Batch Normalization Accelerating Deep Network Training by Reducing Internal Covariate Shift

Amit Patel Sambit Pradhan



- Introduction
- Internal Covariate Shift
- Batch Normalization
- Computational Graph of Batch Normalization
- Real World Application
- Practical Demo
- Discussion



Introduction Batch Normalization

The story Scareey Cat

Baby kitten for Christmas! Yey!



Young, Näve, Innocent, Scared of everything = Not Trained

Təsk – Trəin it to eət















First steppint traintee surpresty a Catzo Eat

Features

Training Data









Model

Output













Training Data with distribution of the features is different







The Pizza Super Duper Man



The Pizza Super Duper Man



The Pizza Super Duper Man - Normalizes







Covariate Shift

- Covariate The Features of the Input Data
- Covariate Shift Formally, is defined as a change in the distribution of a function's domain.
- Feature space of Target is drastically different than Source
- Training Data with distribution of the features is different.
- Informally, when input feature change, and algorithm can't deal with it, thus slowing the training.

Let's say you have a goal to reach, which is easier, a fixed goal a goal that keeps moving about? It is clear that a static goal is much easier to reach than a dynamic goal.

- IID Independent and identically distributed Each random variable has the same probability distribution as the others and all are mutually independent.
- Between the source (S) and target (T) given the same observation X = x conditional distributions of Y is same in both domains.

$$P_{s}(Y|X = x) = P_{t}(Y|X = x) \forall x \in X$$

In real world – In the wild
$$P_{s}(Y|X = x) \neq P_{t}(Y|X = x)$$

- To analyze the issue We consider a Parametric Model Family {P(Y|X, θ)}_{θ∉θ}
- We select a model {P(Y|X, θ*)}, which minimizes the expected classification error.
- If none of the models in the model family can exactly match the true relation between X and Y => there does not exist any θ∉Θ such that P(Y|X = x, θ) = P(Y|X = x) ∀ x∈X => We call this a misspecified model family.
- With a misspecified model family, the optimal model we select depends on P(X) and if Ps(X) ≠ Pt(X) then the optimal model for the target domain will differ from that for the source domain

- The intuitive is that the optimal model performs better in dense regions of than in sparse regions of , because the dense regions dominate the average classification error, which is what we want to minimize. If the dense regions of are different in the source and target then the optimal model for the source domain domain will no longer be optimal for the target domain.
- We re-weigh the $\frac{P_t(x,y)}{P_s(x,y)} = \frac{P_t(x)}{P_s(x)} \frac{P_t(y|x)}{P_s(y|x)}$

$$=\frac{P_t(x)}{P_s(x)}.$$

• We therefore re-weight each training instance with $\frac{P_t(x)}{P_r(x)}$

 Improving predictive inference under covariate shift by weighting the log-likelihood function - *Hidetoshi Shimodaira*

- Deep learning Is parameterized in a hierarchical fashion
- The first layere (Input Layer) looksa at the souce and the output of the first layer feeds the second layer, the second feeds the third, and so on.
- The distribution of the input is important Because we are actually learning a MODEL from the training data and NOT the right model.
- Internal Covariate Shift Small changes to the network get amplified down the network which leads to change in the input distribution to internal layers of the deep network.

Internal covariate shift refers to covariate shift occurring within a neural network, i.e. going from (say) layer 2 to layer 3. This happens because, as the network learns and the weights are updated, the distribution of outputs of a specific layer in the network changes. This forces the higher layers to adapt to that drift, which slows down learning.



Internal Covariate Shift applied to its parts, such as a sub-network or a layer.

 $\ell = F_2(F_1(\mathbf{u}, \Theta_1), \Theta_2)$

- Where F1 and F2 are arbitrary transformations, and the parameters ⊖1, ⊖2 are to be learned so as to minimize the loss ℓ.
- Learning $\Theta 2$ can be viewed as if the inputs $x = F1(u, \Theta 1)$ are fed into the sub-network

$$\ell = F_2(\mathbf{x}, \Theta_2).$$

$$\Theta_2 \leftarrow \Theta_2 - \frac{\alpha}{m} \sum_{i=1}^m \frac{\partial F_2(\mathbf{x}_i, \Theta_2)}{\partial \Theta_2}$$

 (for batch size m and learning rate α) is exactly equivalent to that for a stand-alone network F2 with input x.

