

(8) - convergence

BR: continue with training NN \rightarrow optimisation

- 14.5: key parts

- BR: goes over lecture 4.5 key parts again (due to lack of clarity)

- NN require training to rep. function (given architecture)

\hookrightarrow to minimise loss on training set

\hookrightarrow via gradient descent

\hookrightarrow backpropagation $\frac{\partial E_n}{\partial w_{ij}}$

gradient descent algo

- to minimise $f(x)$ wrt x

(*) initial guess

(*) At each guess, find point where function is decreasing with ✓

If derivative is positive \rightarrow take a step in -ve direction (??)

negative \rightarrow take a step in the +ve.

neural nets:

training error = $\frac{1}{T} \sum_t \text{Div}(y_t, d_t; w_1, w_2, \dots, w_R)$

u.s. vector formulation u(10) = need to review and not proceed until reviewed. ✓

\Rightarrow trivial form of forward propagation

dimensionality

backward propagation

BR: gradient descent update rule

$$\underline{x}^{R+1} = \underline{x}^R - \eta \nabla_{x^R} f^T \quad y = f(\underline{x})$$

(*) this dimensionality

check is a code-sanity check

w(13): check

scalar $\nabla_x y = x^T$ vector or mat

$$\nabla_{y_N} \text{Div} = \nabla_y \text{Div}$$

- $\boxed{ } \text{ Div} = ; \quad \nabla_{y_N} \text{ Div} \quad \boxed{ }$

y_N

- gradient / deriv. of divergence not output Div y_N

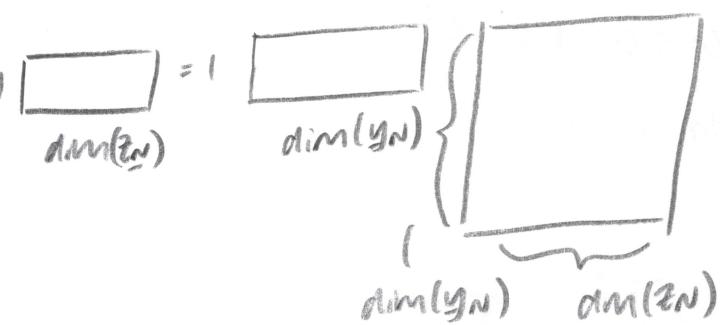
- introduction

Then take a step backward

compute gradient/div of Div w.r.t affine vector \underline{z}_N

$$\nabla_{\underline{z}_N} \text{Div} = \nabla_{\underline{y}_N} \text{Div} \cdot J_{\underline{y}_N}(\underline{z}_N)$$

(*) Take previously computed $\nabla_{\underline{y}_N} \text{Div}$,
postmultiply by Jacobian
 $J_{\underline{y}_N}(\underline{z}_N) \in \mathbb{R}^{\dim(\underline{y}_N) \times \dim(\underline{z}_N)}$



(*) dimensionality illus.

$$\nabla_{\underline{y}_{N-1}} \text{Div} = \nabla_{\underline{y}_N} \text{Div} \cdot J_{\underline{y}_N}(\underline{z}_N) \cdot W_N$$

(*) Take previously computed
 $\nabla_{\underline{z}_N}(\text{Div})$ and postmultiply
by weights matrix W_N (④) - check

Hence

$$\nabla_{\underline{y}_{N-1}} \text{Div} = 1 \begin{bmatrix} \dots \\ \dots \end{bmatrix} / \dim(\underline{y}_N)$$

(*) Each step back
is one matrix
mult. i.e.
 $W_N \in \mathbb{R}^{\dim(z_N) \times \dim(y_N)}$

Each backward step :-

i) Past an activation

→ postmultiply by Jacobian

ii) Past an ...

