

Lecture 13

Overfitting & DD
Weight Decay
Dropout
Normalizations
...

OH \rightarrow 9:30, 10:30 MC 2505

HW5 out - Pb2
needs corrections

Q2 3/10 or 3/12

sol 4 - available
tomorrow

Lecture Notes IV – Neural Networks, Part 2

Marina Meilă
mmp@uwaterloo.ca

With Thanks to Pascal Poupart & Gautam Kamath
Cheriton School of Computer Science
University of Waterloo

February 8, 2026

Training a single unit ✓

Training a 2-layer network ✓

Training a L -layer network ✓

Backpropagation in practice ✓

Practical remedies to I, T, S, L, O

Saturation

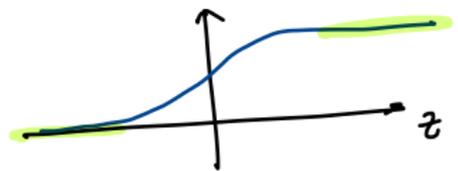
Over fitting

Reading HTF Ch.: 11.3 Neural networks, Murphy Ch.: (16.5 neural nets), Bach Ch.: -, Deep Learning Book (Goodfellow, Bengio, Courville) 6.1-4, ResNet 7.6, ConvNet 9., Autoencoders 14.1, Dive Into Deep Learning 4.1-4.3.

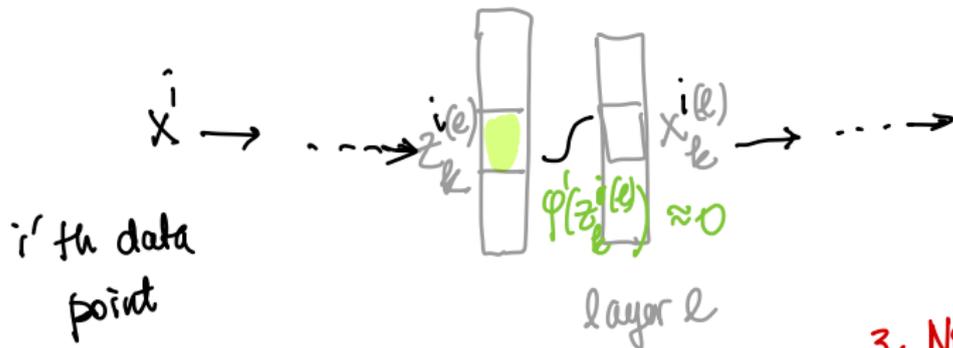
Backpropagation – some issues

- I Computation – how many ops / iteration?
- T Convergence – how many iterations (T) ?
- S Saturation – $\phi'(x) \approx 0$ for large $|z|$
- L Local minima
- 0 Overfitting?

= vanishing gradients



$|z| \text{ large} \Rightarrow \phi'(z) \approx 0$



Fixes:

1. ReLU
2. initialization

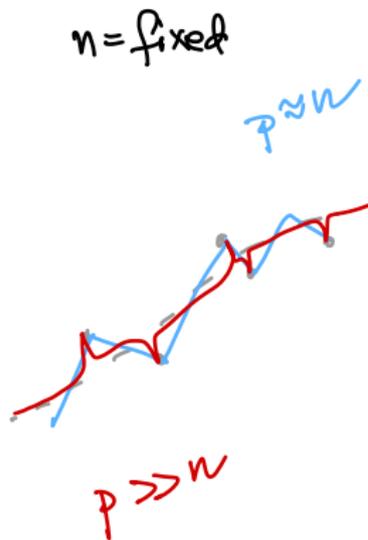
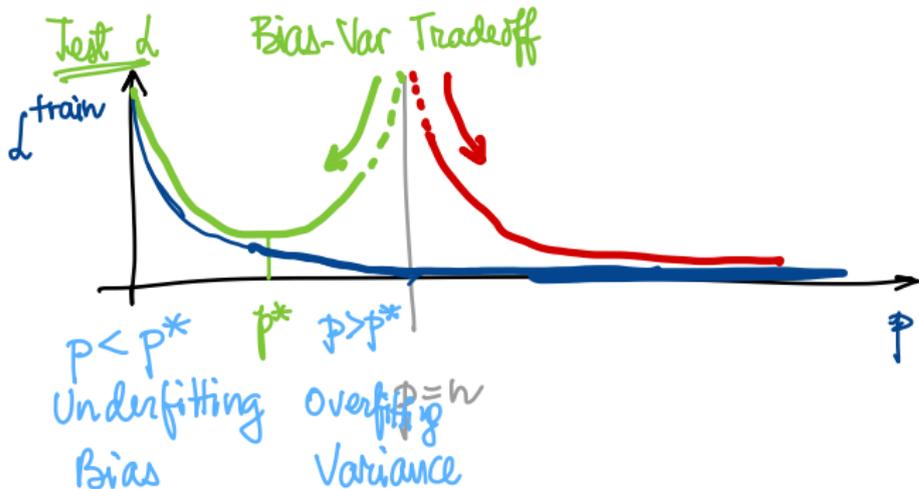
$w_j \approx 0$

