



FRIEDRICH-ALEXANDER-
UNIVERSITÄT
ERLANGEN-NÜRNBERG
SCHOOL OF ENGINEERING

Regularization

K. Breininger, S. Vesal, B. Geissler, N. Maul, L. Reeb, M. Vornehm, Z. Yang, A. Popp, M. Nau,
S. Gündel, F. Denzinger, F. Thamm, C. Bergler, S. Jaganathan, F. Meister, C. Liu, T. Würfl
Pattern Recognition Lab, Friedrich-Alexander University of Erlangen-Nürnberg
June 19, 2020



Tasks in this exercise

L.L.

1. Optimization Constraints: Augmenting the loss function
2. Dropout **Layer**
3. Batch Normalization **Layer**
4. LeNet: Put everything together (**optional**)
5. RNN layer: Elman Unit
6. LSTM layer: Backpropagation at its best! (**optional**)



FRIEDRICH-ALEXANDER-
UNIVERSITÄT
ERLANGEN-NÜRNBERG
SCHOOL OF ENGINEERING

Optimization Constraints: Loss function augmentation



General outline

- Constraints change the total loss ...
- ... and have influence on the weight update of the respective layer!

Goal: weights smaller

But: large weights \rightarrow punish them!

General outline

- Constraints change the total loss ...
 - ... and have influence on the weight update of the respective layer!
 - Implement constraints as separate classes
- **Independent** of loss function

General outline

- Constraints change the total loss ...
- ... and have influence on the weight update of the respective layer!
- Implement constraints as separate classes
- Independent of loss function *Since Loss func only cares about y , \hat{y} labels.*
- Constraints **only need** current weights
- Add constraint objects in the optimizer

General outline

- Constraints change the total loss ...
- ... and have influence on the weight update of the respective layer!
- Implement constraints as separate classes
- **Independent** of loss function
- Constraints **only need** current weights
- Add constraint objects in the optimizer
- Since constraints generate part of the loss:
- Change Neural Network container class (and associated classes) to “channel” and gather **regularization loss** for **all layers**

L₂ regularization

every layer with weights: get $\lambda \|w\|_2^2$

- Forward pass:

$$\tilde{L}(\mathbf{w}) = L(\mathbf{w}) + \underbrace{\lambda \|\mathbf{w}\|_2^2}_{\text{L2 regularization}}$$

- Backward pass:

$$\mathbf{w}^{(k+1)} = \underbrace{(1 - \eta \lambda) \mathbf{w}^{(k)}}_{\text{Shrinkage}} - \eta \frac{\partial L}{\partial \mathbf{w}^{(k)}}$$

Note: The influence of constraints is controlled via λ . Because lambda is a python keyword, you want to use e.g. alpha instead.

L₁ regularization

- Forward pass:

$$\tilde{L}(\mathbf{w}) = L(\mathbf{w}) + \underbrace{\lambda \|\mathbf{w}\|_1}_{\text{L1 regularization}}$$

- Backward pass:

$$\mathbf{w}^{(k+1)} = \underbrace{\mathbf{w}^{(k)} - \eta \lambda \text{sign}(\mathbf{w}^{(k)})}_{\text{Other shrinkage}} - \eta \frac{\partial L}{\partial \mathbf{w}^{(k)}}$$



FRIEDRICH-ALEXANDER-
UNIVERSITÄT
ERLANGEN-NÜRNBERG
SCHOOL OF ENGINEERING

Dropout



Method

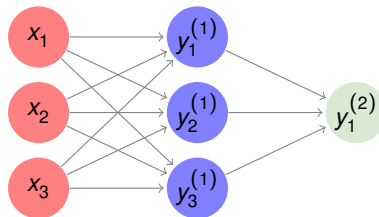


Figure: Dropout

- Implement this as a **fixed-function layer**

Method

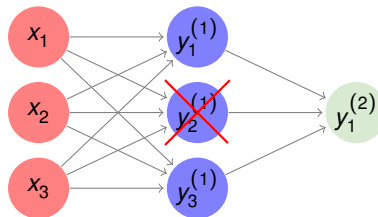


Figure: Dropout

- Implement this as a **fixed-function layer**
- Randomly set **activations** $\mapsto 0$ with probability $1 - p$

