

# Course Notes: Deep Learning for Visual Computing

Peter Wonka

August 30, 2021

# Contents

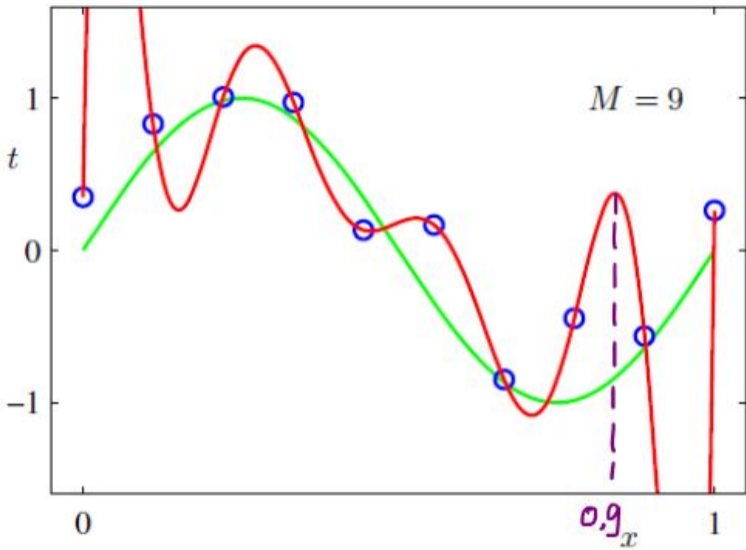
<b>1</b>	<b>Network Regularization</b>	<b>4</b>
1.1	Introduction	5
1.2	Design Choices	8
1.3	L2 regularization	10
1.4	L1 regularization	11
1.5	Elastic Net regularization	12
1.6	What parameters to regularize?	13
1.7	Max Norm Constraint	14
1.8	Dropout	15
1.8.1	Dropout vs. Inverted Dropout	17
1.8.2	Dropout Code Example	18
1.8.3	Survey	21
1.8.4	Dropout Discussion	22
1.9	Early Stopping	23
1.10	Multi-task Learning	24

1.11 Adding Noise to Network Outputs . . . . .	25
1.12 Batch-normalization as Regularizer . . . . .	26
1.13 DropConnect . . . . .	27
1.14 Stochastic Depth . . . . .	31
1.15 Data Augmentation . . . . .	33

# 1 Network Regularization

## 1.1 Introduction

- Review: Story of polynomial fitting



- Blue points: training data
- Red Line: without regularization the fitted polynomial passes through all points
  - over swinging, between training samples values can be very small and very large
- Green Line: using  $L_2$  regularization / weight decay

## 1.2 Design Choices

- Every layer can have a **different regularizer**
- Every layer can have a **different weight** for the regularizer
- Every parameter inside each layer can have a **different weight**, e.g. bias terms
- Regularizer sums up terms for each layer

$$R(\mathbf{W}_1, \mathbf{W}_2, \dots, \mathbf{W}_k) = \sum_{i=1}^k \lambda_i R_i(\mathbf{W}_i) \quad (1.1)$$

- $k$  number of layers,  $\mathbf{W}_i$  parameter tensor of layer  $i$
- Problems:
  - some layers have parameters that cannot be arranged in a tensor
  - not all parameters participate, different weights per parameters



- **Proposed notation:**

$$R(w_1, w_2, \dots, w_k) = \sum_{i \in I} \dots \quad (1.2)$$

- $I$  index set describing all parameters participating in the regularization
- $w_i$  individual scalar parameter
- omit parameters  $\lambda_i$  (can be set per layer or per weight)

## 1.3 L2 regularization

- add an  $L_2$  penalty to encourage small weights

$$R_{L2} = \sum_{i \in I} w_i^2 \quad (1.3)$$

- Comments
  - $L_2$  regularization is popular
  - $L_2$  regularization is also called **weight decay**
  - Check if implementation uses a factor of  $\frac{1}{2}\lambda$  or just  $\lambda$ 
    - $\frac{1}{2}$  is there to cancel the 2 when computing gradients

