

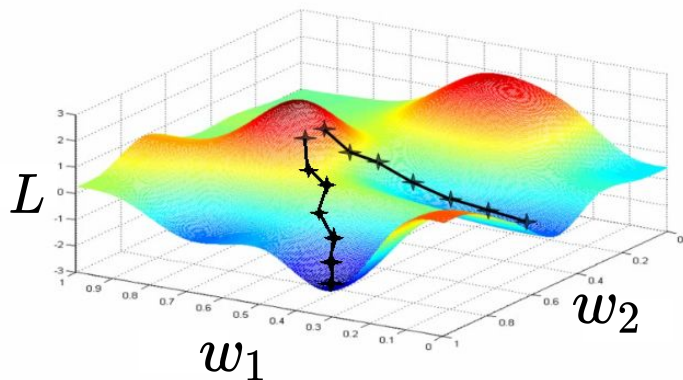
Optimization

Today

- Today is about optimizing a deep model using a training set
- Remember from last time:
 - A model that computes an output using parameters (\mathbf{W}, \mathbf{b})
 - A loss L between the output and the desired outputs (e.g. labels)
 - A method for computing the gradient \mathbf{g} of L w.r.t. the parameters

Gradient descent

- Last time you tried *minibatch GD*
 - Random subsets of training data in each iteration
- But why not the full set, given the loss function we are trying to minimize?
 - MSE loss: $L(\hat{y}, y) = \frac{1}{m} \sum_{i=1}^m \|y_i - \hat{y}_i\|^2$
 - Log loss: $L(\hat{y}, y) = -\frac{1}{m} \sum_{i=1}^m y_i \log \hat{y}_i$
- One loss function \rightarrow one big gradient \rightarrow one step at a time?



<https://infinitygc.com.au/wp-content/uploads/2018/03/88007-02.jpg>

<http://blog.datumbox.com/wp-content/uploads/2013/10/gradient-descent.png>

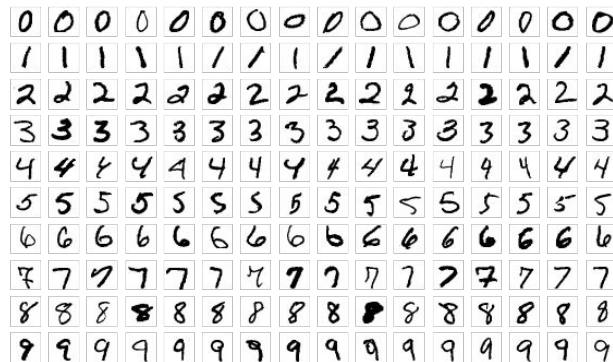
Batch gradient descent

- Infeasible (impossible) for huge training data sets
 - Just think of the forward pass in a small 3-layer network:
- If possible, overall slow
 - Huge efforts for computing only a single, average gradient
- Possible redundancy
 - Training sets often contain many near-duplicate examples
 - Using all examples means a waste of resources
- Determinism
 - May not be preferred
 - We only have an incomplete sample of data

$$\mathbf{h}_1 = f(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1)$$

$$\mathbf{h}_2 = f(\mathbf{W}_2 \mathbf{h}_1 + \mathbf{b}_2)$$

$$\hat{\mathbf{y}} = f(\mathbf{W}_3 \mathbf{h}_2 + \mathbf{b}_3)$$



Additional challenges

- Non-convexity of L
 - Local minima
 - Saddle points
- Steep cliffs
 - Exploding gradients
- The *empirical risk*
 - We may only learn from a training set
 - But we really care about the future performance on a test set
 - We need to use a smooth (differentiable) loss (log loss)
 - But we really care about the (non-differentiable) 0-1 classification loss

Variants of stochastic gradient descent (SGD)

- “Pure” SGD (or online SGD)
 - Select one random training example at a time, compute gradient, update parameters
 - High variance, but many parameter updates
- Minibatch SGD
 - Select a random subset of training examples, compute gradient, update parameters
 - Trade-off between
 - exploration/variance (small batch size),
 - speed (small batch size),
 - stability in gradient estimates (large batch size)

Minibatch SGD

Algorithm 8.1 Stochastic gradient descent (SGD) update

Require: Learning rate schedule $\epsilon_1, \epsilon_2, \dots$

Require: Initial parameter θ

$k \leftarrow 1$

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.

