

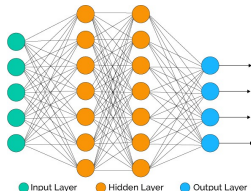
# Machine learning for particle physicists

## III. How to train better networks

Anja Butter

26th Vietnam School of Physics

# Recap - Neural networks basics



DNN:

$$\tilde{\mathbf{y}} = \sigma_n(\mathbf{W}_n \cdot \sigma_{n-1}(\mathbf{W}_{n-1} \cdot \dots \cdot \sigma_1(\mathbf{W}_1 \mathbf{X})))$$

Backpropagation:

$$\nabla_{\mathbf{w}_i} \mathcal{L} = \frac{\partial \mathcal{L}}{\partial \sigma_n} \frac{\partial \sigma_n}{\partial \text{lin}_n} \frac{\partial \text{lin}_n}{\partial \sigma_{n-1}} \cdots \frac{\partial \text{lin}_i}{\partial \mathbf{w}_i}$$

Mini-batch gradient descent: weight update based on batch of data

Overfitting: network learns noise

→ control: validation + test dataset

# Plan for today

How to train better networks

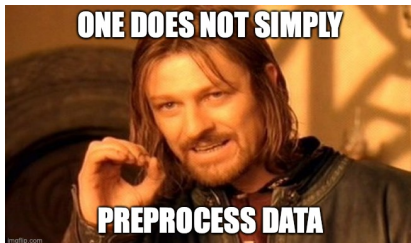
Practical tips to avoid frustration :)

- ① Data preprocessing
- ② Network initialization
- ③ Optimizing the training procedure
- ④ Regularization
- ⑤ Hyperparameter tuning

# Data Preprocessing

Why preprocessing?

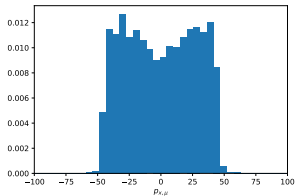
- input features with different scales  
eg.  $\text{jet} = (\text{charge}, n_{\text{particles}}, p_T, M, \eta, \phi)$
- large value with small spread  
eg.  $pp \rightarrow Z \rightarrow ll, m_{ll} \in [80 \text{ GeV} - 100 \text{ GeV}]$
- weights usually initialized to be sensitive in range  $[-1, +1]$
- classification output in range  $[0, 1]$
- training more efficient/stable if features are also in range  $[-1, +1]$



# Rescaling

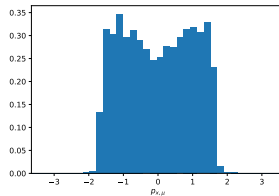
Example:  $pp \rightarrow Z \rightarrow \mu^+ \mu^-$

Rule of thumb: rescale to  $\mu = 0, \sigma = 1$



$$\frac{p_X - \bar{p}_X}{\sigma(p_X)}$$

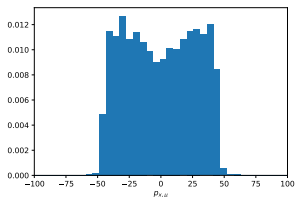
→



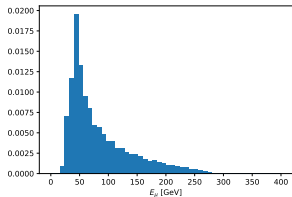
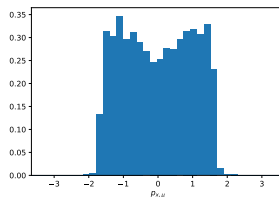
# Rescaling

Example:  $pp \rightarrow Z \rightarrow \mu^+ \mu^-$

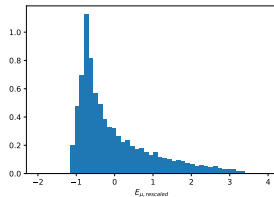
Rule of thumb: rescale to  $\mu = 0, \sigma = 1$



$$\frac{p_X - \bar{p}_X}{\sigma(p_X)}$$



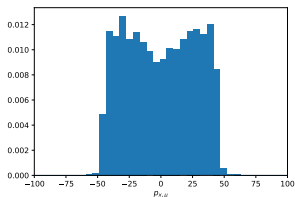
$$\frac{E - \bar{E}}{\sigma(E)}$$



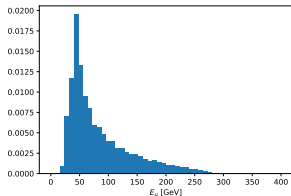
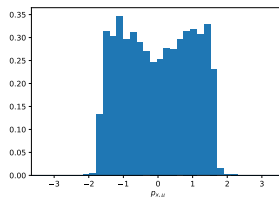
# Rescaling