## 1.4 L1 regularization

- add an  $L_1$  penalty to encourage small weights:

$$R_{L1} = \sum_{i \in I} |w_i| \quad (1.4)$$

## 1.5 Elastic Net regularization

- combine  $L_1$  and  $L_2$  regularization:
  - separately add an  $L_1$  and  $L_2$  regularizer with different weights

$$R_{EL} = \sum_{i \in I} (\lambda_1 |w_i| + \lambda_2 w_i^2) \quad (1.5)$$

## 1.6 What parameters to regularize?

- What weights should be regularized using  $L_2$  or  $L_1$  regularization (What is  $I$ )?
  - common: multiplicative weights of linear layers (FC layers)
  - common: multiplicative weights of convolutional layers (CONV layers)
  - maybe / maybe not: bias terms in FC and CONV layers?
  - ???: batch normalization parameters
- How to set weights  $\lambda_i$ ?
  - Vary as little as possible

## 1.7 Max Norm Constraint

- Ensure that the sum of weights for one output / neuron  $\leq$  threshold *const*:
  - for each output in a linear layer /each conv filter check that:

$$\sum_{i \in NI} (w_i^2) \leq \mathbf{const}^2 \quad (1.6)$$

- $NI$  is the index set describing all the weights for a particular output / neuron
- If the weight exceeds the threshold *const* project onto the ball with radius *const*

## 1.8 Dropout

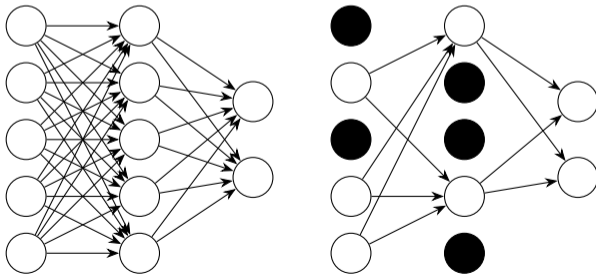
- Literature:
  - original paper: [Srivastava, Dropout: A Simple Way to Prevent Neural Networks from Overfitting](#)
  - survey paper: [Labach et al., Survey of Dropout Methods for Deep Neural Networks](#)
- Idea:
  - **Parameter**: probability  $p$  of keeping a value
  - During **training**:
    - randomly set some values to 0 with probability  $1 - p$
    - Scale other parameters by  $1/p$
  - During **inference**:
    - Eliminate the dropout layer

- Computation:

- Easiest to imagine as a separate layer that receives an input tensor  $\mathbf{X}_{in} \in R^{C,W,H}$  and outputs a modified tensor  $\mathbf{X}_{out}$ .
- For each value in the tensor:

$$\mathbf{X}_{out}(c, w, h) = \frac{1}{p} \times [rand < p] \times \mathbf{X}_{in}(c, w, h) \quad (1.7)$$

- $[rand < p]$  is the Iverson bracket





## 1.8.1 Dropout vs. Inverted Dropout

- Dropout changes the expected value of the output
- Two strategies
  - Dropout: scale the weights by  $p$  during inference
  - Inverted Dropout: scale the weights by  $\frac{1}{p}$  during training after randomly deleting weights with probability  $p$ .
- Discussion:
- Inverted Dropout is better, because the code during inference remains the same

