec} H}{\partial \text{vec} \bar{w}^h}$$

$$y = 1 w^o + H(\bar{w}^h) w^o$$

$$\text{vec}(y) = \text{vec}(1 w^o) + \text{vec}(H(\bar{w}^h) w^o)$$

$$= 1 w^o + (w^o \otimes I_m) \text{vec}(H)$$

$$\frac{\partial \text{vec}(y)}{\partial \text{vec}(H)} = w^o \otimes I_m \quad \} \text{ } m \times pm$$

$$\frac{\partial \text{vec}(H)}{\partial \bar{w}^h} = I_p \otimes \bar{x} \quad \} \text{ } pm \times p \bar{n}$$

$$\begin{aligned} D_{\bar{w}^h} y &= (w^o \otimes I_m)(I_p \otimes \bar{x}) \\ &= (w^o \otimes I_p) \otimes (I_m \bar{x}) \end{aligned} \quad \left[\begin{array}{l} \text{Kronecker} \\ (A \otimes B)(C \otimes D) \\ A \otimes B = C \otimes D \end{array} \right]$$

$$= \underbrace{w^o \otimes \bar{x}}_{(m \times p \bar{n})} \checkmark$$

$$\bar{w}^h : \bar{n} \times p$$

$$\bar{w}^o : \bar{p}$$

$$\bar{H} : m \times \bar{p}$$

Activation functions?

$$y = g(\bar{H}(\bar{w}^h) \bar{w}^o)$$

$$\bar{H} = [1 \ f(\bar{x} \bar{w}^h)] = [1 \ H(\bar{x} \bar{w}^h)]$$

Here f & g act element-wise

→ Need to be careful when
calc. derivatives →

Element-wise functions & derivatives

we have considered 'element-wise' functions by defining for scalar input function $\phi(z)$,

$$\phi(x) = \begin{bmatrix} \phi(x_1) \\ \phi(x_2) \\ \phi(x_3) \end{bmatrix}$$

for vector input x .

More explicitly we could write

$$\tilde{\phi}(x) = \begin{bmatrix} \tilde{\phi}_1(x_1, x_2, x_3) \\ \tilde{\phi}_2(x_1, x_2, x_3) \\ \tilde{\phi}_3(x_1, x_2, x_3) \end{bmatrix}$$

$\tilde{\phi}$ is ϕ for vector input

$$\text{where } \tilde{\phi}_1(x_1, x_2, x_3) := \phi(x_1)$$

$$\tilde{\phi}_2(x_1, x_2, x_3) := \phi(x_2)$$

$$\tilde{\phi}_3(x_1, x_2, x_3) := \phi(x_3) \text{ etc.}$$

$\tilde{\phi}$ is our original function of a scalar var.

Element-wise functions & derivatives

Note then that the derivative of our elementwise function of a vector is

$$D_x \tilde{\phi}(x) = \frac{\partial \tilde{\phi}(x)}{\partial x^T}$$

$$= \begin{bmatrix} \frac{\partial \tilde{\phi}_1}{\partial x_1} & \frac{\partial \tilde{\phi}_1}{\partial x_2} & \frac{\partial \tilde{\phi}_1}{\partial x_3} \\ \frac{\partial \tilde{\phi}_2}{\partial x_1} & \frac{\partial \tilde{\phi}_2}{\partial x_2} & \frac{\partial \tilde{\phi}_2}{\partial x_3} \\ \frac{\partial \tilde{\phi}_3}{\partial x_1} & \frac{\partial \tilde{\phi}_3}{\partial x_2} & \frac{\partial \tilde{\phi}_3}{\partial x_3} \end{bmatrix}$$

$$= \begin{bmatrix} \phi'(x_1) & 0 & 0 \\ 0 & \phi'(x_2) & 0 \\ 0 & 0 & \phi'(x_3) \end{bmatrix}$$

Define $\phi(x) = \begin{bmatrix} \phi'(x_1) \\ \phi'(x_2) \\ \phi'(x_3) \end{bmatrix}$ scalar derivative, i.e. $\frac{d\phi}{dx}$ of $\phi(x)$

↑ vector arg.
scalar deriv

Then $D_x \tilde{\phi}(x) = \text{diag}(\phi'(x))$ { diag(x) = $\begin{bmatrix} x_1 & 0 & 0 \\ 0 & x_2 & 0 \\ 0 & 0 & x_3 \end{bmatrix}$ etc }

& $d\tilde{\phi}(x; dx) = \text{diag}(\phi'(x)) dx$

{ usually just write $\phi(x)$ cf. $\tilde{\phi}(x)$ }

functions of a matrix?

