



# Regularization for Deep Learning

Jie Tang

November 5, 2019

# Overview

- 1 Background: Overfitting
- 2 Parameter Norm Penalties
- 3 Dataset Augmentation
- 4 Noise Robustness
- 5 Early Stopping
- 6 Parameter Sharing
- 7 Ensemble methods
- 8 Dropout
- 9 Other regularization methods

# Outline

- 1 Background: Overfitting
- 2 Parameter Norm Penalties
- 3 Dataset Augmentation
- 4 Noise Robustness
- 5 Early Stopping
- 6 Parameter Sharing
- 7 Ensemble methods
- 8 Dropout
- 9 Other regularization methods

# Background: Overfitting

- In supervised learning settings, overfitting is usually a potential problem.
- Overfitting occurs when a model is excessively complex, such as having too many free parameters.

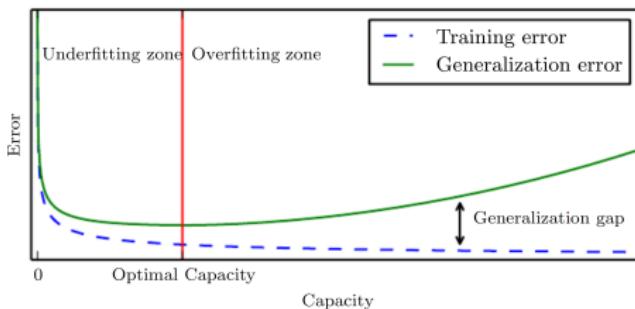


Figure 1: Typical relationship between model capacity and error.  
More complex models tend to overfit.

# Example: Polynomial Regression

- The goal of polynomial regression is to find the best fit polynomial to the observed data points.

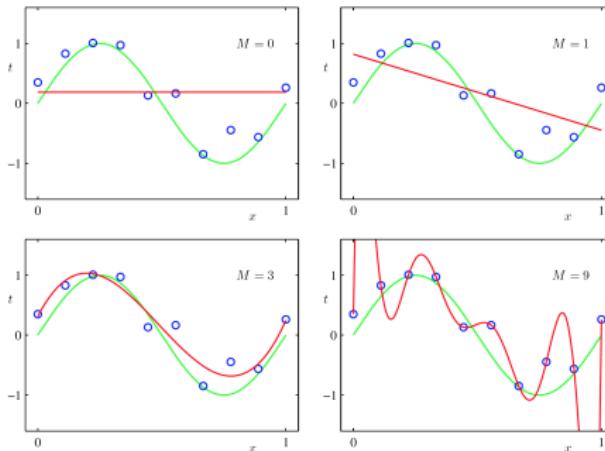


Figure 2: Plots of polynomials having various orders  $M$  (red curves) fitted to the data points (blue circles). A high order polynomial (e.g.  $M = 9$ ) can perfectly fit the observed data, but gives a very poor representation of the function used to generate the data (green curves).

# Regularization for Deep Learning

- Regularization can be motivated as a technique to improve the **generalization** of a learned model.
  - Many regularization strategies are explicitly designed to reduce the test error, possibly at the expense of increased training error.
- Ideas behind the regularization.
  - Encode specific kinds of prior knowledge.
  - Express preference for simpler model.
  - Make an underdetermined problem determined
- In the context of deep learning, most regularization strategies are based on regularizing estimators.
  - Regularization of an estimator works by trading increased bias for reduced variance.
  - In practical scenarios, we almost always do find that the best fitting model (in the sense of minimizing generalization error) is a large model that has been regularized appropriately.

# Outline

- 1 Background: Overfitting
- 2 Parameter Norm Penalties
- 3 Dataset Augmentation
- 4 Noise Robustness
- 5 Early Stopping
- 6 Parameter Sharing
- 7 Ensemble methods
- 8 Dropout
- 9 Other regularization methods

# Regularized Objective Function

- The **regularized objective function**  $\tilde{J}$  is introduced by adding a parameter norm penalty  $\Omega$  to the original objective function  $J$ .

$$\tilde{J}(\theta) = J(\theta) + \alpha\Omega(\theta) \quad (1)$$

- $\theta$ : the set of parameters
- $\alpha \in [0, \infty)$ : a hyperparameter that weights the contribution of the norm penalty term  $\Omega$  relative to the standard objective function  $J$
- When the training algorithm minimizes the regularized objective function  $\tilde{J}$ , it will decrease both the original objective  $J$  and some measure on the parameters  $\theta$  (or some subset of the parameters).
  - Different choices for the parameter norm  $\Omega$  can result in different solutions being preferred.

# $L^2$ Regularization

- term is formulated as  $\Omega(\theta) = \frac{1}{2} \|\theta\|_2^2$ .<sup>1</sup>
- The behavior of  $L^2$  regularization:
  - For simplicity, we assume no bias parameters, so  $\theta$  is just  $w$ .  
The regularized objective function is

$$\tilde{J}(w) = \frac{\alpha}{2} \|w\|_2^2 + J(w) \quad (2)$$

- Take a single gradient step to update the weights:

$$\nabla_w \tilde{J}(w) = \alpha w + \nabla_w J(w) \quad (3)$$

$$w \leftarrow w - \epsilon (\alpha w + \nabla_w J(w)) \quad (4)$$

$$\Leftrightarrow w \leftarrow (1 - \epsilon \alpha) w - \epsilon \nabla_w J(w) \quad (5)$$

- $L^2$  regularization: **shrink** the weight by a constant factor on each step.

---

<sup>1</sup> $L^2$  regularization is commonly known as *weight decay*. It is also known as *ridge regression* or *Tikhonov regularization*.

## Taylor series approximation–review

- Taylor expansion: if  $f(x)$  has second-order derivatives, then we have

$$f(x) \approx f(x_0) + (x - x_0)f'(x) + \frac{1}{2}(x - x_0)^2 f''(x) \quad (6)$$

- Gradient:

$$g = \nabla_{\theta} J(\theta) \quad (7)$$

- Hessian matrix:

$$H_{i,j} = \frac{\partial}{\partial \theta_j} \frac{\partial J(\theta)}{\partial \theta_i} \quad (8)$$

- $H$  is symmetric, so there exists

$$Q = [v_1, \dots, v_n] \quad (9)$$

$$H = Q \Lambda Q^{\top} \quad (10)$$

- A second derivative in direction  $d$  is given by:

$$d^{\top} H d \quad (11) \quad 10 / 66$$

# $L^2$ Regularization

- What happens over the entire training step?
  - Make a quadratic approximation to  $J(\mathbf{w})$  at the optimal weights  $\mathbf{w}^* = \arg \min_{\mathbf{w}} J(\mathbf{w})$ .<sup>2</sup>

$$\hat{J}(\mathbf{w}) = J(\mathbf{w}^*) + \frac{1}{2}(\mathbf{w} - \mathbf{w}^*)^T \mathbf{H}(\mathbf{w} - \mathbf{w}^*) \quad (12)$$

where  $\mathbf{H}$  is the Hessian matrix of  $J$  with respect to  $\mathbf{w}$  evaluated at  $\mathbf{w}^*$ .  $\mathbf{H}$  is positive semidefinite because  $\mathbf{w}^*$  is a minimum of  $J$ .<sup>3</sup>

- The minimum of  $\hat{J}$  occurs where its gradient is equal to  $\mathbf{0}$ .

$$\nabla_{\mathbf{w}} \hat{J}(\mathbf{w}) = \mathbf{H}(\mathbf{w} - \mathbf{w}^*) = \mathbf{0} \quad (13)$$

---

<sup>2</sup>If the objective function is truly quadratic, then the approximation is perfect.

