

17: SGD, overfitting, regular., tricksRecap

- Rprop, quickprop \rightarrow dynamics lost \rightarrow efficiency issues
- Momentum smooth fluctuations via using running averages
- Polyak \rightarrow rad
- Nesterov's method
- Nesterov paper (in Russian)
- Incremental methods: size of training set presented.
 - potentially leads to high variance of estimates
- (*) Convergence for incremental updates depends on learning rate (schedule)
 - sufficient conditions on learning rate
- (*) Still require learning rate schedule

(*) Rand algos (topic)

- every time you update, doing it on different loss (in the sense that it is own completely different set of samples)
- variance in this case refers to differences of loss functions and recommended updates across sample draws.

(*) Momentum helps (in this case?) - using running averages

- dampen/smooth oscillations (of what?) ; explores over error surface?

(*) Incremental SGD, mini-batch \rightarrow momentum may perform variance cleanup(*) minibatch with momentum pseudocode(*) Nesterov's accelerated gradient

- Superior to standard momentum
- provably faster (in Russian paper)

(*) ALSO can be combined with incremental methods

(B) momentum style methods using running averages
shrink variance of estimated gradients

(A) ^(A) you need
to clarify variance
of what exactly

(BR): what is being smoothed using running average

BR: only using 1st moment of derivatives (here may be other moments using leveraging other statistical info)
momentum

(W) (A): compute expected value of derivative
at the next step; take step along expected value \rightarrow sufficient?

- expectation may be 0; but variance may be high (~~do not~~ ^{do not} first 1st mom)
or variance near 0 (1st mom ok)

smoothing trajectory:

(*) horizontal direction: - steps always positive

vertical -||- :- steps swinging sign

(*) this is after momentum is applied (?)

- make use of 2nd moments of corrections

(*) momentum \rightarrow running average; but still have learning rate

BR: Account for 2nd moment of deriv by adjusting learning rate (?)

BR: dampen step size in directions of high variance

Increase -||- in -||- low variance.

(*) normalise steps with variance (high-level)

- BR: compare net movement vs total movement
(arbitrary) direct.)

i) if in past, total movement in y -comp high } shrink step size
net movement -||- small } η_j

ii) -||- total movement large } increase step size
net movement large } η_j

(*) modify gradient based update with relative scaling in axis dir; according to variation.

RMS prop (Hinton)

- Notation:- $\partial_w D$ is seen as $\frac{\partial D}{\partial w}$ (derivative)

$\partial_w^2 D$ is not $\frac{\partial^2 D}{\partial w^2}$ but $\left[\frac{\partial D}{\partial w} \right]^2$ (i.e. squared deriv.; not 2nd.)

$$\mathbb{E}[\partial_w^2 D] = \mathbb{E}\left[\left(\frac{\partial D}{\partial w}\right)^2\right] \quad (\text{mean squared deriv.})$$

RMS Prop

(*) focus is mean squared derivative:

use
(*) 2nd moment
of derivative

$$(*) \quad \mathbb{E}[\partial_w^2 D]_K = \gamma \mathbb{E}[\partial_w^2 D]_{K-1} + (1-\gamma) (\partial_w^2 D)_K$$

$$(*) \quad w_{K+1} = w_K - \frac{\eta}{\sqrt{\mathbb{E}[\partial_w^2 D]_K + \epsilon}} \partial_w D$$

- standard gradient descent (*)

- normalising learning rate by RMS value (*)

(*) similarity to Rprop.

RMSprop pseudocode

(*) maintain running average of 2nd moment of the derivative $\mathbb{E}[\partial_w^2 D]_K$

(*) normalise by square root of 2nd moment.

Q: where is the 1st moment?

- momentum / nesterov - 1st moment; RMSprop - 2nd moment

- BR: try both?

ADAM: RMSprop with momentum

- ADAM utilises smoothed version of momentum-augmented gradient with both 1st and 2nd moments

- (*) - maintain running estimate of mean derivative for each param.
- maintain " " of mean squared value of derivatives
- scale update of param by inverse of root mean squared deriv.

$$(6) \quad m_R = \delta m_{R-1} + (1-\delta)(\partial_w D)_K$$

②② - surely divide by 0?

$$v_R = \gamma v_{R-1} + (1-\gamma)(\partial_w^2 D)_K$$

$$\hat{m}_R = \frac{m_R}{1-\delta^R} \quad \hat{v}_R = \frac{v_R}{1-\gamma^R}$$

$$w_{R+1} = w_R - \frac{\eta}{\sqrt{\hat{v}_R + \epsilon}} \hat{m}_R$$

- (*) $(1-\delta^R)$ - ^{typ.} delta values close to 1 → initialisation of m_{R-1} and v_{R-1} at 0
→ δ^R close to 0 as δ close to 1
- slow update at early iterations of means/variances; 1st, 2nd moments

