

6 Review

(*) Note the key point → learning rate/step-size is same for every component of the parameter update. (I)

(*) Essence of normalised update rule is to amend the error/loss function; which is a function defined over a multi-variate parameter space; by effectively rescaling/normalising parameters.

(*) For the multivariate case, we have :-
convex loss

$$(6) \quad w^{(k+1)} = w^{(k)} - \eta H_E(w^{(k)})^{-1} \nabla_w E(w^{(k)})^T$$

(*) Normalisation constant is the Hessian; and we set $\eta = 1$

(*) Note similarity with setting of $\eta = 1$ in scalar case

(*) In this case, we alter the parameters/gradient of error and keep learning rate $\eta = 1$ to resolve issue (I)

(*) Derivative-information inspired algorithms - RProp and Quickprop

(*) Both make use of information associated with the derivative/2nd deriv.
(i.e. some kind of compression/approx.)

(*) Both make component-wise updates to the parameters to improve convergence of gradient descent

RPROP
(A2), (A3)

• Note Rprop and quickprop
are modifications to
how the parameters are
updated. (indep.)

(*) Resilient propagation

(*) Steps in different dir. not coupled
• Accounts for overshooting a minimum

Rprop (at each time-step)

- At each time:

→ If derivative at current location recommends continuing in
some dir as before (i.e. no change of sign)
- increase the step size; continue in same dir.

→ If the derivative has changed sign (i.e. minimum overshoot)

- reduce the step and reverse direction

• Rprop → see slides (main here) (*) tractability

• The former and latter are governed by parameters α and β respectively

• Select initial \hat{w} and compute derivative

- Step Δw against the derivative (direction that reduces)

$$\rightarrow \Delta w = \text{sign} \left(\frac{dE(\hat{w})}{dw} \right) \Delta w$$

$$\rightarrow \hat{w} = \hat{w} - \Delta w$$

• Compute deriv in new location (governed by α)

- No change in sign from previous location; increase step size and take large step.

$$*\alpha > 1$$

$$\rightarrow \Delta w = \alpha \Delta w$$

$$\rightarrow \hat{w} = \hat{w} - \Delta w$$

• Compute deriv in new location

- If derivative has changed sign; return to previous location

- shrink step size (governed by β)

- Take a smaller step forward

$$\beta < 1$$

$$\rightarrow \hat{w} = \hat{w} + \Delta w$$

$$\rightarrow \Delta w = \beta \Delta w$$

$$\rightarrow \hat{w} = \hat{w} - \Delta w$$

(*) Note that Rprop is still used with backpropagation; which is used to compute derivatives

Rprop (pseudocode)

► set $\alpha = 1.2$ $\beta = 0.5$

► For each layer l , and each i, j :-

parameters
 $w_{l,i,j}$ updated
independently

- Init $w_{l,i,j}, \Delta w_{l,i,j} > 0$
- $\text{prev } D(l,i,j) = \frac{\text{derr}(w_{l,i,j})}{\Delta w_{l,i,j}}$ (*) obtained via backprop
- set $\Delta w_{l,i,j} = \text{sign}(\text{prev } D(l,i,j)) \Delta w_{l,i,j}$
- WHILE not converged:-

○ • $w_{l,i,j} = w_{l,i,j} - \Delta w_{l,i,j}$

• $D(l,i,j) = \frac{\text{derr}(w_{l,i,j})}{\Delta w_{l,i,j}}$ (†) as above

• IF $\text{sign}(\text{prev } D(l,i,j)) = \text{sign}(D(l,i,j))$:-

• $\Delta w_{l,i,j} = \alpha \Delta w_{l,i,j}$

• $\text{prev } D(l,i,j) = D(l,i,j)$

○ ELSE:

• $w_{l,i,j} = w_{l,i,j} + \Delta w_{l,i,j}$

• $\Delta w_{l,i,j} = \beta \Delta w_{l,i,j}$

(*) Quickprop (Scott Fahlman)

- similar to Rprop; updates parameters independently, component-wise
- quickprop can be viewed as employing Newton updates (i.e. 2nd order method involving Hessian) except
 - i) dimensions componentwise treated indep.
 - ii) uses approximation of 2nd deriv via finite differences

quickprop-key equations

(for each dimension of the weight:-)