<sup>3</sup>There is no first-order term in this quadratic approximation, because  $\mathbf{w}^*$  is defined to be a minimum where the gradient vanishes.

## $L^2$ Regularization

- Adding the  $L^2$  regularization gradient, the regularized version of  $\hat{J}$  reaches the minimum at  $\tilde{\mathbf{w}}$ :

$$\alpha \tilde{\mathbf{w}} + \mathbf{H}(\tilde{\mathbf{w}} - \mathbf{w}^*) = 0 \quad (14)$$

$$\Leftrightarrow \tilde{\mathbf{w}} = (\mathbf{H} + \alpha \mathbf{I})^{-1} \mathbf{H} \mathbf{w}^* \quad (15)$$

- As  $\alpha$  approaches 0, the regularized solution  $\tilde{\mathbf{w}}$  approaches  $\mathbf{w}^*$ .
- As  $\alpha$  grows, decompose  $\mathbf{H}$  into a diagonal matrix  $\Lambda$  and an orthonormal basis of eigenvectors  $\mathbf{Q}$  such that  $\mathbf{H} = \mathbf{Q} \Lambda \mathbf{Q}^T$ .<sup>4</sup>

$$\tilde{\mathbf{w}} = (\mathbf{Q} \Lambda \mathbf{Q}^T + \alpha \mathbf{I})^{-1} \mathbf{Q} \Lambda \mathbf{Q}^T \mathbf{w}^* \quad (16)$$

$$= (\mathbf{Q}(\Lambda + \alpha \mathbf{I}) \mathbf{Q}^T)^{-1} \mathbf{Q} \Lambda \mathbf{Q}^T \mathbf{w}^* \quad (17)$$

$$= \mathbf{Q}(\Lambda + \alpha \mathbf{I})^{-1} \Lambda \mathbf{Q}^T \mathbf{w}^* \quad (18)$$

---

<sup>4</sup> $\mathbf{H}$  is real and symmetric.

# $L^2$ Regularization

$$\tilde{\mathbf{w}} = \mathbf{Q}(\Lambda + \alpha \mathbf{I})^{-1} \Lambda \mathbf{Q}^T \mathbf{w}^*$$

- The effect  $L^2$  Regularization is to rescale  $\mathbf{w}^*$  along the axes defined by the eigenvectors of  $\mathbf{H}$ .
- The component of  $\mathbf{w}^*$  that is aligned with the  $i$ -th eigenvector of  $\mathbf{H}$  is rescaled by a factor  $\frac{\lambda_i}{\lambda_i + \alpha}$ .
- **Bayesian prior interpretation:**  $L^2$  regularization is equivalent to MAP Bayesian inference with a **Gaussian prior** on the weights (as we discussed before).

# $L^2$ Regularization

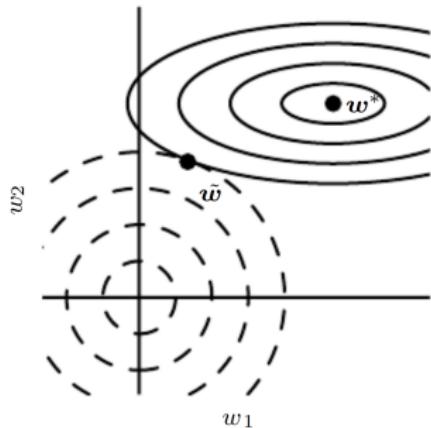


Figure 3: An illustration of the effect of  $L^2$  regularization. The solid ellipses represent contours of equal value of the original objective  $J$ . The dotted circles represent contours of equal value of the  $L^2$  regularizer.

- At the point  $\tilde{w}$ , these competing objectives reach an equilibrium.
- In the first dimension, the objective function does not increase much when moving horizontally away from  $w^*$ . Thus, the regularizer has a strong effect on this axis.
- In the second dimension, the objective function is very sensitive to movements away from  $w^*$ . As a result, the regularizer affects the position of  $w_2$  relatively little.

## $L^2$ Regularization on Linear Regression

- How does these effects related to machine learning in particular?
- Take linear regression for example, the objective function is the sum of squared errors:

$$J(\mathbf{w}; \mathbf{X}, \mathbf{y}) = (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) \quad (19)$$

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} J(\mathbf{w}; \mathbf{X}, \mathbf{y}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (20)$$

- Add  $L^2$  regularization on the objective function:

$$\tilde{J}(\mathbf{w}; \mathbf{X}, \mathbf{y}) = (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) + \frac{1}{2} \alpha \mathbf{w}^T \mathbf{w} \quad (21)$$

$$\tilde{\mathbf{w}}^* = \arg \min_{\mathbf{w}} \tilde{J}(\mathbf{w}; \mathbf{X}, \mathbf{y}) = (\mathbf{X}^T \mathbf{X} + \alpha \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \quad (22)$$

## $L^2$ Regularization on Linear Regression

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} J(\mathbf{w}; \mathbf{X}, \mathbf{y}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\tilde{\mathbf{w}}^* = \arg \min_{\mathbf{w}} \tilde{J}(\mathbf{w}; \mathbf{X}, \mathbf{y}) = (\mathbf{X}^T \mathbf{X} + \alpha \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

- Using  $L^2$  Regularization replaces  $\mathbf{X}^T \mathbf{X}$  with  $(\mathbf{X}^T \mathbf{X} + \alpha \mathbf{I})$ .
- The original matrix is proportional to the covariance matrix with the diagonal entries corresponding to the variance of each input feature.
- The new matrix is the same as the origin one with the addition of  $\alpha$  to the diagonal.
- $L^2$  Regularization causes the learning algorithm to “perceive” the input  $\mathbf{X}$  as having **higher variance**, which makes it shrink the weights on features whose covariance with the output target is low compared to this added variance.

# $L^1$ Regularization

- $L^1$  regularization term is formulated as  $\Omega(\theta) = \|\theta\|_1$ .
- Again, we assume no bias parameter ( $\theta = w$ ).
- The regularized objective function and its gradient:

$$\tilde{J}(w) = \alpha \|w\|_1 + J(w) \quad (23)$$

$$\nabla_w \tilde{J}(w) = \alpha \text{sign}(w) + \nabla_w J(w) \quad (24)$$

- The regularization contribution to the gradient is a **constant** factor with the sign of  $w$ .

# $L^1$ Regularization

- Make a quadratic approximation to the objective function via its Taylor series:

$$\nabla_{\mathbf{w}} \hat{J}(\mathbf{w}) = \mathbf{H}(\mathbf{w} - \mathbf{w}^*) \quad (25)$$

- The  $L^1$  penalty does not admit clean algebraic expressions in the case of a fully general Hessian, we will further **assume** that the Hessian is **diagonal**  $\mathbf{H} = \text{diag}(H_{1,1}, \dots, H_{n,n})$ , where each  $H_{i,i} > 0$ .<sup>5</sup>
- The quadratic approximation of the  $L^1$  regularized objective function decomposes into a sum over the parameters:

$$\tilde{J}(\mathbf{w}) = J(\mathbf{w}^*) + \sum_i \frac{1}{2} H_{i,i} (w_i - w_i^*)^2 + \alpha |w_i| \quad (26)$$

---

<sup>5</sup>This assumption holds if the data has been preprocessed to remove all correlation between the input features, which may be accomplished using PCA.