AD-HOC ext/correction

⇒ early iterations → don't prioritise $(1-\delta)$
As $R \rightarrow \infty$, δ^R and γ^R will become large; (0)

④④: A little conf.
here.
- Review this

- through this lens; can generate more ideas/tricks

- Alex Radford animations ④④ - ④ - animations

long valley

- RMSprop - steps become fixed when slope is constant

- derivative $= 1 \Rightarrow$ fixed step size
RMS value of deriv

(*) Recall

$$\text{loss} = \frac{1}{T} \sum_t \text{div}(y_t, d_t; w_1, w_2, \dots, w_K)$$

(error)

BR: ability to minimise loss ; and hence divergence depends on nature of divergence.

- results in good divergence

- (*) Hill climbing is a really good analogy for this

(*) more 'bowly' than quadratic \rightarrow flat near min.

(*) divergence depends on activation (pre-act)

- L₁ divergence and KL divergence

quad ?

(*) Think that L₂ is bowl-shaped X wrong?

- KL divergence same as cross-entropy, with extra term.

L₂ or KL single (logistic fn)

- consider perceptron with 2 inputs, no bias

- compute output over training samples

- compute output of logistic function

- compute L₂ div / KL div on logistic output and target outputs.

- Parameters are the 2 weights of logistic function (?)

- 2 dimensions

- (*) (*) - It's not the shape wrt output of classifier

- rather the shape of loss wrt to parameters of the classifier

↓

KL divergence \rightarrow convex bowl ; L₂ - non convex
(although global minimum)

(*) We don't know what loss fn looks like

- examine in case near end of network

→ pedagogical
heuristic for
intuition

- how do methods leave scaddle points
- (*) local behavior depends on nature of error surface (where gradients point)
- (*) (P): difficult to say which is optimal → we don't know what error surface in general looks like at any point
- But with loss function of a particular kind; we can recommend technique as we know how they behave.

(*) story so far
- see slides.

(*) convergence can be improved with smoothed updates
e.g. RMSprop, ADAM etc.

(*) distinction between ^{variance from} randomly sampling data differently
AND natural variance over entire batch

(P) need assurance of convergence to 'right spot' over entire training data.
(learning rate)

(Q) using normalisation of step size via 2^{nd} moments
- do we reduce learning rate according to a schedule as before?

BR: learning rate adjusts

- 2^{nd} moment techniques - do not exogenously adjust learning rate η
as 2^{nd} moment adjusts it automatically; excludes need for finding out
optimal learning rate / schedule

BUT: we have other hyperparams to tune.

(related to initial step size; or how running avg. maintained)
- no such thing as a free lunch

(P) generalisation/tricks of the trade:

- divergence
- dropout
- batch norm
- gradient clip, data aug, hecks

See papers

backprop → error backprop

$$\nabla_{\theta} \frac{1}{2} \|y - d\|_2^2 = (y - d) \nabla y(z)$$

- propagate $(y - d)$ backwards

(*) choice of divergence affects both learned network and results

(*) covariate shifts

- Assumption using mini-batch :-

(*) each mini-batch has more or less same statistics as every other mini-batch
so that minimising for one mini-batch \Rightarrow close to minimum for other
mini-batches; and that overall mini-batches; correct soln!

(*) when minibatches are far apart,
minimisation for one does not give info about others.

(*) (A4): inconfortable \rightarrow presentation too informal review

- is this to do with uniform convergence?

- is there a better way of presenting using this?

- what is 'space' of the training samples shown

Q. don't understand his heuristic presentation of moving mini-batches

\rightarrow standard mean 0, unit variance ; then move
pile to where they think they lie if considered all of them at once

Batch normalization

(*) covariate adjustment after the weighted addition of inputs (affine trans.)
and before activation

Training: Adjustment
our mini-batches.

- done independently for each unit

(*) Purpose of batch norm \rightarrow cancel out covariate shifts

batch norm:

- standardise mini-batch; 'move it to correct pos'; to be determined via a parameter setting

⑦

A6: review + 1st diagram: → ⑧ ⑨ ⑩

- compute mean of mini-batch μ_B and σ_B^2 (estimators?)
- normalise instances of mini-batch (every)
- then conduct shift to a unit-specific location (per neuron)
 - as param for each neuron lie in a specific loc in param space

- (*) comprised of 2 operations (key)

derivatives: exploits additivity

- batch norm complicates additivity.

$$\text{div}(MB) = \frac{1}{B} \sum_t \text{div}(y_t(x_t, \mu_B, \sigma_B^2), d_t(x_t))$$

- divergence for all y_t depends on all x_t within minibatch.

Batch norm - vector fn. over minibatch

⑧: Review

- (*) shows not dimensions of input; but B instances of mini-batch
- (*) each line rep one training instance of minibatch at one node in the network

(*) note dependence / influence diag style arg.