Method

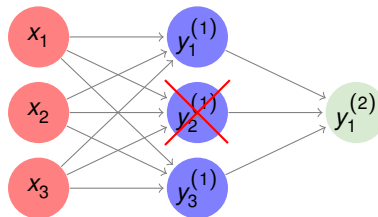


Figure: Dropout

- Implement this as a **fixed-function layer**
- Randomly set **activations** $\mapsto 0$ with probability $1 - p$ *\Rightarrow reduce the energy /*
- **Test-time**: multiply activations with p *\Leftarrow decrease energy lose some info
(since energy is high)*

Inverted Dropout

- Can we get rid of the dropout layer at test-time?

Inverted Dropout

- Can we get rid of the dropout layer at test-time?
- Change the behavior during training
- Multiply activations in forward-pass only during training by $\frac{1}{p}$ \Rightarrow energy \uparrow
dropout with $(1-p) \dots / p$
- Note: the backward pass has to be adapted as well!
derivative !



FRIEDRICH-ALEXANDER-
UNIVERSITÄT
ERLANGEN-NÜRNBERG
SCHOOL OF ENGINEERING

Batch normalization



Forward pass

→ Normalization as a new layer with 2 parameters, γ and β

Forward pass

→ Normalization as a new layer with 2 parameters, γ and β vector (length = channels)

input matrix

$$\tilde{X} = \frac{\underline{X} - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

μ_B and σ_B from **batch**

batch

$$\begin{cases} x_1 & x_2 & x_3 & \dots \\ x_1 & x_2 & x_3 & \dots \end{cases}$$

input dimension

Forward pass

→ Normalization as a new layer with 2 parameters, γ and β

$$\tilde{\mathbf{X}} = \frac{\mathbf{X} - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

μ_B and σ_B from **batch**

$$\hat{\mathbf{Y}} = \underbrace{\gamma}_{\text{element-wise}} \tilde{\mathbf{X}} + \beta$$

Forward pass

→ Normalization as a new layer with 2 parameters, γ and β

$$\tilde{\mathbf{X}} = \frac{\mathbf{X} - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

μ_B and σ_B from **batch**

$$\hat{\mathbf{Y}} = \gamma \tilde{\mathbf{X}} + \beta$$

- μ, σ have the **same dimension** as the **input vectors**

Forward pass

→ Normalization as a new layer with 2 parameters, γ and β

$$\tilde{\mathbf{X}} = \frac{\mathbf{X} - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

μ_B and σ_B from **batch**

$$\hat{\mathbf{Y}} = \gamma \tilde{\mathbf{X}} + \beta$$

- μ, σ have the same dimension as the **input vectors**
- β, γ and μ_B, σ_B have same **dimension** to be able to preserve **identity**

Forward pass

→ Normalization as a new layer with 2 parameters, γ and β

$$\tilde{\mathbf{X}} = \frac{\mathbf{X} - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

μ_B and σ_B from **batch**

$$\hat{\mathbf{Y}} = \gamma \tilde{\mathbf{X}} + \beta \quad \rightsquigarrow \text{FC Layer!}$$

- μ, σ have the **same dimension** as the **input vectors**
need to train in optimization
- β, γ and μ_B, σ_B have same **dimension** to be able to preserve **identity**
- Notice that β is a **bias**

Test time

- Test-time: replace μ_B and σ_B with μ and σ of the **training set**

Test time

- Test-time: replace μ_B and σ_B with μ and σ of the **training set**
- It's **expensive** to calculate the true training set mean and variance