- for $i=1:N$

$$w_i^{(k+1)} = w_i^{(k)} - E''(w_i^{(k)} | w_j^{(k)}, j \neq i)^{-1} E'(w_i^{(k)} | w_j^{(k)}, j \neq i)$$

$$w_i^{(k+1)} = w_i^{(k)} - \underbrace{D(w_i^{(k)}, w_i^{(k-1)})^{-1}}_{\text{finite diff. approx to 2nd derivative}} E'(w_i^{(k)} | w_j^{(k)}, j \neq i)$$

assuming quad $E()$

Yielding: $w_i^{(k+1)} = w_i^{(k)} - \left(\frac{E'(w^{(k)}) - E'(w^{(k-1)})}{\Delta w^{(k-1)}} \right)^{-1} E'(w^{(k)})$

- for every layer l , every connection from node i in the $(l-1)^{th}$ layer
to node j in the l^{th} layer:-

$$\Delta w_{l,i,j}^{(k)} = \frac{\Delta w_{l,i,j}^{(k-1)}}{\frac{E'(w_{l,i,j}^{(k)}) - E'(w_{l,i,j}^{(k-1)})}{\underbrace{E'(w_{l,i,j}^{(k)})}_{\text{- compute using backprop}}}}$$

$$w_{l,i,j}^{(k+1)} = w_{l,i,j}^{(k)} - \Delta w_{l,i,j}^{(k)}$$

(*) issues with treating dimensions independently/making component-wise
updates to the parameters

- solution converges ^{smoothly} in some dir; may oscillate/diverge in others

(*) motivation for momentum/nestcov's acc. gradient methods

- keep track of oscillations/div/convergence man indicator

- for smooth convergence; increase step size

- oscillation/divergence \rightarrow shrink.

Momentum methods → ① distill pub (Polyak 1964)

- key equations already covered in notes

- β is set at around 0.9

② note the order of operation that the equations specify

$$\Delta \underline{w}^{(k)} = \beta \Delta \underline{w}^{(k-1)} - \eta \nabla_{\underline{w}} E(\underline{w}^{(k-1)}) \quad (\text{omitting layer subscripts})$$

③: $\nabla_{\underline{w}^{(k)}} E(\underline{w}^{(k-1)})$

- gradient of error eval at $\underline{w}^{(k-1)}$
(previous it.)

Nesterov's accel. grad.

(*) Similar to momentum, except change in order of operations

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

(*) At each iteration:-

- i) extend previous step
- ii) compute gradient step at resultant position
- iii) Add the two

(*) Principle behind momentum (and Nesterov's accel. gradients?)

- ideally - a component specific step size
→ too many more param (maintain sep step for every weight)

④: Adaptive solutions: (in learning rate/step-size)

- start with common learning rate
- shrink step size when weight oscillates
- expand step size when weight moves directionally consistently

⑤ Recap: improving convergence

- i) 2nd order methods (normalise variance across dim.)
- ii) Adaptive/decaying learning rates
- iii) Decouple dimensions (Rprop, Quickprop)
- iv) Momentum methods → emphasize steady improv; decouple unstable dir.

- (*) incremental updates → Q ch 8
- full batch parameter updates → computationally costly
(memory, time wise also)
 - note that in the 'gradient descent algorithm' previously pres;
recall: (BR notation): $E(w) = \frac{1}{T} \sum_{t=1}^T \text{div}(f(x_t; w), d_t)$
 - $\nabla_w E = \frac{1}{T} \sum_{t=1}^T \nabla_w \text{div}(x_t, d_t)$
- overall error as an additive sum of per instance divergences
- gradient distributes our sums nicely
- Q 5.9

- (*) There are 2 key design choices here:-
- i) the size of the training set we decide to use to compute the gradient of the error over.
 - ii) the frequency with which we update parameters.
- (*) At this stage of the lecture, a distinction is being made between computing the per-instance divergence and the overall error over the entire training set (batch); and making parameter updates accordingly.
- (*) incremental updates occur and refer to instance by instance computation of per instance divergence; additive error; and follow by parameter update after every instance.
- (*) I have accidentally included other issues raised here; which will come up later: so (I) is useful to have for late exposition.
- (*) Pathological behaviour with incremental updating.
- (*) The key point is that cyclic behaviour as a pathology of algorithm behaviour can occur if we decide to do incremental updates without injecting stochasticity

(*) We can get stuck in cycles and algorithm fails to converge.