# $L^1$ Regularization

- The problem of minimizing the approximate objective function has an analytical solution for each  $i$ :

$$w_i = \text{sign}(w_i^*) \max\{|w_i^*| - \frac{\alpha}{H_{i,i}}, 0\} \quad (27)$$

- Consider the situation where  $w_i^* > 0$  for all  $i$ . There are two possible outcomes:

- The case where  $w_i^* \leq \frac{\alpha}{H_{i,i}}$ . Here the optimal value of  $w_i$  under the regularized objective is simply  $w_i = 0$ .
- The case where  $w_i^* > \frac{\alpha}{H_{i,i}}$ . In this case, the regularization just shifts the optimal value by a distance equal to  $\frac{\alpha}{H_{i,i}}$ .

- A similar process happens when  $w_i^* < 0$ , but with the  $L^1$  penalty making  $w_i$  less negative by  $\frac{\alpha}{H_{i,i}}$ , or 0.
- $L^1$  regularization results in a solution that is more **sparse**. It has been used extensively as a **feature selection** mechanism.
- Bayesian prior interpretation:**  $L^1$  regularization can be explained as a Bayesian inference with an Laplace distribution over  $\mathbf{w}$ :

$$\log p(\mathbf{w}) = \sum_i \log \text{Laplace}(w_i; 0, \frac{1}{\alpha}) = -\alpha \|\mathbf{w}\|_1 + n \log \alpha - n \log 2$$

## Norm Penalties as Constrained Optimization

- Consider a constrained optimization problem:

$$\begin{aligned} & \min J(\theta) \\ \text{s.t. } & \Omega(\theta) - k \leq 0 \end{aligned}$$

we can solve the constrained problem by constructing a generalized Lagrange function:

$$\mathcal{L}(\theta, \alpha) = J(\theta) + \alpha(\Omega(\theta) - k) \quad (28)$$

$$\theta^* = \arg \min_{\theta} \max_{\alpha, \alpha \geq 0} \mathcal{L}(\theta, \alpha) \quad (29)$$

- If the value of  $\alpha$  is fixed, we can get the solution:

$$\theta^* = \arg \min_{\theta} \mathcal{L}(\theta, \alpha^*) = \arg \min_{\theta} J(\theta) + \alpha^* \Omega(\theta) \quad (30)$$

- This is exactly the same as the regularized objective function.

# Norm Penalties as Constrained Optimization

- We can thus think of a parameter norm penalty as imposing a constraint on the weights.
  - The parameter  $\alpha$  encourages  $\Omega(\theta)$  to **shrink**: larger  $\alpha$  will result in a smaller constraint region; smaller  $\alpha$  will result in a larger one.
- Sometimes we may prefer explicit constraints rather than penalties.
  - If we have an idea of what value of  $k$  is appropriate and do not want to spend time searching for the value of  $\alpha$  that corresponds to this  $k$ .
  - Penalties can cause non-convex optimization procedures to get stuck in **local minimal**.
  - Explicit constraints can be useful because they impose some stability on the optimization procedure.
- Hinton et al. (2012c) recommend using constraints combined with a high learning rate to allow rapid exploration of parameter space while maintaining some stability.

# Regularization and Under-Constrained Problems

- In some cases, regularization is necessary for machine learning problems to be properly defined. Many linear models in machine learning depend on inverting the matrix  $\mathbf{X}^T \mathbf{X}$ . This is not possible whenever  $\mathbf{X}^T \mathbf{X}$  is **singular**. In this case, many forms of regularization correspond to inverting  $\mathbf{X}^T \mathbf{X} + \alpha \mathbf{I}$  instead.
- Most forms of regularization are able to guarantee the convergence of iterative methods applied to **underdetermined problems**. For example, weight decay will cause gradient descent to quit increasing the magnitude of the weights when the slope of the likelihood is equal to the weight decay coefficient.

# Outline

- 1 Background: Overfitting
- 2 Parameter Norm Penalties
- 3 Dataset Augmentation**
- 4 Noise Robustness
- 5 Early Stopping
- 6 Parameter Sharing
- 7 Ensemble methods
- 8 Dropout
- 9 Other regularization methods

# Dataset Augmentation

- The best way to make a machine learning model generalize better is to train it on **more data**.
- For some machine learning tasks, the algorithm is expected to be robust to noise or be **invariant** to some transformations.
- It is reasonably straightforward to generate new fake data by transforming the inputs.
- Dataset augmentation has been a particularly effective technique for tasks like object recognition, speech recognition.

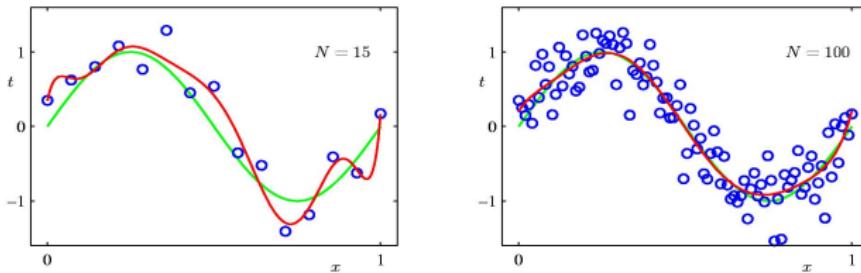
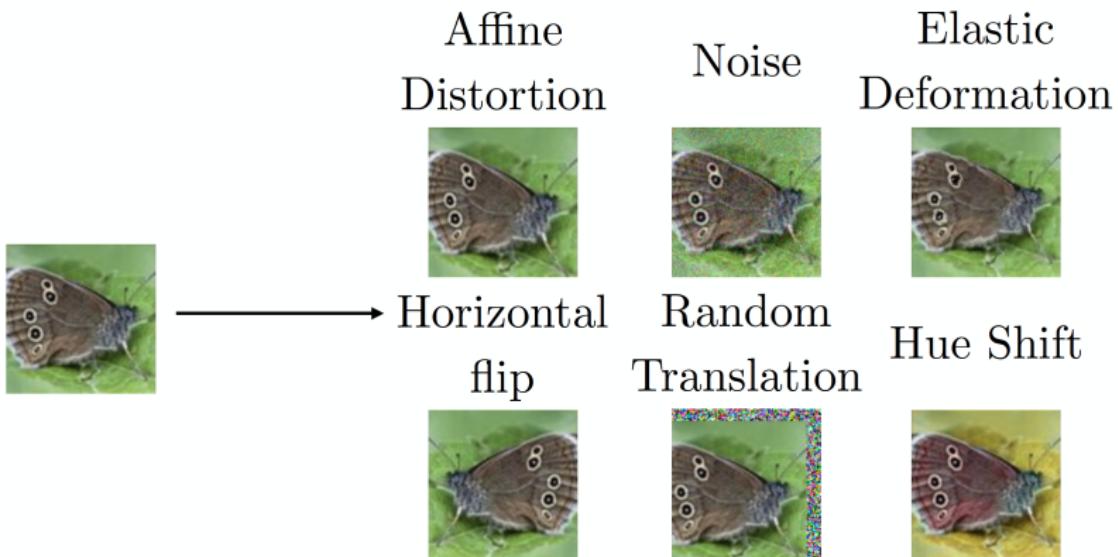


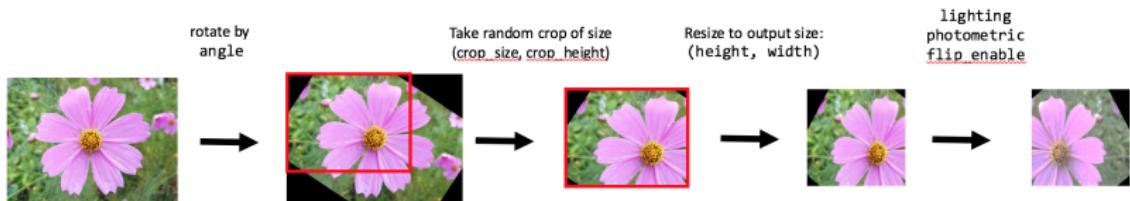
Figure 4: Recall the data-augment examples for fitting the polynomial functions.

