



Software



# Regularization Techniques for Deep Learning

Materials from

- Intel Deep Learning <https://www.intel.com/content/www/us/en/developer/learn/course-deep-learning.html>
- Improving Deep Neural Networks <https://www.deeplearning.ai/>

# Legal Notices and Disclaimers

This presentation is for informational purposes only. INTEL MAKES NO WARRANTIES, EXPRESS OR IMPLIED, IN THIS SUMMARY.

Intel technologies' features and benefits depend on system configuration and may require enabled hardware, software or service activation. Performance varies depending on system configuration. Check with your system manufacturer or retailer or learn more at [intel.com](https://www.intel.com).

This sample source code is released under the [Intel Sample Source Code License Agreement](#).

Intel and the Intel logo are trademarks of Intel Corporation in the U.S. and/or other countries.

\*Other names and brands may be claimed as the property of others.

Copyright © 2017, Intel Corporation. All rights reserved.

# Regularizing Neural Networks

We have several means by which to help “regularize” neural networks  
– that is, to prevent overfitting

- Regularization penalty in cost function
- Dropout
- Early stopping "หยุดตรงจุดที่ point ที่เหมาะสม"
- Stochastic / Mini-batch Gradient descent (to some degree)

# Penalized Cost function

- One option is to explicitly add a penalty to the loss function for having high weights.
- This is a similar approach to Ridge Regression

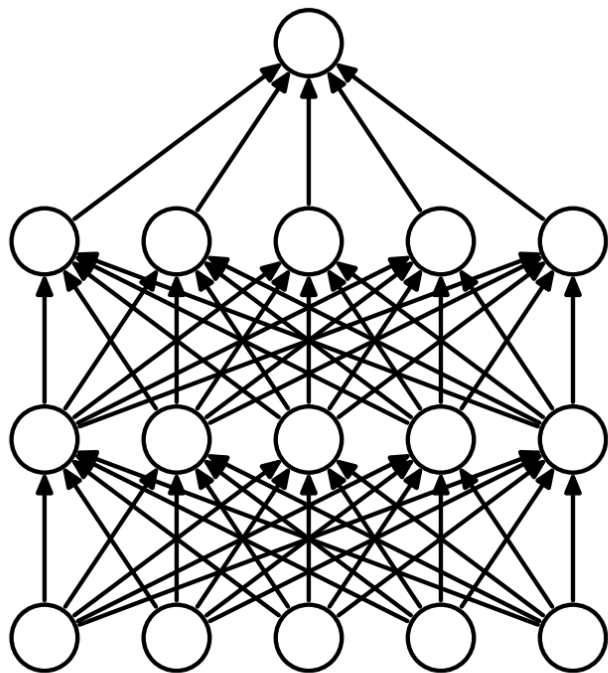
$$J = \frac{1}{2n} \sum_{i=1}^n (\hat{y}_i - y_i)^2 + \lambda \sum_{j=1}^m W_j^2$$

- Can have an analogous expression for Categorical Cross Entropy

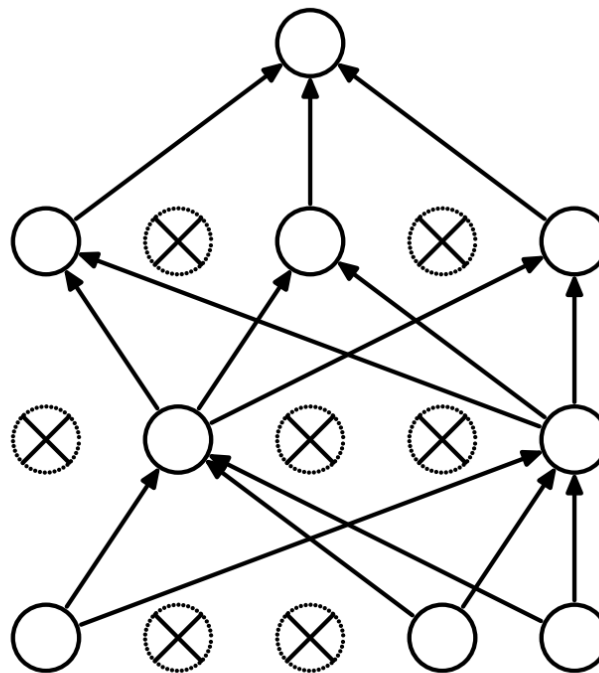
# Dropout

- Dropout is a mechanism where at each training iteration (batch) we randomly remove a subset of neurons
- This prevents the neural network from relying too much on individual pathways, making it more “robust”
- At test time we “rescale” the weight of the neuron to reflect the percentage of the time it was active