→ postmultiply by weights matrix (**) Review + check

In between, compute :- (of most int.)

$$\nabla_{\underline{w}_R} \text{Div} = \underline{y}_{R-1} \nabla_{\underline{z}_R} \text{Div}$$

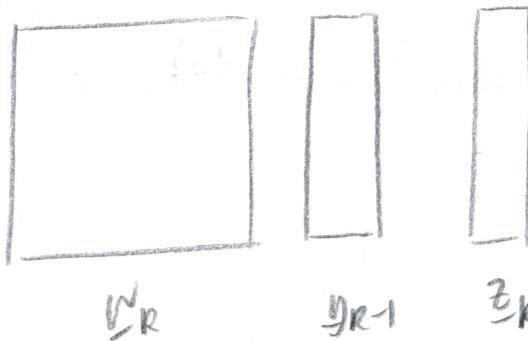
$$\nabla_{\underline{b}_R} \text{Div} = \nabla_{\underline{z}_R} \text{Div}$$

BR: $\dim(\underline{b}_R) = \dim(\underline{z}_R)$ as $\underline{z}_R = \underline{w}_R \underline{y}_{R-1} + \underline{b}_R$

but how does that
 $\Rightarrow \nabla_{\underline{b}_R} \text{Div} = \nabla_{\underline{z}_R} \text{Div} (?)$ w(15) → review

- vector y_{R-1}
- matrix w_R

(?)



scalar
 $\text{div}()$ from z_R

$$\nabla_{w_R} \text{div}() = y_{R-1} \nabla_{z_R} \text{div}()$$

- $\nabla_{z_R} \text{div}()$ already computed. (row vector $\in \mathbb{R}^{\text{dim}(z_R) \times 1}$)

- y_R - column vector $\in \mathbb{R}^{\text{dim}(y_R) \times 1}$

- Hence $\nabla_{w_R} \text{div} = y_{R-1} \nabla_{z_R} \text{div}() \in \mathbb{R}^{\text{dim}(y_R) \times \text{dim}(z_R)}$

$$w_R \in \mathbb{R}^{\text{dim}(z_R) \times \text{dim}(y_R)}$$

Recursive formulation.

(*)

then this

(*) entire backprop pseudocode:-

(BETTER VERSION)

Init: $y_N = Y$ $g_0 = z$ $\nabla_{y_N} \text{div} = \nabla_z \text{div}$

(*) supplement/summarise/
check understand

For $k = N, \dots, 1$:-

compute $\nabla_{y_k} (z_k)$

$$\nabla_{z_k} \text{div} = \nabla_{y_k} \text{div} \nabla_{y_k} (z_k) \quad \nabla_{y_{k-1}} \text{div} = \nabla_{z_k} \text{div} w_R \quad (\text{recursion})$$

$$\text{and } \nabla_{w_R} \text{div} = y_{k-1} \nabla_{z_R} \text{div}$$

(gradient
comp.)

$$\nabla_{w_R} \text{div} = \nabla_{z_R} \text{div}$$

(*) entire NN training algo \rightarrow check your understand logic

- KL divergence rather than L_2 divergence \rightarrow faster optim.

convergence: - will we get to location? } \rightarrow focus
- How long will it take

- network of MLP with thresholds \rightarrow not poss. to train
.. replace with diff. activation; replace error 'count' with divergence
- Allows differentiability, parameter tuning
- Does this work?
- Questions about convergence of gd.

②: Does backprop yield desirable results?

- 2 issues

i) changed definition of error

- loss is no longer error (in perptr. sense); but average of divergences
- see slides
- minimising loss does not necessarily mean minimising
classification error

③ - Quiz: missed slides (W) A6

backprop fails to sep...

(*) red line is global minimum; where $\text{loss}_{\text{fn}}/\text{has}$ minimum

(*) Add a further training instance;
will perceptron find linear class.

(*) consider add. of one more point and influence on loss for

i) perceptron

ii) backpropagation

(*) uncomfortable with hard-hairy exposition here \rightarrow ② solve?

(*) see slides

(*) consider sensitivity of estimated decision boundary to add. training data

(W): an imposture - not a drug,

{ overly simplified it's a
→ presentation feature
of bias-variance?

