# Today's Agenda

1) Backprop Recap

2) Batch Norm

3) Optimizers

## Batch Normalization
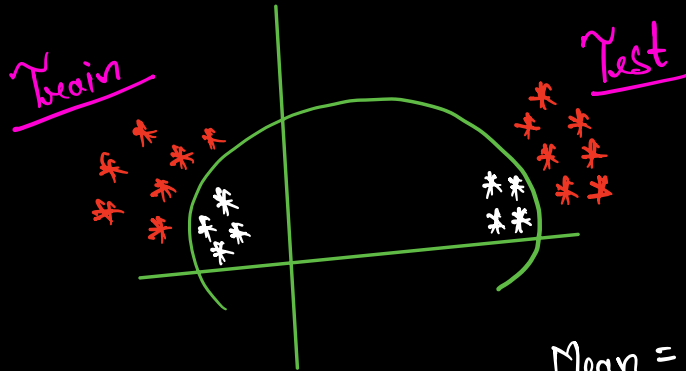
1) Internal covariate shift



| | |
|---|---|
| Training Distribution | Testing Distribution |
| | seen classes     unseen classes |
| dogs in water   cats in grass | dogs in grass   cats in water   bicycle, boat... |

{ Object }- Dog
Background — Water

train
test
↳ unseen
overfitting

Train    Test

Retrain the
network

Normalization

Mean = 0
Std = 1

$$X_i = \frac{X_i - Mean_i}{Std\ Dev_i}$$

Mini batch = hyperparameter → Batch Size

| 10 | 12 | Placed |
|----|----|--------|
| 6  | 8  | 0      |
| 9  | 9  | 1      |
| 8  | 8  | 1      |
| 5  | 6  | 1      |
| .  | .  | .      |
| .  | .  | .      |
| .  | .  | .      |

$h_1$   $h_2$   $h^3$   $h^4$

Independent Networks

1. Training Time
$\hookrightarrow$ low LR

$Z_{11} = wx + b$

$= w_1 \, 10 + w_2 \, 12 + b$

$g(z) = a_{11}$

$Z_{12}$
$13$
$14$

✓ $Z_{11} \longrightarrow Z_{11}^N \longrightarrow g\left(Z_{11}^N\right) \longrightarrow a_{11}$

OR

✳ $Z_{11} \longrightarrow g\left(Z_{11}\right) \longrightarrow a_{11} \longrightarrow a_{11}^N$ ✳

$$\frac{Z_{11} - \mu}{\sigma}$$

Mini Batch
Batch Size
$= h$

Mean $= \frac{1}{4} \cdot \overset{2}{\cdot}$
$\mu$

$\underline{\mu = \frac{1}{m} \sum z_{11}^i}$

**Input:** Values of $x$ over a mini-batch: $\mathcal{B} = \{x_{1...m}\}$;
Parameters to be learned: $\gamma, \beta$
**Output:** $\{y_i = \text{BN}_{\gamma,\beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m}\sum_{i=1}^{m} x_i \qquad \text{// mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m}\sum_{i=1}^{m}(x_i - \mu_{\mathcal{B}})^2 \qquad \text{// mini-batch variance}$$

$$\widehat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \qquad \text{// normalize}$$

$$y_i \leftarrow \gamma \widehat{x}_i + \beta \equiv \text{BN}_{\gamma,\beta}(x_i) \qquad \text{// scale and shift}$$

$$\sigma = \sqrt{\frac{1}{m} \sum_{i=0}^{m} (z_{11}^{i} - \mu)}$$

## Normalize

$$z_{11}^{N} = \frac{z_{11} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$= \frac{z_{11} - \mu}{\sigma + \sqrt{\epsilon}} \nearrow \underline{\text{Error Term}}$$

## Scale & Shift

$$z_{11}^{BN} = \gamma \, z_{11}^{N} + \beta$$

## Batch Normalization Parameters

$$\rightarrow \gamma, \beta \leftarrow \quad \rightarrow \text{Learnable Parameters}$$

$$\gamma\beta \leftarrow 0$$
$$0 \rightarrow \gamma\beta$$
$$Q$$

$\overset{\gamma'}{\gamma \beta}$

<u>Test / Inference</u>

Moving        Average

Mean

$\mu$

Std. dev

$\sigma$

$$\frac{\mu_1 + \mu_2 + \mu_3 + \mu_4 + \mu_5 \dots}{10} = \underline{\underline{\mu_{10}}}$$

$$= \overset{2}{\overset{.}{.}}$$

During Training

data points = 1000

Batch Size = 100

Iterations = $\underline{10}$

<u>Non Learnable Parameters</u>

1) Mov. Avg Mean          *

2) Mov. Avg Std dev