# Dropout - Visualization



(a) Standard Neural Net

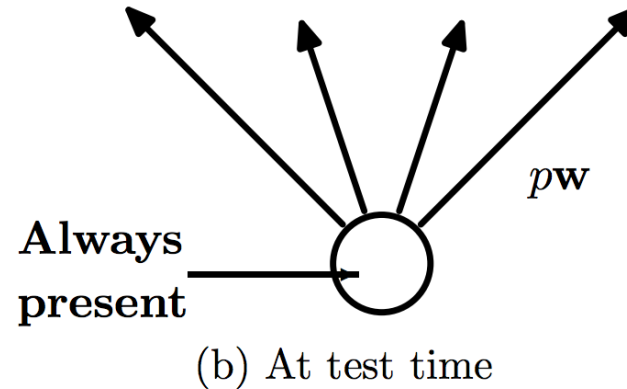
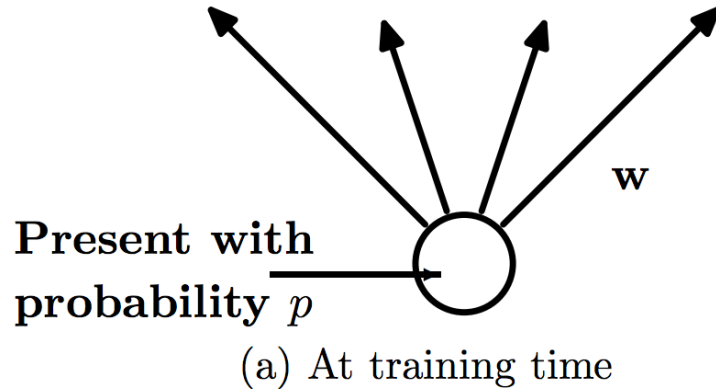


(b) After applying dropout.

probability  
 $p$  to keep : 0.4  
 $p$  to drop : 0.6

# Dropout - Visualization

- If the neuron was present with probability  $p$ , at test time we scale the outbound weights by a factor of  $p$ .



# Early Stopping

"Call back func"

called when some ... was triggered

- Another, more heuristical approach to regularization is early stopping.
- This refers to choosing some rules after which to stop training.
- Example:
  - Check the validation log-loss every 10 epochs.
  - If it is higher than it was last time, stop and use the previous model (i.e. from 10 epochs previous)



# Optimizers

- We have considered approaches to gradient descent which vary the number of data points involved in a step.
- However, they have all used the standard update formula:

$$W := W - \alpha \cdot \nabla J$$

- There are several variants to updating the weights which give better performance in practice.
- These successive “tweaks” each attempt to improve on the previous idea.
- The resulting (often complicated) methods are referred to as “optimizers”.

# Exponentially Weighted Averages

$$\theta_1 = 40^\circ\text{F} \quad 4^\circ\text{C} \leftarrow$$

$$\theta_2 = 49^\circ\text{F} \quad 9^\circ\text{C}$$

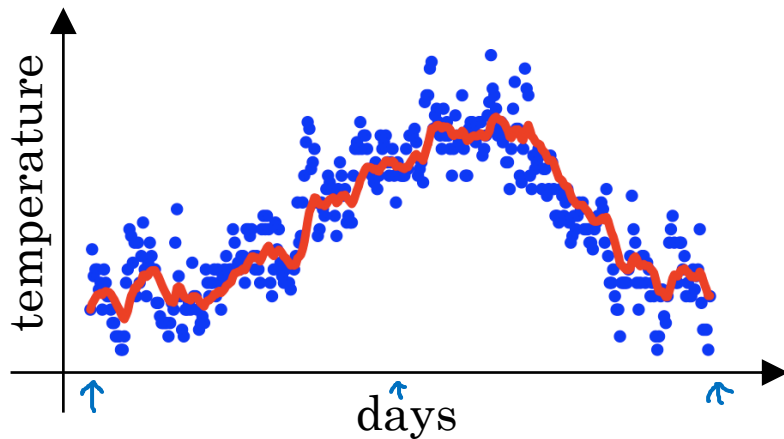
$$\theta_3 = 45^\circ\text{F} \quad \vdots$$