 Compute gradient estimate: $\hat{\mathbf{g}} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

 Apply update: $\theta \leftarrow \theta - \epsilon_k \hat{\mathbf{g}}$

$k \leftarrow k + 1$

end while

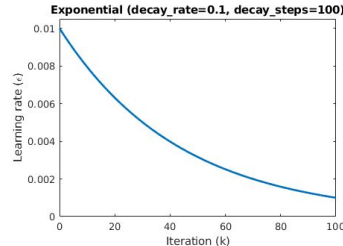
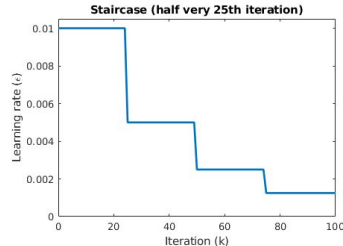
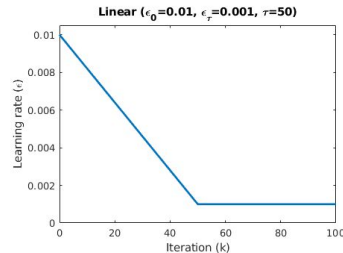
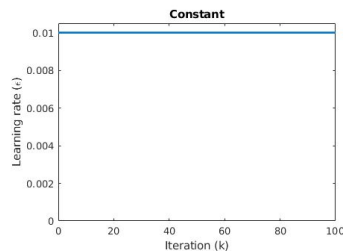
Learning rate schedule

- Common approaches for SGD

- Constant: $\epsilon_k = \epsilon_0$
- Linear decay: $\epsilon_k = (1 - \frac{k}{\tau})\epsilon_0 + \frac{k}{\tau}\epsilon_\tau$
- Piecewise constant or “staircase”
- Exponential: $\epsilon_k = \epsilon_0 \cdot \text{decay_rate}^{k/\text{decay_steps}}$

- Resources

- <https://pytorch.org/docs/stable/optim.html#how-to-adjust-learning-rate>
- A very common approach is the StepLR schedule

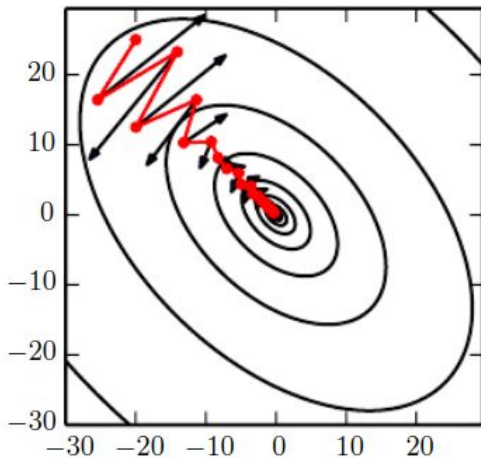
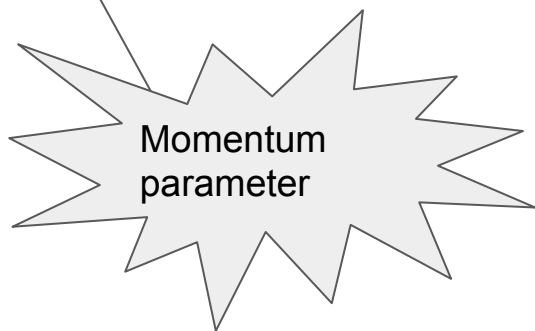


Momentum

- Here's the update rule for (S)GD: $\mathbf{W} \leftarrow \mathbf{W} - \epsilon \mathbf{g}$
- Momentum tries to accelerate learning by including previous gradients
- Define a velocity vector \mathbf{v} , initialized to zero
- The momentum update rule:

$$\mathbf{v} \leftarrow \alpha \mathbf{v} - \epsilon \mathbf{g}$$

$$\mathbf{W} \leftarrow \mathbf{W} + \mathbf{v}$$



Momentum

Algorithm 8.2 Stochastic gradient descent (SGD) with momentum

Require: Learning rate ϵ , momentum parameter α

Require: Initial parameter θ , initial velocity v

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set $\{x^{(1)}, \dots, x^{(m)}\}$ with corresponding targets $y^{(i)}$.

 Compute gradient estimate: $g \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(x^{(i)}; \theta), y^{(i)})$.

 Compute velocity update: $v \leftarrow \alpha v - \epsilon g$.