Element-wise function of a matrix A ($m \times n$)

$$\tilde{f}(A) := \begin{bmatrix} f(A_{11}) & f(A_{12}) & \dots \\ \vdots & & \\ f(A_{m1}) & \dots & f(A_{mn}) \end{bmatrix}$$

$$D_A \tilde{f}(A) = \frac{\partial \text{vec}(\tilde{f}(A))}{\partial \text{vec}(A)^T} \quad \left. \right\} \text{dim: } mn \times mn$$

$$= \frac{\partial (\text{vec}(\tilde{f}(A)))}{\partial \text{vec}(A)^T}$$

$$= \text{diag}(\text{vec}(f'(A)))$$

where $f'(A)$ $\left. \right\}$ matrix $\xrightarrow{\text{vec}}$ vector
 Scalar $f'(z) = \frac{df(z)}{dz}$ applied elementwise to A

$$\left[\text{diag}(mn) \rightarrow mn \times mn \right]$$

defn: vectorised form of elementwise function of matrix input:

$$\boxed{\tilde{F}(\text{vec}(A)) = \text{vec}(\tilde{f}(A))}$$

$$= \begin{bmatrix} f(A_{11}) \\ f(A_{21}) \\ \vdots \\ f(A_{12}) \\ \vdots \\ f(A_{mn}) \end{bmatrix}$$

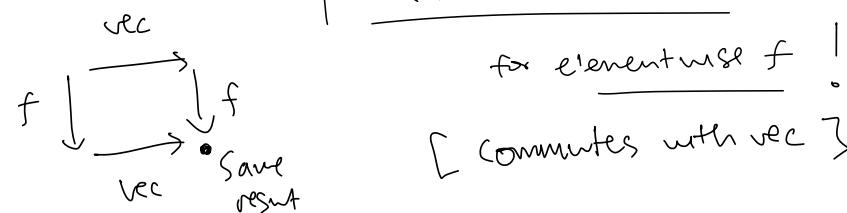
$$\boxed{= \tilde{f}(\text{vec}(A))} \quad \begin{array}{l} \text{def for} \\ \text{vector} \\ \text{input} \end{array}$$

for element-wise f !

usually write $\tilde{f}(\text{vec}(A))$

$$\text{as } \boxed{f(\text{vec}(A))}$$

$$\text{so } \boxed{\text{vec}(f(A)) = f(\text{vec}(A))} \neq$$



for elementwise f !

$\left[\text{commutes with vec} \right]$

So, for element-wise, we have:

$$\begin{aligned}
 d\tilde{f}(A) &= d f(\text{vec}(A)) = d \text{vec}(f(A)) \\
 &= \underbrace{\text{diag}\left(f'(\text{vec}(A))\right)}_{\substack{\text{Scalar} \\ \text{derivative}}} d \text{vec } A \underbrace{\quad}_{\substack{\text{applied element-wise} \\ \text{to vec}(A)}}
 \end{aligned}$$

$$\nabla_A \tilde{f}(A) = \text{diag}(\text{vec}(f'(A)))$$

$$= \text{diag}(f'(\text{vec}(A)))$$

↗ ↗
 scalar derivative applied
 derivative to
 vec A

$$\text{So...} \quad \text{vector} \quad \bar{w}^o = \begin{bmatrix} \bar{p} \\ \vdots \end{bmatrix}, \quad \bar{w}^h = \begin{bmatrix} p \\ \vdots \end{bmatrix}$$