- derivative computed globally over entire minibatch must consider all z_j 's wrt z_i

- not locally for every instance in mini-batch

- one-to-one interactions.

Batchnorm dependency diagram ⊗ ⊕

(A8) Review

(A8) Review
(A) This should be sufficient to write derivatives. (next slides)

-Batchnorm - backprop

(*) Complexity in backprop happens after $\frac{d\text{Div}}{d\bar{z}} = f'(\bar{z}) \frac{d\text{Div}}{dy}$ step.

(k) For each neuron have γ and β parameters

-Have to compute driv not ρ or β also.

Params are multi-batch specific during forward pass

- running mean-separate

BR: $\Sigma = \mu n + \beta$ *n-computed
during forced pass*

div

$$\frac{d \text{Div}}{d\beta} = \frac{d \text{Div}}{d\tilde{\ell}} \cdot \frac{d \tilde{\ell}}{d\beta}$$

\approx
 $= 1$

$$= \frac{d\text{Div}}{d\hat{z}}$$

$$\frac{d\text{Div}}{d\gamma} = \frac{d\text{Div}}{dZ} \frac{dZ}{d\gamma} = \frac{d\text{Div}}{dZ} u$$

- Have to keep backpropagating

$$\frac{d\text{Div}}{du} = \frac{d\text{Div}}{dz} \frac{dz}{du} = \frac{d\text{Div}}{dz} \gamma$$

BR: issue is taking $\frac{\partial \text{Div}}{\partial u}$ back to unnormalised z

Final step:- $\frac{\partial \text{Div}}{\partial z_i}$

$$\text{Div} = \text{function}(u_i, \mu_B, \sigma_B^2) \quad u_i = f(z_i, \mu_B, \sigma_B^2)$$

$$\frac{\partial \text{Div}}{\partial z_i} = \frac{\partial \text{Div}}{\partial u_i} \frac{\partial u_i}{\partial z_i} + \frac{\partial \text{Div}}{\partial \sigma_B^2} \frac{\partial \sigma_B^2}{\partial z_i} + \frac{\partial \text{Div}}{\partial \mu_B} \frac{\partial \mu_B}{\partial z_i}$$

(A7): Review

(A8): Use influence diagram

(to account for all paths)

- remember route of influence topologically through network

(*) derivation of $\frac{\partial \text{Div}}{\partial \sigma_B^2}, \frac{\partial \sigma_B^2}{\partial z_i}$

(*) derivation of $\frac{\partial \text{Div}}{\partial \mu_B}, \frac{\partial \mu_B}{\partial z_i}$

(*) derivation of $\frac{\partial \text{Div}}{\partial u_i}, \frac{\partial u_i}{\partial z_i}$

(A9)

BR: $\frac{\partial \text{Div}}{\partial z_i}$ - simple formula \rightarrow work though it

Batch normalis: inference

(A10): A lot to review \rightarrow (A)

Batch normalis.

(*) no need to apply to all layers

(*) Anecdotal evidence \rightarrow no need for dropout

(*) Batch norm: within each batch for each neuron

- for sigmoids, moves data to where 'it's most linear'

i.e. 0 position (?)

- for ReLUs: using 0 as standard mean, not good idea

- due to $\frac{1}{2}$ data ending up in '0 area' where there is no derivative

batchnorm-empirical diff

greatly

- Ioffe & Szegedy (2015) - reduce training time

Data under-specification

- 2D plot; dense training points (see divergence diagrams)
- overfitting
- BR(x): get a sense of how little data we have (esp. in high dimensions)
e.g. 100 dim, binary inputs
- inputs on corners of 100-dim hypercube
 - How many corners? 10^{30} locations where data can exist
- e.g. 1 quadrillion training instances - 10^{15}

⑩ High dimensionality

$$\frac{10^{15}}{10^{30}} = \frac{1}{10^{15}}$$

(*) Require additional constraints

Weight manipulation

- (*) indiv neurons capable of arbitrary changes in steepness
 - (A8) How does sigmoid activation change in response to inputs / affine input weights changing?

(*) \rightarrow Heaviside step rather than smooth fn.

(*) a) Steep changes \rightarrow overfitted responses \rightarrow perceptions with large w

(*) constraint weights to be low \downarrow
smooth output

(regularisation)

- place constraints on weights

- Add regularisation via modifying loss via weight regulariser

Regularising weights

- combinatorial regulariser

Network structure

- (*) to impose smoothness

(*) for a given no. param deeper nets impose more smoothness than shallow ones.

- each layer places constraints on what can be represented

- next layer treats what it already has

- (*) depth and training data for deeper nets.

III - Review insight

- L₁ regularisation
 - Noise added reg.
 - Dropout
- } ^(a12) → go through + review