$\vdots$

$$\theta_{180} = 60^\circ\text{F} \quad 15^\circ\text{C}$$

$$\theta_{181} = 56^\circ\text{F} \quad \vdots$$

$\vdots$



$$V_0 = 0 \quad \text{beta}$$

$$V_1 = 0.9 V_0 + 0.1 \theta_1 \quad \text{1-beta} \quad \text{data}$$

$$V_2 = 0.9 V_1 + 0.1 \theta_2$$

$$V_3 = 0.9 V_2 + 0.1 \theta_3$$

$\vdots$

$$V_t = 0.9 V_{t-1} + 0.1 \theta_t$$

# Exponentially weighted averages <sup>moving</sup>

$$\underline{V_t} = \underline{\beta} \underline{V_{t-1}} + \underline{(1-\beta)} \underline{\Theta_t} \leftarrow$$

$\beta = 0.9$  :  $\approx 10$  days' temper.

$\beta = 0.98$  :  $\approx 50$  days

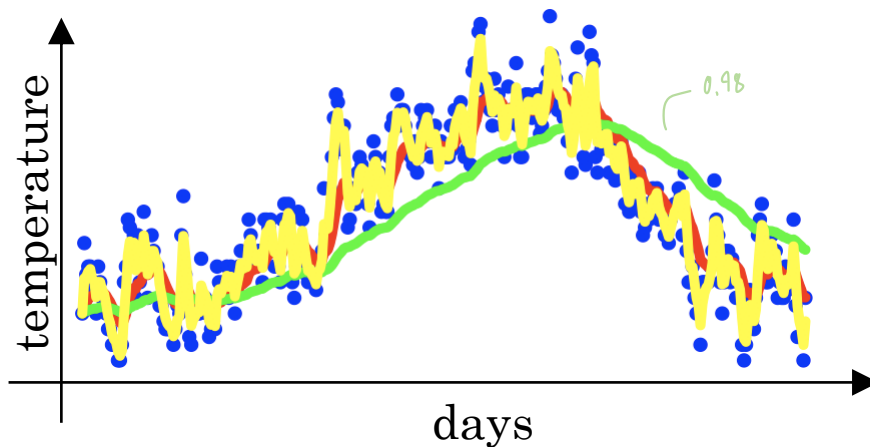
$\beta = 0.5$  :  $\approx 2$  days

$V_t$  is approximately

average over

$\rightarrow \approx \frac{1}{1-\beta}$  days' temperature.

$$\frac{1}{1-0.98} = 50$$



# Momentum

- Idea, only change direction by a little bit each time.
- Keeps a “running average” of the step directions, smoothing out the variation of the individual points.

$$v_t = \beta \cdot v_{t-1} + (1 - \beta) \cdot \nabla J$$

Handwritten annotations for the first equation:

- $\beta$ : MOMENTUM!!
- $v_{t-1}$ : previous one
- $\nabla J$ : gradient
- Curly brace under  $\beta \cdot v_{t-1}$ : average value over time

$$W = W - \alpha \cdot v_t$$

Handwritten annotations for the second equation:

- $\alpha$ : learning rate

- Here,  $\beta$  is referred to as the “momentum”. It is generally given a value  $< 1$

# Momentum

- Idea, only change direction by a little bit each time.
- Keeps a “running average” of the step directions, smoothing out the variation of the individual points.