$$y = g(\overline{H}(\bar{w}^h) \bar{w}^o) = g([1 \ H] \bar{w}^o)$$

$$\Rightarrow dy = \text{diag}(g'(\bar{H}(\bar{w}^h) \bar{w}^o)) [d\bar{H} \bar{w}^o + \bar{H} d\bar{w}^o]$$

$$= \left\{ \begin{array}{l} \text{diag}\left(g'(\bar{H}(\bar{w}^h) \bar{w}^o)\right) d\bar{H} \bar{w}^o \\ + \\ \underbrace{\text{diag}\left(g'(\bar{H}(\bar{w}^h) \bar{w}^o)\right) \bar{H} d\bar{w}^o}_{D_{\bar{w}^o}} \end{array} \right\} ? \rightarrow D_{\bar{w}^h} \bar{w}^h$$

Now:

$$\text{vec} \left(\text{diag} \left(g'(\bar{H}(\bar{w}^n)\bar{w}^o) \right) d\bar{H}\bar{w}^o \right)$$

$$= \text{diag} \left(g'(\bar{H}(\bar{w}^u) \bar{w}^o) \right) d\bar{H}\bar{w}^o$$

(as already m-dim vector)

$$\text{Also } d\bar{H} = [0 \quad dH], \quad \bar{w}^o = [w^o_0 \quad w^o]^T$$

$$\Rightarrow d\bar{H}\bar{W}^o = dHW^o \quad \Rightarrow$$

$$\bar{H} = m \left[\begin{array}{c} \\ \end{array} \right] \quad \bar{\omega}^o = \bar{p} \left[\begin{array}{c} \\ \end{array} \right]$$

$$H = \bar{x} \bar{w}^h$$

$$\Rightarrow \text{diag}(g'(\bar{H}(\bar{w}^h) \bar{w}^o)) d\bar{H} \bar{w}^o$$

$$= \text{diag}(g'(\bar{H}(\bar{w}^h) \bar{w}^o)) dH w^o$$

$$= \left[(\bar{w}^o)^T \otimes \text{diag}(g'(\bar{H}(\bar{w}^h) \bar{w}^o)) \right] d\text{vec}(H)$$

{ by Kronecker prod rule }.

Now:

$$\left[(\bar{w}^o)^T \otimes \text{diag}(g'(\bar{H}(\bar{w}^h) \bar{w}^o)) \right] d\text{vec}(H)$$

$$\overbrace{d\text{vec}(H)} = d\text{vec}(f(\bar{x} \bar{w}^h))$$

$$= \text{diag}(f'(\text{vec}(\bar{x} \bar{w}^h)) d(\text{vec}(\bar{x} \bar{w}^h)))$$

$$\& \text{vec}(\bar{x} \bar{w}^h) = \text{vec}(\bar{x} \bar{w}^h \mathbb{I}_P)$$

$$= (\mathbb{I}_P \otimes \bar{x}) d\text{vec}(\bar{w}^h)$$

$$\therefore d\text{vec}(H) =$$

$$\underbrace{\text{diag}(f'(\text{vec}(\bar{x} \bar{w}^h)) (\mathbb{I}_P \otimes \bar{x}) d\text{vec}(\bar{w}^h))}_{\rightarrow}$$

So,

$$\begin{aligned} & \text{diag}(g'(\bar{H}(\bar{w}^h) \bar{w}^o)) d\bar{H} \bar{w}^o \\ &= \left[(\bar{w}^o)^T \otimes \text{diag}(g'(\bar{H}(\bar{w}^h) \bar{w}^o)) \right] d\text{vec}(H) \\ &= \end{aligned}$$

$$\left\{ \left[(\bar{w}^o)^T \otimes \text{diag}(g'(\bar{H}(\bar{w}^h) \bar{w}^o)) \right] \text{diag}(f'(\text{vec}(\bar{x} \bar{w}^h)) (\mathbb{I}_P \otimes \bar{x})) \cdot d\text{vec}(\bar{w}^h) \right\}$$

$$\Rightarrow D_{\bar{w}^h}(y) = \left[(\bar{w}^o)^T \otimes \text{diag}(g'(\bar{H}(\bar{w}^h) \bar{w}^o)) \right] - \underbrace{\left[\text{diag}(f'(\text{vec}(\bar{x} \bar{w}^h)) (\mathbb{I}_P \otimes \bar{x})) \right]}_{\text{phen!}} \quad \text{---} \quad \text{times}$$

$$\bar{x} \in \mathbb{R}^p$$

sense

check?