# Dataset Augmentation: Example

- In objective recognition, we would like our classifiers to be invariant to translations such as translating, rotating, scaling and etc.

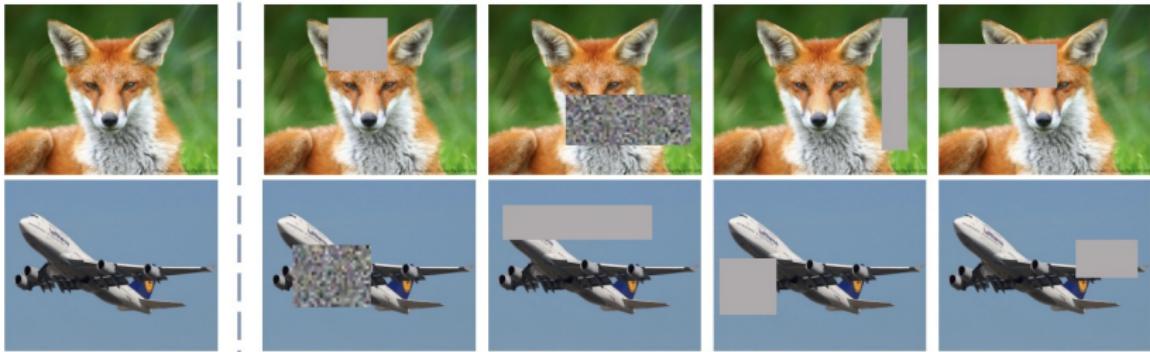


# Random Crop



- Randomly crops a fix-sized rectangle of the image, and resize it to the input scale.
- Random Cropping is a very useful data augmentation technique that forces the deep learning model to learn robustly and ignore positional bias in the training set.

# DropBlock



- Randomly selects a random-sized rectangle of the image, fill it with grey or random color.
- DropBlock(Ghiasi et al. 2018) forces the model not to rely on subtle local patterns of the image and guides it to learn the general concept.

# Outline

- 1 Background: Overfitting
- 2 Parameter Norm Penalties
- 3 Dataset Augmentation
- 4 Noise Robustness
- 5 Early Stopping
- 6 Parameter Sharing
- 7 Ensemble methods
- 8 Dropout
- 9 Other regularization methods

## Noise Injection at the Inputs

- One way to improve the robustness of neural networks is simply to train them with random noise applied to their inputs. It can be motivated as a **data augmentation** strategy.
- For some models, the addition of noise with infinitesimal variance at the input of the model is **equivalent** to imposing a penalty on the norm of the weights (Bishop, 1995a,b).
- Input noise injection is part of some unsupervised learning algorithms such as the **denoising autoencoder**.
- Noise injection also works when the noise is applied to the hidden units, which can be seen as doing dataset augmentation at multiple levels of abstraction.

# Noise Injection at the Weights

- Another way that noise has been used in the service of regularizing models is by adding it to the **weights**. This technique has been used primarily in the context of recurrent neural networks.
- This can be interpreted as a stochastic implementation of a Bayesian inference over the weights.
  - **The Bayesian treatment** of learning would consider the model weights to be uncertain and representable via a probability distribution.
  - Adding noise to the weights is a practical, stochastic way to reflect this **uncertainty**.

Example:

$$\mu_N = \frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2} \cdot \left( \frac{1}{N} \sum_{n=1}^N x_n \right) + \frac{\sigma^2}{N\sigma_0^2 + \sigma^2} \mu_0 \quad (31)$$

- This can also be interpreted as equivalent (under some assumptions) to a more traditional form of regularization.

# Noise Injection at the Weights

- Effect of weight noise in the regression setting:
  - The object is to minimize the least-square cost between the model predictions  $\hat{y}(x)$  and the true values  $y$ :

$$J = \mathbb{E}_{p(x,y)}[(\hat{y}(x) - y)^2] \quad (32)$$

- Assume that with each input presentation we also include a random perturbation  $\epsilon_w \sim \mathcal{N}(\epsilon; \mathbf{0}, \eta I)$ . We denote the perturbed model as  $\hat{y}_{\epsilon_w}(x)$ . The objective function thus become:

$$\tilde{J}_w = \mathbb{E}_{p(x,y,\epsilon_w)}[(\hat{y}_{\epsilon_w}(x) - y)^2] \quad (33)$$

$$= \mathbb{E}_{p(x,y,\epsilon_w)}[\hat{y}_{\epsilon_w}^2(x) - 2y\hat{y}_{\epsilon_w}(x) + y^2] \quad (34)$$

- For small  $\eta$ , the minimization of  $J$  with added weight noise (with covariance  $\eta I$ ) is equivalent to minimization of  $J$  with an additional regularization term:  $\eta \mathbb{E}_{p(x,y)}[\|\nabla_w \hat{y}(x)\|^2]$
- This form of regularization encourages the parameters to go to regions of parameter space where small perturbations of the weights have a relatively small influence on the output.

## Noise Injection at the Outputs

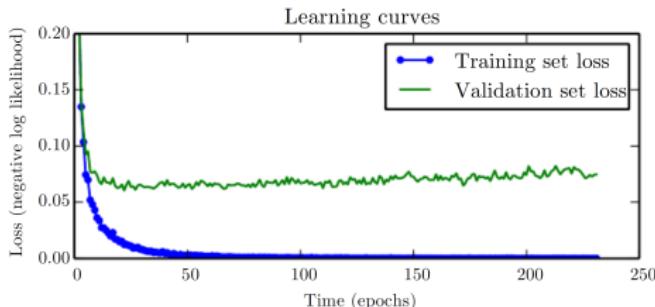
- Most datasets have some amount of mistakes in the  $y$  labels. It can be harmful to maximize  $\log p(y|x)$  when  $y$  is a mistake.
  - For example, we can assume that for some small constant  $\epsilon$ , the training set label  $y$  is correct with probability  $1 - \epsilon$ , and otherwise any of the other possible labels might be correct.
  - This assumption is easy to incorporate into the cost function analytically, rather than by explicitly drawing noise samples.
- **Label smoothing** regularizes a model based on a softmax with  $k$  output values by replacing the hard 0 and 1 classification targets with targets of  $\frac{\epsilon}{k}$  and  $\frac{1-\epsilon}{k}$ , respectively.

# Outline

- 1 Background: Overfitting
- 2 Parameter Norm Penalties
- 3 Dataset Augmentation
- 4 Noise Robustness
- 5 Early Stopping
- 6 Parameter Sharing
- 7 Ensemble methods
- 8 Dropout
- 9 Other regularization methods

# Early Stopping

- When training large models with sufficient representational capacity to overfit the task, we often observe that training error decreases steadily over time, but validation set error begins to **rise** again.



- We can obtain a model with better validation set error (and thus, hopefully better test set error) by returning to the parameter setting at the point in time with the lowest validation set error.
- This strategy is known as **early stopping**. It is probably the most commonly used form of regularization in deep learning.

# Early Stopping

---

**Algorithm 7.1** The early stopping meta-algorithm for determining the best amount of time to train. This meta-algorithm is a general strategy that works well with a variety of training algorithms and ways of quantifying error on the validation set.

---

Let  $n$  be the number of steps between evaluations.