## 1.8.2 Dropout Code Example

---

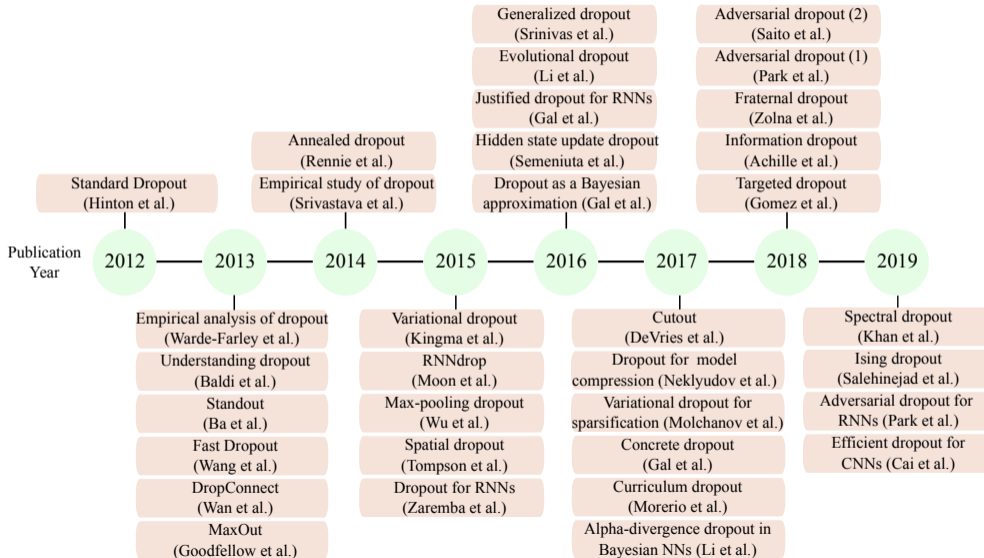
```
1  """
2  Inverted Dropout: recommended implementation
3  Drop and scale at train time and don't do anything at test time
4  """
5
6  import numpy as np
7
8  # probability of keeping a unit active. higher = less dropout
9  p = 0.5
10
11 # simulated input tensor of size 3 x 4
12 Input = np.random.randn(3,4)
13
14 # dropout computation. Notice /p!
15 Mask = (np.random.rand(*Input.shape) < p) / p
16 Output = Input * Mask # drop!
17
18 print(Input)
```

```
19 print (Mask)
20 print (Output)
```

---



## 1.8.3 Survey



## 1.8.4 Dropout Discussion

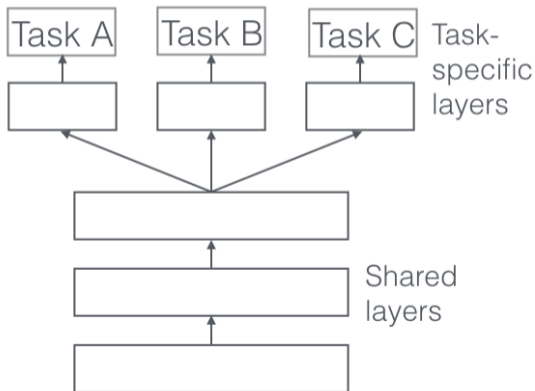
- Dropout can be interpreted as adding noise to the data
  - regularization by making the network robust to noise
- Dropout can be interpreted as simultaneously training an ensemble of networks (disputed)
- Dropout is typically only applied to **certain layers**
  - probability  $p$  changes per layer
  - more common in later layers of image classification
- Dropout was very popular before 2015
  - Some networks don't use dropout

## 1.9 Early Stopping

- When to stop training?
  - One Idea: plot the error on **training set** and **validation set**
  - stop training when error on validation set starts to increase
- Early stopping has a similar effect to  $L_2$  regularization

## 1.10 Multi-task Learning

- Develop a network to solve **multiple tasks** / problems at once
- Example: predict **segmentation**, **normals**, **depth** for each pixel in an image





## 1.11 Adding Noise to Network Outputs

- Example: **Label Smoothing**
- Assume a classification problem with one hot encoding of ground truth labels

$$\begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \end{pmatrix} \quad (1.8)$$

- Use smoothed labels,  $\epsilon$  smoothing parameter,  $k = \text{number of classes} - 1$

$$\begin{pmatrix} \epsilon/k \\ \epsilon/k \\ 1 - \epsilon \\ \epsilon/k \end{pmatrix} \quad (1.9)$$

