



國立交通大學  
National Chiao Tung University

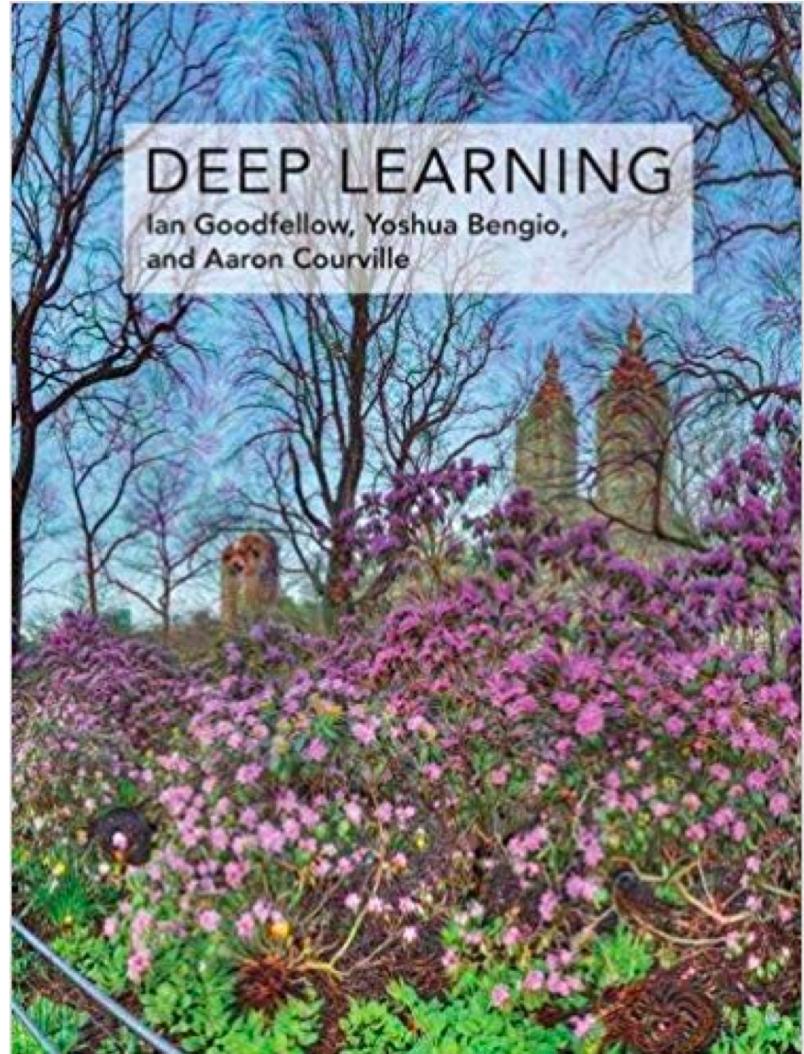
# Deep Learning

## 深度學習

### Fall 2018

*Optimization for  
Training  
(Chapter 8.6-8.7)*

Prof. Chia-Han Lee  
李佳翰 副教授





# Approximate second-order methods

- In this lecture we discuss the application of **second-order methods** to training deep networks.
- For simplicity of exposition, the only objective function we examine is the empirical risk:

$$J(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x}, y \sim \hat{p}_{\text{data}}(\mathbf{x}, y)} [L(f(\mathbf{x}; \boldsymbol{\theta}), y)] = \frac{1}{m} \sum_{i=1}^m L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}), y^{(i)}). \quad (8.25)$$

- The methods we discuss here extend readily, however, to more general objective functions, such as those that include parameter regularization terms.



# Newton's method

- Newton's method is an optimization scheme based on using a **second-order Taylor series expansion to approximate  $J(\theta)$**  near some point  $\theta_0$ , ignoring derivatives of higher order:

$$J(\theta) \approx J(\theta_0) + (\theta - \theta_0)^\top \nabla_{\theta} J(\theta_0) + \frac{1}{2}(\theta - \theta_0)^\top H(\theta - \theta_0), \quad (8.26)$$

where  $H$  is the Hessian of  $J$  with respect to  $\theta$  evaluated at  $\theta_0$ .

- If we then solve for the critical point of this function, we obtain the Newton parameter update rule:

$$\theta^* = \theta_0 - H^{-1} \nabla_{\theta} J(\theta_0). \quad (8.27)$$



# Newton's method

- Thus for a **locally quadratic function** (with positive definite  $\mathbf{H}$ ), by rescaling the gradient by  $\mathbf{H}^{-1}$ , Newton's method **jumps directly to the minimum**.
- If the **objective function is convex but not quadratic** (there are higher-order terms), this update can be **iterated**, yielding the training algorithm associated with Newton's method.
- For surfaces that are **not quadratic**, as long as the **Hessian remains positive definite**, Newton's method can be applied iteratively. This implies a **two-step iterative procedure**. First, update or compute the inverse Hessian. Second, update the parameters according to

$$\boldsymbol{\theta}^* = \boldsymbol{\theta}_0 - \mathbf{H}^{-1} \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}_0). \quad (8.27)$$



# Newton's method

---

**Algorithm 8.8** Newton's method with objective  $J(\boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^m L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}), y^{(i)})$

---

**Require:** Initial parameter  $\boldsymbol{\theta}_0$

**Require:** Training set of  $m$  examples

**while** stopping criterion not met **do**

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

    Compute Hessian:  $\mathbf{H} \leftarrow \frac{1}{m} \nabla_{\boldsymbol{\theta}}^2 \sum_i L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}), \mathbf{y}^{(i)})$