- perceptron: low bias, high variance

- Backprop: low variance

for MLPs other than single perceptron \rightarrow same issue.

(*) Backprop BR: Backprop will often not find a separating solution even though the soln. is within the class of functions learnable by the network.

(*) As separating soln. is not a feasible optimum for the loss function

- Backprop trained NN class. has a lower variance than an optimal classifier for training data

(W) (A8): - review how this claim is formalised

the error/loss surface

(assumed)

- find minimum point of this surface : single global optimum

- what does loss surface actually look like?

BR: open question \rightarrow contentious.

Saddle point \rightarrow see diagram

(*) surface increases in some dir., increases in others.

(*) - some Hessian eigenvalues +ve; some -ve

Notably: ~~but~~ see literature: (W) (A9):

Story so far:

- not guaranteed to find 'true soln' ?? (W) (A10)

- why / now have gone from assuming global minima \rightarrow local minima

convergence \rightarrow local minima

1. converge?

2. How long?

(*) Hard to analyse for MLP with complex loss fn (surface?)

use/make convex functions (best-case for quadratic loss)

BR: informal, intuitive def of convex: connecting line + intersection. ✓
convex function / set ✓

(*) see contour plots of surface for illustration of convexity

i) converging a(i): theoretically $\rightarrow ?$

ii) jittering

iii) diverging

convergence rate: (theoretical const.)

$$R = \frac{|f(x^{(k+1)}) - f(x^*)|}{|f(x^{(k)}) - f(x^*)|}$$

- $x^{(k+1)}$ - k^{th} iteration
- x^* - optimal value of x

- If R is constant

④ Review slides, supp here

quadratic surfaces (convergence)

- a determines shape of quadratic (ad.: min or max)

- 2nd derivative also

- computing optimal step size. BR claims it is one step/line via Taylor's expansion

$$f(x) = f(x_0) + f'(x)(x - x_0) + f''(x)(x - x_0)^2$$

✓ BR: check

(*) first K derivatives of resulting expansion are equivalent to first K derivatives of function at that point

(*) how many terms for Taylor series of quadratic?

0th, 1st, 2nd order term \rightarrow exact for quadratic

(*) final Taylor exp. at $w^{(K)}$

$$E(w^{(K)}) + E'(w^{(K)})(w - w^{(K)}) + \frac{1}{2} E''(w^{(K)})(w - w^{(K)})^2$$

minimise wrt w :

$$\frac{dE(w^{(K)})}{dw} = E'(w^{(K)}) + \frac{1}{2} E''(w^{(K)})(w - w^{(K)}) = 0$$

Hence $E''(w^{(k)})^{-1} = g^{-1}$ (?) ⑥ ⑦ check ✓ (*)

Set $\eta_{opt} = E''(w^{(k)})^{-1}$

- set stepsize to be inverse of 2nd derivative of loss (?)

- get to optimum in 1 step

- new and interesting analysis

(*) consider $\eta < \eta_{opt}$ (monotonic)

$2\eta_{opt} > \eta > \eta_{opt}$ (oscillating convergence)

⑧: bound on step size before it gets bad?

$\eta > 2\eta_{opt}$ (divergence)

- ⑨ results for quadratic

- For generic diff. convex opti objective functions

- same logic as above:-

$$\eta_{opt} = \left(\frac{d^2 E(w^{(k)})}{dw^2} \right)^{-1}$$

- functions of multivariate inputs

Obj.

$$E = g(w)$$

$$w = \begin{bmatrix} w_1 \\ \vdots \\ w_N \end{bmatrix}$$

E is a scalar $g()$ is scalar fm
of vector w .