$$v_t = \beta \cdot v_{t-1} + (1 - \beta) \cdot \nabla J$$

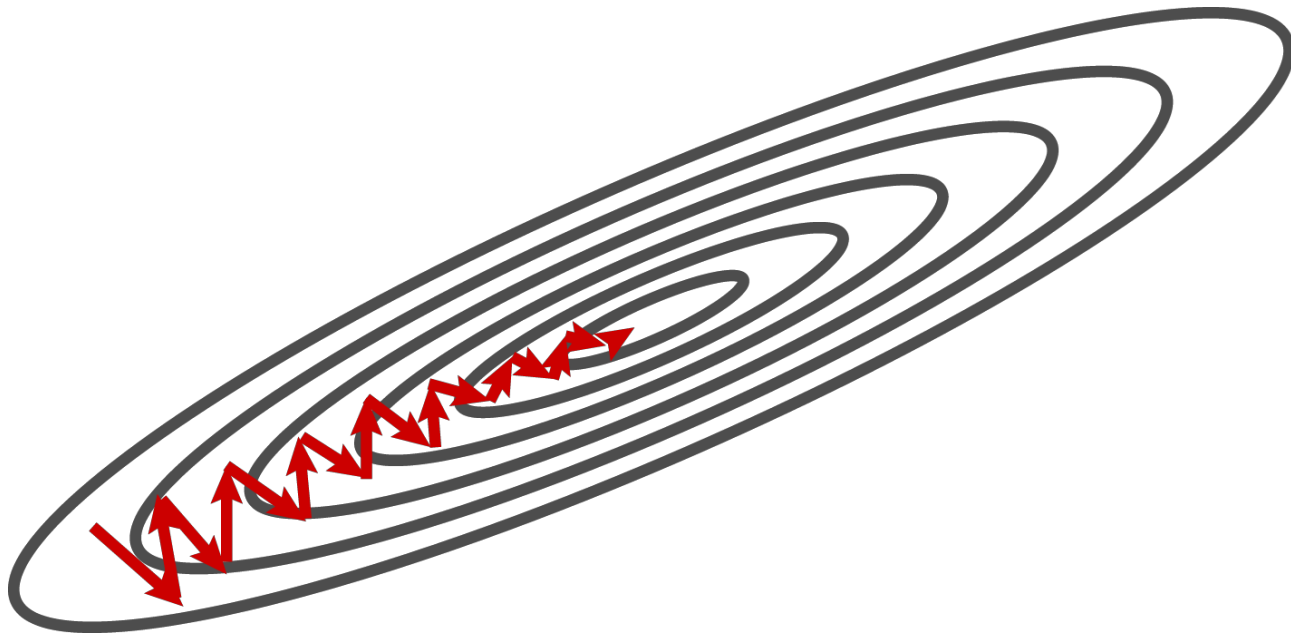
*↗ default: 0.9*

$$W = W - \alpha \cdot v_t$$

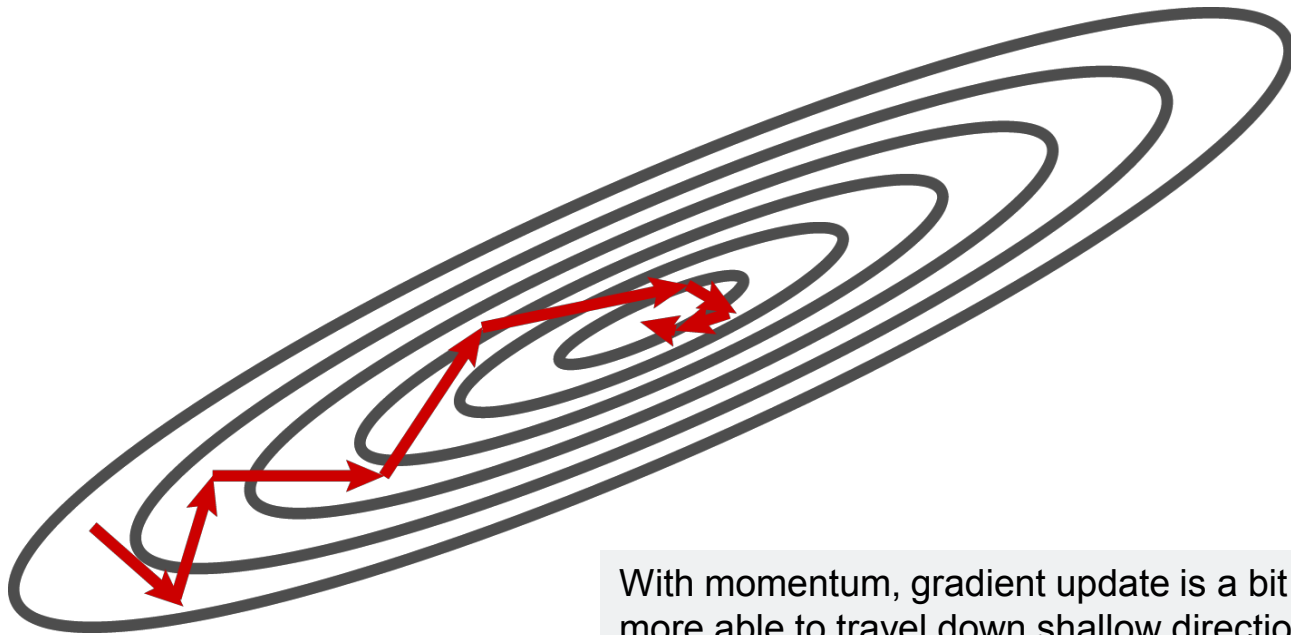
*(learning rate)* often omitted in literatures  
*not general in some paper*