3. Normalization $j=1:p$

Overfitting

$$p = \# \text{parameters}$$

- ▶ p can be very large deep networks
- ▶ Overfitting possible for everyday/in-house/not too large networks
- ▶ (Remedies later on in this lecture)
- ▶ Modern (very large) nn regime
 - ▶ Very large nn, with $p \gg n$ are not subject to the classical bias-variance tradeoff
 - ▶ This was a surprising discovery! The phenomenon is called double descent
 - ▶ (A separate lecture on it if time)



Normalizations – Controlling saturation S,O

- ▶ These refer to **data/inputs** to neurons, not to weights
- ▶ **Classic standardization**

1) Recenter data around mean $x^i \leftarrow x^i - \underline{\mu}$

$$\mu = \frac{1}{n} \sum_{i=1}^n x^i \quad (33)$$

2) Rescale each data coordinate $x_j^i \leftarrow x_j^i / \underline{\sigma_j}$

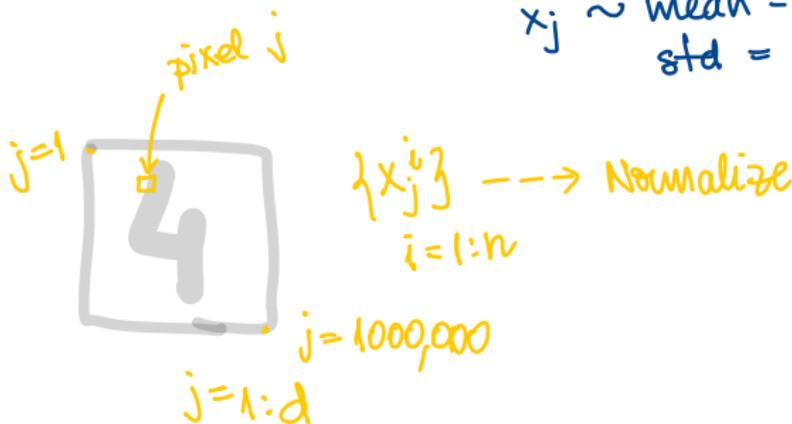
$j=1:d$

$$\sigma_j^2 = \frac{1}{n} \sum_{i=1}^n (x_j^i)^2 \quad (34)$$

$\sigma_j = \sqrt{\sigma_j^2}$

- ▶ Batch normalization – standardization within a batch
- ▶ Layer normalization – standardization within a layer

$$x_j \sim \begin{matrix} \text{mean} = 0 \\ \text{std} = 1 \end{matrix} \quad j=1:d$$



Batch normalization **S** $i \in \mathcal{B}$

1. ▶ Normalize the inputs in each \mathcal{B}
 3. ▶ Normalize the inputs to the sigmoids in each batch and layer $\{z^{i(l)}, i \in \mathcal{B}\}$
 - ↳ Normalize the inputs to the sigmoids in each batch and neuron $\{z_j^{i(l)}, i \in \mathcal{B}\}, j = 1 : m_l,$
 $l = 1 : L - 1$
 2. ▶ Normalize the inputs to each layer, i.e. $x^{i(l)}$ for $x^i \in \mathcal{B}, l = 0 : L - 1$
- ▶ Then add **learnable** scale and shift (same for all batches), e.g. $x^i \leftarrow \sigma_0 \frac{x^i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} + \mu_0$

$$x^{i(0)} \equiv x^i$$

$x^{i(l)}$ input to layer $l+1$

Layer normalization

 x_i

- ▶ Standardize $\{z_j^{(l)}, j = 1 : m_l\}$ for each data point and each layer



Avoiding overfitting \mathcal{O}

- ▶ Regularization, e.g. **weight decay**
- ▶ Dropout
 - ▶ Data augmentation
 - ▶ Early stopping
- ▶ Model averaging (also called **bagging**) or Bayesian Neural Networks

Weight decay - 0

train for min

Bias

- ▶ **Idea** Prevent overfitting by penalizing weights that are too large (Regularization)
- ▶ New "loss"

$$\mathcal{L}_\lambda(\mathbb{W}) = \mathcal{L}(\mathbb{W}) + \frac{\lambda}{2} \|\mathbb{W}\|^2 \quad (35)$$

and $\|\mathbb{W}\|^2 = \sum_{j=1}^p w_j^2$

Regularization term

- ▶ Effect on gradient $g_\lambda = \frac{\partial \mathcal{L}_\lambda}{\partial \mathbb{W}}$

 $\lambda > 0$ reg. param.