Batch Normalization

Accelerating Deep Network Training by Reducing Internal Covariate Shift



Batch Normalization

Batch Normalization – Is a process normalize each scalar feature independently, by making it have the mean of zero and the variance of 1 and then scale and shift the normalized value *for each training mini-batch thus* reducing internal covariate shift fixing the distribution of the layer inputs x as the training progresses.

Using mini-batches of examples, as opposed to one example at a time, is helpful in several ways.

- First, the gradient of the loss over a mini-batch is an estimate of the gradient over the training set, whose quality improves as the batch size increases.
- Second, computation over a batch can be much more efficient than m computations for individual examples, due to the parallelism afforded by the modern computing platforms.

To remedy internal covariate shift, the solution proposed in the paper is to normalize each batch by both mean and variance

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

$$\begin{split} \mu_{\mathcal{B}} &\leftarrow \frac{1}{m} \sum_{i=1}^{m} x_{i} & // \text{ mini-batch mean} \\ \sigma_{\mathcal{B}}^{2} &\leftarrow \frac{1}{m} \sum_{i=1}^{m} (x_{i} - \mu_{\mathcal{B}})^{2} & // \text{ mini-batch variance} \\ \widehat{x}_{i} &\leftarrow \frac{x_{i} - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^{2} + \epsilon}} & // \text{ normalize} \\ y_{i} &\leftarrow \gamma \widehat{x}_{i} + \beta \equiv \text{BN}_{\gamma,\beta}(x_{i}) & // \text{ scale and shift} \end{split}$$

 $\gamma^{(k)} = \sqrt{\operatorname{Var}[x^{(k)}]}$

 $\beta^{(k)} = \mathbf{E}[x^{(k)}]$

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

Let us say that the layer we want to normalize has *d* dimensions $\mathbf{x} = (x_1, ..., x_d)$. Then, we can normalize the kth dimension as follows:

$$\hat{x}^k = \frac{x^k - E[x^k]}{\sqrt{Var[x^k]}}$$

- the parameters γ and β are to be learned, but it should be noted that the BN transform does not independently process the activation in each training example.
- Rather, BNγ,β(x) depends both on the training example and the other examples in the mini-batch. The scaled and shifted values y are passed to other network layers.

During training we need to backpropagate the gradient of loss & through this transformation, as well as compute the gradients with respect to the parameters of the BN transform.

$$\begin{split} \frac{\partial \ell}{\partial \hat{x}_{i}} &= \frac{\partial \ell}{\partial y_{i}} \cdot \gamma \\ \frac{\partial \ell}{\partial \sigma_{\mathcal{B}}^{2}} &= \sum_{i=1}^{m} \frac{\partial \ell}{\partial \hat{x}_{i}} \cdot (x_{i} - \mu_{\mathcal{B}}) \cdot \frac{-1}{2} (\sigma_{\mathcal{B}}^{2} + \epsilon)^{-3/2} \\ \frac{\partial \ell}{\partial \mu_{\mathcal{B}}} &= \left(\sum_{i=1}^{m} \frac{\partial \ell}{\partial \hat{x}_{i}} \cdot \frac{-1}{\sqrt{\sigma_{\mathcal{B}}^{2} + \epsilon}} \right) + \frac{\partial \ell}{\partial \sigma_{\mathcal{B}}^{2}} \cdot \frac{\sum_{i=1}^{m} -2(x_{i} - \mu_{\mathcal{B}})}{m} \\ \frac{\partial \ell}{\partial x_{i}} &= \frac{\partial \ell}{\partial \hat{x}_{i}} \cdot \frac{1}{\sqrt{\sigma_{\mathcal{B}}^{2} + \epsilon}} + \frac{\partial \ell}{\partial \sigma_{\mathcal{B}}^{2}} \cdot \frac{2(x_{i} - \mu_{\mathcal{B}})}{m} + \frac{\partial \ell}{\partial \mu_{\mathcal{B}}} \cdot \frac{1}{m} \\ \frac{\partial \ell}{\partial \gamma} &= \sum_{i=1}^{m} \frac{\partial \ell}{\partial y_{i}} \cdot \hat{x}_{i} \\ \frac{\partial \ell}{\partial \beta} &= \sum_{i=1}^{m} \frac{\partial \ell}{\partial y_{i}} \end{split}$$