—

$$\boxed{\text{if}} \quad g(z) = z \\ g'(z) = 1, \quad \}$$

then e.g.

$$\text{diag}(g'(\mathbf{x})) = \begin{bmatrix} 1 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 \end{bmatrix} = I_{\text{len}(\mathbf{x})}$$

etc

$$\boxed{\text{So}} \quad \bar{H}\bar{w}^o \quad \} \quad g'(\bar{H}\bar{w}^o) = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \}_{\mathbf{m}}$$

dim: $(m \times p)(p) \ni m$

$$\Rightarrow \text{diag}(g'(\bar{H}(\bar{w}^o)\bar{w}^o)) = I_m$$

$$\Rightarrow (\bar{w}^o)^T \otimes \text{diag}(g'(\bar{H}(\bar{w}^o)\bar{w}^o))$$

$$= (\bar{w}^o)^T \otimes I_m$$

$$\text{Also } \boxed{\text{if}} \quad f(z) = z, \quad f'(z) = 1$$

$$\Rightarrow f'(\underbrace{\text{vec}(\bar{x}\bar{w}^h)}_{mp}) = I_{mp} \quad \} \text{ one vector length } mp$$

$$\Rightarrow \text{diag}(f'(\text{vec}(\bar{x}\bar{w}^h)))$$

$$= I_{mp}$$

so ...

$$D_{\bar{w}^h} = ((\bar{w}^o)^T \otimes I_m) I_{mp} (I_p \otimes \bar{x})$$

$$= ((\bar{w}^o)^T \otimes I_m) (I_p \otimes \bar{x})$$

$$= (\bar{w}^o)^T I_p \otimes I_m \bar{x}$$

$$= (\bar{w}^o)^T \otimes \bar{x} \quad / \quad \begin{array}{l} \text{(no activation} \\ \text{function} \\ \text{result)} \end{array}$$

Now we can compute

$$D_{\bar{w}} \lambda(y(\bar{w}), d) =$$

& use (mini-batch) gradient descent

$$\bar{w}(k+1) = \bar{w}(k) - \eta D_{\bar{w}} d$$

to find weights!

First, note:

$$D_y^1 = (y - d)^T \quad \} (1 \times m)$$

$$D_{\bar{W}^0}^{(k)} = \text{diag} \left(g'(\bar{H}(\bar{W}^k)\bar{W}^0) \right) \bar{H} \Big\}^{(m \times m)(m \times p)}$$

$$\mathbb{D}_{\bar{w}^h} = \left[\begin{array}{c} \left(\mathbb{W}^\circ \right)^\top \otimes \text{diag} \left(g'(\bar{H}(\bar{w}^h)\bar{w}^v) \right) \\ \vdots \\ \text{diag} \left(f'(\text{vec}(\bar{X}\bar{w}^h)) \right) (\mathbb{I}_p \otimes \bar{X}) \end{array} \right] \underbrace{\text{times}}_{\text{times}} \left. \begin{array}{c} (\text{mxmp}) \\ \\ (\text{mp} \times \text{mp}) \\ \text{times} \end{array} \right.$$

Backpropagation & adjoint/reverse mode diff.

To compute the derivatives via the chain rule we end up needing to compute:

$A B C$ or $A B C D$

for various size arrays (vectors & matrices)

→ while the final answer doesn't depend on the order we multiply these eg $A(BC)$ or $(AB)C$, the computational cost does!

