

Today's Agenda

1) Loss Functions

Regression

Mean Squared Error

$$MSE = \frac{1}{n} \sum_{i=1}^n (y - \hat{y})^2$$

Squares of Errors \longrightarrow large errors will be penalized

Outlier Sensitivity \longrightarrow Sensitive

2) MAE (Mean Absolute Error)

$$MAE = \frac{1}{n} \sum_{i=1}^n |y - \hat{y}|$$

Measure the mean absolute difference

1) Robust to outliers

* Not differentiable at zero

2) Linear Penalty

* gradient is constant *
Slower Convergence

3) Huber Loss (Smooth MAE)

$$\text{Huber} = \begin{cases} \frac{1}{2} (y - \hat{y})^2 & \text{if } |y - \hat{y}| \leq \delta \\ \delta |y - \hat{y}| - \frac{1}{2} \delta^2 & \text{otherwise} \end{cases}$$

Tuning $\leftarrow \delta \rightarrow$ threshold (0, 1)

MSE \rightarrow Small Errors
MAE \rightarrow Large Errors

Classification

1) Log Loss (Binary Cross Entropy)

* Binary Classification *

$$BCE = -\frac{1}{n} \sum_{i=1}^n [y \cdot \log(\hat{y}) + (1-y) \log(1-\hat{y})]$$

Predicted Probability

Penalize \rightarrow wrong predictions

Activation \rightarrow Sigmoid

Strong gradients

* Sensitive to mislabeled data *

Outlier Sensitivity \rightarrow High

Multiclass Classification

1) Categorical Cross Entropy

$$CCE :- -\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^K y_{ij} \log(\hat{y}_{ij})$$

n = data points

K = no of classes

Outlier Sensitivity \rightarrow High

Activation \Rightarrow Softmax

* Probabilistic *

3) Hinge Loss

$$L = \frac{1}{n} \sum_{i=1}^n$$

$$\max(0, 1 - y \cdot \hat{y})$$

$$y_{\text{true}} = \{-1, 1\}$$

Outlier Sensitivity :- Medium

* Binary Classification *

Weight Initialization

1) Zero Weight Ini

$$W = 0, \quad b = 0$$

All neurons learn the same features

* No learning *

* Same gradient values *

Never, $W = 0$

2) Random Initialization

$$W \sim N(0, 0.1)$$

$$(-0.01, 0.01)$$

✓ shallow networks

* Vanishing gradient *

Xavier / Glorot Initialization

n_{in} = No of input neurons

n_{out} = No of output neurons

Uniform

$$W \sim U\left(-\sqrt{\frac{6}{n_{in} + n_{out}}}, \sqrt{\frac{6}{n_{in} + n_{out}}}\right)$$

Normal

$$W \sim N\left(0, \frac{2}{n_{in} + n_{out}}\right)$$

Prevents Vanishing gradient

Fast convergence

Use :- Sigmoid, tanh

Not use :- ReLU

layer :- Fully Connected / Hidden / Dense / Linear

* ANN *

CNN \rightarrow ReLU ✓

2) He Initialization / Kaiming Initialization

Uniform

$$W = U \left(-\frac{\sqrt{6}}{n_{in}}, \frac{\sqrt{6}}{n_{in}} \right)$$

Normal

$$W = N \left(0, \frac{2}{n_{in}} \right)$$

Best \rightarrow Relu Family works very good.

Cons:- Sigmoid

Pros:-
Prevent V. Gr. P
Faster Convergence

Leccun Initialization LECUN

$$W \sim N\left(0, \frac{1}{n_{in}}\right)$$

* SELU *

1) Dropout

2) Batch Normalize

Dropout

Regularization Technique

prevent overfitting

Layer \rightarrow Dropout

10 Neurons

2 - Neurons

deactivated
turn-off

$$p = (0 - 1)$$

$$p = (0.2)$$

20% deactivate

Shuffle

* loss of Info *

Overfitting

Small-medium

$$(0.1 - 0.5)$$

$$(0.2 \cdot 25 \cdot 3)$$

1 \rightarrow Destructive

* Normalize * \rightarrow dataset

Pros:-

- 1) Improve generalization
- 2) Simple to implement
- 3) Works good with deep networks

