# Lecture 13

Overfitting & DD
Normalizations
Dropout
Weight decay
. . .

Q2 — WP 20% 3/5
              80% 3/10

L$\overline{\text{IV}}$-2   updated

# Lecture Notes IV – Neural Networks, Part 2

Marina Meilă

mmp@uwaterloo.ca

February 8, 2026

Training a single unit ✔

Training a 2-layer network ✔
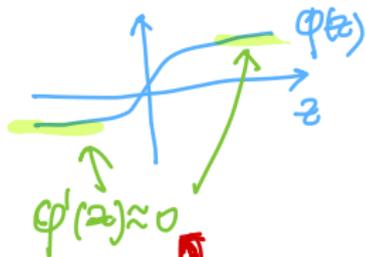
Training a $L$-layer network ✔

Backpropagation in practice ✔

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

**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 ($T$) ? ✓  # epochs = T

S  Saturation – $\phi'(x) \approx 0$ for large $|z|$ ←

L  Local minima ✓

0  Overfitting? ←

$\phi(z)$

$z$

$\phi'(z) \approx 0$

VANISHING gradients

1 epoch = 1 pass over $\mathcal{D}$
       = 1 iteration GD
       = $\frac{n}{n'}$ iterations of SGD ($n'$)

want: $\frac{\partial \mathcal{L}}{\partial W} = 0 = \sum_i \frac{\partial \mathcal{L}}{\partial W}(y^i, f(x^i; W))$
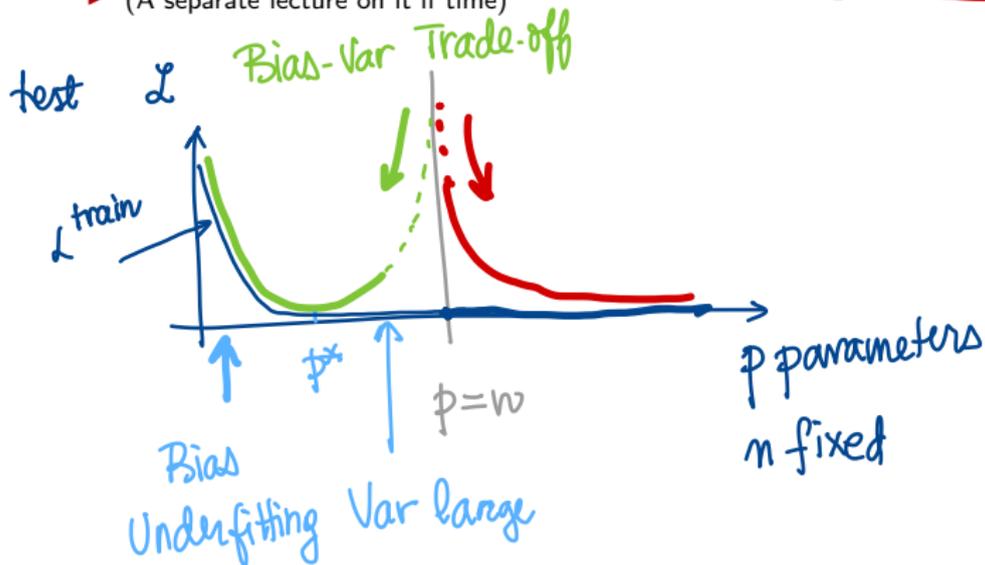
$W \in \mathbb{R}^P$

layer $\ell$

$z_k^{(\ell)}$   $x_k^{(\ell)}$

$|z_k^{(\ell)}(x^i)|$ large

# Overfitting

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

- ▶ $p$ 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
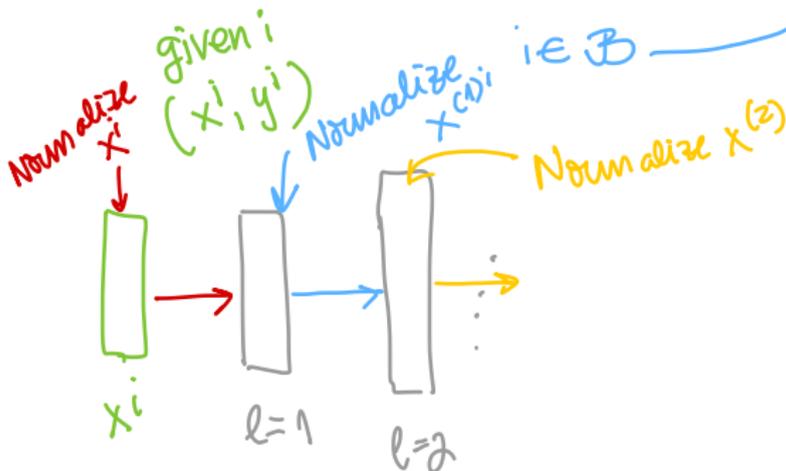  - ▶ (A separate lecture on it if time)

