

16 - Convergence in neural networks

- HW - educational, not punitive exercises (by working together)

Training a network LS slides

- continues LS slides
S. 2019

$$U(W) = \frac{1}{N} \sum_X \text{div}(f(X; W), D(X))$$

$$\hat{W} = \underset{W}{\operatorname{argmin}} U(W)$$

$$W_k = W_{k-1} - \eta \nabla_W U(W)^T$$

→ you've seen this already

(*) Gradient of total loss → average of gradients of loss for indiv instances

$$U(W) = \frac{1}{N} \sum_X \text{div}(f(X; W), D(X))$$

$$\nabla_W U(W) = \frac{1}{N} \sum_X \nabla_W \text{div}(f(X; W), D(X))$$

computed via backprop

(*) Backpropagation → an algorithm for computing derivatives

(*) Differentiable divergence is a proxy for class. error

- Minimising loss is expected, but not guaranteed, to minimise class. error

(36705 analogies?)

Issues with gradient descent

(*) Multivariate function of many param - gradient descent takes step 'back' against derivative (①②)

① Size of step / learning rate ^{same} in every component of the param

(Blunt tool)

(*) Loss function has eccentricities in different directions

2nd order methods

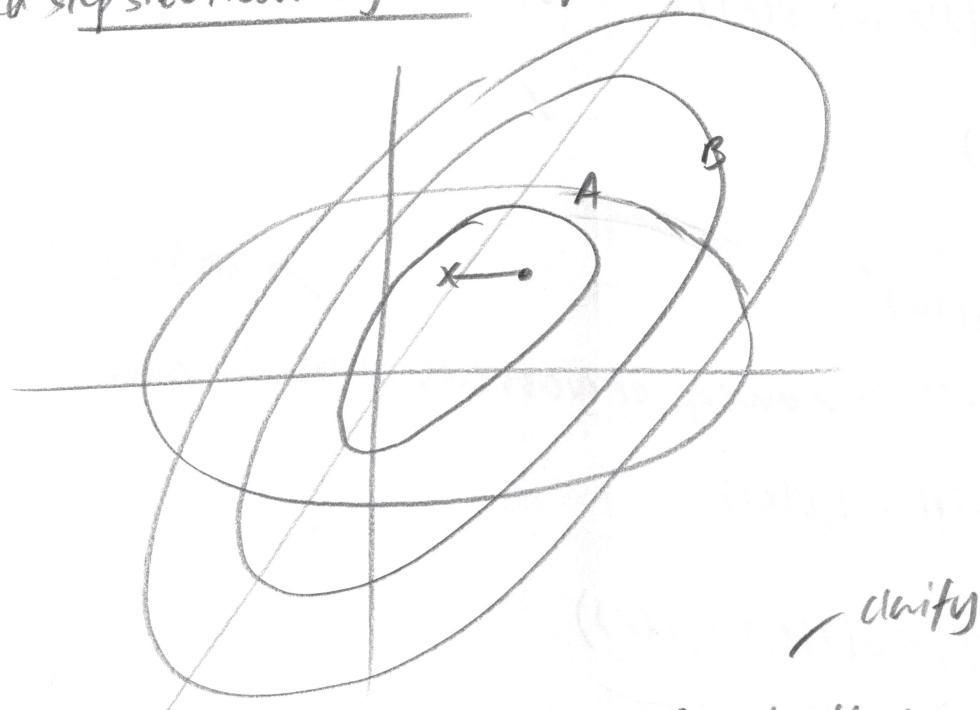
- Normalise variation along comp to mitigate diff opt. learning rates for diff comp.

- import: infeasible
- 2nd order approx + step - away from minimum (not-convex)

(*) see slides → ① ②: Review

- release fixed step size / learning rate requirement.

