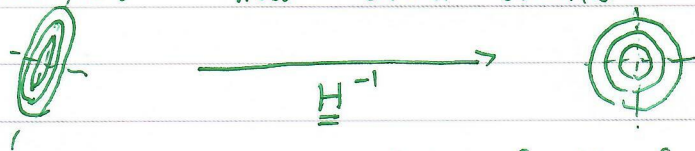* difficulty in DL Optimization
  - local min : not as bad as imagined
  - saddle / local max : easily escapable (esp. w/ SGD)
  - plateau : need larger $\eta$ (learning rate) ;
  - ravines : need avoid oscillation



$$\underline{\underline{H}}^{-1}$$

(Quasi-)Newton, but infeasible for DL usually

- slow computation of gradient : SGD on minibatch

$$\Downarrow$$

"instable" estimate of gradient

* ┃running┃ average estimate ┃ of gradient from ┃SG┃

$$\boxed{\underline{V_t}} = \beta \cdot \boxed{\underline{V_{t-1}}} + (1-\beta)\boxed{\underline{\Delta}_t}$$

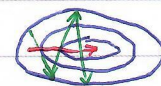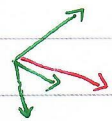current average    previous average    $\searrow$ stochastic gradient at $t$-th iteration

update : $\underline{W_t} = \underline{W}_{t-1} - \eta \, \underline{V}_t$

(SGD with) momentum
- stochastic instability cancels out

- pass ravine faster



* ┃oscillation┃ prevention by per-component learning rate

larger gradient component
want: smaller step

update : $\underline{W_t} = \underline{W}_{t-1} - \eta \; \underline{\Delta}_t \cdot / \; "magnitude(\underline{\Delta}_t)"$

$$\frac{\underline{ill}}{\sqrt{\underline{u}_t + \epsilon}}$$

running average

$$\underline{u}_t = \beta \, \underline{u}_{t-1} + (1-\beta)(\underline{\Delta}_t \cdot * \underline{\Delta}_t)$$

RMSprop

\* Adam = SGD + momentum + RMSProp (+ other tricks)

$$\underline{V}_t = \beta_1 \underline{V}_{t-1} + (1-\beta_1) \underline{\Delta}_t$$

$$\downarrow$$
$$0.9$$

$$\underline{u}_t = \beta_2 \underline{u}_{t-1} + (1-\beta_2)(\underline{\Delta}_t \cdot {*}\ \underline{\Delta}_t)$$

$$\downarrow \qquad\qquad 0.001 \leftarrow \dfrac{\eta}{\sqrt{t/N}} \text{ (decaying)}$$
$$0.999$$

$$\underline{W}_t = \underline{W}_{t-1} - \dfrac{\eta_t}{\sqrt{\underline{u}_t} + \epsilon} \cdot {*}\ \underline{V}_t$$

$$\downarrow$$
$$10^{-8}$$

\* difficulty in DL generalization

- co-adaptation :

(consistent mistakes from some neurons) corrected by fitting other neurons

$$\Downarrow \qquad\qquad\qquad\qquad\qquad\qquad\qquad \Downarrow$$

like noise                                like overfit

and vice versa

(dependence)

\* break the dependence : shut down neurons (randomly)



mini-batch-1                    mini-batch-2

dropout : drop P
         keep (1-p)  $\Rightarrow$ many thin networks (combined)

$$\Downarrow$$

slows down convergence          noise as regularization

but faster per-iteration                    like DAE

**\* dropout during testing**

- full-net prediction w/o changing

$$E(S_{i,}^{(l)}) = (1-p) S_i^{(l)}$$

in training           in testing

- test-time "pseudo-" **dropout**

$$x_i^{(l)} = \theta\left((1-p)\, S_i^{(l)}\right)$$

need to record $p$, less flexibility
for changing $p$ per neuron
or dynamically

- inverted **dropout**

training :   **dropout**   &   $x_i^l = \theta\left(S_i^{(l)} / (1-p)\right)$

testing :      $x_i^{(l)} = \theta\left((1-p')\, S_i^{(l)} / (1-p')\right)$

$$= \theta\left(S_i^{(l)}\right) \quad\quad \text{unchanged} \quad (\text{usually preferred})$$