

⑦

Training Neural Network (II)

- Fancy optimization
- Regularization
- Transfer Learning

* Optimization

→ Core strategy in training neural network is an optimization problem.

① # Problems with SGD

- ① → If loss changes quickly in one direction and slowly in another.



{ Loss function has
high condition number }

- ⇒ Very slow progress along shallow dimension, jitter along steep direction.

- ⇒ This problem is much more common in higher dimension.

- ② → ~~If the loss function has a local minimum or saddle point.~~

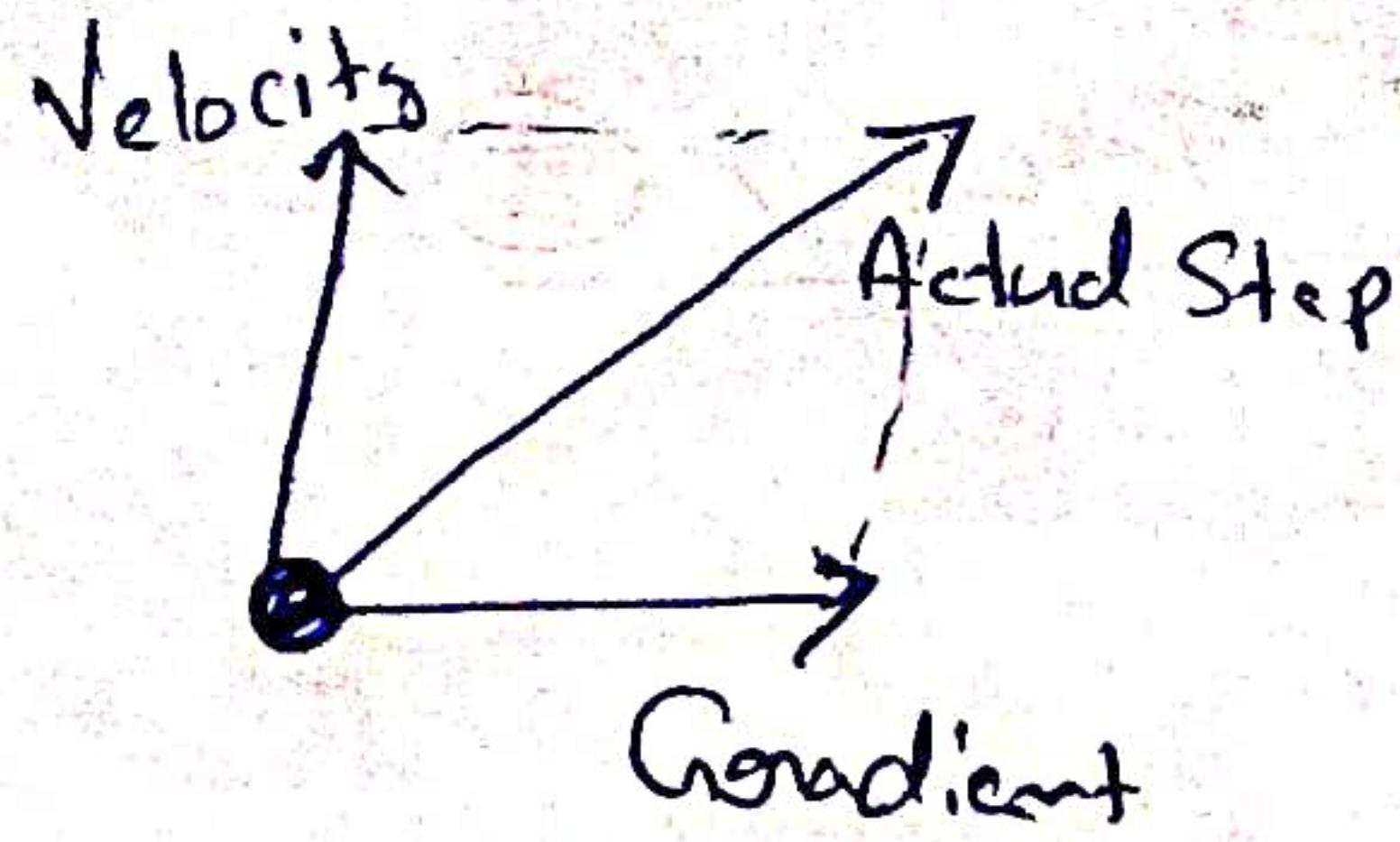
① SGD + Momentum

$$V_{t+1} = \gamma V_t + \nabla f(x_t)$$

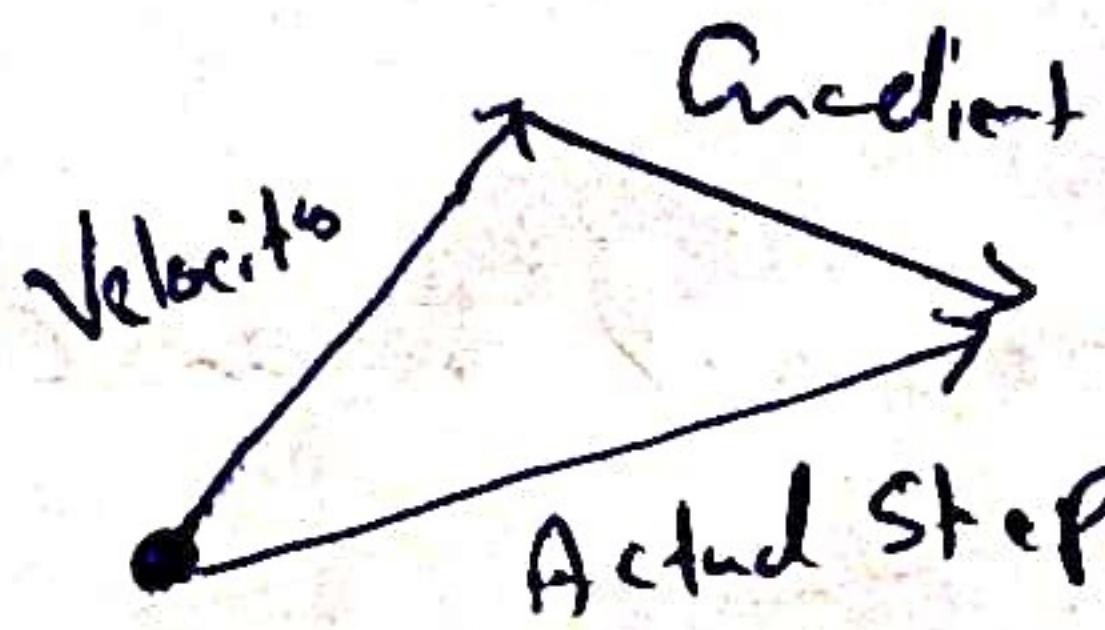
$$x_{t+1} = x_t - \alpha V_{t+1}$$

⇒ This simple strategy solves

- Local minima
- Saddle points
- & poor conditioning problems.



② Nesterov Momentum



$$V_{t+1} = \gamma V_t - \alpha \nabla f(x_t + \gamma V_t)$$

$$x_{t+1} = x_t + V_{t+1}$$

⇒ Change of variable $\tilde{x}_t = x_t + \gamma V_t$ and rearrange;

$$V_{t+1} = \gamma V_t - \alpha \nabla f(\tilde{x}_t)$$

$$\tilde{x}_{t+1} = \tilde{x}_t + V_{t+1} + \gamma (\nabla V_{t+1} - \nabla V_t)$$

⇒ Adding momentum has an effect of overshooting the minimum.

③ AdaGrad

$$\text{grad_squared} = 0$$

while True:

$$dx = \text{Compute-Gradient}(x)$$

$$\text{grad_squared} += dx * dx$$