    Compute Hessian inverse:  $\mathbf{H}^{-1}$

    Compute update:  $\Delta \boldsymbol{\theta} = -\mathbf{H}^{-1} \mathbf{g}$

    Apply update:  $\boldsymbol{\theta} = \boldsymbol{\theta} + \Delta \boldsymbol{\theta}$

**end while**

---



# Newton's method

- If the eigenvalues of the Hessian are not all positive, for example, near a saddle point, then Newton's method can actually cause updates to move in the wrong direction.
- This situation can be avoided by **regularizing the Hessian**. Common regularization strategies include adding a constant,  $\alpha$ , along the diagonal of the Hessian. The regularized update becomes

$$\boldsymbol{\theta}^* = \boldsymbol{\theta}_0 - [H(f(\boldsymbol{\theta}_0)) + \alpha \mathbf{I}]^{-1} \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_0). \quad (8.28)$$

- This regularization strategy is used in **approximations to Newton's method**, such as the **Levenberg-Marquardt algorithm**.



# Newton's method

- When there are more extreme directions of curvature, the value of  $\alpha$  would have to be sufficiently large to offset the negative eigenvalues. As  $\alpha$  increases in size, however, the Hessian becomes dominated by the  $\alpha I$  diagonal, and the direction chosen by Newton's method converges to the standard gradient divided by  $\alpha$ .
- When strong negative curvature is present,  $\alpha$  may need to be so large that Newton's method would make smaller steps than gradient descent with a properly chosen learning rate.



# Newton's method

- The application of Newton's method for training large NNs is limited by the **significant computational burden**.
- The number of elements in the Hessian is squared in the number of parameters, so with  $k$  parameters (and for even very small neural networks, the number of parameters  $k$  can be in the millions), **Newton's method would require the inversion of a  $k \times k$  matrix**—with computational complexity of  $O(k^3)$ .
- Since the parameters will change with every update, the inverse Hessian **has to be computed at every training iteration**. As a consequence, only networks with a very small number of parameters can be practically trained via Newton's method.



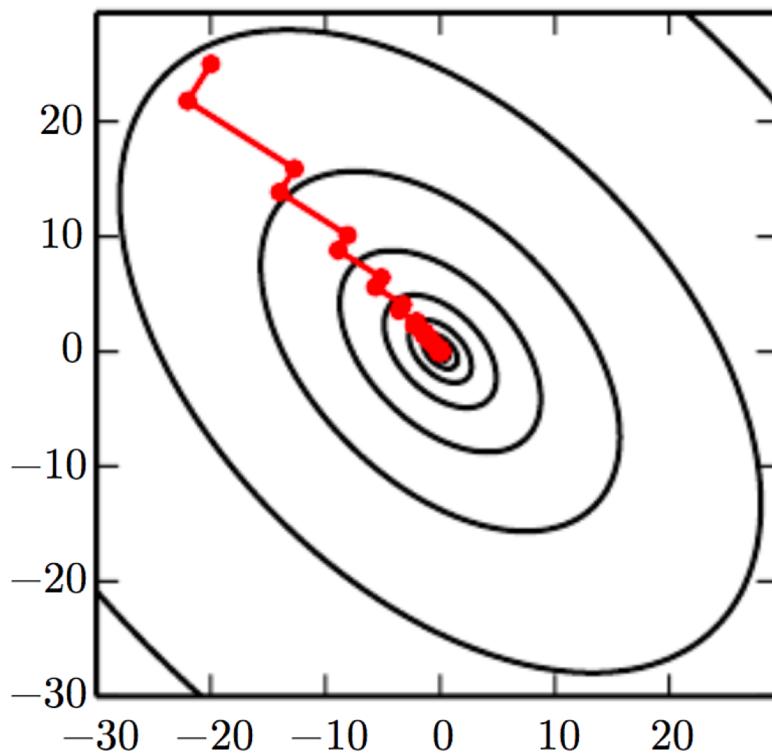
# Conjugate gradients

- Conjugate gradients is a method to efficiently avoid the calculation of the inverse Hessian by iteratively descending conjugate directions.
- The inspiration for this approach follows from the study of the weakness of the method of steepest descent, where line searches (to evaluate  $f(x - \epsilon \nabla_x f(x))$  for several values of  $\epsilon$  and choose the one that results in the smallest objective function value) are applied iteratively in the direction associated with the gradient.
- Figure 8.6 illustrates how the method of steepest descent progresses in a rather ineffective back-and-forth zig-zag pattern. This happens because each line search direction, when given by the gradient, is guaranteed to be orthogonal to the previous line search direction.*DL Fall 2018*



# Conjugate gradients

The method of steepest descent involves jumping to the point of lowest cost along the line defined by the gradient at the initial point on each step. This resolves some of the problems seen with using a fixed learning rate in figure 4.6, but even with the optimal step size the algorithm still makes back-and-forth progress toward the optimum.