A-diagonal

- quadratic convex paraboloid:-

$$E = \frac{1}{2} w^T A w + b^T w + c$$

$$= \frac{1}{2} [w_1 \dots w_N] \begin{bmatrix} a_{11} & & \\ & \ddots & \\ & & a_{NN} \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_N \end{bmatrix} + [w_1 \dots w_N] \begin{bmatrix} b_1 \\ \vdots \\ b_N \end{bmatrix} + c$$

add

If A diagonal:- - sum of indep. quadratics

$$E = \frac{1}{2} \sum_i a_{ii} w_i^2 + \sum_i b_i w_i + c = \sum_i \left(\frac{1}{2} a_{ii} w_i^2 + b_i w_i \right) + c$$

for 2d:

A - diagonal

$$E = \left(\frac{1}{2} a_{11} w_1^2 + b_1 w_1 + c_1 \right) + \left(\frac{1}{2} a_{22} w_2^2 + b_2 w_2 + c_2 \right)$$

- This is a bowl in 3 dimensions i.e. $z = E(w_1, w_2)$

- Take level sets / slices of bowl

- Concentric ellipses when viewed from top

- What is shape of horizontal slices e.g.

- Fix w_1 and vary w_2 (W) AII: Review argument here.
- Vary w_1 and fix w_2

For each quadratic, check step size prescription

$$M_{1,\text{opt}} = a_{11}^{-1} \quad M_{2,\text{opt}} = a_{22}^{-1}$$

vector update rule (W) AII review

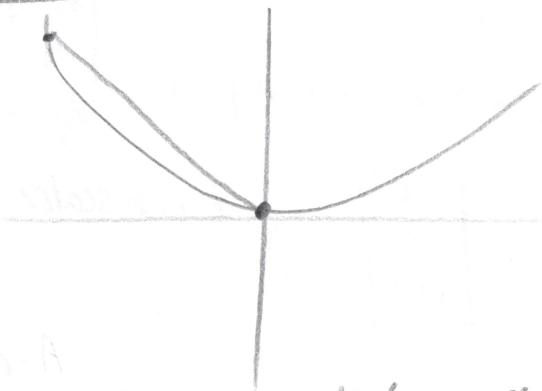
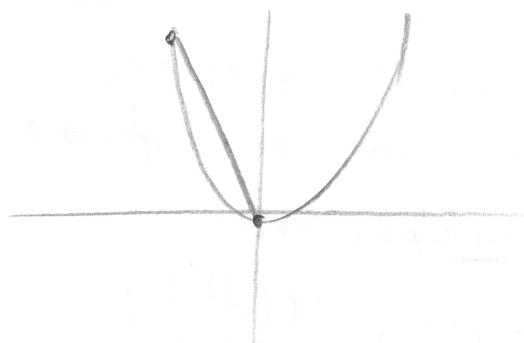
- Gradient orthogonal to level set

- Apply to quadratics you have decoupled

problem with vector update rule

$$\underline{w}^{(k+1)} \leftarrow \underline{w}^{(k)} - \eta \nabla_{\underline{w}} E$$

$$w_i^{(k+1)} = w_i^{(k)} - \eta \frac{dE(w_i^{(k)})}{dw}$$



- Optimal step-sizes for each of independently uncoupled quadratics
intuition:

$$\text{and } M_{i,\text{opt}} = \left(\frac{d^2 E(w_i^{(k)})}{dw_i^2} \right)^{-1} = a_{ii}^{-1}$$

(W) AII
④ conflict between optimal size solns. makes learning slow.

$$-\eta < 2 \min_i M_{i,\text{opt}} \quad M_{i,\text{opt}} < \eta < 2 M_{i,\text{opt}}$$

(W) AIV: Review the panels in light of above presentation.

- Dependence on learning rate

⑩ convergence behavior becomes increasingly unpredictable as dimensions increase.

⑪ fastest convergence; based on intuition with quadratic surfaces, learning rate η must be close to both largest η_i, opt and smallest η_i, opt .

- to ensure converge in every direction/dimension

- BR: generally infeasible

⑫ BR: generally slow if $\frac{\max_i \eta_i, \text{opt}}{\min_i \eta_i, \text{opt}}$ large

- gradient descent - step size η is fixed for every direction

⑬ ⑭ do you misunderstand the formulation of the problem?