BR:



- in example A: moving in x -direction; does not affect y -direction
(in axis aligned case)

- in example B: Have to move along recalibrated x - y axes for these contours for independence of dimensions x only

- Heuristic illustration of fixed step size; at when it may not be a good thing

i.e.: when there is 'dependence in direction'

- derivative inspired algos

- Rprop

- quickprop

- resilient propagation

- Adaptive to sign of derivative

- not considering value of deriv

(*) Accounts for overshooting
- introduce scaling factor

- (*) Algorithmic explanation of Rprop
- independently for every direction
-
- (*) Rprop pseudocode

(i) Review logic

(ii) Algo

- use intuition to us.
- ensure you understand now formalised.

@ does this require backprop?

- Yes; for $\frac{\partial \text{err}(w_{i,j})}{\partial w_{i,j}}$

(BR): Much more efficient than gradient descent

- (*) NO convexity

quickprop - Scott Fahlman

- Newton updates (comptd of 2nd deriv.) (empirically derived)
- instab. for non-convex objectives
- Better than Rprop (faster)
- Newton: Multiply by inverse Hessian (normalise)

(i) Review

(*) Tying step sizes for all dimension is subopt. (57)

- treat each dimension separately

- Rprop, quickprop issues:-

- ignores dependence between dimensions
- unexpected behaviour
- slow

(ii) suggests we should find a more nuanced way of addressing (iv)

↗ (i) Use 16.5.2019 slides here

(*) Close look at convergence

(*) ^{BR} in deep learning rates

↳ converges in some dir; diverges in others

- steps; in directions of convergence (smoothly)

- steps all pointing in same direction

- steps; in directions of divergence/oscillation

- steps not pointing in same dir (net distance travelled given a direction is low)

vector quantity with dir

(*) Proposal - keep track of oscillations

- Bouncing around (encoded as frequent sign changes?) \rightarrow decrease step size
- smooth convergence - increase step size.
- Another way of seeing it is changing the nature of the feedback w.r.t. for an adaptive learning rate - time series analogy (e.g.)

(*) Momentum methods

- maintain running av. of all past steps

(@) - can encode differences between smooth conv / oscillation / divergence

- smooth conv \rightarrow running average will have large value

- oscillation/divergence \rightarrow running average will be small

(cancellation of the +ve and -ve swings)

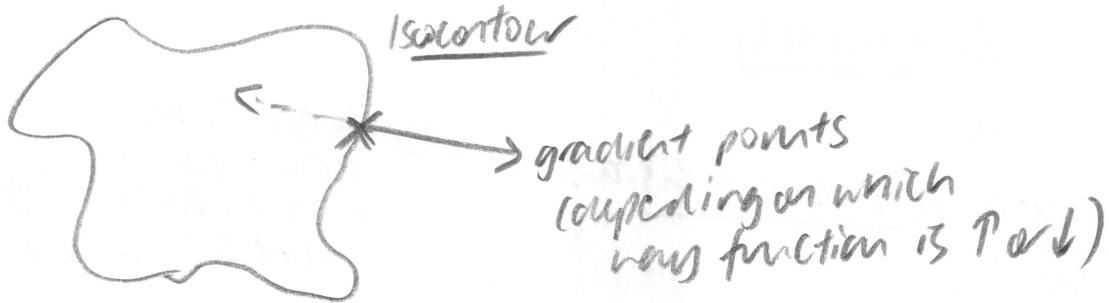
(elimin fluctuations,

reduce step size)

update with running average, rather than curv gradient.

(*) Momentum update

- BR: remember gradient is orthogonal to level set



(*) Plain gradient update

- BR: "gradient always pointing immediately downwards rather than towards centre of ellipse"

Final update
is average
- momentum
 \nearrow intuit.

- BR: "top-bottom (y-direction) - oscillating back and forth
left-right (x-direction) - moving gradually rightwards"

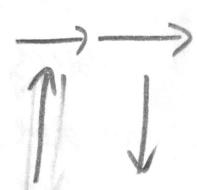
(*) Example (using running av.)