# Conjugate gradients

- Let the previous search direction be  $\mathbf{d}_{t-1}$ . At the minimum, where the line search terminates, the directional derivative is zero in direction  $\mathbf{d}_{t-1}$ :  $\nabla_{\theta}J(\boldsymbol{\theta}) \cdot \mathbf{d}_{t-1} = 0$ . Since the gradient at this point defines the current search direction,  $\mathbf{d}_t = \nabla_{\theta}J(\boldsymbol{\theta})$  will have no contribution in the direction  $\mathbf{d}_{t-1}$ . Thus  $\mathbf{d}_t$  is orthogonal to  $\mathbf{d}_{t-1}$ .
- The choice of orthogonal directions of descent do not preserve the minimum along the previous search directions. This gives rise to the zig-zag pattern of progress, where by following the gradient at the end of each line search we undo progress we have already made in the direction of the previous line search.  
**Conjugate gradients** seeks to address this problem.



# Conjugate gradients

- In the method of conjugate gradients, we seek to find a search direction that is **conjugate to the previous line search direction**; that is, it will not undo progress made in that direction.
- At training iteration  $t$ , **the next search direction**  $\mathbf{d}_t$  takes the form:

$$\mathbf{d}_t = \nabla_{\theta} J(\theta) + \beta_t \mathbf{d}_{t-1}, \quad (8.29)$$

where  $\beta_t$  is a coefficient whose magnitude controls how much of the direction,  $\mathbf{d}_{t-1}$ , we should add back to the current search direction.



# Conjugate gradients

- Two directions,  $d_t$  and  $d_{t-1}$ , are defined as conjugate if  $d_t^T H d_{t-1} = 0$ , where  $H$  is the Hessian matrix.
- The straightforward way to impose conjugacy would involve calculation of the eigenvectors of  $H$  to choose  $\beta_t$ , which would not satisfy our goal of developing a method that is more computationally viable than Newton's method for large problems.
- Can we calculate the conjugate directions without resorting to the calculations of the eigenvectors of  $H$  ? Fortunately, the answer to that is yes.



# Conjugate gradients

- Two popular methods for computing  $\beta_t$  are

1. Fletcher-Reeves:

$$\beta_t = \frac{\nabla_{\theta} J(\theta_t)^T \nabla_{\theta} J(\theta_t)}{\nabla_{\theta} J(\theta_{t-1})^T \nabla_{\theta} J(\theta_{t-1})} \quad (8.30)$$

2. Polak-Ribière:

$$\beta_t = \frac{(\nabla_{\theta} J(\theta_t) - \nabla_{\theta} J(\theta_{t-1}))^T \nabla_{\theta} J(\theta_t)}{\nabla_{\theta} J(\theta_{t-1})^T \nabla_{\theta} J(\theta_{t-1})} \quad (8.31)$$

- For a quadratic surface, the conjugate directions ensure that **the gradient along the previous direction does not increase in magnitude**. We therefore **stay at the minimum along the previous directions**. Thus, in a  $k$ -dimensional parameter space, the conjugate gradient method **requires at most  $k$  line searches to achieve the minimum**.



# Conjugate gradients

---

**Algorithm 8.9** The conjugate gradient method

---

**Require:** Initial parameters  $\theta_0$

**Require:** Training set of  $m$  examples

Initialize  $\rho_0 = \mathbf{0}$

Initialize  $g_0 = 0$

Initialize  $t = 1$

**while** stopping criterion not met **do**

    Initialize the gradient  $\mathbf{g}_t = \mathbf{0}$

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

    Compute  $\beta_t = \frac{(\mathbf{g}_t - \mathbf{g}_{t-1})^\top \mathbf{g}_t}{\mathbf{g}_{t-1}^\top \mathbf{g}_{t-1}}$  (Polak-Ribière)

    (Nonlinear conjugate gradient: optionally reset  $\beta_t$  to zero, for example if  $t$  is a multiple of some constant  $k$ , such as  $k = 5$ )

    Compute search direction:  $\rho_t = -\mathbf{g}_t + \beta_t \rho_{t-1}$

    Perform line search to find:  $\epsilon^* = \operatorname{argmin}_{\epsilon} \frac{1}{m} \sum_{i=1}^m L(f(\mathbf{x}^{(i)}; \theta_t + \epsilon \rho_t), \mathbf{y}^{(i)})$

    (On a truly quadratic cost function, analytically solve for  $\epsilon^*$  rather than explicitly searching for it)

    Apply update:  $\theta_{t+1} = \theta_t + \epsilon^* \rho_t$

$t \leftarrow t + 1$

**end while**

---



# Nonlinear conjugate gradients

- Without any assurance that the objective is quadratic, the conjugate directions are no longer assured to remain at the minimum of the objective for previous directions. As a result, the **nonlinear conjugate gradients** algorithm includes occasional resets where the method of conjugate gradients is restarted with line search along the unaltered gradient.
- It is often beneficial to **initialize** the optimization with a few iterations of stochastic gradient descent before commencing nonlinear conjugate gradients.
- While the (nonlinear) conjugate gradients algorithm has traditionally been cast as a batch method, minibatch versions have been used successfully for training NNs.



- The Broyden–Fletcher–Goldfarb–Shanno (BFGS) algorithm attempts to bring some of the advantages of Newton’s method without the computational burden.
- BFGS is similar to the conjugate gradient method. However, BFGS takes a more direct approach to the approximation of Newton’s update.
- Recall that Newton’s update is given by

$$\theta^* = \theta_0 - H^{-1} \nabla_{\theta} J(\theta_0), \quad (8.32)$$

where  $H$  is the Hessian of  $J$  with respect to  $\theta$  evaluated at  $\theta_0$ . The primary computational difficulty in applying Newton’s update is the calculation of the inverse Hessian  $H^{-1}$ .