Let  $p$  be the “patience,” the number of times to observe worsening validation set error before giving up.

Let  $\theta_o$  be the initial parameters.

$\theta \leftarrow \theta_o$

$i \leftarrow 0$

$j \leftarrow 0$

$v \leftarrow \infty$

$\theta^* \leftarrow \theta$

$i^* \leftarrow i$

**while**  $j < p$  **do**

    Update  $\theta$  by running the training algorithm for  $n$  steps.

$i \leftarrow i + n$

$v' \leftarrow \text{ValidationSetError}(\theta)$

**if**  $v' < v$  **then**

$j \leftarrow 0$

$\theta^* \leftarrow \theta$

$i^* \leftarrow i$

$v \leftarrow v'$

**else**

$j \leftarrow j + 1$

**end if**

**end while**

---

Best parameters are  $\theta^*$ , best number of training steps is  $i^*$

---

# Early Stopping

- Early stopping requires a **validation** set, which means some training data is not fed to the model.
- To best exploit this extra data, one can perform extra training after the initial training with early stopping has completed.
- There are two basic strategies one can use for this extra training procedure.

# Early Stopping

---

**Algorithm 7.2** A meta-algorithm for using early stopping to determine how long to train, then retraining on all the data.

---

Let  $\mathbf{X}^{(\text{train})}$  and  $\mathbf{y}^{(\text{train})}$  be the training set.

Split  $\mathbf{X}^{(\text{train})}$  and  $\mathbf{y}^{(\text{train})}$  into  $(\mathbf{X}^{(\text{subtrain})}, \mathbf{X}^{(\text{valid})})$  and  $(\mathbf{y}^{(\text{subtrain})}, \mathbf{y}^{(\text{valid})})$  respectively.

Run early stopping (Algorithm 7.1) starting from random  $\boldsymbol{\theta}$  using  $\mathbf{X}^{(\text{subtrain})}$  and  $\mathbf{y}^{(\text{subtrain})}$  for training data and  $\mathbf{X}^{(\text{valid})}$  and  $\mathbf{y}^{(\text{valid})}$  for validation data. This returns  $i^*$ , the optimal number of steps.

Set  $\boldsymbol{\theta}$  to random values again.

Train on  $\mathbf{X}^{(\text{train})}$  and  $\mathbf{y}^{(\text{train})}$  for  $i^*$  steps.

---

---

**Algorithm 7.3** Meta-algorithm using early stopping to determine at what objective value we start to overfit, then continue training until that value is reached.

---

Let  $\mathbf{X}^{(\text{train})}$  and  $\mathbf{y}^{(\text{train})}$  be the training set.

Split  $\mathbf{X}^{(\text{train})}$  and  $\mathbf{y}^{(\text{train})}$  into  $(\mathbf{X}^{(\text{subtrain})}, \mathbf{X}^{(\text{valid})})$  and  $(\mathbf{y}^{(\text{subtrain})}, \mathbf{y}^{(\text{valid})})$  respectively.

Run early stopping (Algorithm 7.1) starting from random  $\boldsymbol{\theta}$  using  $\mathbf{X}^{(\text{subtrain})}$  and  $\mathbf{y}^{(\text{subtrain})}$  for training data and  $\mathbf{X}^{(\text{valid})}$  and  $\mathbf{y}^{(\text{valid})}$  for validation data. This updates  $\boldsymbol{\theta}$ .

$\epsilon \leftarrow J(\boldsymbol{\theta}, \mathbf{X}^{(\text{subtrain})}, \mathbf{y}^{(\text{subtrain})})$

**while**  $J(\boldsymbol{\theta}, \mathbf{X}^{(\text{valid})}, \mathbf{y}^{(\text{valid})}) > \epsilon$  **do**

    Train on  $\mathbf{X}^{(\text{train})}$  and  $\mathbf{y}^{(\text{train})}$  for  $n$  steps.

**end while**

---

# Early Stopping

- How does early stopping act as a regularizer?
- Bishop (1995a) and Sjoberg and Ljung (1995) argued that early stopping has the effect of restricting the optimization procedure to a relatively small volume of parameter space in the neighborhood of the initial parameter value  $\theta_0$ .
- We can show, in the case of a simple linear model with a quadratic error function and simple gradient descent, early stopping is **equivalent** to  $L^2$  regularization.

# Early Stopping

- We examine a simple setting where the only parameters are linear weights ( $\theta = \mathbf{w}$ ).
- Model the cost function  $J$  with a quadratic approximation in the neighborhood of the empirically optimal value of the weights  $\mathbf{w}^*$ :

$$\hat{J}(\mathbf{w}) = J(\mathbf{w}^*) + \frac{1}{2}(\mathbf{w} - \mathbf{w}^*)^T \mathbf{H}(\mathbf{w} - \mathbf{w}^*) \quad (35)$$

$$\nabla_{\mathbf{w}} \hat{J}(\mathbf{w}) = \mathbf{H}(\mathbf{w} - \mathbf{w}^*) \quad (36)$$

- For simplicity, initialize the parameter vector to the origin ( $\mathbf{w}^{(0)} = \mathbf{0}$ ).
- The trajectory followed by the parameter vector during training:

$$\mathbf{w}^{(\tau)} = \mathbf{w}^{(\tau-1)} - \epsilon \nabla_{\mathbf{w}} \hat{J}(\mathbf{w}^{(\tau-1)}) \quad (37)$$

$$= \mathbf{w}^{(\tau-1)} - \epsilon \mathbf{H}(\mathbf{w}^{(\tau-1)} - \mathbf{w}^*) \quad (38)$$

$$\mathbf{w}^{(\tau)} - \mathbf{w}^* = (\mathbf{I} - \epsilon \mathbf{H})(\mathbf{w}^{(\tau-1)} - \mathbf{w}^*) \quad (39)$$

# Early Stopping

- Exploiting the eigen-decomposition of  $\mathbf{H}$ :  $\mathbf{H} = \mathbf{Q}\Lambda\mathbf{Q}^T$ , where  $\Lambda$  is a diagonal matrix and  $\mathbf{Q}$  is an orthonormal basis of eigenvectors.

$$\mathbf{w}^{(\tau)} - \mathbf{w}^* = (\mathbf{I} - \epsilon \mathbf{Q}\Lambda\mathbf{Q}^T)(\mathbf{w}^{(\tau-1)} - \mathbf{w}^*) \quad (40)$$

$$\mathbf{Q}^T(\mathbf{w}^{(\tau)} - \mathbf{w}^*) = (\mathbf{I} - \epsilon \Lambda)\mathbf{Q}^T(\mathbf{w}^{(\tau-1)} - \mathbf{w}^*) \quad (41)$$

$$\mathbf{Q}^T \mathbf{w}^{(\tau)} = (\mathbf{I} - (\mathbf{I} - \epsilon \Lambda)^r) \mathbf{Q}^T \mathbf{w}^* \quad (42)$$

- Recall Eq. 18 in  $L^2$  regularization:

$$\tilde{\mathbf{w}} = \mathbf{Q}(\Lambda + \alpha \mathbf{I})^{-1} \Lambda \mathbf{Q}^T \mathbf{w}^* \quad (43)$$

$$\mathbf{Q}^T \tilde{\mathbf{w}} = (\Lambda + \alpha \mathbf{I})^{-1} \Lambda \mathbf{Q}^T \mathbf{w}^* \quad (44)$$

$$\mathbf{Q}^T \tilde{\mathbf{w}} = (\mathbf{I} - (\Lambda + \alpha \mathbf{I})^{-1} \alpha) \mathbf{Q}^T \mathbf{w}^* \quad (45)$$