- ⑮ rescale axes in response to different optimal learning rates for different directions.

- scale w_1 and w_2 axes $\hat{w} = S w$ $w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$ $S = \begin{bmatrix} s_1 & 0 \\ 0 & s_2 \end{bmatrix}$ $\hat{w} = \begin{bmatrix} \hat{w}_1 \\ \hat{w}_2 \end{bmatrix}$

⑯ ⑰ check

yielding:- $E = \frac{1}{2} \hat{w}^T \hat{Q} \hat{w} + \hat{b}^T \hat{w} + c$ (as A is
(2nd?)

optimal step size $\eta=1$ (check; as 'derivative' one in every direction)

- how to find S ?

S

(*) Review
these derivations

⑱ solve for scaling matrix by inspection.

⑲ solve quadratic (?) in scaled space yielding:-

$$\hat{w}^{(k+1)} = w^{(k)} - \eta \nabla_{\hat{w}} E(\hat{w}^{(k)})^T \quad (\text{make subst.})$$

$$w^{(k+1)} = w^{(k)} - \eta A^{-1} \nabla_w E(w^{(k)})^T \quad - \text{gets you to optimum}\\ \text{in single step if } \eta=1$$

⑩ generic differentiable multivariate convex fns:-

- Taylor exp:- check dimensionality $(*) \underline{w}^{(k+1)} = \underline{w}^{(k)} - \eta \nabla_{\underline{w}} E(\underline{w}^{(k)})^T \nabla_{\underline{w}} E(\underline{w}^{(k)})^{-1}$

⑪ what justification for using only 1st two terms of Taylor expansion for generic multivariates ^{2nd} (quadratic approx.)

- notice role of Hessian here (derivative of scalar fn not good)
vector input) due to convexity?

(*) Heated localised optim. with quad. approx.

⑫ why $\eta=1$ by Newton's method.

Issues - Hessian - normalising by Hessian

⑬ 100,000 parameters \rightarrow difficult to inv / compute

⑭ Non convex functions; Hessian may not be positive semi-definite
⑤ (surely positive definite)

- algorithm diverges - away from optimum

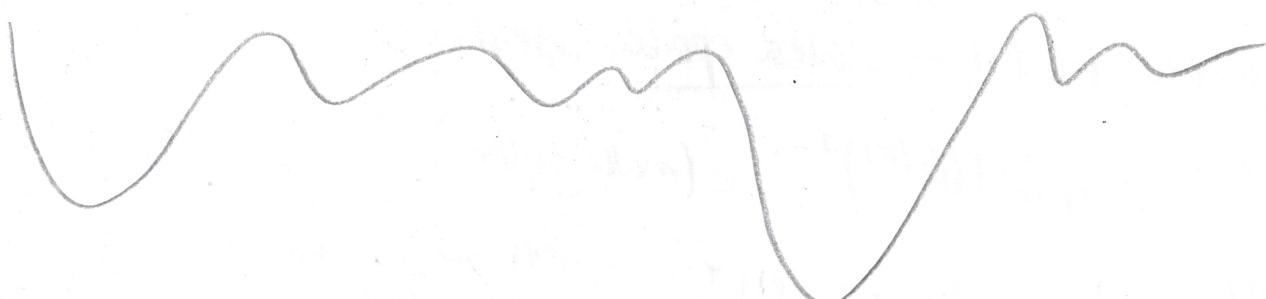
- eigenvalues \rightarrow what are the consequences?

(*) Quantitative techniques for approximating Hessian, improve positive-definiteness

BR: analysis was focused on ensuring step size was not so large as to cause divergence within a convex region $\eta < 2\eta_{opt}$

(*) But for non-convex functions, is divergence a bad thing?

⑯



(*) $\eta > 2m_{\text{opt}}$ may help escape local optima; but always having $\eta > 2m_{\text{opt}}$ will ensure that we never actually find a sol.