Layer 1 \rightarrow Input (opening)

⋮

mid - closing

Layer 2 \rightarrow Output (closing)

10 \rightarrow

$\left\{ \begin{array}{c} 6 \\ 7 \\ 8 \\ 9 \end{array} \right\}$

layer 1 $\rightarrow \gamma$
 2 $\rightarrow \gamma$ 0.1
 3 $\rightarrow \gamma$ 0.2
 4 $\rightarrow \gamma$ 0.3
 5 $\rightarrow \gamma$ 0.3
 6 $\rightarrow \gamma$ 0.4

Batch Normalization

Samples, Features

$$x_i = \frac{x_i - \text{Mean}}{\text{Std. Dev}}$$

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_1 \dots x_m\}$;

Parameters to be learned: γ, β

Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

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

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

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

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

Algorithm 1: Batch Normalizing Transform, applied to activation x over a mini-batch.

Mean, $\mu_B \rightarrow \frac{1}{m} \sum_{i=1}^m x_i$

Variance, $\sigma_B^2 \rightarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2$

Normalization, $\tilde{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma^2 + \epsilon}}$

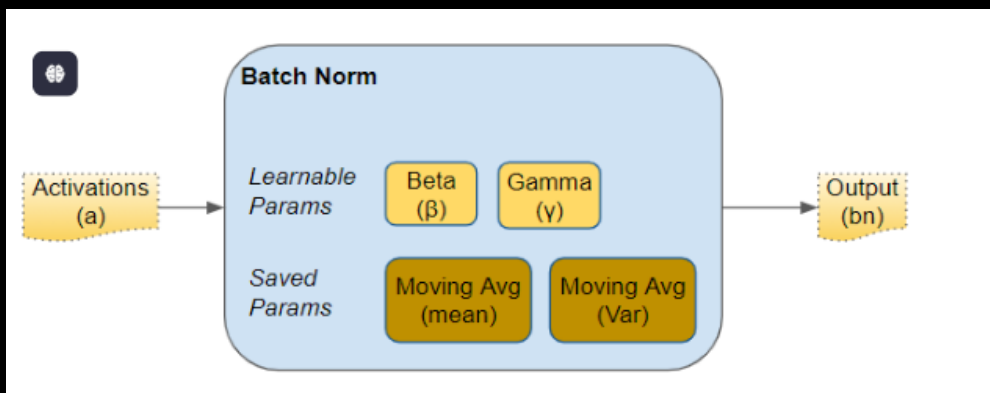
Scaling & Shifting

$$y_i \rightarrow \gamma \tilde{x}_i + \beta$$

β / $\gamma \rightarrow$ learnable Parameters
 Shift / scaling

$\beta, \gamma \rightarrow$ optimal activation

$w, b, \text{BN} \rightarrow$ learnable
 non-learnable



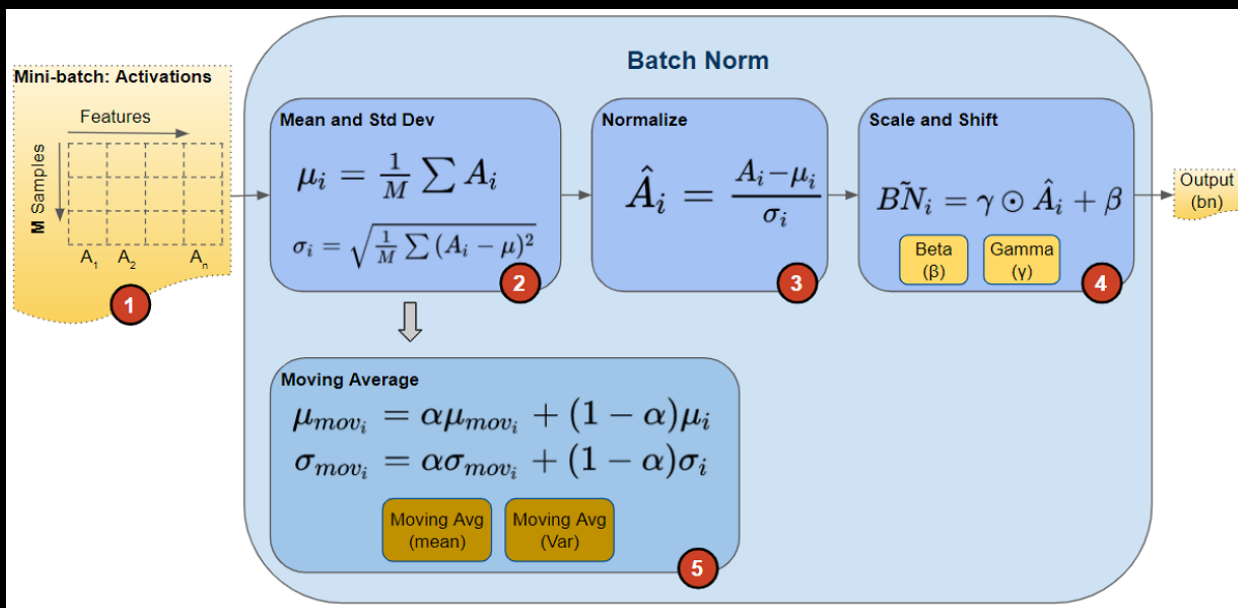
Training

b.s $n = 32$

* Saved * Model State

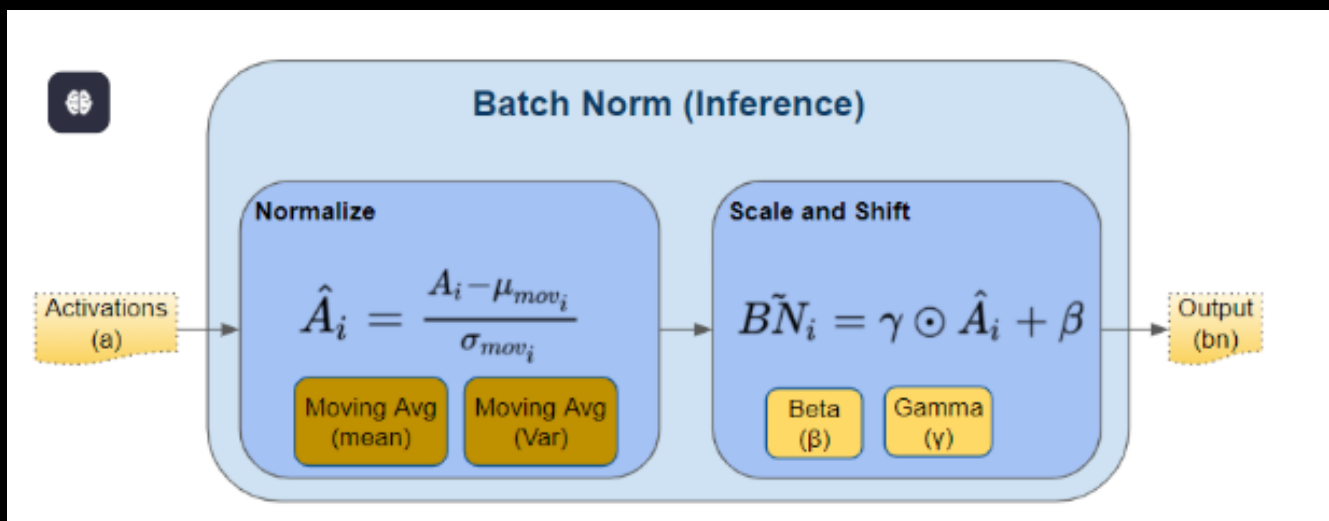
Inference

* Mov Avg { Mean } *
* { Var } *



$\alpha + (1 - \alpha)$

$$\beta = 0.9 = \begin{cases} \beta x + (1 - \beta) x \\ 0.9 + (1 - 0.9) \\ 0.9 + 0.1 \end{cases}$$



Placement of BN layer

- 1) Before Activation layer
- 2) After Activation layer

* Hidden layer \rightarrow Activation \rightarrow BN

