

CS convergence - review

(A1): see old book (B): moving in opposite direction of gradient

- $y = f(x)$; $\frac{dy}{dx} = f'(x)$
 - $f(x+\epsilon) \approx f(x) + \epsilon f'(x)$ (from definition of derivative as $\lim_{\epsilon \rightarrow 0} \frac{f(x+\epsilon) - f(x)}{\epsilon}$)
 - $f(x - \epsilon \text{sign}(f'(x))) < f(x)$ for small enough ϵ (?)
 - reduce $f(x)$ by moving x in small steps with opposite sign derivative
 - Best understood graphically; formal results (?)
- (A1): The directional derivative $u^T \nabla f$ (a unit vector) is the slope of function f in direction u .
- formally; it is the derivative of $f(x + \alpha u)$ with respect to α , eval at $\alpha=0$.
- via chain rule: - $\frac{\partial}{\partial \alpha} f(x + \alpha u) \Big|_{\alpha=0} = u^T \nabla_x f(x)$ f(.) - scalar fn

- (*)
- minimise f ; find direction in which f decreases fastest
 - use the directional derivative

$$\min_{u, u^T u = 1} u^T \nabla_x f(x) = \min_{u, u^T u = 1} \|u\|_2 \|\nabla_x f(x)\|_2 \cos \theta$$

$$\left. \begin{array}{l} \text{Note } u^T u = 1 \rightarrow \text{basis vectors} \\ \text{using dot product cosine rule} \end{array} \right\} = \min_u \cos \theta$$

- minimised when u points in opposite direction of gradient i.e.
- $\cos \theta = 180^\circ = \pi$

- gradient points uphill; negative gradient points downhill
- decrease f by moving in direction of negative gradient
- gradient descent / steepest descent:

$$x' = x - \epsilon \nabla_x f(x)$$

ε-learning rate

- ε selected in different ways
 - fix constant
 - evaluate $f(x - \epsilon \nabla_x f(x))$ for several values of ε, select one with smallest obj fn value → line search

(B) dimensionality checks

- introduction of transpose seems like a convention
- recall $\nabla_{y_N} \text{Div}$ is the gradient of the per instance loss/divergence with respect to the post output activation vector from the N^{th} layer, y_N .
- $\nabla_{z_N} \text{Div} = \nabla_{y_N} \text{Div} \cdot J_{y_N}(z_N)$
- $\nabla_{z_N} \text{Div}$ is the gradient of Div not pre-activation affine combo vector from N^{th} layer, z_N .
- note gradient of a scalar fn w/ vector is itself a vector

(*) in vector formulation;

- i) each backward step past an activation

$$\text{i.e. } \nabla_{z_N} \text{Div} = \nabla_{y_N} \text{Div} \cdot J_{y_N}(z_N)$$

involves post multiplication by a Jacobian

- ii) each backward step past an affine combo

$$\text{i.e. } \nabla_{y_{N-1}} = \nabla_{z_N} \text{Div} \cdot J_{y_{N-1}}(z_N) = \underbrace{\nabla_{y_N} \text{Div} \cdot J_{y_N}(z_N)}_{\text{from earlier (cached)}} \cdot \underbrace{w_N}_{\text{yields a weights matrix } w_N}$$

involves post multiplication by a Jacobian (which for an affine function)

from earlier (cached)

yields a weights matrix w_N

$$\nabla_{\mathbf{D}\mathbf{R}} \text{Div} = \nabla_{\mathbf{Z}\mathbf{R}} \text{Div} \cdot " \nabla_{\mathbf{D}\mathbf{R}} \mathbf{z}_{\mathbf{R}} " = \nabla_{\mathbf{Z}\mathbf{R}} \text{Div} \cdot \mathbf{I} = \nabla_{\mathbf{Z}\mathbf{R}} \text{Div}$$

⑦ Dimensionality → revisit

(*) On comparison of perceptron and backprop in terms of bias-variance

- OK with the intuitive sense in which the perceptron algo is more 'sensitive' in its estimate of $\hat{\mathbf{w}}$ to new training examples.
- AND the sense in which backpropagation for a linear classifier is not as 'sensitive' to addition of new training examples.
- But I am unsure as to how to formulate this formally
- ↗ intuitively variance is related to error from small sensitivity to fluctuations in training

→ ⑧-revisit

⑧ Claims regarding backprop generally

→ see paper
 'Backprop fails to sep...'
 - Brady et al. (1989)

⑨ Structure of the error surface

- lots of excellent papers
- At this stage, jury's out

⑩ Pedagogical

(*) convergence rate → analysis definition of convergence

(*) definition of R is OK

(*) But explanation of linear convergence; and exponentially fast is a little unclear.

(*) There are some more rigorous definitions in Nocedal & Wright (1999) and Tibshirani's 10-725 notes

(*) Quadratic surfaces

- (*) Yes; Taylor series can be seen as encoding all of derivatives of a function f about the point at which it is being expanded i.e. at $\hat{w}^{(k)}$.

$$(*) \text{ (Aii)} - E'(w^{(k)}) + E''(w^{(k)})(w - w^{(k)}) = 0$$

$$\Rightarrow E'(w^{(k)}) + wE''(w^{(k)}) - w^{(k)}E''(w^{(k)}) = 0$$

$$\Rightarrow w = \frac{w^{(k)}E''(w^{(k)}) - E'(w^{(k)})}{E''(w^{(k)})}$$

→ much clearer

$$\Rightarrow w_{\min} = w^{(k)} - \frac{E'(w^{(k)})}{E''(w^{(k)})} \quad \text{as required}$$

(*) Key step is to compute

i) Analytic optimisation of w^* via Taylor approx of quadratic error about an arbitrary w_k ; calculus

ii) gradient descent.