→ in general, want e.g. vector times matrix instead of matrix times matrix

Note eg

$$\begin{aligned}
 \text{Note e.g. } & \quad \text{smaller} \\
 \underline{\underline{D_1}}_{\text{no}} & = \underbrace{\text{vector} \times \text{matrix} \times \text{matrix}}_{\text{vector}} \\
 & = \left\{ \underbrace{\text{vector} \times \text{matrix}}_{\text{vector} \times (\text{matrix} \times \text{matrix})} \right\} \times \text{matrix}
 \end{aligned}$$

\Rightarrow If output is small, use "backwards order"
 \rightarrow "backpropagation"!
 \rightarrow also called adjoint / reverse mode differentiation

Exercises (brave!)

- Implement what we derived (I might later...)

- Compare to eg Pytorch

↳ does at 'lower' operator level!

↳ automates for you!

- Read 'The matrix calculus you need for machine learning'

by Parr & Howard (2018)

on canvas &

'explained.ai/matrix-calculus/'

online

- Other refs

- Maynuss & Nendecker (1999, 2019)
(Matrix differential calculus)

- Leedfellow, Bengio, Courville (2016)

'Deep learning'

Physics-informed neural nets (PINNs)

→ combine neural nets + ϵg differential eqn penalty } or in assignment?
(NN instead of spline)

Scientific Machine Learning Through Physics-Informed Neural Networks: Where we are and What's Next

Salvatore Cuomo¹ · Vincenzo Schiano Di Cola² · Fabio Giampaolo¹ ·
Gianluigi Rozza³ · Maziar Raissi⁴ · Francesco Piccialli¹ 

Received: 14 January 2022 / Revised: 6 July 2022 / Accepted: 6 July 2022 /
Published online: 26 July 2022
© The Author(s) 2022

Eg:



Abstract
Physics-Informed Neural Networks (PINNs) are neural networks (NNs) that encode model equations, like Partial Differential Equations (PDE), as a component of the neural network itself. PINNs are nowadays used to solve PDEs, fractional equations, integral-differential equations, and stochastic PDEs. This novel methodology has arisen as a multi-task learning framework in which a NN must fit observed data while reducing a PDE residual. This article provides a comprehensive review of the literature on PINNs: while the primary goal of the study was to characterize these networks and their related advantages and disadvantages, The review also attempts to incorporate publications on a broader range of collocation-based physics informed neural networks, which stars form the vanilla PINN, as well as many other variants, such as physics-constrained neural networks (PCNN), variational hp-VPINN, and conservative PINN (CPINN). The study indicates that most research has focused on customizing the PINN through different activation functions, gradient optimization techniques, neural network structures, and loss function structures. Despite the wide range of applications for which PINNs have been used, by demonstrating their ability to be more feasible in some contexts than classical numerical techniques like Finite Element Method (FEM), advancements are still possible, most notably theoretical issues that remain unresolved.

Keywords Physics-Informed Neural Networks · Scientific Machine Learning · Deep Neural Networks · Nonlinear equations · Numerical methods · Partial Differential Equations · Uncertainty

Physics-informed neural networks (PINNs) for fluid mechanics: a review

Shengze Cai¹ · Zhiping Mao² · Zhicheng Wang³ · Minglang Yin^{4,5} · George Em Karniadakis^{1,4}

Received: 21 May 2021 / Accepted: 10 September 2021 / Published online: 23 January 2022
© The Chinese Society of Theoretical and Applied Mechanics and Springer-Verlag GmbH Germany, part of Springer Nature 2021



Abstract

Despite the significant progress over the last 50 years in simulating flow problems using numerical discretization of the Navier-Stokes equations (NSE), we still cannot incorporate seamlessly noisy data into existing algorithms, mesh-generation is complex, and we cannot tackle high-dimensional problems governed by parametrized NSE. Moreover, solving *inverse flow problems* is often prohibitively expensive and requires complex and expensive formulations and new computer codes. Here, we review *flow physics-informed learning*, integrating seamlessly data and mathematical models, and implement them using physics-informed neural networks (PINNs). We demonstrate the effectiveness of PINNs for inverse problems related to three-dimensional wake flows, supersonic flows, and biomedical flows.

Keywords Physics-informed learning · PINNs · Inverse problems · Supersonic flows · Biomedical flows