(*) Incremental update + SGD

- In context of incremental updates (param. update after one training instance) we have to randomly permute training data between epochs.

- definitions:
- (I)- Iterations can make multiple passes over training data
 - (II)- One pass through entire training set - epoch
 - (III)- Epoch over a training set of T samples results in T updates of the parameters
(more incremental updates)

(*) Issues with incremental updates

- Heuristic example to understand when, under what conditions incm. updates work.
- Extreme case (data clotting):-

$$\frac{dE}{dw_{ij}^{(k)}} = \frac{1}{T} \sum_{i=1}^T \frac{\partial \text{div}(y(x_i), d_i)}{\partial w_{ij}^{(k)}} = \frac{\partial \text{div}(y(x_i), d_i)}{\partial w_{ij}^{(k)}}$$

- When all training instances are same, average of derivatives over all per-instance divergences is same as derivative of single instance.

(*) Batch vs incremental updates (SGD) (A12)

- (i)- Batch \rightarrow operates over T training instances to get a single parameter update

(ii): Incremental/online SGD:- Same computation yields T parameter updates

- (iii): Incremental/online SGD only reveals one sense of 'stochastic'; in the sense of random permutation of training instances for which parameter updates are conducted.

(iv): There is another (vitamin) sense in which stochastic gradient descent has statistical connections (in that we are computing an estimator of a comp. infeasible quantity)

- (*) converges → learning rate η with mean updates / SGD A14
- Intuitive explanation offered by BR is ~~contained in~~
is captured in formal results.
- The decreasing learning rate η is part of a package of assumptions that ensure that SGD converges:-
 - i) SGD converges almost surely to a global minimum for a convex loss
 - ii) SGD converges otherwise almost surely to a local minimum.

- A consequence of the Robbins-Siegmund theorem

(*) reference material for SGD:-

- also cited in DL book.

i) Papers by Bottou → placed in folder.

→ in particular 2003 SS lectures

→ Cn8

ii) more info in DL book (will supplement with formalism)

(*) There are further sufficient conditions on learning rate η (other than the decay) mentioned above to guarantee SGD convergence.

- Already noted

- Also correspondence with DL book concerning heuristic details.

(*) convergence results are a little unclear to me

(a) Can you link this to your work on 36-705; that is
conventional convergence/boundedness?

- Is the convergence pointwise or uniform?

(*) Following results are placed here, but require further scrutiny: -

- convergence: No. Iterations taken to get within ϵ of optimal sol.

$$|f(w^{(k)}) - f(w^*)| < \epsilon$$

(*) SGD convergence results

• Strongly convex fn; optimal learning rate:-

$$|\underline{w}^{(k)} - \underline{w}^*| < \frac{1}{R} |\underline{w}^{(0)} - \underline{w}^*| \quad \text{"iterations to } \epsilon\text{-convergence } O\left(\frac{1}{\epsilon}\right)" \quad (?)$$

• Generically convex

- " ϵ -convergence" of $\frac{1}{\sqrt{R}}$ using learning rate of $\frac{1}{\sqrt{R}}$

DL Book Ch 8:

SGD for convex, excess error $J(\underline{\theta}) - \min_{\underline{\theta}} J(\underline{\theta}) = O\left(\frac{1}{\sqrt{R}}\right)$ after k ite.

- SGD for strongly convex $\underline{\theta}$ —————— $= O\left(\frac{1}{R}\right)$

(*) Batch gradient convergence

• Strongly convex functions

$$|\underline{w}^{(k)} - \underline{w}^*| < c^R |\underline{w}^{(0)} - \underline{w}^*| \quad \text{"ite to } \epsilon\text{-convergence } O\left(\log\left(\frac{1}{\epsilon}\right)\right)"$$

• Generically convex

"Iterations to ϵ -convergence $O\left(\log\left(\frac{1}{\epsilon}\right)\right)" \quad (?)$

(*) Batch gradient descent is faster (corroborates DL) in theory

- SGD performs 1 update for every batch update

DL (Goodfellow) \rightarrow some discussion on Game-Rao bound, boundedness of generalisation error, practical adv. of SGD.

(*) See 10-725 notes; ~~but SGD not fully worked out yet~~ from a formal classical / convex opt. perspective.