- The approach adopted by **quasi-Newton methods** (of which the BFGS algorithm is the most prominent) is to **approximate the inverse with a matrix  $M_t$**  that is **iteratively refined by low-rank updates** to become a better approximation of  $H^{-1}$ .
- Once the inverse Hessian approximation  $M_t$  is updated, the direction of descent  $\rho_t$  is determined by  $\rho_t = M_t g_t$ . A **line search is performed in this direction to determine the size of the step**,  $\epsilon^*$ , taken in this direction. The final update to the parameters is given by

$$\theta_{t+1} = \theta_t + \epsilon^* \rho_t. \quad (8.33)$$



- Like the method of conjugate gradients, the BFGS algorithm **iterates line searches with the direction incorporating second-order information**. Unlike conjugate gradients, however, the success of the approach is **not heavily dependent on the line search finding a point very close to the true minimum along the line**.
- Thus, relative to conjugate gradients, BFGS has the advantage that **it can spend less time refining each line search**.
- However, the BFGS algorithm **must store the inverse Hessian matrix**,  $M$ , that requires  $O(n^2)$  memory, making BFGS impractical for most modern deep learning models that typically have millions of parameters.



# Limited memory BFGS (L-BFGS)

- The memory costs of the BFGS algorithm can be significantly decreased by **avoiding storing the complete inverse Hessian approximation  $M$ .**
- The L-BFGS algorithm begins with **the assumption that  $M^{(t-1)}$  is the identity matrix**, rather than storing the approximation from one step to the next.
- If used with exact line searches, the directions defined by L-BFGS are mutually conjugate. However, this procedure remains well behaved when the minimum of the line search is reached only approximately.
- L-BFGS can be generalized to include more information about the Hessian by storing some of the vectors used to update  $M$  at each time step, with costs  $O(n)$ .



# Optimization strategies and meta-algorithms

- Many optimization techniques are not exactly algorithms but rather general templates that can be specialized to yield algorithms, or subroutines that can be incorporated into many different algorithms.



# Batch normalization

- Batch normalization is a method of adaptive reparametrization, motivated by the difficulty of training very deep models.
- The gradient tells how to update each parameter, under the assumption that the other layers do not change. In practice, we update all the layers simultaneously.
- When we make the update, unexpected results can happen because many functions composed together are changed simultaneously, using updates that were computed under the assumption that the other functions remain constant.



# Batch normalization

- As a simple example, suppose we have a deep neural network that has only one unit per layer and does not use an activation function at each hidden layer:  $\hat{y} = x\omega_1\omega_2\omega_3 \dots \omega_l$ . Here,  $\omega_i$  provides the weight used by layer  $i$ . The output of layer  $i$  is  $h_i = h_{i-1}\omega_i$ . The output  $\hat{y}$  is a linear function of the input  $x$  but a nonlinear function of the weights  $\omega_i$ .
- Suppose our cost function has put a gradient of 1 on  $\hat{y}$ , so we wish to decrease  $\hat{y}$  slightly. The back-propagation algorithm can then compute a gradient  $g = \nabla_{\omega}\hat{y}$ .



# Batch normalization

- Consider what happens when we make an update  $\omega \leftarrow \omega - \epsilon g$ . The first-order Taylor series approximation of  $\hat{y}$  predicts that the value of  $\hat{y}$  will decrease by  $\epsilon g^T g$ . If we wanted to decrease  $\hat{y}$  by 0.1, this first-order information available in the gradient suggests we could set the learning rate  $\epsilon$  to  $\frac{0.1}{g^T g}$ .
- The actual update will include second/third-order effects, up to effects of order  $l$ . The new value of  $\hat{y}$  is given by

$$x(w_1 - \epsilon g_1)(w_2 - \epsilon g_2) \dots (w_l - \epsilon g_l). \quad (8.34)$$

An example of one second-order term arising from this update is  $\epsilon g_1 g_2 \prod_{i=3}^l \omega_i$ . This term might be negligible if  $\prod_{i=3}^l \omega_i$  is small, or might be exponentially large if the weights on layers 3 through  $l$  are greater than 1.



# Batch normalization

- This makes it very hard to choose an appropriate learning rate, because the effects of an update to the parameters for one layer depend so strongly on all the other layers.
- Second-order optimization algorithms address this issue by computing an update that takes these second-order interactions into account, but in very deep networks, even higher-order interactions can be significant.
- Even second-order optimization algorithms are expensive and usually require numerous approximations that prevent them from truly accounting for all significant second-order interactions. Building an  $n$ -th order optimization algorithm for  $n > 2$  thus seems hopeless.



# Batch normalization

- Batch normalization provides an elegant way of reparametrizing almost any deep network. The reparametrization significantly reduces the problem of coordinating updates across many layers. Batch normalization can be applied to any input or hidden layer.
- Let  $H$  be a minibatch of activations of the layer to normalize, arranged as a design matrix, with the activations for each example appearing in a row of the matrix. To normalize  $H$ , we replace it with

$$H' = \frac{H - \mu}{\sigma}, \quad (8.35)$$

where  $\mu$  is a vector containing the mean of each unit and  $\sigma$  is a vector containing the standard deviation of each unit.