Test time



- Test-time: replace μ_B and σ_B with μ and σ of the **training set**
- It's **expensive** to calculate the true training set mean and variance
- Therefore a moving average is common:

$$\tilde{\mu}^{(k)} \approx \alpha \tilde{\mu}^{(k-1)} + (1 - \alpha) \mu_B^{(k)}$$

$$\tilde{\sigma}^{(k)} \approx \alpha \tilde{\sigma}^{(k-1)} + (1 - \alpha) \sigma_B^{(k)}$$

Test time

- Test-time: replace μ_B and σ_B with μ and σ of the **training set**
- It's **expensive** to calculate the true training set mean and variance
- Therefore a **moving average** is common:

$$\tilde{\mu}^{(k)} \approx \alpha \tilde{\mu}^{(k-1)} + (1 - \alpha) \mu_B^{(k)}$$

$$\tilde{\sigma}^{(k)} \approx \alpha \tilde{\sigma}^{(k-1)} + (1 - \alpha) \sigma_B^{(k)}$$

- Moving average **decay** α (e.g. 0.8)

Backward pass

- Gradient **with respect to weights** is simply:

$$\frac{\partial L}{\partial \gamma} = \sum_{b=1}^B \frac{\partial L}{\partial \hat{\mathbf{Y}}_b} \tilde{\mathbf{X}}_b = \sum_{b=1}^B \mathbf{E}_b \tilde{\mathbf{X}}_b$$

- For the **bias** likewise we have:

$$\frac{\partial L}{\partial \beta} = \sum_{b=1}^B \frac{\partial L}{\partial \hat{\mathbf{Y}}_b} = \sum_{b=1}^B \mathbf{E}_b$$

Backward pass

The **gradient with respect to the input** is more complicated, but here it is:

error-tensor

$$\frac{\partial L}{\partial \tilde{\mathbf{X}}} = \frac{\partial L}{\partial \hat{\mathbf{Y}}} \odot \gamma$$

$$\frac{\partial L}{\partial \sigma_B^2} = \sum_{b=1}^B \frac{\partial L}{\partial \tilde{\mathbf{X}}_b} \odot (\mathbf{x}_b - \mu_B) \odot \frac{-1}{2} (\sigma_B^2 + \epsilon)^{\frac{-3}{2}}$$

$$\frac{\partial L}{\partial \mu_B} = \left(\sum_{b=1}^B \frac{\partial L}{\partial \tilde{\mathbf{X}}_b} \odot \frac{-1}{\sqrt{\sigma_B^2 + \epsilon}} \right) + \underbrace{\frac{\partial L}{\partial \sigma_B^2} \odot \frac{\sum_{b=1}^B -2(\mathbf{x}_b - \mu_B)}{B}}_0$$

$$\frac{\partial L}{\partial \mathbf{X}} = \frac{\partial L}{\partial \tilde{\mathbf{X}}} \odot \frac{1}{\sqrt{\sigma_B^2 + \epsilon}} + \frac{\partial L}{\partial \sigma_B^2} \odot \frac{2(\mathbf{X} - \mu_B)}{B} + \frac{\partial L}{\partial \mu_B} \odot \frac{1}{B}$$

Backward pass

- \odot denotes an element-wise multiplication. Always check the dimensionality of your matrices!

Backward pass

- \odot denotes an element-wise multiplication. Always check the dimensionality of your matrices!
- To make life easier, we will provide the code for the computation of the gradient with respect to the input:

Backward pass

- \odot denotes an element-wise multiplication. Always check the dimensionality of your matrices!
- To make life easier, we will provide the code for the computation of the gradient with respect to the input:
- `compute_bn_gradients`



Convolutional Batch Normalization

- In CNNs batch normalization is adjusted to work **similar to convolution**

Convolutional Batch Normalization

- In CNNs batch normalization is adjusted to work **similar to convolution**
- A scalar μ, σ is calculated for the H **channels**