 Apply update: $\theta \leftarrow \theta + v$.

end while

Nesterov momentum

- Instead of directly updating the velocity using the current gradient, *predict* a new gradient:

$$\tilde{\mathbf{W}} \leftarrow \mathbf{W} + \alpha \mathbf{v}$$

estimate $\hat{\mathbf{y}}$ using $\tilde{\mathbf{W}}$ as model parameters...

$$\mathbf{g} \leftarrow \nabla L(\hat{\mathbf{y}}, \mathbf{y})$$

- After then, it's standard momentum:

$$\mathbf{v} \leftarrow \alpha \mathbf{v} - \epsilon \mathbf{g}$$

$$\mathbf{W} \leftarrow \mathbf{W} + \mathbf{v}$$

Nesterov momentum

Algorithm 8.3 Stochastic gradient descent (SGD) with Nesterov momentum

Require: Learning rate ϵ , momentum parameter α

Require: Initial parameter θ , initial velocity v

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding labels $\mathbf{y}^{(i)}$.

 Apply interim update: $\tilde{\theta} \leftarrow \theta + \alpha v$.

 Compute gradient (at interim point): $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\tilde{\theta}} \sum_i L(f(\mathbf{x}^{(i)}; \tilde{\theta}), \mathbf{y}^{(i)})$.

 Compute velocity update: $v \leftarrow \alpha v - \epsilon \mathbf{g}$.

 Apply update: $\theta \leftarrow \theta + v$.

end while

Adaptive learning rates

- The learning rate is a *hyperparameter*, which is hard to tune
- In a high-dimensional error landscape, it could be helpful to:
 - Decrease the step in steep directions
 - Increase the step in almost flat directions

Adaptive learning rates

- Three variants in the book
 - AdaGrad
 - RMSProp
 - Adam
- And many more...

AdaGrad

- Divide current gradient by magnitude of entire past history of gradients
- Rapid decay in learning rate

Algorithm 8.4 The AdaGrad algorithm

Require: Global learning rate ϵ

Require: Initial parameter θ

Require: Small constant δ , perhaps 10^{-7} , for numerical stability

Initialize gradient accumulation variable $\mathbf{r} = \mathbf{0}$

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.

 Compute gradient: $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$.

 Accumulate squared gradient: $\mathbf{r} \leftarrow \mathbf{r} + \mathbf{g} \odot \mathbf{g}$.

 Compute update: $\Delta \theta \leftarrow -\frac{\epsilon}{\delta + \sqrt{\mathbf{r}}} \odot \mathbf{g}$. (Division and square root applied element-wise)

 Apply update: $\theta \leftarrow \theta + \Delta \theta$.

end while

RMSProp

- Add a momentum-like exponential decay to the gradient history
- Gradients in extreme past have less influence
- Less rapid decay in learning rate

Algorithm 8.5 The RMSProp algorithm

Require: Global learning rate ϵ , decay rate ρ

Require: Initial parameter θ

Require: Small constant δ , usually 10^{-6} , used to stabilize division by small numbers

Initialize accumulation variables $\mathbf{r} = 0$

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.

 Compute gradient: $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$.

 Accumulate squared gradient: $\mathbf{r} \leftarrow \rho \mathbf{r} + (1 - \rho) \mathbf{g} \odot \mathbf{g}$.

 Compute parameter update: $\Delta \theta = -\frac{\epsilon}{\sqrt{\delta + \mathbf{r}}} \odot \mathbf{g}$. ($\frac{1}{\sqrt{\delta + \mathbf{r}}}$ applied element-wise)

 Apply update: $\theta \leftarrow \theta + \Delta \theta$.

end while

Adam

- Call the 1st and 2nd power of the gradient *moments*
- Adam uses a momentum-like decay on both moments
- In addition, it uses *bias correction* to avoid initial instabilities of the moments

Algorithm 8.7 The Adam algorithm

Require: Step size ϵ (Suggested default: 0.001)

Require: Exponential decay rates for moment estimates, ρ_1 and ρ_2 in $[0, 1)$.
(Suggested defaults: 0.9 and 0.999 respectively)

Require: Small constant δ used for numerical stabilization (Suggested default: 10^{-8})

Require: Initial parameters θ

Initialize 1st and 2nd moment variables $\mathbf{s} = \mathbf{0}$, $\mathbf{r} = \mathbf{0}$

Initialize time step $t = 0$

while stopping criterion not met **do**

Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.