# Batch normalization

- The arithmetic is based on broadcasting the vector  $\mu$  and the vector  $\sigma$  to be applied to every row of the matrix  $H$ . Within each row, the arithmetic is element-wise, so  $H_{i,j}$  is normalized by subtracting  $\mu_j$  and dividing by  $\sigma_j$ . The rest of the network then operates on  $H'$  in exactly the same way that the original network operated on  $H$ .
- At training time,

$$\mu = \frac{1}{m} \sum_i H_{i,:} \quad (8.36)$$

and

$$\sigma = \sqrt{\delta + \frac{1}{m} \sum_i (\mathbf{H} - \boldsymbol{\mu})_i^2}, \quad (8.37)$$

where  $\delta$  is a small positive value such as  $10^{-8}$  to avoid encountering the undefined gradient of  $\sqrt{z}$  at  $z = 0$ .



# Batch normalization

- Crucially, we back-propagate through these operations for computing the mean and the standard deviation, and for applying them to normalize  $H$ .
- This means that the gradient will never propose an operation that acts simply to increase the standard deviation or mean of  $h_i$ ; the normalization operations remove the effect of such an action and zero out its component in the gradient. This was a major innovation of the batch normalization approach.



# Batch normalization

- The previous approach of adding penalties to the cost function to encourage units to have normalized activation statistics usually resulted in imperfect normalization.
- The previous approach of intervening to renormalize unit statistics after each gradient descent step usually resulted in significant wasted time, as the learning algorithm repeatedly proposed changing the mean and variance, and the normalization step repeatedly undid this change.
- Batch normalization reparametrizes the model to make some units always be standardized by definition, deftly sidestepping both problems.



# Batch normalization

- At test time,  $\mu$  and  $\sigma$  may be replaced by running averages that were collected during training time. This allows the model to be evaluated on a single example, without needing to use definitions of  $\mu$  and  $\sigma$  that depend on an entire minibatch.



# Batch normalization

- Revisiting the  $\hat{y} = x\omega_1\omega_2 \dots \omega_l$  example, we can resolve the difficulties in learning this model by normalizing  $h_{l-1}$ . Suppose that  $x$  is drawn from a unit Gaussian. Then  $h_{l-1}$  will also come from a Gaussian, because the transformation from  $x$  to  $h_l$  is linear. However,  $h_{l-1}$  will no longer have zero mean and unit variance. After applying batch normalization, we obtain the normalized  $\hat{h}_{l-1}$  that restores the zero mean and unit variance properties.
- For almost any update to the lower layers,  $\hat{h}_{l-1}$  will remain a unit Gaussian. The output  $\hat{y}$  may then be learned as a simple linear function  $\hat{y} = \omega_l \hat{h}_{l-1}$ . Learning in this model is now very simple because the parameters at the lower layers do not have an effect in most cases; their output is always renormalized to a unit Gaussian.



# Batch normalization

- In some rare cases, the lower layers can have an effect. Changing one of the lower layer weights to 0 can make the output become degenerate, and changing the sign of one of the lower weights can flip the relationship between  $\hat{h}_{l-1}$  and  $y$ .
- Without normalization, nearly every update would have an extreme effect on the statistics of  $h_{l-1}$ . Batch normalization has thus made this model significantly easier to learn.



# Batch normalization

- In this example, the ease of learning came at the cost of making the lower layers useless. In our linear example, the lower layers no longer have any harmful effect, but they also no longer have any beneficial effect. This is because we have normalized out the first- and second-order statistics, which is all that a linear network can influence. In a deep neural network with nonlinear activation functions, the lower layers can perform nonlinear transformations of the data, so they remain useful.
- Batch normalization acts to standardize only the mean and variance of each unit in order to stabilize learning, but it allows the relationships between units and the nonlinear statistics of a single unit to change.



# Batch normalization

- Because the final layer of the network is able to learn a linear transformation, we may actually **wish to remove all linear relationships between units within a layer.**
- Unfortunately, **eliminating all linear interactions is much more expensive** than standardizing the mean and standard deviation of each individual unit, and so far batch normalization remains the most practical approach.



# Batch normalization

- Normalizing the mean and standard deviation of a unit can **reduce the expressive power of the neural network** containing that unit.
- To maintain the expressive power of the network, it is common to replace the batch of hidden unit activations  $H$  with  $\gamma H' + \beta$  rather than simply the normalized  $H'$ . The variables  $\gamma$  and  $\beta$  are learned parameters that allow the new variable to have any mean and standard deviation.
- At first glance, this may seem useless—why did we set the mean to 0, and then introduce a parameter that allows it to be set back to any arbitrary value  $\beta$ ?



# Batch normalization

- The answer is that the new parametrization can represent the same family of functions of the input as the old parametrization, but **the new parametrization has different learning dynamics.**
- In the old parametrization, the mean of  $H$  was determined by a complicated interaction between the parameters in the layers below  $H$ .
- In the new parametrization, **the mean of  $\gamma H' + \beta$  is determined solely by  $\beta$ .** The new parametrization is much easier to learn with gradient descent.



# Batch normalization

- Most neural network layers take the form of  $\phi(\mathbf{X}\mathbf{W} + \mathbf{b})$ . It is natural to wonder whether we should apply batch normalization to the input  $\mathbf{X}$ , or to the transformed value  $\mathbf{X}\mathbf{W} + \mathbf{b}$ .
- Batch normalization is recommended to be applied to  $\mathbf{X}\mathbf{W} + \mathbf{b}$ . More specifically,  $\mathbf{X}\mathbf{W} + \mathbf{b}$  should be replaced by a normalized version of  $\mathbf{X}\mathbf{W}$  because the bias term becomes redundant with the  $\beta$  parameter.
- In convolutional networks, it is important to apply the same normalizing  $\mu$  and  $\sigma$  at every spatial location within a feature map, so that the statistics of the feature map remain the same regardless of spatial location.



