\* Initialization for Deep Learning

activation $\Rightarrow$ bad $W_{ij}^{(\ell)}$ properties

e.g. tanh : $\boxed{|W_{ij}^{(\ell)}| \text{ too big}}$ $\longrightarrow$ saturation

$\boxed{W_{ij}^{(\ell)} = 0, \text{ or constant}}$ $\rightarrow$ saddle , symmetry

ReLU : $\boxed{W_{ij}^{(\ell)} \text{ too negative}}$ $\longrightarrow$ saturation ( signal shutdown)

\* want :, random small

usually 0-mean

uniform $[-L, +L]$

(truncated) Gaussian $(0, \sigma^2)$

ReLU max$(s, 0)$

weight variance $\dfrac{2}{d^{(\ell-1)}}$

tanh

weight variance $\dfrac{1}{d^{(\ell-1)}}$

weight variance $\dfrac{1}{d^{(\ell)}}$

$\approx \boxed{\begin{array}{c} \text{variance of } x_i^{(\ell-1)} \\ \hline \text{variance of } x_j^{(\ell)} \end{array}}$

w/ assumptions

$\Longrightarrow$ "

tanh$(s) \approx s$

variance of $\delta^{(\ell-1)}$ $\approx$ variance of $\delta^{(\ell)}$ when tanh$(s) \approx s$

He init.

weight variance $\dfrac{2}{d^{(\ell-1)} + d^{(\ell)}}$

$[-U, U]$ uniform

$\dfrac{6}{\sqrt{1 + d^{(\ell-1)} + d^{(\ell)}}}$ $\leftarrow$

$- - - - - - - -$ Glorot init.

\* difficulty in DL optimization

surface
- local min: not as bad as imagined
- saddle / local max: easily escapable (esp. w/ SGD)
- plateau : ___ need [larger $\eta$] (learning rate)
- ravine , avoid oscillation



hessian $\underline{\underline{H}}^{-1}$ &w

(Quasi-)Newton , not feasible for DL
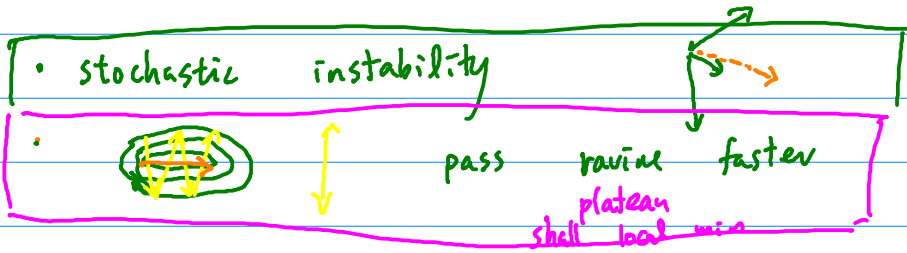
- slow computation of gradient : SGD on minibatch
  ⇓
  "instable" estimate of gradient

---

\* [running] average estimate of SG        (SGD) with __momentum__

$$\underline{V}_t = \beta \cdot \underbrace{[\underline{V}_{t-1}}_{\text{previous average}} + (1-\beta)\underline{\Delta}_t$$

$\underline{V}_t$ → Current average

→ SG at t-th iteration

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

- stochastic instability



pass ravine faster
     plateau
shall local min



larger gradient component
want: smaller step          )     per-component learning rate

update :   $\underline{W}_t = \underline{W}_{t-1} - \eta \cdot \underline{\Delta}_t \boxed{\cdot /}$  "magnitude$(\underline{\Delta}_t)$"
                                                          |||
running average

$$\underline{U}_t = \beta \cdot \underline{U}_{t-1} + (1-\beta)\underbrace{\underline{\Delta}_t \boxed{\cdot *} \underline{\Delta}_t}_{\text{RMS prop}}$$

$\sqrt{\underline{U}_t + \epsilon}$

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

$$V_t = \boxed{\beta_1} V_{t-1} + (1-\beta_1) \Delta_t$$

(0.9 above $\beta_1$)

$$U_t = \boxed{\beta_2} U_{t-1} + (1-\beta_2)(\Delta_t \cdot * \Delta_t)$$

(0.999 below $\beta_2$)

$$W_t = W_{t-1} - \frac{\eta_t}{\sqrt{U_t} + \epsilon} \cdot * V_t$$

0.001

$\dfrac{\eta}{\sqrt{t/N}}$    decaying

$\rightarrow 10^{-8}$