## Early Stopping

- Compare Eq. 42 and Eq. 45, we see that  $L^2$  regularization and early stopping can be seen to be equivalent if the hyperparameters  $\epsilon$ ,  $\alpha$  and  $\tau$  are chosen such that:

$$\mathbf{I} - (\mathbf{I} - \epsilon \boldsymbol{\Lambda})^r = \mathbf{I} - (\boldsymbol{\Lambda} + \alpha \mathbf{I})^{-1} \alpha \quad (46)$$

- Going even further, by taking logarithms and using the series expansion for  $\log(1 + x)$ , we can conclude that if all  $\lambda_i$  are small (that is,  $\epsilon \lambda_i \ll 1$  and  $\lambda_i/\alpha \ll 1$ ) then:

$$\tau \approx \frac{1}{\epsilon \alpha} \quad (47)$$

$$\alpha \approx \frac{1}{\tau \epsilon} \quad (48)$$

- Parameter values corresponding to directions of significant curvature are regularized less than directions of less curvature. In the context of early stopping, this means that parameters that correspond to directions of significant curvature tend to learn early relative to parameters corresponding to directions of less curvature.

## Early Stopping

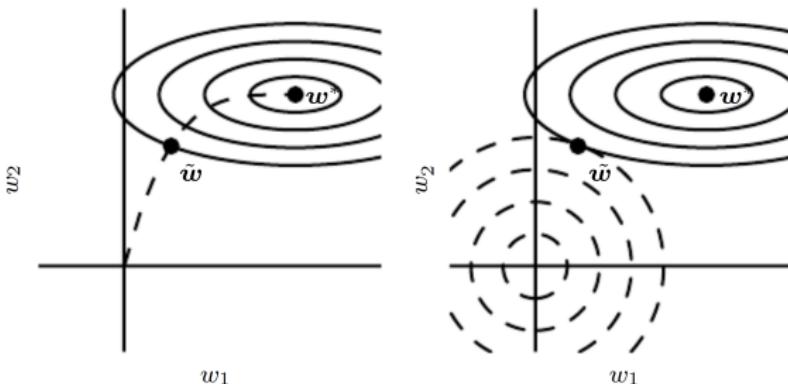


Figure 5: An illustration of the effect of early stopping. (Left) The solid contour lines indicate the contours of the negative log-likelihood. The dashed line indicates the trajectory taken by SGD beginning from the origin. Rather than stopping at the point  $w$  that minimizes the cost, early stopping results in the trajectory stopping at an earlier point  $\tilde{w}$ . (Right) An illustration of the effect of  $L^2$  regularization for comparison. The dashed circles indicate the contours of the  $L^2$  penalty, which causes the minimum of the total cost to lie nearer the origin than the minimum of the unregularized cost.

# Outline

- 1 Background: Overfitting
- 2 Parameter Norm Penalties
- 3 Dataset Augmentation
- 4 Noise Robustness
- 5 Early Stopping
- 6 **Parameter Sharing**
- 7 Ensemble methods
- 8 Dropout
- 9 Other regularization methods

# Parameter Sharing

- Sometimes we might not know precisely what values the parameters should take, but we know, from knowledge of the domain and model architecture, that there should be some dependencies between the model parameters.
  - Consider that there are two models performing the same classification task but with different input distributions. Formally, we have model  $A$  with parameters  $\mathbf{w}^{(A)}$  and model  $B$  with parameters  $\mathbf{w}^{(B)}$ .
  - Imagine that the tasks are similar enough that we believe the model parameters should be close to each other:  $w_i^{(A)}$  should be close to  $w_i^{(B)}$  for all  $i$ .
  - While we can use a parameter norm penalty  $\|\mathbf{w}^{(A)} - \mathbf{w}^{(B)}\|_2^2$  to regularize parameters to be close to one another, the more popular way is to use constraints: to force sets of parameters to be equal.
- This method of regularization is often referred to as **parameter sharing**.
- By far the most popular and extensive use of parameter sharing occurs in **convolutional neural networks** (CNNs) applied to computer vision.

# Semi-Supervised Learning

- In the context of deep learning, semi-supervised learning usually refers to learning a representation  $\mathbf{h} = f(\mathbf{x})$ .
  - The goal is to learn a representation so that examples from the same class have similar representations.
  - A linear classifier in the new space may achieve better generalization in many cases (Belkin and Niyogi, 2002; Chapelle et al., 2003).
- Instead of having separate unsupervised and supervised components in the model, one can construct models in which a **generative model** of either  $P(\mathbf{x})$  or  $P(\mathbf{x}, \mathbf{y})$  shares parameters with a discriminative model of  $P(\mathbf{y}|\mathbf{x})$ .
  - One can then trade-off the supervised criterion  $-\log P(\mathbf{y}|\mathbf{x})$  with the unsupervised or generative one (such as  $-\log P(\mathbf{x})$  or  $-\log P(\mathbf{x}, \mathbf{y})$ ).
  - The generative criterion then expresses a particular form of prior belief that the structure of  $P(\mathbf{x})$  is connected to the structure of  $P(\mathbf{y}|\mathbf{x})$  in a way that is captured by the shared parametrization.

# Multi-Task Learning

- Multi-task learning (Caruana, 1993) is a way to improve generalization by pooling the examples arising out several tasks.
- In the same way that additional training examples put more pressure on the parameters towards values that generalize well, when part of a model is shared across tasks, that part of the model is more constrained towards good values, often yielding better generalization.
- From the point of view of deep learning, the underlying prior belief is the following: *among the factors that explain the variations observed in the data associated with the different tasks, some are shared across two or more tasks.*

# Multi-Task Learning

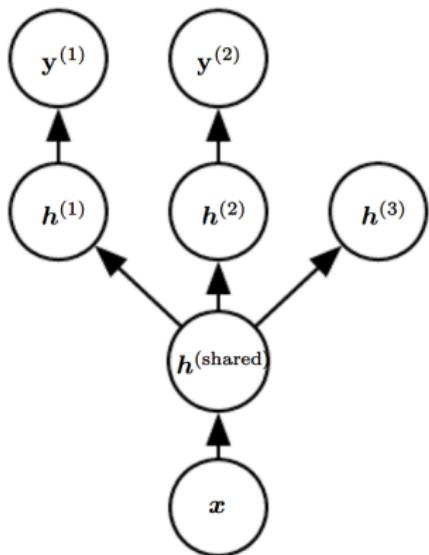


Figure 6: Multi-task learning

- Left figure illustrates a very common form of multi-task learning, in which different supervised tasks share the same input  $x$ , as well as some intermediate-level representation  $h^{(\text{shared})}$ .
- Improved generalization can be achieved because of the shared parameters, for which statistical strength can greatly improved.
- This will happen only if some assumptions about the statistical relationship between the different tasks are valid.

# Outline

- 1 Background: Overfitting
- 2 Parameter Norm Penalties
- 3 Dataset Augmentation
- 4 Noise Robustness
- 5 Early Stopping
- 6 Parameter Sharing
- 7 Ensemble methods**
- 8 Dropout
- 9 Other regularization methods

# Bagging and Ensemble Methods

- Bagging (short for **bootstrap aggregating**) is a technique for reducing generalization error by combining several models.
- The idea is to train several different models separately, then have all of the models vote on the output for test examples. This is an example of a general strategy in machine learning called **model averaging**. Techniques employing this strategy are known as **ensemble methods**.
- The reason that model averaging works is that different models will usually not make all the same errors on the test set.
  - Neural networks can often benefit from model averaging even if all of the models are trained on the same dataset. Differences in random initialization, random selection of minibatches, hyperparameters, or different outcomes of non-deterministic implementations of neural networks are often enough to cause different members of the ensemble to make partially independent errors.