Example:  $pp \rightarrow Z \rightarrow \mu^+ \mu^-$

Rule of thumb: rescale to  $\mu = 0, \sigma = 1$

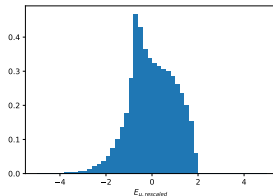


$$\frac{p_X - \bar{p}_X}{\sigma(p_X)}$$



$$\frac{E' - \bar{E}'}{\sigma(E')}$$

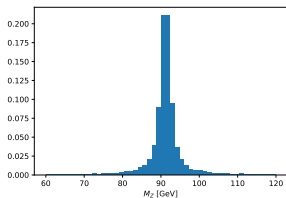
$E' = \log(E - 20)$



# Rescaling

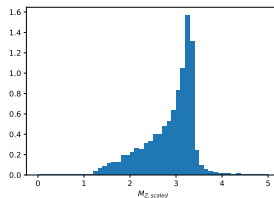
Example:  $pp \rightarrow Z \rightarrow \mu^+ \mu^-$

Exception: Correlated observables



$$\frac{p_{i,\mu} - \bar{p}_{i,\mu}}{\sigma(p_{i,\mu})}$$

→

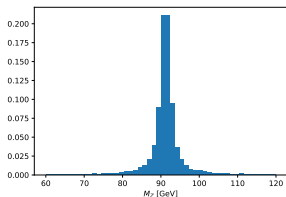




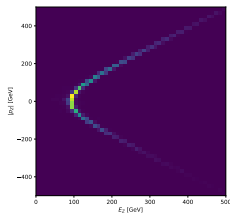
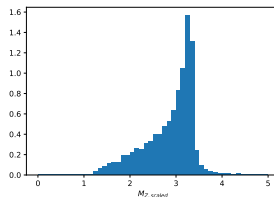
# Rescaling

Example:  $pp \rightarrow Z \rightarrow \mu^+ \mu^-$

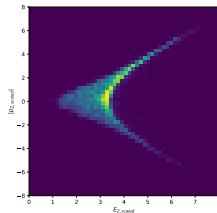
Exception: Correlated observables



$$\frac{p_{i,\mu} - \bar{p}_{i,\mu}}{\sigma(p_{i,\mu})}$$



$$\frac{p_{i,\mu} - \bar{p}_{i,\mu}}{\sigma(p_{i,\mu})}$$

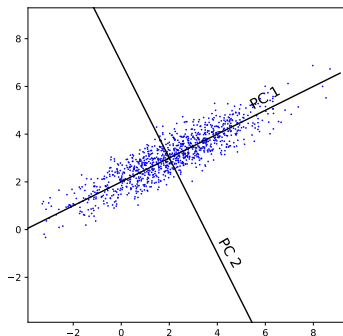


$\Rightarrow$  Use same scale for  $p_{i,\mu}$

# PCA

## Principal component analysis

- directions maximizing variance
  - eigenvector of covariance matrix
  - $\text{cov}(\mathbf{X}) = \mathbf{X}^T \mathbf{X}$
- 
- + facilitates training
  - + useful for interpretation
  - + can reduce data dimension



## ② Network initialization

# Network initialization

We know how to update weights. But how do we start?

①  $w_i = 1$ ?

# Network initialization

We know how to update weights. But how do we start?

①  $w_i = 1$ ?

symmetric initialization  $\Rightarrow$  symmetric updates  $\Rightarrow$  identical weights  $\nexists$

# Network initialization

We know how to update weights. But how do we start?