$$\frac{\partial \mathcal{L}_\lambda}{\partial w_j} = \frac{\partial}{\partial w_j} \left(\mathcal{L}(\mathbb{W}) + \frac{\lambda}{2} \|\mathbb{W}\|^2 \right) = \frac{\partial \mathcal{L}(\mathbb{W})}{\partial w_j} + \lambda w_j \quad (36)$$

$$w_j^{t+1} \leftarrow w_j^t - \eta \left(\frac{\partial \mathcal{L}(\mathbb{W})}{\partial w_j} + \lambda w_j^t \right) = \underbrace{w_j^t - \eta \frac{\partial \mathcal{L}(\mathbb{W})}{\partial w_j}}_{\text{old update}} - \eta \lambda w_j^t \quad (37)$$

$$= w_j^t (1 - \eta \lambda) - \eta \frac{\partial \mathcal{L}(\mathbb{W})}{\partial w_j} \quad (38)$$

$$g_j = \frac{\partial \mathcal{L}}{\partial w_j}$$

\uparrow
 ≤ 1

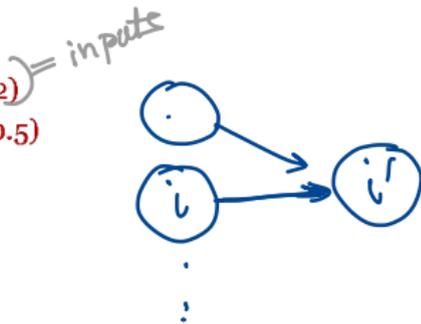
$$\frac{\partial}{\partial w_j} \left(\frac{1}{2} \sum_{j=1}^p w_j^2 \right) = w_j$$

Dropout — specific nn

Dropout

- Idea: randomly “drop” some units from the network when training
- Training: at each iteration of gradient descent
 - Each input unit is dropped with probability p_1 (e.g., 0.2)
 - Each hidden unit is dropped with probability p_2 (e.g., 0.5)
- Prediction (testing):
 - Multiply each input unit by $1 - p_1$
 - Multiply each hidden unit by $1 - p_2$

no drop out



<https://www.cs.toronto.edu/~rsalakhu/papers/srivastava14a.pdf> “Dropout: A simple way to prevent neural networks from overfitting” by Srivastava, Hinton, Krizhevsky, Sutskever, Salkhutinov, JMLR 2014.

Training with dropout

Train

- ▶ For each iteration
 - ▶ For each training example (x^i, y^i)
 1. Sample $r^{(l)} \in \{0, 1\}^{m_l}$ from Bernoulli($p_{\text{drop}}^{(l)}$), for $l = 1 : L$
 2. Backward propagation: skip/ignore dropped out units (do not compute gradient contribution from them)
 3. Update: only $w_i^{(l)}$ with $r_i^{(l)} = 1$

Test/Predict

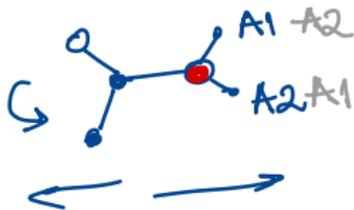
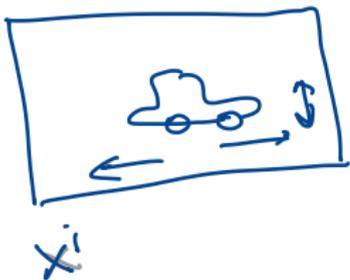
- ▶ Forward propagation $x^{(l)} = \phi(W^{(l)}z^{(l)}(1 - p_{\text{drop}}^{(l)}))$

Intuition

- Dropout can be viewed as an approximate form of ensemble learning
- In each training iteration, a different subnetwork is trained
- At test time, these subnetworks are “merged” by averaging their weights

Data augmentation - 0

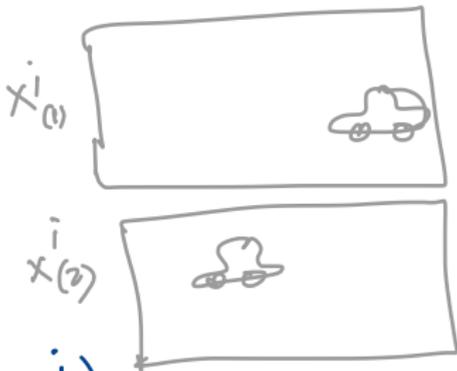
x has symmetries
 invariances



create synthetic

$x^i_{(1,2,\dots)}$ from x^i

to "teach" nn
 about invariance



$$(x^i, y^i) \rightarrow (x^i_{(1,2,\dots)}, y^i)$$

Drawbacks: $n \rightarrow 5n \gg n$

Better: Invariant / Equivariant nn.

Early stopping (Regularization)