Convolutional Batch Normalization

- In CNNs batch normalization is adjusted to work **similar to convolution**
- A scalar μ, σ is calculated for the H **channels**
- Implementation can be reused, by observing
 - that **spatial dimensions** M, N can be treated **like the batch dimension** B

Convolutional Batch Normalization

- In CNNs batch normalization is adjusted to work **similar to convolution**
- A scalar μ, σ is calculated for the H **channels**
- Implementation can be reused, by observing
 - that **spatial dimensions** M, N can be treated **like the batch dimension** B
 - we can **reshape** the $B \times H \times M \times N$ tensor to $B \times H \times M \cdot N$
 - because of our format we have to **transpose** from $B \times H \times M \cdot N$ to $B \times M \cdot N \times H$
 - and afterwards **reshape again** to have a $B \cdot M \cdot N \times H$ tensor

Convolutional Batch Normalization

- In CNNs batch normalization is adjusted to work **similar to convolution**
- A scalar μ, σ is calculated for the H **channels**
- Implementation can be reused, by observing
 - that **spatial dimensions** M, N can be treated **like the batch dimension** B
 - we can **reshape** the $B \times H \times M \times N$ tensor to $B \times H \times M \cdot N$
 - because of our format we have to **transpose** from $B \times H \times M \cdot N$ to $B \times M \cdot N \times H$
 - and afterwards **reshape again** to have a $B \cdot M \cdot N \times H$ tensor
- Consequently we have to **reverse this** before returning the **output**

Convolutional Batch Normalization

- In CNNs batch normalization is adjusted to work **similar to convolution**
- A scalar μ, σ is calculated for the **H channels**
- Implementation can be reused, by observing
 - that **spatial dimensions** M, N can be treated **like the batch dimension** B
 - we can **reshape** the $B \times \textcircled{H} \times M \times N$ tensor to $B \times H \times M \cdot N$
 - because of our format we have to **transpose** from $B \times H \times M \cdot N$ to $B \times M \cdot N \times H$
 - and afterwards **reshape again** to have a $B \cdot M \cdot N \times H$ tensor
- Consequently we have to **reverse this** before returning the **output**
- ... and do the **same** in the **backward pass**



FRIEDRICH-ALEXANDER-
UNIVERSITÄT
ERLANGEN-NÜRNBERG
SCHOOL OF ENGINEERING

LeNet (optional)



LeNet architecture

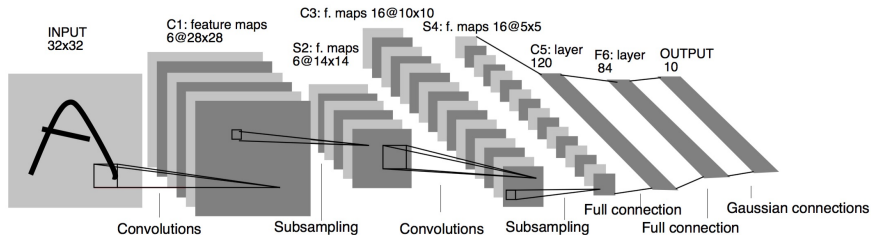


Figure: LeNet

Modified LeNet architecture

Deviations

- Input is 28×28
- Our conv only supports “same” padding - so C3 has **larger activation maps**
- Input to **C5** is also **larger**
- We only implemented ReLUs, so **no** TanH
- We also use the implemented SoftMax **instead of** RBF units

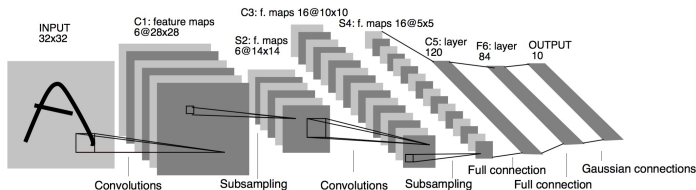


Figure: LeNet



Thanks for listening.
Any questions?