(i) occurs in this controlled setting of quadratic error

(*) compute: gradient descent $w^{(k+1)} = w^{(k)} - \eta \frac{dE(w^{(k)})}{dw}$

analytic: $w_{\min} = w^{(k)} - \frac{E'(w^{(k)})}{E''(w^{(k)})}$

(*) notationally; $\frac{dE(w^{(k)})}{dw} = E'(w^{(k)})$

(*) Hence: setting $\eta = \frac{1}{E''(w^{(k)})}$ in gradient descent algo

(*) essentially showing how learning rate could be set analytically and associated intuition.

- Q1 - I can't see formal reasons why:-
- $\eta < \eta_{\text{opt}} \Rightarrow$ monotonic convergence
 - $\eta_{\text{opt}} < \eta < 2\eta_{\text{opt}} \Rightarrow$ oscillating convergence
 - $\eta > 2\eta_{\text{opt}} \Rightarrow$ divergence.
- see PL book
to do with
rel. of gradient
and Hessian
I suspect

Q2 - add to queries

(*) Generic diff. convex objectives:-

$$E = E(w^{(k)}) + (w - w^{(k)}) \frac{dE(w^{(k)})}{dw} + \frac{1}{2} (w - w^{(k)})^T \frac{d^2 E(w^{(k)})}{dw^2} + \dots$$

(*) Taylor expand error, $E(w)$ about $w^{(k)}$

(*) Using some logic; set

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

Q3 - relation between Taylor series and Newton's method \rightarrow refresh

Multivariate $E(w)$

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

- for $w \in \mathbb{R}^n$ i.e. $w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$, diagonal A

- Plot contours $z = E(w_1, w_2)$

- point is that you can decouple $E(w)$ into 2 independent quadratics in w_1 and w_2 (i.e. components of w); each with their own coefficients.

- can use earlier expressions η_{opt} to find a learning rate for each quadratic

Q4 - Is η actually a vector? in this case?

- currently assuming that η is a scalar; but wouldn't issue BR be noting; i.e. learning rate η is same for all components of $\frac{dE(w^{(k)})}{dw}$ or $\eta \nabla_{w_i} E$ just an artefact of his notation.

- what is what prevents us from using:-

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

(W): Think
clarity would
be gained by
better ind. directional
deriv.

(*) For now note that for MV $E(\underline{w})$ with decoupled quadratics, we have $M_{1,\text{opt}}$ and $M_{2,\text{opt}}$ and that it will not in general be the case that $M_{1,\text{opt}} = M_{2,\text{opt}}$. ; while η is scalar \rightarrow issues

(*) Note the above is heuristic illustration and that convex. beh. gets unpredictable

(*) Formal results on convergence of gradient descent in language of convex optimisation, in part:-

- i) quadratic strongly convex fns.
- ii) lipschitz-smooth convex fns.

(*) Axis rescaling

- (*) formally; may be a form of weight normalisation (see paper by Salimans, Kingma (2016))

(*) BR shows this heuristically;

(*) essentially, amounts to rescaling weights $\hat{\underline{w}} = S \underline{w}$

→ rewriting $E(\underline{w})$ as $E(\hat{\underline{w}})$; solving for diagonal scaling mat S ; relating $\nabla_{\underline{w}} E$ to $\nabla_{\hat{\underline{w}}} E$.

② ensit this \rightarrow involved calcul. (not involved, but I need to make on)

③ :- key eq. in notes

(*) generic differentiable multivariate convex fns

$$E(\underline{w}) \approx E(\underline{w}^{(k)}) + \nabla_{\underline{w}} E(\underline{w}^{(k)}) (\underline{w} - \underline{w}^{(k)}) + \frac{1}{2} (\underline{w} - \underline{w}^{(k)})^T H_E(\underline{w}^{(k)}) (\underline{w} - \underline{w}^{(k)}) + \dots$$

$$\underline{w}^{(k+1)} = \underline{w}^{(k)} - \eta H_E(\underline{w}^{(k)})^{-1} \nabla_{\underline{w}} E(\underline{w}^{(k)})^T \rightarrow \textcircled{2} \quad \text{OLS}$$

(*) ILD with Q.A. - O/S

(?) confused by distinction between applying normalisation and use of Hessian → O/S.

DL - 1st order opt - gradient

2nd order opt - Hessian

and Nocedal & Wright (1999)

(*) some issues regarding Hessian:

i) computationally infeasible → matrix inves.

ii) non convex fn ⇒ Hessian not PSD or PD : divergence

(*) cluster of methods aimed at approximating (i) → i.e. approx Hessian.

(*) learning rate / step size η

- previous discussion → $\eta < 2\eta_{\text{opt}}$ to prevent divergence for convex fns
(key)

(*) looking forward; there are 2 main methods other than fixing learning rate η :

i) use a decay schedule

ii) Adaptive methods

(*) decay schedules for learning rates often used with validation set

- escaping local minima vs convergence

(?) - unclear whether the sufficient conditions specified in this lecture for convergence hold under which assumptions.

- results directly lifted from SGD Ch 8.3.1. DL.

(*) in any case; covered in bb in more detail → details.