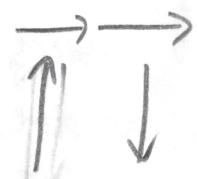
- Average of 2

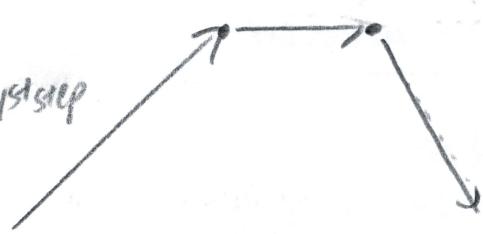
step 1 step 2

- in x-direction



- in y-direction





under momentum
(using running average
of gradient)

(W) (B): review - check you
understand
heuristically

$$(*) \Delta \underline{w}^{(k)} = \beta \Delta \underline{w}^{(k-1)} - \eta \nabla_{\underline{w}} \text{err}(\underline{w}^{(k-1)}) \rightarrow (2) \text{ au}$$

$$\underline{w}^{(k)} = \underline{w}^{(k-1)} + \Delta \underline{w}^{(k)} \Rightarrow \underline{w}^{(k)} - \underline{w}^{(k-1)} = \Delta \underline{w}^{(k)}$$

(W): distinguish
layer subscripts
and iteration
no. superscripts

(*) Gradient descent / training with momentum

- initialise all weights: w_1, w_2, \dots, w_R

- \forall layers R , mit $\nabla_{w_R} \text{err} = 0 \quad \Delta w_R = 0$
- $\forall t=1:T$ (training instances)

- \forall every layer R
 - compute gradient $\nabla_{w_R} \text{div}(y_t, d_t)$
 - $\nabla_{w_R} \text{err} \leftarrow \frac{1}{T} \nabla_{w_R} \text{div}(y_t, d_t)$

- \forall every layer R ,
- $$\Delta \underline{w}_R^{(t)} = \beta \Delta \underline{w}_R^{(t-1)} - \eta \nabla_{\underline{w}_R} \text{err}$$

$$w_R = w_R + \Delta w_R$$

(*) refer to
previous;
as no iteration
superscripts.

Momentum update

- (W) - Analyse intuitively the eq:-

$$\Delta \underline{w}^{(k)} = \beta \Delta \underline{w}^{(k-1)} - \eta \nabla_{\underline{w}} \text{err} \quad (\text{update})$$

(*) (B): check you
understand this
in the hill-climbing
language ✓

- At any iteration: to compute current step:-

1. compute gradient at current loc $(\nabla_{\underline{w}} \text{err}(\underline{w}^{(k-1)}))$
2. Add m scaled previous step (running average) $(\beta \Delta \underline{w}^{(k-1)})$

- use hill-climbing analogy :-
- reverse order of operation \rightarrow Nesterov's accelerated gradient

(x) Nesterov's accelerated gradient

- At any iteration; to complete current step:-
- Extend previous step
- Compute gradient step at result. pos.
- Add the 2 to obtain final
- BR: Still maintaining running average; but not computing grad. at current loc; rather, ...

(W) (A) Review -

check you

selected

not hill-climbing;

vectors

(W) (A) : - see

- probably faster than
simple grad. descent

(x) Nesterov's Accelerated gradient update

$$\Delta w^{(k)} = \beta \Delta w^{(k-1)} - \eta J_{w_k} \text{Err}(w^{(k-1)} + \beta \Delta w^{(k-1)})$$

$$w^{(k)} = w^{(k-1)} + \Delta w^{(k)}$$

(D)

BR: Nesterov's method is faster than momentum

- for convex fns; optimum

(x) Gradient descent pseudocode with Nesterov method.

- initialise all weights w_1, w_2, \dots

- Do:

- \forall layers K , mit $J_{w_K} \text{Err} = 0$, $\Delta w_K = 0$
- \forall layer K :-

$$w_K = w_K + \beta \Delta w_K$$

- $\forall t=1, \dots, T$ (training instance)

- \forall every layer k :

- compute gradient $J_{w_K} \text{Div}(y_t, d_t)$

- $J_{w_K} \text{Err} \leftarrow \frac{1}{t} J_{w_K} \text{Div}(y_t, d_t)$

• for every layer k :-

$$\underline{W}_R = \underline{W}_R - \eta \nabla_{\underline{W}_R} Err$$

$$\Delta \underline{W}_R = \beta \Delta \underline{W}_R - \eta \nabla_{\underline{W}_R} Err$$

① Review pseudocode
check untested.