- relation of empirical and expected error, per training instance error

- going back to DL equations:-

$$J^*(\theta) = \mathbb{E}_{(x,y) \sim p_{\text{data}}} L(f(x;\theta), y)$$

$$J(\theta) = \mathbb{E}_{(x,y) \sim \hat{p}_{\text{data}}} \hat{L}(f(x;\theta), y)$$

ERM:

$$\min_{\theta} J(\theta) = \min_{\theta} \mathbb{E}_{(x,y) \sim \hat{p}_{\text{data}}(x,y)} [L(f(x;\theta); y)] = \frac{1}{m} \sum_{i=1}^m L(f(x^{(i)}; \theta), y^{(i)})$$

(empirical risk)

Q: Having difficulty squaring BR's presentations symbolically with DL book.

- list claims verbally \rightarrow re-review \rightarrow O/S

- very confused by: - (BR's presentation)

i) we minimise empirical risk ✓

ii) Empirical risk is an unbiased estimate of expected loss

iii) NO guarantee minimising empirical risk will minimise expected loss

iv) variance of empirical risk $\text{Var}(\text{loss}) = \frac{1}{N} \text{Var}(\text{distr})$
(derivation given)

v) OK large variance is great like that the w that minimises the empirical risk will differ significantly from w that minimises the expected loss.

W: It will help if this was placed in a formal pointwise vs uniform convergence setting \rightarrow O/S using std. notation

- BR very unclear here

O/S ②

(*) Huristic point :-

- see graph; variance of true fn - estimated is higher over incremental update SGD than over entire batch

(*) Mini-batch updates

- Again, some / similar arguments which are very unclear

DL Ch8: Next piece clarifies what is going on.

$$\text{MLE: } \theta_M = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^m \log p_{\text{model}}(x^{(i)}, y^{(i)}; \theta)$$

• MLE estimation: $\max \sum (\text{sum}) \Leftrightarrow \text{max. exp. our empirical distri defn by training set.}$

$$\text{i.e. } J(\theta) = \mathbb{E}_{(x,y) \sim \hat{p}_{\text{data}}} \log p_{\text{model}}(x, y; \theta)$$

• obj fn J can be viewed as expectations over training set

$$\textcircled{R} \quad \nabla_{\theta} J(\theta) = \nabla_{\theta} \mathbb{E}_{(x,y) \sim \hat{p}_{\text{data}}} \log p_{\text{model}}(x, y; \theta)$$

$$= \mathbb{E}_{(x,y) \sim \hat{p}_{\text{data}}} [\nabla_{\theta} \log p_{\text{model}}(x, y; \theta)]$$

• computing $\mathbb{E}(\cdot)$ very expensive; eval model on every example in dataset

• compute these expectations using a random sample of training data points and taking an average.

⇒ mini-batch

(*) motivation for mini-batch gradient descent :-

- follows gradient of generalisation error $J^*(\theta) = \mathbb{E}_{(x,y) \sim \hat{p}_{\text{data}}} L(f(x; \theta), y)$
(expected)

If no examples are repeated

(*) discrete x, y ; abstract gen error (as a sum of arbitrary L)

$$J^*(\theta) = \sum_x \sum_y p_{\text{data}}(x, y) L(f(x; \theta), y)$$

$$g = \nabla_{\theta} J^*(\theta) = \sum_x \sum_y p_{\text{data}}(x, y) \nabla_{\theta} L(f(x; \theta), y)$$

exact grad.
of gen. error

* (6) Obtain an unbiased estimator of exact gradient of $\text{gen. } g$ by sampling a minibatch $\{x^{(1)}, \dots, x^{(m)}\}$ of training examples with w.r.t. targets $y^{(i)}$ from original distri. of data.

- compute grad. of loss wrt param on mini-batch

$$\hat{g} = \frac{1}{m} \nabla_{\theta} \sum_i L(f(x^{(i)}, \theta), y^{(i)})$$

est. of grad of
loss wrt param
over minibatch

- updating θ in direction of $\hat{g} \rightarrow \text{SGD or general. error}$

(*) If examples not reused

(*) Minibatch conv.

- convex fn; conv. for SGD $O(\frac{1}{\sqrt{k}})$

- minibatch updates with size b batches: - $O\left(\frac{1}{\sqrt{bk}} + \frac{1}{k}\right)$

- improvement of \sqrt{b} over SGD

- But batch size $b \rightarrow b$ times as many comp. lit. as SGD

- degradation of \sqrt{b} .

In practice

- Non conv. obj.

- mini batches more effective with right η .

- vector processing