# Normalizations – Controlling saturation **S,O**

Std in very large nn

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

1) Recenter data around mean    $x^i \leftarrow x^i - \mu$    $\mu = \dfrac{1}{n}\sum_{i=1}^{n} x^i$    (33)

2) Rescale each data coordinate $j$: $x_j^i \leftarrow x_i^i/\sigma_j$    $\sigma_j^2 = \dfrac{1}{n}\sum_{i=1}^{n}(x_j^i)^2$    (34)

$\Rightarrow \{x_j^i\}_{i=1:n}$    0 mean
1 = variance

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



given i
$(x^i, y^i)$    $i \in B$

Normalize $x^i$    Normalize $x^{(1)i}$    Normalize $x^{(2)}$

$x^i$    $\ell = 1$    $\ell = 2$

# Batch normalization

- ▶ Normalize the inputs in each $\mathcal{B}$
- ▶ 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 $\{z_j^{i(l)}, i \in \mathcal{B}\}, j = 1 : m_l$, $l = 1 : L - 1$
- ▶ 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. $x^i \leftarrow \sigma_0 \dfrac{x^i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} + \mu_o$

1) $\{x^i\}_{i \in \mathcal{B}}$ $\xrightarrow{\text{Normalize}}$

2) $\{z_i^{(l)}\}_{i \in \mathcal{B}}$ $\xrightarrow{\text{Normalize}}$ $l = 1 : L$ recursively

Store all $\mu$'s, $\sigma$'s

At Test time

For Prediction : $x \leftarrow \dfrac{x - \mu}{\sigma}$

# Layer normalization

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

# Avoiding overfitting O

▶ Regularization, e.g. weight decay  $\underline{\text{Bias}}$   $w_i's$ not too large !!
▶ Dropout
▶ Data augmentation
▶ Early stopping
▶ Model averaging (also called **bagging**) or Bayesian Neural Networks

# Weight decay – O

$\mathcal{L} = \text{TRAINING LOSS}$

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

$$\text{min} \to \mathcal{L}_\lambda(\mathbb{W}) = \underset{\text{no regularization}}{\underbrace{\mathcal{L}(\mathbb{W})}} + \underset{\substack{\text{regularization} \\ \text{term} > 0}}{\underbrace{\frac{\lambda}{2}\|\mathbb{W}\|^2}} \tag{35}$$

$\text{min}$

$\lambda > 0$ reg. parameter

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

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

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

$$w_j^{t+1} \leftarrow w_j^t - \eta\left(\frac{\partial \mathcal{L}(\mathbb{W})}{\partial w_j} + \lambda w_j^t\right) = w_j^t - \eta\frac{\partial \mathcal{L}(\mathbb{W})}{\partial w_j} - \eta\lambda w_j^t \tag{37}$$

$$= w_j^t\underset{<1}{\underbrace{(1 - \eta\lambda)}} - \eta\underset{g_j}{\underbrace{\frac{\partial \mathcal{L}(\mathbb{W})}{\partial w_j}}} \tag{38}$$

*decay*

*Same as GD before*

$$\|\mathbb{W}\|^2 = \sum_{j=1}^p w_j^2 \Rightarrow \frac{\partial}{\partial w_j}\sum_{j'=1}^n w_{j'}^2 = 2w_j$$

# 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 $1 - p_2$

UNIVERSITY OF
WATERLOO

# 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 $x^{(l)} = \phi(W^{(l)} z^{(l)} (1 - p_{\text{drop}}^{(l)})$
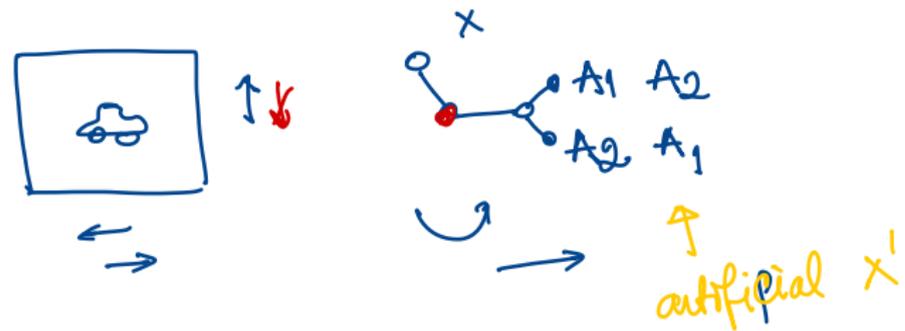
use all units

# 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

UNIVERSITY OF
**WATERLOO**

Data augmentation – **O** $\longrightarrow$ Better solution — incorporate symmetries in model

If there are symmetries in $X$, then create synthetic examples



$X$

$A_1$  $A_2$

$A_2$  $A_1$

artificial $X'$

$\mathcal{D} \longrightarrow \mathcal{D}^{\wedge S}$

and all synthetic examples