# Ensemble Methods

- Consider for example a set of  $k$  regression models.
- Suppose that each model makes an error  $\epsilon_i$  on each example. The errors are drawn from a multivariate normal distribution with  $\mathbb{E}(\epsilon_i^2) = v$  and  $\mathbb{E}(\epsilon_i \epsilon_j) = c$ . Then the error made by the average prediction is  $\frac{1}{k} \sum_i \epsilon_i$ . The expected squared error of the ensemble predictor is

$$\mathbb{E}\left((\frac{1}{k} \sum_i \epsilon_i)^2\right) = \frac{1}{k^2} \mathbb{E}\left(\sum_i (\epsilon_i^2 + \sum_{j \neq i} \epsilon_i \epsilon_j)\right) = \frac{1}{k} v + \frac{k-1}{k} c \quad (49)$$

- If the errors are perfectly correlated and  $c == v$ , the mean error reduces to  $v$ , the model averaging does not help at all. If the errors are perfectly uncorrelated and  $c = 0$ , the expected squared error of the ensemble is only  $\frac{1}{k} v$ .
- On average, **the ensemble will perform at least as well as any of its members**, and if the members make independent errors, the ensemble will perform significantly better than its members.

# Bagging

- Different ensemble methods construct the ensemble of models in different ways.
  - For example, each member of the ensemble could be formed by training a completely different kind of model using a different algorithm.
- Bagging is a method that allows the same kind of model, training algorithm and objective function to be reused several times.
- Bagging first constructs  $k$  different datasets by sampling with replacement from the original dataset.
- Model  $i$  is then trained on dataset  $i$ . The differences between which examples are included in each dataset result in differences between the trained models.

# Bagging

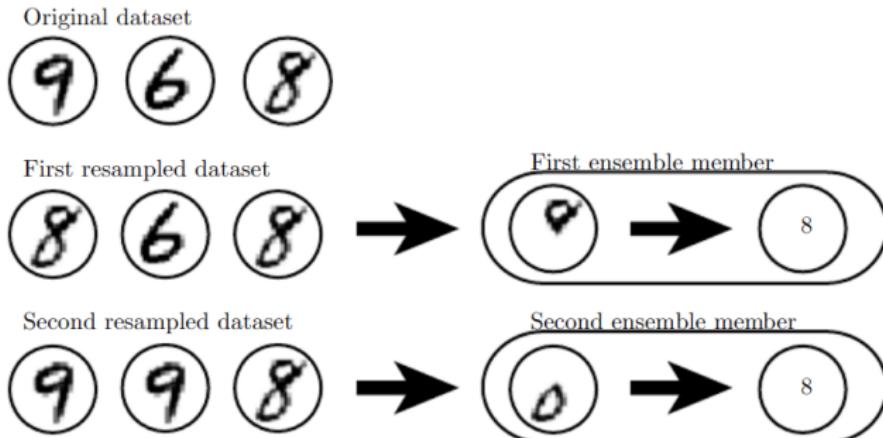


Figure 7: An example of training an 8 detector on the dataset containing an 8, a 6 and a 9. On the first dataset, the detector learns that a loop on top of the digit corresponds to an 8. On the second dataset, the detector learns that a loop on the bottom of the digit corresponds to an 8. Each of these individual classification rules is brittle, but if we average their output then the detector is robust.

# Outline

- 1 Background: Overfitting
- 2 Parameter Norm Penalties
- 3 Dataset Augmentation
- 4 Noise Robustness
- 5 Early Stopping
- 6 Parameter Sharing
- 7 Ensemble methods
- 8 Dropout**
- 9 Other regularization methods

# Dropout

- **Dropout** (Srivastava et al., 2014) provides a computationally inexpensive but powerful method of regularizing a broad family of models.
- Dropout can be thought of as a method of making bagging practical for ensembles of very many large neural networks.
  - Bagging involves training and evaluating multiple models, which seems impractical when each model is a large neural network.
  - Dropout trains the ensemble consisting of all sub-networks that can be formed by removing non-output unit from an underlying base network.
  - In most modern neural networks, we can effectively remove a unit from a network by multiplying its output value by zero.

# Dropout

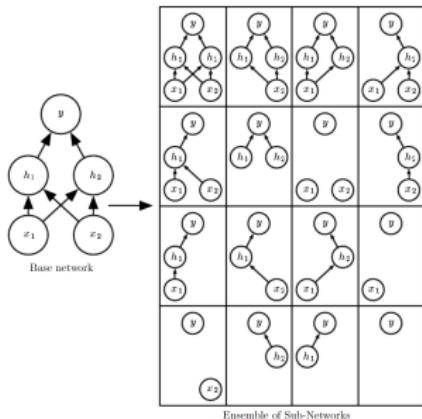


Figure 8: There are 16 possible subsets of a base network with 2 visible units and 2 hidden units.

- Dropout aims to approximate the bagging process, but with an *exponentially* large number of neural networks.
  - To train with dropout, we use a minibatch-based algorithm.
  - Each time we load an example into a minibatch, we randomly sample a binary mask for each of the input and hidden units.<sup>a</sup>
  - We then run forward propagation, back-propagation, and the learning update as usual.

<sup>a</sup>The mask probability is a set hyperparameter .

# Dropout

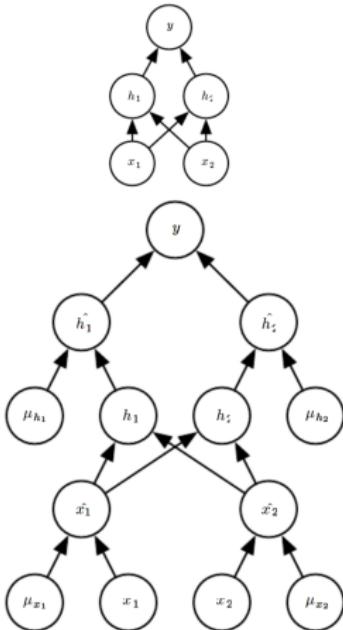


Figure 9: An example of forward propagation with dropout.

- To perform forward propagation with dropout, we randomly sample a vector  $\mu$  with one entry for each input or hidden unit int the network.
- Each unit in the network is multiplied by the corresponding mask, and then forward propagation continues through the rest of the network as usual.
- This is equivalent to randomly selecting one of the sub-networks from the base network and running forward propagation through it.

# Dropout

- Dropout training is not quite the same as bagging training.
- Bagging:
  - The models are all independent.
  - Each model is trained to convergence on its respective training set.
- Dropout:
  - The models share parameters, with each model inheriting a different subset of parameters from the base network.
  - Typically most models are not explicitly trained at all—the model is large enough that it would be infeasible to sample all possible sub-networks. Instead, a tiny fraction of the possible sub-networks are each trained for a single step, and the parameter sharing causes the remaining sub-networks to arrive at good settings of the parameters.