$$x = \text{learning_rate} * \frac{dx}{\sqrt{\text{grad_squared}} + 1e-7}$$

④ RMSProp

Grad-Squared = 0

While True:

$dx = \text{Compute-gradient}(x)$

$$\text{Grad-Squared} = (\text{decay-rate} * \text{Grad-Squared}) + (1 - \text{decay-rate}) * dx * dx$$

$$x -= \text{learning-rate} * \frac{dx}{\sqrt{\text{Grad-Squared}} + 1e-7}$$

⑤ Adam (almost)

first-moment = 0

Second-moment = 0

While True:

$dx = \text{Compute-gradient}(x)$

first-moment = $\beta_1 * \text{first-moment} + (1 - \beta_1) * dx$

Second-moment = $\beta_2 * \text{Second-moment} + (1 - \beta_2) * dx * dx$

$$x -= \text{learning-rate} * \frac{\text{first-moment}}{\sqrt{\text{Second-moment}} + 1e-7}$$

\Rightarrow Sort of like RMSProp with momentum,

⑥ Adam (full form)

$$\text{first-moment} = 0$$

$$\text{Second-moment} = 0$$

For t in range (num-iterations):

$$dx = \text{Compute-Gradient}(x)$$

$$\text{first-moment} = \text{beta1} * \text{first-moment} + (1 - \text{beta1}) * dx$$

$$\text{Second-moment} = \text{beta2} * \text{Second-moment} + (1 - \text{beta2}) * dx^2$$

$$\text{First-bias} = \frac{\text{first-moment}}{1 - \text{beta1}^{**t}}$$

$$\text{Second-bias} = \frac{\text{Second-moment}}{1 - \text{beta2}^{**t}}$$

$$x' = \text{learning-rate} * \frac{\text{First-bias}}{\sqrt{\text{Second-bias}} + 1e-7}$$

⇒ Bias correction for the fact that first & second moment estimates starts at zero.