- Here,  $\beta$  is referred to as the “momentum”. It is generally given a value  $< 1$

# Gradient Descent vs Momentum



# Gradient Descent vs Momentum



With momentum, gradient update is a bit smoother and more able to travel down shallow directions.

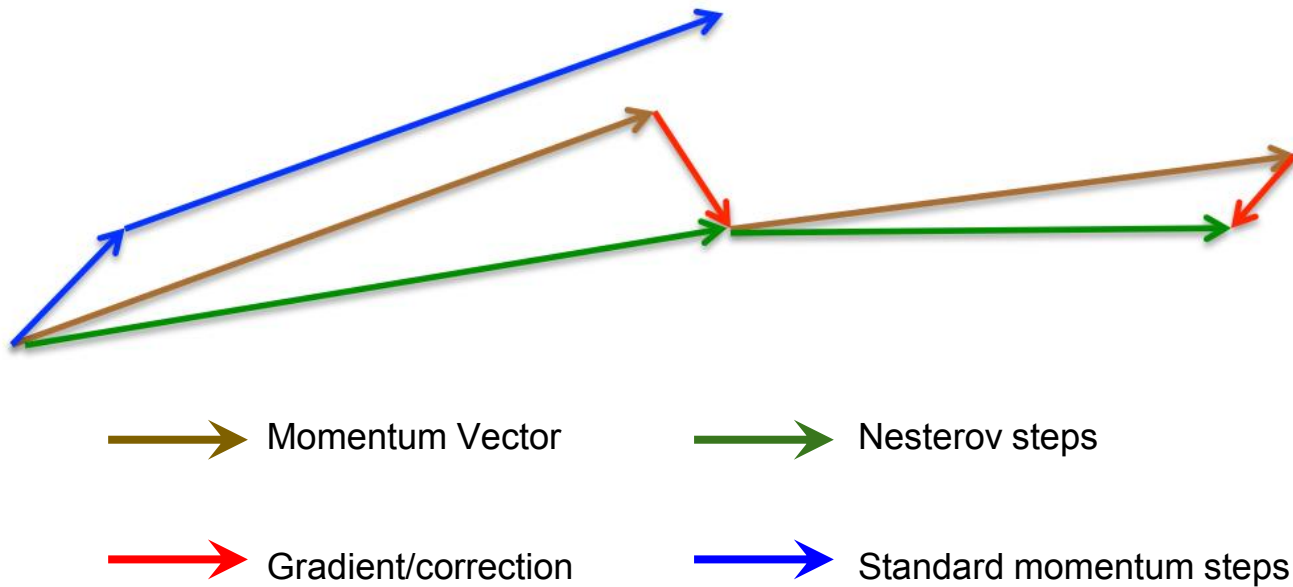
# Nesterov Momentum

- Idea: Control “overshooting” by looking ahead.
- Apply gradient only to the “non-momentum” component.

$$v_t = \beta \cdot v_{t-1} + \alpha \nabla J(W - \beta \cdot v_{t-1})$$
$$W = W - v_t$$



# Nesterov Momentum



# AdaGrad (Adaptive Gradient Optimizer)

- Idea: scale the update for each weight separately.
- Keep running sum of previous updates
- Divide new updates by factor of previous sum

$$W = W - \frac{\alpha}{\sqrt{G_t + \epsilon}} \nabla J \quad G_t = G_{t-1} + (\nabla J)^2$$

*dynamic learning rate*  
*high  $\rightarrow$  low but it'll be too low at further point*

- Instead of using constant learning rate, the learning rate is then divided by the square root of the sum of each component separately.
- As a result of this normalization, those weights that are associated with high gradients will have their learning rates suppressed more.
- This aggressive decay of the rate, however, turns out to be too strong.