{ (ii) issue here (particularly notation)}

• until error converged

(*) see story so far - convergence

(*) incremental updates

(*) training formulation

② graph illustration

- \rightarrow all input-output pairs

- Adjust parameters only at input-output pairs } gradient descent

issue: (*) Process all training points before making a single adjustment

- batch update

- compute loss over all training points
(average)

(*) Process entire training set in one batch; getting errors, adjusting param.

(*) 'handkerchief grabbing'

(*) incremental update

- Adjust fraction at one training point at at a time

- Processed all training points \rightarrow adjusted entire fraction.

(*) incremental updates: ~~soak~~ (pseudocode)

① A10: How is this diff vs gradient descent (batch)

• Given $(X_1, d_1), (X_2, d_2), \dots, (X_T, d_T)$

• Initialise all weights $\underline{W}_1, \underline{W}_2, \dots, \underline{W}_R$ (via backprop)

• Do:- (until loss converged) (I)

(I) • For all $t = 1, 2, \dots, T$:-

• For every layer K :-

- compute $\nabla_{\underline{W}_K} \text{Div}(Y_t, d_t)$

- update : $\underline{W}_R = \underline{W}_R - \eta \nabla_{\underline{W}_R} \text{Div}(Y_t, d_t)^T$ (II)

(*) Adjust param. after every training inst.

- (*) Iterations \rightarrow multiple passes over training data
- single pass through entire training data \rightarrow epoch

- (I) \rightarrow over multiple epochs
(II) - one epoch
(III) \rightarrow one update

(*) A heuristic illustration of cyclic behavior

- insert stochasticity i.e. randomly

(*) SGD - insert random permutation of training data between epochs
~~intrinsic~~ i.e. change order in which you present
training instances

Q: why does process of incremental updates work?

Q: under what conditions

- Heuristic, simplistic exp.

(*) use extreme example

- use expected behavior of gradient

- Batch vs SGD

Q: review argument for batch vs SGD i

Q: under what conditions do incremental updates work?

- caveats - learning rate

Q: review heuristic argument for why
learning rate must decrease with iterations

(*) incremental update, SGD, changing step size (with iteration)

(*) supplement to existing $\textcircled{W}\textcircled{A}14$: review

SGD convergence

(*) SGD converges 'almost surely' to a global/local minimum for most functions.
sufficient condit.

$$\sum_k \eta_k = \infty, \quad \sum_k \eta_k^2 < \infty$$

- eventually entire param space can be searched

- steps shrink (finite sum of squared step sizes finite)

- \textcircled{W} test converging series that satisfies both above:-

$$\eta_k \propto \frac{1}{k}$$

- optimal rate of shrinking step size for convex fns.

- learning rates determined heuristically

- convex loss; SGD \rightarrow optimum sol

- non-convex loss; SGD \rightarrow local minimum.

$\textcircled{W}\textcircled{A}16$: Harmonic series

SGD convergence

(*) some key results here $\textcircled{W}\textcircled{A}17$ - review + digest converge results

- strong convexity assumption

- generically / weakly convex

Batch gradient convergence

- some distinctions as above

SGD example

- k-means -

(*) vertical bars - variation of error across runs

(*) SGD larger variation, fast converge; higher error converging points? @ UTTY?

(*) Batch gradient descent

- modelling a function

- distinction between

empirical est. of error vs expected gen. error
(over the distn)
(over empirical distn)

② - empirical error is an unbiased estimate of expected error
(loss)

BR: Hence we can compute / minimise loss as we have done

③ How about variance

$$(*) \text{Var}\left(\frac{1}{N} \sum_i \text{div}(x_i)\right) = \frac{1}{N^2} \sum_i \text{Var}(\text{div}) = \frac{N \text{Var}(\text{div})}{N^2} = \frac{\text{Var}(\text{div})}{N}$$

(variance of empirical risk)

(*) compare variance of empirical / sample error
mols batch and SGD

(W) Q16: review
mathematically

(+) SGD - examines error/divergence for a single instance

$$\hookrightarrow l = \text{div}(x_i) \quad E[l] = E[\text{div}(x_i)]$$

- SGD minimises an unbiased estimator (loss) of actual function
we want to minimise

- How about variance for SGD?