# Coordinate descent

- In some cases, it may be possible to **solve an optimization problem quickly by breaking it into separate pieces**. If we minimize  $f(\mathbf{x})$  with respect to a single variable  $x_i$ , then minimize it with respect to another variable  $x_j$ , and so on, repeatedly cycling through all variables, we are guaranteed to arrive at a (local) minimum. This practice is known as **coordinate descent**, because we **optimize one coordinate at a time**.
- More generally, block coordinate descent refers to **minimizing with respect to a subset of the variables simultaneously**. The term “coordinate descent” is often used to refer to **block coordinate descent** as well as the strictly individual coordinate descent.



# Coordinate descent

- Coordinate descent makes the most sense when the different variables in the optimization problem can be clearly separated into groups that play relatively isolated roles, or when optimization with respect to one group of variables is significantly more efficient than optimization with respect to all of the variables.
- For example, consider the cost function

$$J(\mathbf{H}, \mathbf{W}) = \sum_{i,j} |H_{i,j}| + \sum_{i,j} (\mathbf{x} - \mathbf{w}^\top \mathbf{h})_{i,j}^2. \quad (8.38)$$

This function describes a learning problem called sparse coding, where the goal is to find a weight matrix  $\mathbf{W}$  that can linearly decode a matrix of activation values  $\mathbf{H}$  to reconstruct the training set  $\mathbf{X}$ .



# Coordinate descent

- Most applications of sparse coding also involve weight decay or a constraint on the norms of the columns of  $W$ , to prevent extremely small  $H$  and large  $W$ .
- The function  $J$  is **not convex**. However, we can **divide the inputs to the training algorithm into two sets**: the **dictionary parameters  $W$**  and **the code representations  $H$** . Minimizing the objective function with respect to either one of these sets of variables is a **convex problem**.
- **Block coordinate descent** thus gives us an optimization strategy that allows us to use **efficient convex optimization algorithms**, by **alternating between optimizing  $W$  with  $H$  fixed**, then **optimizing  $H$  with  $W$  fixed**.



# Coordinate descent

- Coordinate descent is **not a very good strategy when the value of one variable strongly influences the optimal value of another variable**, as in the function  $f(x) = (x_1 - x_2)^2 + \alpha(x_1^2 + x_2^2)$  where  $\alpha$  is a positive constant. The first term encourages the two variables to have similar value, while the second term encourages them to be near zero. The solution is to set both to zero.
- **Newton's method can solve the problem in a single step** because it is a positive definite quadratic problem. For small  $\alpha$ , however, **coordinate descent will make very slow progress** because the first term does not allow a single variable to be changed to a value that differs significantly from the current value of the other variable.



# Polyak averaging

- Polyak averaging consists of averaging several points in the trajectory through parameter space visited by an optimization algorithm.
- If  $t$  iterations of gradient descent visit points  $\theta^{(1)}, \dots, \theta^{(t)}$ , then the output of the Polyak averaging algorithm is  $\hat{\theta}^{(t)} = \frac{1}{t} \sum_i \theta^{(i)}$ . On some problem classes, such as gradient descent applied to convex problems, this approach has strong convergence guarantees.
- When applied to neural networks, the idea is that the optimization algorithm may leap back and forth across a valley several times without ever visiting a point near the bottom of the valley. The average of all the locations on either side should be close to the bottom of the valley.



# Polyak averaging

- In nonconvex problems, the path taken by the optimization trajectory can be very complicated and visit many different regions. Including points in parameter space from the distant past that may be separated from the current point by large barriers in the cost function does not seem like a useful behavior.
- As a result, when applying Polyak averaging to nonconvex problems, it is typical to use an exponentially decaying running average:

$$\hat{\theta}^{(t)} = \alpha \hat{\theta}^{(t-1)} + (1 - \alpha) \theta^{(t)}. \quad (8.39)$$



# Supervised pretraining

- Sometimes, directly training a model to solve a specific task can be too ambitious if the model is complex and hard to optimize or if the task is very difficult.
- It is sometimes **more effective** to train a simpler model to solve the task, then make the model more complex.
- It can also be **more effective** to train the model to solve a simpler task, then move on to confront the final task.
- These strategies that involve **training simple models on simple tasks before confronting the challenge of training the desired model to perform the desired task** are collectively known as **pretraining**.



# Supervised pretraining

- Greedy algorithms break a problem into many components, then solve for the optimal version of each component in isolation. Unfortunately, combining the individually optimal components is not guaranteed to yield an optimal complete solution.
- Nonetheless, greedy algorithms can be computationally much cheaper than algorithms that solve for the best joint solution, and the quality of a greedy solution is often acceptable if not optimal.
- Greedy algorithms may also be followed by a fine-tuning stage in which a joint optimization algorithm searches for an optimal solution to the full problem. Initializing the joint optimization algorithm with a greedy solution can greatly speed it up and improve the quality of the solution.