⇒ Adam with $\text{beta1} = 0.9$

$$\text{beta2} = 0.999$$

$$\text{learning-rate} = 1e-3 \text{ or } 1e-5$$

is a great starting point for many models!

* Learning rate decay over time!

④ Step decay

⇒ Decay learning rate by half every few epochs.

⑤ Exponential decay

$$\alpha = \alpha_0 e^{-kt}$$

⑥ Yt decay

$$\alpha = \frac{\alpha_0}{1+kt}$$

⇒ More critical with SGD + Momentum
↳ Less common with Adam.

* Second-Order Optimization

- (1) Use gradient & Hessian to form quadratic approximation
- (2) Step to the minima of the approximation.

⇒ Second-order Taylor expansion:

$$J(\theta) \approx J(\theta_0) + (\theta - \theta_0)^T \nabla_{\theta} J(\theta_0) + \frac{1}{2} (\theta - \theta_0)^T H(\theta - \theta_0)$$

⇒ Solving for the critical point we obtain the Newton Parameter update:

$$\theta^* = \theta_0 - H^{-1} \nabla_{\theta} J(\theta_0)$$

⇒ Hessian has $\mathcal{O}(N^2)$ elements, invert it take $\mathcal{O}(N^3)$

↳ For deep learning this is not a good idea
 ↳ $N = (\text{Tens or Hundreds of}) \text{ Millions.}$

* Quasi-Newton methods (BFGS most popular)
↳ This works with approximation of Hessian.
⇒ L-BFGS (Limited memory BFGS)

- Does not form/store the full inverse Hessian.
- Usually works very well in full batch deterministic mode
- Does not transfer very well to mini-batch setting.

⇒ In practice:

- Adam is a good default choice in most cases
- If you can afford to do full batch updates then try out L-BFGS.

* Model Ensembles

1. Train multiple independent models
 2. At test time average their results
Enjoy 2% extra performance
- ⇒ Instead of training independent models, use multiple snapshots of a single model during training.
- ⇒ Instead of using actual parameter vector, keep a moving average of the parameter vector and use that at test time. (Polyak averaging)

* Regularization: Dropout

- ⇒ In each forward pass, randomly set some neurons to zero. Probability of dropping is a hyperparameter (0.5 is common)
- ⇒ Forces the network to have a redundant representation;
 - ↳ Prevents co-adaptation of features.
- ⇒ Dropout makes an output random!

$$y = f_w(x, z)$$

Random mask

→ Want to "average out" the randomness at test time

$$y = f(x) = E_z[f(x, z)] = \int p(z) f(x, z) dz$$

But this integral seems hard...

- ⇒ At test time all neurons are active always.
 - ↳ We must scale the activations so that for each neuron:
- Output at test time = expected output at training time

⇒ More common inverted dropout

↓

In the train time divide by P
, to make test time fast

* Regularization: A common pattern

Training: Add some kind of randomness

$$y = f_w(x, z)$$

Testing: Average out randomness
(Sometime approximate)

$$y = f(x) = E_z [f(x, z)] = \int p(z) f(x, z) dz$$

* Data Augmentation

⇒ Random crops & scales

Training: Sample random crops / scales

Testing: Average a fixed set of crops

⇒ Color Jitter, translation, rotation,
, stretching, shearing, lens distortion etc..

⇒ Other regularizations:

- Dropout
- Fractional Max Pooling
- Stochastic Depth

* Transfer Learning

65 You need a lot of data if you want to train / use CNNs²²

	Very similar dataset	Very different dataset
Very little data	Use Linear Classifier on top layers	You're in trouble. Try linear classifier from different stages
Quite a lot of data	Fine tune a few layers	Fine tune a large number of layers

⇒ Takeaway for your projects & beyond:

Have some dataset of interest but it has $< \sim 1M$ images?

① Find a very large dataset that has similar data, train a big ConvNet there

② Transfer learn to your dataset.

⇒ Deep learning frameworks provide a "Model Zoo" of pretrained models so you don't need to train your own.