- graphical illustration of variance issues/charact

- sample estimate / empirical loss estimates area using average
length of distances

- this average distance depends on where training samples
are taken

- increase no. of samples

- BR: "large variance depending on n" - variation between epoch?

will be large if you have a small no. of training samples, but increase
no. of training samples, variation ↓

(*) refers to estimated \hat{w} ?
⑩⑪: find the relevant formalism to bind these intuitions (which
are well pres.)

OR: creating a tension between SGD and batch updating to
motivate mini-batch

(*) Minibatch updates

- not one instance at a time (SGD) nor
- whole training set at a time (batch)
- use subset of training samples at a time

, very clear motivation
for mini-batch

(*) Incremental update: minibatch update (pseudocode)

- see this as building on gradient descent (stochastic) with shrinking learning rate.

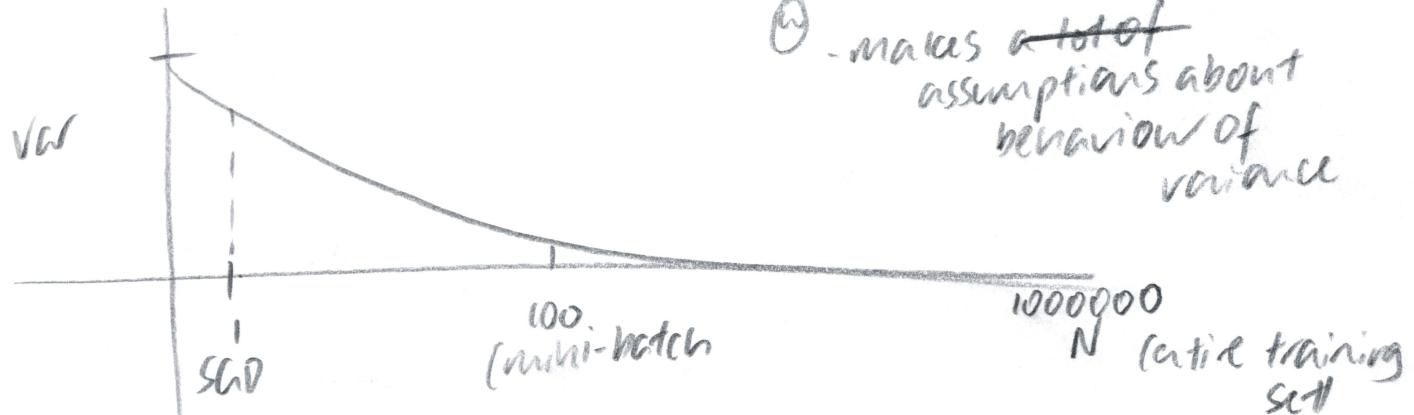
⑩⑪⑫ - check you understand algorithmically

- Randomly select subset of training samples
- permute instances between epochs
- After permutation, run over mini-batch
- Compute average derivative of over mini-batch; before making updates

(*) Batch loss is an unbiased estimate of expected loss
variance of batch error smaller than sample error variance in SGD.

- variance against
batch size

⑩ - makes a lot of
assumptions about
behavior of
variance



- minibatch convergence
- GPU becomes important here
- time taken to process instance similar to processing batch

(*) Empirical comparison of SGD, mini-batch, batch.

(*) Measuring loss

- relies on estimating loss over entire training set (figures)

training and mini-batches

- mini-batch size in practice \rightarrow hyperparameter
- convergence depends on
learning rate

- simple technique \rightarrow start from fixed/exogenous parameters
(e.g. Nesterov and 1st moment of gradient)

- other methods
 - learning rate part of estimation