# RMSProp (Root Mean Square Propagation)

- Quite similar to AdaGrad.
- Rather than using the sum of previous gradients, decay older gradients more than more recent ones.
- More adaptive to recent updates
- It provides an adaptive learning rate which suppresses learning rates for weights with large frequent gradient updates.

$$S_{dW} = \beta S_{dW_{prev}} + (1 - \beta) \left( \frac{\partial J}{\partial W} \right)^2$$

$$W = W - \alpha \frac{dW}{\sqrt{S_{dW}} + \epsilon}$$

prevent divided by 0

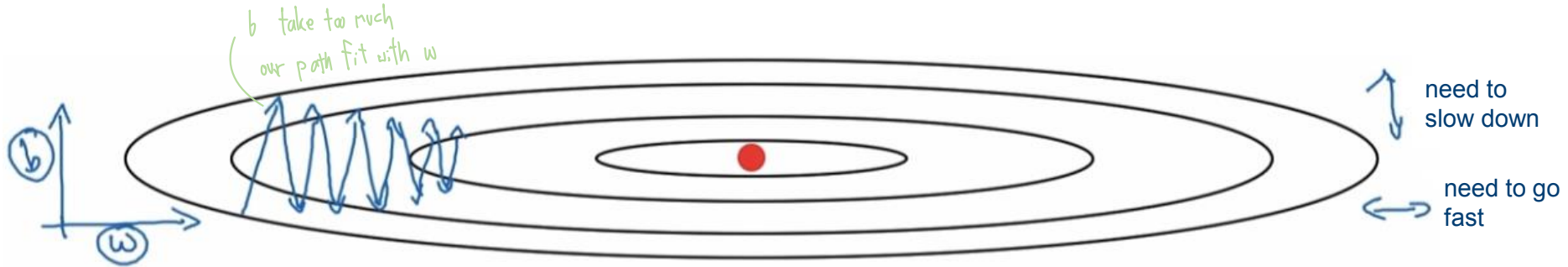
# RMSProp (Case Study)

$$S_{dW} = \beta S_{dW_{prev}} + (1 - \beta)(dW)^2$$

$$S_{db} = \beta S_{db_{prev}} + (1 - \beta)(db)^2$$

$$W = W - \alpha \frac{dW}{\sqrt{S_{dW}} + \epsilon}$$

$$b = b - \alpha \frac{db}{\sqrt{S_{db}} + \epsilon}$$



# Adam (Adaptive Moment Estimation)

- Idea: blending between momentum and RMSprop.
- For iteration  $t$ :

$$V_{dW} = \beta_1 V_{dW_{prev}} + (1 - \beta_1) dW$$

$$\hat{V}_{dW} = \frac{V_{dW}}{1 - \beta_1^t}$$

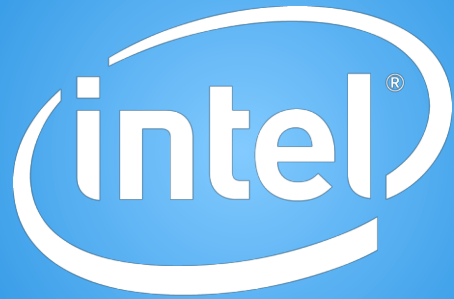
$$S_{dW} = \beta_2 S_{dW_{prev}} + (1 - \beta_2) (dW)^2$$

$$\hat{S}_{dW} = \frac{S_{dW}}{1 - \beta_2^t}$$

$$W := W - \alpha \frac{\hat{V}_{dW}}{\sqrt{\hat{S}_{dW} + \epsilon}}$$

# Which one should I use?!

- RMSProp and Adam seem to be quite popular now.
- Difficult to predict in advance which will be best for a particular problem.
- Still an active area of inquiry.



Software