⑩ Adaptive step sizes?

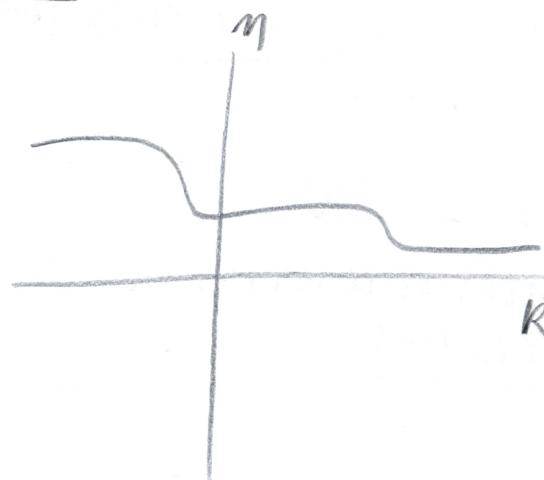
- decaying learning rate - post-normalisation $\eta > 2$
(reduce with iterations)

decay schedules

$$\text{linear } \eta_R = \frac{\eta_0}{R+1}$$

common approach
- see slides

$$\text{quadratic } \eta_R = \frac{\eta_0}{(R+1)^2}$$

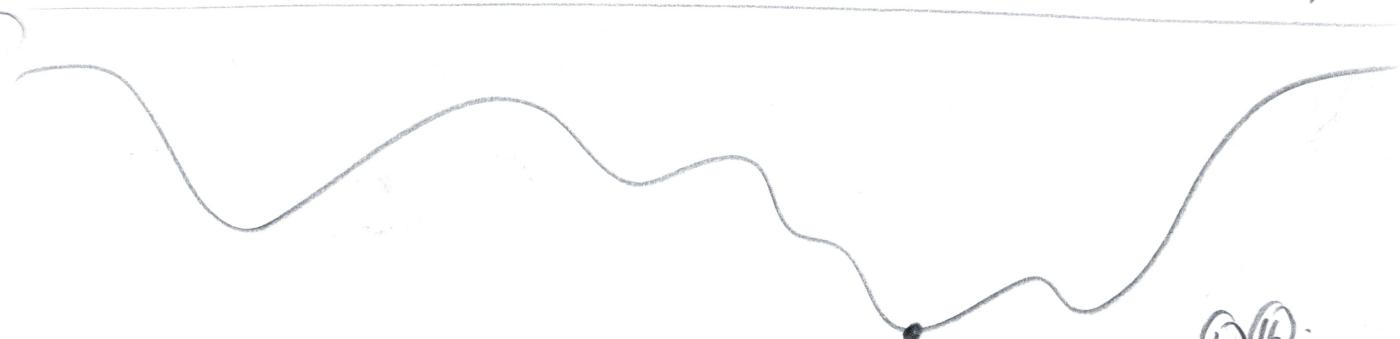


$$\text{exponential } \eta_R = \eta_0 e^{-\beta R} \quad \beta > 0$$

- note in ⑩ that function derivatives are bounded

$$|f'(x)| < B$$

- what criterion to ensure I can find global minimum?



require:-

$$\sum_{R=1}^K \eta_R |f'(x)| > D$$

so

$$\sum_{R=1}^K \eta_R = \infty$$

(?)

- make sure
Step sizes sum to
infinity

- otherwise if start in wrong place, never able to reach optimum.
- want step size to also shrink a little \rightarrow another critera.

• $\sum_{R=1}^K m_R^2 < C$ (i.e. sum of squared step sizes is bounded).

(*) Use 2 criteria allow exploration of entire space; convergence to solution.

candidate function - $\frac{1}{R}$ boundary where sum of step sizes is infinite
sum of square step sizes is finite

$$\sum_{n=0}^{\infty} \frac{1}{R} - \text{harmonic series (wtf)}$$

(WTF) Review this thinking

(*) Review summary

(*) many convergence issues arise because we force learning rate on all parameters