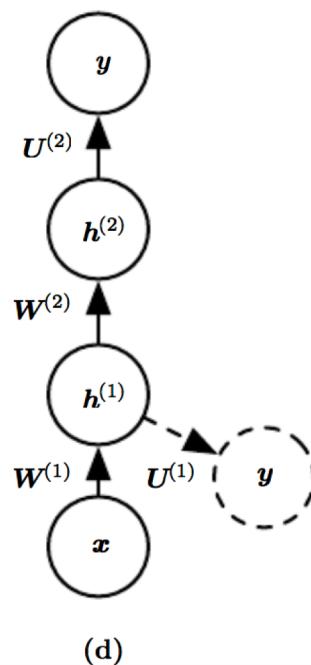
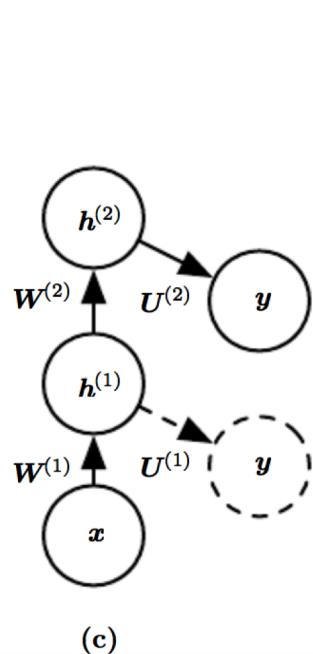
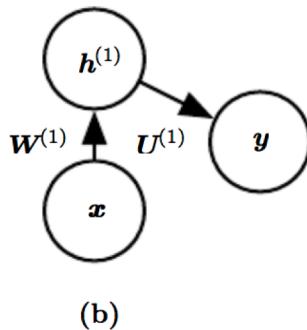
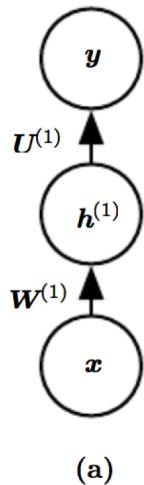


# Supervised pretraining

- Those **pretraining algorithms that break supervised learning problems into other simpler supervised learning problems** are known as **greedy supervised pretraining**.
- In the original version of greedy supervised pretraining, each stage consists of a supervised learning training task **involving only a subset of the layers** in the final neural network. Each added hidden layer is pretrained as part of a shallow supervised MLP, taking as input the output of the previously trained hidden layer.



# Supervised pretraining



(a) We start by training a sufficiently shallow architecture. (b) Another drawing of the same architecture. (c) We keep only the input-to-hidden layer of the original network and discard the hidden-to-output layer. We send the output of the first hidden layer as input to another supervised single hidden layer MLP that is trained with the same objective as the first network was, thus adding a second hidden layer. This can be repeated for as many layers as desired. (d) Another drawing of the result, viewed as a feedforward network.



# Supervised pretraining

- Instead of pretraining one layer at a time, we may **pretrain a deep convolutional network** (eleven weight layers) and then **use the first four and last three layers** from this network **to initialize even deeper networks** (with up to nineteen layers of weights). The **middle layers** of the new, very deep network are **initialized randomly**. The new network is then **jointly trained**.
- Another option is to **use the outputs** of the previously trained MLPs, as well as the raw input, as inputs for each added stage.
- Why would greedy supervised pretraining help? The hypothesis is that **it helps** to provide better guidance to **the intermediate levels** of a deep hierarchy.



# Supervised pretraining

An approach related to supervised pretraining extends the idea to the context of **transfer learning**:

- Pretrain a deep convolutional net with eight layers of weights on a set of tasks (a subset of the 1,000 ImageNet object categories).
- Then initialize a same-size network with the first  $k$  layers of the first net.
- All the layers of the second network (with the upper layers initialized randomly) are then jointly trained to perform a different set of tasks (another subset of the 1,000 ImageNet object categories), with fewer training examples than for the first set of tasks.



# Supervised pretraining

- Another related line of work begins by training a network that has low enough depth and great enough width to be easy to train.
- This network then becomes a teacher for a second network, designated the student. The student network is much deeper and thinner (eleven to nineteen layers) and would be difficult to train with SGD under normal circumstances.
- The training of the student network is made easier by training the student network not only to predict the output for the original task, but also to predict the value of the middle layer of the teacher network.



# Supervised pretraining

- This extra task provides a set of hints about how the hidden layers should be used and can simplify the optimization problem.
- Additional parameters are introduced to regress the middle layer of the five layer teacher network from the middle layer of the deeper student network. Instead of predicting the final classification target, the objective is to predict the middle hidden layer of the teacher network.
- The lower layers of the student networks thus have two objectives: to help the outputs of the student network accomplish their task, as well as to predict the intermediate layer of the teacher network.



# Supervised pretraining

- Although a thin and deep network appears to be more difficult to train than a wide and shallow network, **the thin and deep network may generalize better** and certainly has **lower computational cost** if it is thin enough to have far fewer parameters.
- Without the hints on the hidden layer, the student network performs very poorly in the experiments, on both the training and the test set.
- Hints on middle layers may thus be one of the tools to help train neural networks that otherwise seem difficult to train, but other optimization techniques or changes in the architecture may also solve the problem.



# Designing models to aid optimization

- In practice, it is more important to choose a model family that is easy to optimize than to use a powerful optimization algorithm.
- Most of the advances in neural network learning over the past thirty years have been obtained by changing the model family rather than changing the optimization procedure.
- Stochastic gradient descent with momentum, which was used to train neural networks in the 1980s, remains in use in modern state-of-the-art neural network applications.



# Designing models to aid optimization

- Specifically, modern neural networks reflect a design choice to **use linear transformations between layers and activation functions that are differentiable almost everywhere, with significant slope in large portions of their domain.**
- In particular, model innovations like the LSTM, rectified linear units and maxout units have all **moved toward using more linear functions** than previous models like deep networks based on sigmoidal units.
- These models have nice properties that **make optimization easier**. The gradient flows through many layers provided that the Jacobian of the linear transformation has reasonable singular values.