To remedy internal covariate shift, the solution proposed in the paper is to normalize each batch by both mean and variance





Backpropagation



Chain Rule in Backpropagation

$$egin{aligned} f(x,y) &= xy & o & rac{\partial f}{\partial x} = y & rac{\partial f}{\partial y} = x \ f(x,y) &= x+y & o & rac{\partial f}{\partial x} = 1 & rac{\partial f}{\partial y} = 1 \ f(x,y) &= \max(x,y) & o & rac{\partial f}{\partial x} = 1(x > = y) & rac{\partial f}{\partial y} = 1(y > = x) \ f(x) &= rac{1}{x} & o & rac{d f}{d x} = -1/x^2 \ f_c(x) &= c+x & o & rac{d f}{d x} = 1 \ f(x) &= e^x & o & rac{d f}{d x} = e^x \ f_a(x) &= ax & o & rac{d f}{d x} = a \end{aligned}$$



STEP 9 – Summation gate





STEP 8 – First Multiplication gate





STEP 7 – Second Multiplication gate





STEP 6 – Inverse gate





STEP 5 – Square root gate















STEP 2 – Subtraction gate





STEP 1 – Mean gate









Real World Applications



Advantages of Batch Normalization

- Increased learning rate
- Remove Dropout
- Increased accuracy
- Allow use of saturating nonlinearities

Disadvantages of Batch Normalization

- Difficult to estimate mean and standard deviation of input during testing
- Cannot use batch size of 1 during training
- Computational overhead during training

Practical Demo

CIFAR-10

- 4 convolutional layers.
- ReLu activation function.
- Mini-batch size to be 32.
- Batch Normalization added before each activation.
- 50,000 train samples
- 10,00 test samples

```
def train model(useBN):
   model = Sequential()
   model.add(Convolution2D(32, 3, 3, border mode='same',
                            input shape=X train.shape[1:]))
   if useBN:
       model.add(BatchNormalization())
   model.add(Activation('relu'))
   model.add(Convolution2D(32, 3, 3))
   if useBN:
       model.add(BatchNormalization())
   model.add(Activation('relu'))
   model.add(MaxPooling2D(pool size=(2, 2)))
   model.add(Dropout(0.25))
   model.add(Convolution2D(64, 3, 3, border mode='same'))
   if useBN:
        model.add(BatchNormalization())
   model.add(Activation('relu'))
   model.add(Convolution2D(64, 3, 3))
   if useBN:
       model.add(BatchNormalization())
   model.add(Activation('relu'))
   model.add(MaxPooling2D(pool size=(2, 2)))
   model.add(Dropout(0.25))
   model.add(Flatten())
   model.add(Dense(512))
   if useBN:
       model.add(BatchNormalization())
```

```
model.add(Activation('relu'))
model.add(Dropout(0.5))
model.add(Dense(nb_classes))
if useBN:
    model.add(BatchNormalization())
model.add(Activation('softmax'))
```

Learning curve on CIFAR-10



Related Work

- Recurrent Batch Normalization Batch-normalize the hidden-tohidden transition, thereby reducing internal covariate shift between time steps.
- Weight Normalization: A Simple Reparameterization to Accelerate Training of Deep Neural Networks
- Normalization Propagation: A Parametric Technique for Removing Internal Covariate Shift in Deep Networks
- Layer Normalization

Citations

 <u>http://adsabs.harvard.edu/cgi-bin/nph-</u> ref_query?bibcode=2015arXiv150203167I&refs=CITATIONS&db_key=PRE