Compute gradient: $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

$t \leftarrow t + 1$

Update biased first moment estimate: $\mathbf{s} \leftarrow \rho_1 \mathbf{s} + (1 - \rho_1) \mathbf{g}$

Update biased second moment estimate: $\mathbf{r} \leftarrow \rho_2 \mathbf{r} + (1 - \rho_2) \mathbf{g} \odot \mathbf{g}$

Correct bias in first moment: $\hat{\mathbf{s}} \leftarrow \frac{\mathbf{s}}{1 - \rho_1^t}$

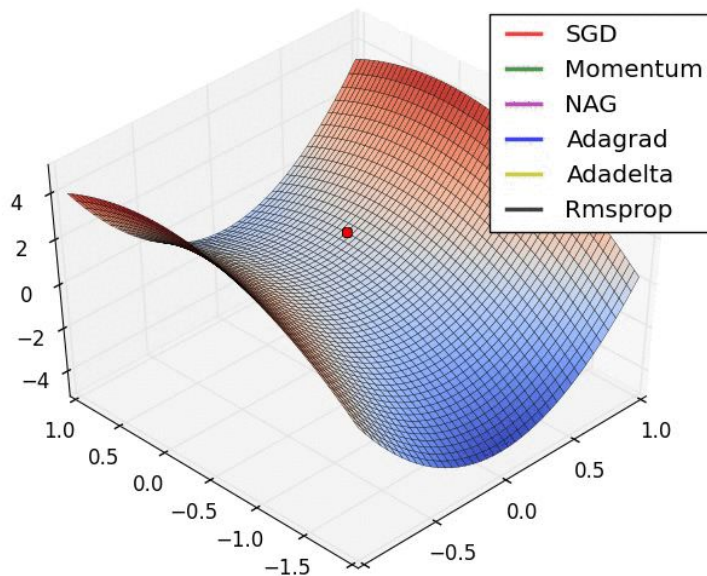
Correct bias in second moment: $\hat{\mathbf{r}} \leftarrow \frac{\mathbf{r}}{1 - \rho_2^t}$

Compute update: $\Delta \theta = -\epsilon \frac{\hat{\mathbf{s}}}{\sqrt{\hat{\mathbf{r}} + \delta}}$ (operations applied element-wise)

Apply update: $\theta \leftarrow \theta + \Delta \theta$

end while

Adaptive learning rates



<http://cs231n.github.io/neural-networks-3/>

Batch normalization

Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift

Sergey Ioffe
Christian Szegedy

Google, 1600 Amphitheatre Pkwy, Mountain View, CA 94043

SIOFFE@GOOGLE.COM
SZEGEDY@GOOGLE.COM

Abstract

Training Deep Neural Networks is complicated by the fact that the distribution of each layer's inputs changes during training, as the parameters of the previous layers change. This slows down the training by requiring lower learning rates and careful parameter initialization, and makes it notoriously hard to train models with saturating nonlinearities. We refer to this phenomenon as *internal covariate shift*, and address the problem by normalizing layer inputs. This greatly improves the training speed, with a

minimize the loss

$$\Theta = \arg \min_{\Theta} \frac{1}{N} \sum_{i=1}^N \ell(x_i, \Theta)$$

where $x_{1...N}$ is the training data set. With SGD, the training proceeds in steps, at each step considering a *mini-batch* $x_{1...m}$ of size m . Using mini-batches of examples, as opposed to one example at a time, is helpful in several ways. First, the gradient of the loss over a mini-batch $\frac{1}{m} \sum_{i=1}^m \frac{\partial \ell(x_i, \Theta)}{\partial \Theta}$ is an estimate of the gradient over the whole data set $\frac{1}{N} \sum_{i=1}^N \frac{\partial \ell(x_i, \Theta)}{\partial \Theta}$. This is an estimate of the gradient of the loss over the whole data set, which is useful for training. Second, the gradient of the loss over a mini-batch is a better estimate of the gradient of the loss over the whole data set than the gradient of the loss over a single example. This is because the gradient of the loss over a single example is a noisy estimate of the gradient of the loss over the whole data set, while the gradient of the loss over a mini-batch is a more accurate estimate.