# Designing models to aid optimization

- Moreover, linear functions consistently increase in a single direction, so even if the model's output is very far from correct, it is clear simply from computing the gradient which direction its output should move to reduce the loss function.
- In other words, modern neural nets have been designed so that their **local gradient information corresponds reasonably well to moving toward a distant solution**.
- Other model design strategies can help to make optimization easier. For example, **linear paths or skip connections between layers** reduce the length of the shortest path from the lower layer's parameters to the output, and thus mitigate the vanishing gradient problem.



# Designing models to aid optimization

- A related idea to skip connections is adding extra copies of the output that are attached to the intermediate hidden layers of the network, as in GoogLeNet and deeply supervised nets.
- These “auxiliary heads” are trained to perform the same task as the primary output at the top of the network to ensure that the lower layers receive a large gradient. When training is complete, the auxiliary heads may be discarded.
- This is an alternative to the pretraining strategies. In this way, one can train jointly all the layers in a single phase but change the architecture, so that intermediate layers (especially the lower ones) can get some hints about what they should do, via a shorter path.



# Continuation methods and curriculum learning

- Many of the challenges in optimization arise from the global structure of the cost function and cannot be resolved merely by making better estimates of local update directions.
- The predominant strategy for overcoming this problem is to attempt to initialize the parameters in a region connected to the solution by a short path through parameter space that local descent can discover.
- Continuation methods are a family of strategies that can make optimization easier by choosing initial points to ensure that local optimization spends most of its time in well-behaved regions of space.



# Continuation methods and curriculum learning

- The idea behind continuation methods is to **construct a series of objective functions over the same parameters.**
- To minimize a cost function  $J(\theta)$ , we **construct new cost functions**  $\{J^{(0)}, \dots, J^{(n)}\}$ . These cost functions are designed to be **increasingly difficult**, with  $J^{(0)}$  being fairly easy to minimize, and  $J^{(n)}$ , the most difficult, being  $J(\theta)$ , the true cost function motivating the entire process. When we say that  $J^{(i)}$  is easier than  $J^{(i+1)}$ , we mean that it is well behaved over more of  $\theta$  space.
- The series of cost functions are designed so that **a solution to one is a good initial point of the next**. We thus begin by solving an easy problem, then refine the solution to solve incrementally harder problems until we arrive at a solution to the true underlying problem. *DL Fall 2018*



# Continuation methods and curriculum learning

- Continuation methods traditionally were mostly designed with the goal of overcoming the challenge of local minima. They were designed to reach a global minimum despite the presence of many local minima.
- These continuation methods would construct easier cost functions by “blurring” the original cost function. This blurring operation can be done by approximating

$$J^{(i)}(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{\theta}' \sim \mathcal{N}(\boldsymbol{\theta}'; \boldsymbol{\theta}, \sigma^{(i)2})} J(\boldsymbol{\theta}') \quad (8.40)$$

via sampling.



# Continuation methods and curriculum learning

- The intuition is that some nonconvex functions become approximately convex when blurred. This blurring preserves enough information about the location of a global minimum that we can find the global minimum by solving progressively less-blurred versions of the problem.

This approach can break down in three different ways.

- First, it might successfully define a series of cost functions where the first is convex and the optimum tracks from one function to the next, arriving at the global minimum, but it might require so many incremental cost functions that the cost of the entire procedure remains high. NP-hard optimization problems remain NP-hard, even when continuation methods are applicable.



# Continuation methods and curriculum learning

The other two ways continuation methods fail both correspond to the method not being applicable.

- First, the function might not become convex, no matter how much it is blurred. Consider, for example, the function  $J(\theta) = -\theta^T \theta$ .
- Second, the function may become convex as a result of blurring, but the minimum of this blurred function may track to a local rather than a global minimum of the original cost function.



# Continuation methods and curriculum learning

- Though continuation methods were mostly originally designed to deal with the problem of local minima, local minima are no longer believed to be the primary problem for neural network optimization. Fortunately, continuation methods can still help.
- The easier objective functions introduced by the continuation method can eliminate flat regions, decrease variance in gradient estimates, improve conditioning of the Hessian matrix, or do anything else that will either make local updates easier to compute or improve the correspondence between local update directions and progress toward a global solution.



# Continuation methods and curriculum learning

- An approach called **curriculum learning**, or shaping, can be interpreted as a continuation method.
- Curriculum learning is based on the idea of planning a learning process to begin by learning simple concepts and progress to learning more complex concepts that depend on these simpler concepts.
- Curriculum learning is justified as a continuation method, where earlier  $J^{(i)}$  are made easier by increasing the influence of simpler examples (either by assigning their contributions to the cost function larger coefficients, or by sampling them more frequently)
- Curriculum learning was also verified as being consistent with the way in which humans teach.



# Continuation methods and curriculum learning

- Another important contribution to research on curriculum learning arose in the context of **training recurrent neural networks to capture long-term dependencies**.
- Much better results were obtained with a **stochastic curriculum**, in which a random mix of easy and difficult examples is always presented to the learner, but where the average proportion of the more difficult examples (here, those with longer-term dependencies) is gradually **increased**. With a deterministic curriculum, no improvement over the baseline (ordinary training from the full training set) was observed.