# Outline

- 1 Background: Overfitting
- 2 Parameter Norm Penalties
- 3 Dataset Augmentation
- 4 Noise Robustness
- 5 Early Stopping
- 6 Parameter Sharing
- 7 Ensemble methods
- 8 Dropout
- 9 Other regularization methods

# Sparse Representations

- We have already discussed how  $L^1$  penalization induces a sparse parametrization—meaning that many of the parameters become zero.
- Another sparsity regularization, representational sparsity, describes a presentation where many of the elements of the representation are 0 (or close to 0).

$$\begin{bmatrix} 18 \\ 5 \\ 15 \\ -9 \\ -3 \end{bmatrix} = \begin{bmatrix} 4 & 0 & 0 & -2 & 0 & 0 \\ 0 & 0 & -1 & 0 & 3 & 0 \\ 0 & 5 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & -1 & 0 & -4 \\ 1 & 0 & 0 & 0 & -5 & 0 \end{bmatrix} \begin{bmatrix} 2 \\ 3 \\ -2 \\ -5 \\ 1 \\ 4 \end{bmatrix}$$
$$\begin{bmatrix} -14 \\ 19 \\ 2 \\ 23 \end{bmatrix} = \begin{bmatrix} 3 & -1 & 2 & -5 & 4 & 1 \\ 4 & 2 & -3 & -1 & 1 & 3 \\ -1 & 5 & 4 & 2 & -3 & -2 \\ 3 & 1 & 2 & -3 & 0 & -3 \\ -5 & 4 & -2 & 2 & -5 & -1 \end{bmatrix} \begin{bmatrix} 0 \\ 2 \\ 0 \\ 0 \\ -3 \\ 0 \end{bmatrix}$$
$$\mathbf{y} \in \mathbb{R}^m \qquad \mathbf{A} \in \mathbb{R}^{m \times n} \qquad \mathbf{x} \in \mathbb{R}^n \qquad \mathbf{y} \in \mathbb{R}^m \qquad \mathbf{B} \in \mathbb{R}^{m \times n} \qquad \mathbf{h} \in \mathbb{R}^n$$

- (a) A linear regression model with parameter sparsity.      (b) A linear regression model with representational sparsity

Figure 10: An illustration of the distinction between **sparse parametrization** and **sparse representation**. In the right expression,  $\mathbf{h}$  is a sparse representation of the data  $\mathbf{x}$ . That is,  $\mathbf{h}$  is a function of  $\mathbf{x}$  that represents the information present in  $\mathbf{x}$  with a sparse vector.

# Sparse Representations

- Norm penalty regularization of representations is performed by adding to the loss function  $J$  a norm penalty on the representation.

$$\tilde{J}(\theta; \mathbf{X}, \mathbf{y}) = J(\theta; \mathbf{X}, \mathbf{y}) + \alpha \Omega(\mathbf{h}) \quad (50)$$

- Like on the parameters, an  $L^1$  penalty on the representation induces representational sparsity as well.
- Other approaches can result in a sparse representation:
  - Penalty derived from a Student-t prior of the representation.
  - KL divergence penalties that constrain the representations to lie on the unit interval.
  - Regularize the average activation across several examples to be near some target value.
  - Orthogonal matching pursuit (OMP) encodes an input  $\mathbf{x}$  with the representation  $\mathbf{h}$  that solves the constrained optimization problem:

$$\arg \min_{\mathbf{h}, \|\mathbf{h}\|_0 \leq k} \|\mathbf{x} - \mathbf{W}\mathbf{h}\|^2$$

where  $\|\mathbf{h}\|_0$  is the number of non-zero entries of  $\mathbf{h}$ .

# Adversarial Training

- In many cases, neural networks have begun to reach human performance, but it is unclear whether these models have obtained a true human-level understanding of these tasks.
- Szegedy et al. (2014b) found that even neural networks that perform at human level accuracy have a nearly 100% error rate on examples that are intentionally constructed.
- An optimization procedure can be used to search for an input  $x'$  near a data point  $x$  such that the model output is very different at  $x'$ .
- In many cases,  $x'$  can be so similar to  $x$  that a human observer cannot tell the difference between the original example and the **adversarial example**, but the network can make highly different predictions.

# Adversarial Training

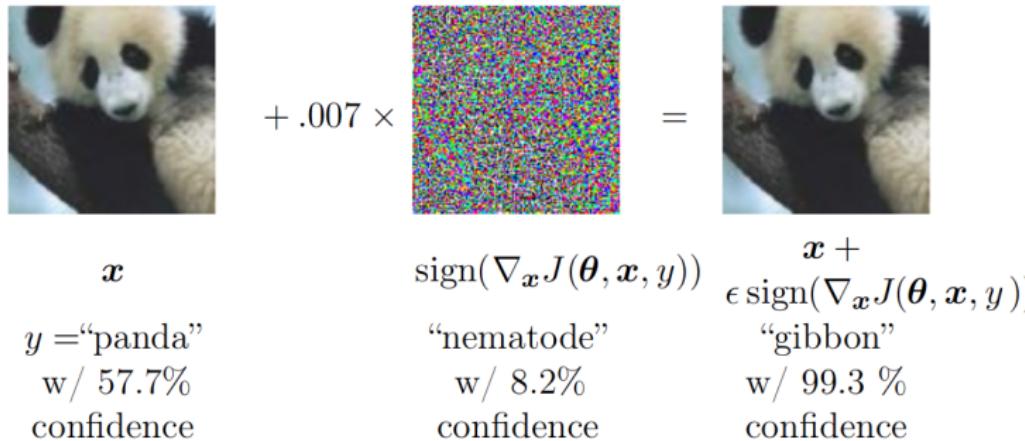


Figure 11: A demonstration of adversarial example generation applied to GoogLeNet on ImageNet. By adding an imperceptibly small vector whose elements are equal to the sign of the elements of the gradient of the cost function with respect to the input, we can change GoogLeNet’s classification of the image.

# Adversarial Training

- Adversarial training is to train on adversarially perturbed examples from the training set.
- Goodfellow et al. (2014b) showed that one of the primary causes of these adversarial examples is excessive linearity.
  - Neural networks are built out of primarily linear building blocks, which results in the overall function being highly linear.
  - The value of a linear function can change very rapidly if it has numerous inputs. If we change each input by  $\epsilon$ , then a linear function with weights  $w$  can change by as much as  $\epsilon\|w\|_1$ , which can be very large if  $w$  is high-dimensional.
- Adversarial training discourages this highly sensitive locally linear behavior by encouraging the network to be locally constant in the neighborhood of the training data.

# Adversarial Training

- Adversarial training helps to illustrate the power of using a large function family in combination with aggressive regularization.
  - Purely linear models are not able to resist adversarial examples.
  - Neural networks are able to represent all sorts of functions, thus can capture the linear trends in the training data while still learning to resist local perturbation.
- Adversarial examples also provide a meaning of accomplishing semi-supervised learning.
  - Given a point  $\mathbf{x}$  where the model outputs  $\hat{y}$  equals to the real label  $y$ , we can seek an adversarial example  $\mathbf{x}'$  close to  $\mathbf{x}$  that causes the classifier to output a label  $\hat{y}' \neq \hat{y}$ .
  - The classifier may then be trained to assign the same label to  $\mathbf{x}$  and  $\mathbf{x}'$ . This encourages the classifier to learn a function robust to small changes anywhere along the manifold.
  - The assumption motivating this approach is that different classes usually lie on disconnected manifolds, and a small perturbation should not be able to jump from one class manifold to another class manifold.(So the real label  $y' = y$  still holds for  $\mathbf{x}$  and  $\mathbf{x}'$ .)

# Overview

- 1 Background: Overfitting
- 2 Parameter Norm Penalties
- 3 Dataset Augmentation
- 4 Noise Robustness
- 5 Early Stopping
- 6 Parameter Sharing
- 7 Ensemble methods
- 8 Dropout
- 9 Other regularization methods

# Thanks.

**HP:** <http://keg.cs.tsinghua.edu.cn/jietang/>

**Email:** jietang@tsinghua.edu.cn