BatchNorm

- Recall the output (the neuronal activity) at any layer i :

$$\mathbf{h}_i = \phi(\mathbf{W}_i \mathbf{h}_{i-1} + \mathbf{b}_i) = \phi(\mathbf{a}_i)$$

Activation function (e.g. ReLu, sigmoid)

- If we are using a (mini)batch of m vectors, we get a data matrix \mathbf{H} instead:

$$\mathbf{H}_i = \phi(\mathbf{W}_i \mathbf{H}_{i-1} + \mathbf{b}_i)$$

- BatchNorm (or BN) normalizes the i 'th layer using the m samples:

$$\mathbf{H}'_i = \frac{\mathbf{H}_i - \boldsymbol{\mu}_i}{\boldsymbol{\sigma}_i}$$

A column vector with the means of all rows of \mathbf{H}

A column vector with the standard deviations of all rows of \mathbf{H}

BatchNorm $\mathbf{H}'_i = \frac{\mathbf{H}_i - \boldsymbol{\mu}_i}{\sigma_i}$

- The computations of mean/std:

$$\boldsymbol{\mu}_i = \frac{1}{m} \sum_{j=1}^m \mathbf{H}_{i,j} \quad \sigma_i = \sqrt{\delta + \frac{1}{m} \sum_{j=1}^m (\mathbf{H}_{i,j} - \boldsymbol{\mu}_{i,j})^2}$$

are part of the feedforward chain (you can see them as a pseudo-layer)

- This means that backprop will “see” these operations and take them into account when computing gradients
- This in return means that learning can focus on other stuff than trying to bring the activations to a “good” range

BatchNorm $\mathbf{H}'_i = \frac{\mathbf{H}_i - \mu_i}{\sigma_i}$

- All units in the network are now normalized - this stabilizes training
 - Read the section in the book for more in-depth of this
- The final stage of BN scales/shifts the layer responses by two trainable parameters:

$$\mathbf{H}'_i \leftarrow \gamma_i \mathbf{H}'_i + \beta_i$$

- During *evaluation/test*, these learned parameters are fixed
 - For the mean/std, running averages are simply accumulated during training

BatchNorm

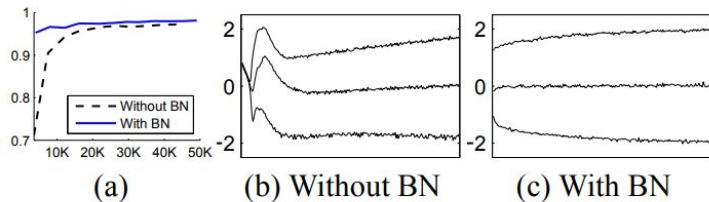


Figure 1: (a) *The test accuracy of the MNIST network trained with and without Batch Normalization, vs. the number of training steps. Batch Normalization helps the network train faster and achieve higher accuracy.* (b, c) *The evolution of input distributions to a typical sigmoid, over the course of training, shown as {15, 50, 85}th percentiles. Batch Normalization makes the distribution more stable and reduces the internal covariate shift.*

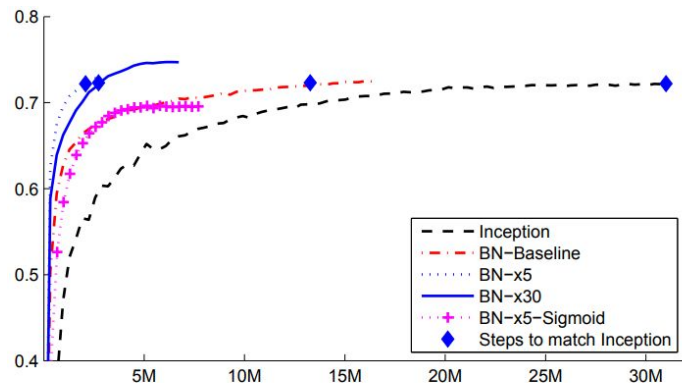


Figure 2: *Single crop validation accuracy of Inception and its batch-normalized variants, vs. the number of training steps.*

Challenge

- Use solution from lecture 3 (MLP) as a starting point
 - Remove the normalization to $[-0.5, 0.5]$ pixel values in the Cifar10 class, i.e. use the full pixel range $[0, 255]$ - remember to still convert to float32
 - Use 1 hidden layer with 100 neurons
 - Now train on the unnormalized pixels without and with BatchNorm (applied after ReLU)
 - When using BatchNorm, use `model.train()/model.eval()` when training/validating (because the layer behaves differently under training and inference)
 - Try a learning rate schedule from `torch.optim.lr_scheduler`
 - Try replacing the SGD optimizer with the Adam optimizer from `torch.optim`
- Bonus:
 - Put back the $[-0.5, 0.5]$ pixel normalization
 - Try to do batch GD instead of minibatch SGD

It's expensive to calculate the gradient, but the gradient has less variance so you can increase the learning rate. Is it enough to compensate?