## 1.12 Batch-normalization as Regularizer

- Batch normalization can be interpreted as follows
  - We estimate mean  $\mu$  and standard deviation  $\sigma$  from a mini-batch
  - $\mu$  and  $\sigma$  are an estimate of the **true data + noise**
  - Noise: difference between true  $\mu$  and  $\sigma$  and mini-batch  $\mu$  and  $\sigma$
  - Batch normalization is a regularizer

## 1.13 DropConnect

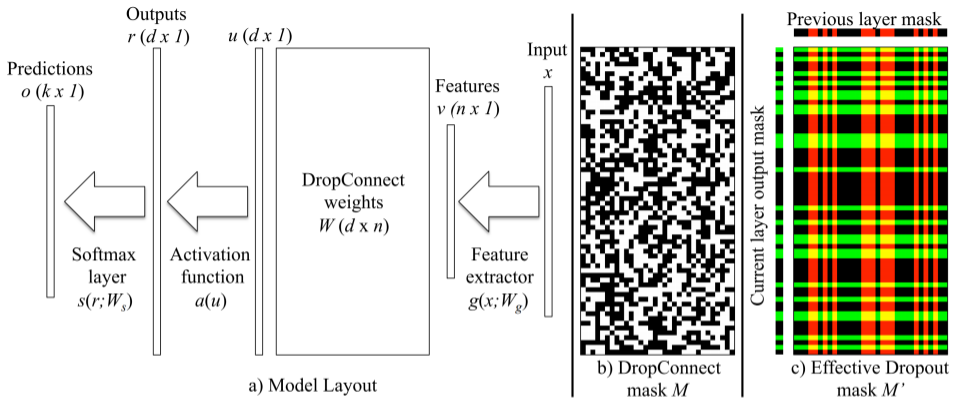
- Dropout: randomly set tensor values to zero
- DropConnect: randomly set network weights to zero
  - Results in a better random pattern
- Dropout example for a linear layer,  $\mathbf{W}$  parameters, input  $\mathbf{x}_{in}$ , output  $\mathbf{x}_{out}$

$$\mathbf{x}_{out} = \mathbf{W}\mathbf{x}_{in} \tag{1.10}$$

$$\begin{pmatrix} 1 & 3 & 4 & 7 \\ 2 & -2 & 5 & 3 \\ 2 & 9 & -1 & -9 \end{pmatrix} \begin{pmatrix} 2 \\ 5 \\ 6 \\ 8 \end{pmatrix} \tag{1.11}$$

$$\begin{pmatrix} 1 & 3 & 4 & 7 \\ 2 & -2 & 5 & 3 \\ 2 & 9 & -1 & -9 \end{pmatrix} \begin{pmatrix} 4 \\ \mathbf{0} \\ 12 \\ \mathbf{0} \end{pmatrix} \tag{1.12}$$

$$\begin{pmatrix} 2 & \mathbf{0} & 8 & \mathbf{0} \\ 4 & \mathbf{0} & 10 & \mathbf{0} \\ 4 & \mathbf{0} & -2 & \mathbf{0} \end{pmatrix} \begin{pmatrix} 2 \\ 5 \\ 6 \\ 8 \end{pmatrix} \tag{1.13}$$

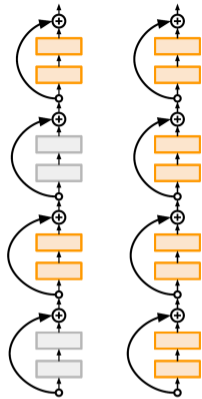


- Figure shows a single DropConnect layer.
- After running feature extractor  $g()$  on input  $x$ , a random mask  $M$  masks out the weight matrix  $W$ .
- Dropout comparison:

- effective weight mask for dropout applied to the previous layer's output (red columns)
- effective weight mask for dropout applied to this layer's output (green rows)

## 1.14 Stochastic Depth

- Literature: [Huang et al., Deep Networks with Stochastic Depth](#)
- Training: randomly skip layers
- Testing: use all layers





## 1.15 Data Augmentation

- Discussed Later