- ①  $w_i = 1?$   $\nexists$
- ②  $w_i \sim \mathcal{N}(\mu = 0, \sigma = 1)?$

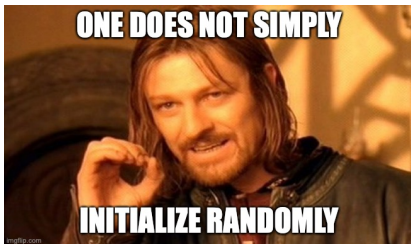
Check for single neuron  $y = w_i x_i$  with  $w_i, x_i$  independent:

$$\begin{aligned}\langle y^2 \rangle &= \sum_i \langle w_i^2 x_i^2 \rangle \\ &= \sum_i \langle w_i^2 \rangle \langle x_i^2 \rangle + \langle x_i^2 \rangle \langle w_i^2 \rangle + \langle w_i^2 \rangle \langle x_i^2 \rangle \\ &= \sum_i \langle w_i^2 \rangle \langle x_i^2 \rangle \quad \leftarrow \langle w_i \rangle = \langle x_i \rangle = 0 \\ &= n_{incoming} \langle w_i^2 \rangle \langle x_i^2 \rangle \quad \text{diverges!}\end{aligned}$$

# Network initialization

We know how to update weights. But how do we start?

- ①  $w_i = 1?$     ⚡
- ②  $w_i \sim \mathcal{N}(\mu = 0, \sigma = 1)?$     ⚡



$\rightarrow \langle w_i^2 \rangle = \frac{1}{n_{incoming}}$  to preserve variance through network

# Network initialization

We know how to update weights. But how do we start?

- ①  $w_i = 1$ ?  $\nexists$
- ②  $w_i \sim \mathcal{N}(\mu = 0, \sigma = 1)$ ?  $\nexists$
- ③ Xavier/Glorot initialization  $w_i \sim \mathcal{N}\left(\mu = 0, \sigma = \sqrt{2/(n_{in} + n_{out})}\right)$ 
  - caveat 1: Same argument for backpropagation  $\rightarrow$  average  $(n_{in} + n_{out})/2$
  - caveat 2: only for  $\approx$  linear activation function eg. tanh



# Network initialization

We know how to update weights. But how do we start?

- ①  $w_i = 1$ ?  $\nexists$
- ②  $w_i \sim \mathcal{N}(\mu = 0, \sigma = 1)$ ?  $\nexists$
- ③ Xavier/Glorot initialization  $w_i \sim \mathcal{N}\left(\mu = 0, \sigma = \sqrt{2/(n_{in} + n_{out})}\right)$
- ④ ReLU  $\rightarrow$  50% of outputs = 0  $\rightarrow$  additional factor 2  
 $\Rightarrow$  He initialization  $\sigma = \sqrt{2/n_{in}}$

# Network initialization

We know how to update weights. But how do we start?

- ①  $w_i = 1$ ?  $\nexists$
- ②  $w_i \sim \mathcal{N}(\mu = 0, \sigma = 1)$ ?  $\nexists$
- ③ Xavier/Glorot initialization  $w_i \sim \mathcal{N}\left(\mu = 0, \sigma = \sqrt{2/(n_{in} + n_{out})}\right)$
- ④ ReLU  $\rightarrow$  50% of outputs = 0  $\rightarrow$  additional factor 2  
 $\Rightarrow$  He initialization  $\sigma = \sqrt{2/n_{in}}$
- ⑤ Glorot & He initialization also available for uniform distributions

# Pretraining

For some tasks we can use pretrained networks

→ trained on large dataset to extract image features

Google's InceptionResNetV2 to identify anomalies in QCD jet images

Best image class to identify QCD jet images:

# Pretraining

For some tasks we can use pretrained networks

→ trained on large dataset to extract image features

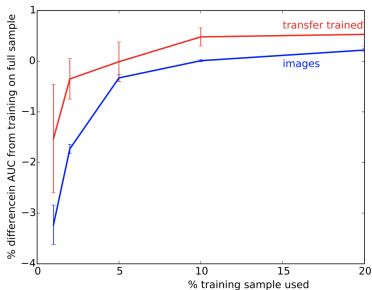
Google's InceptionResNetV2 to identify anomalies in QCD jet images

Best image class to identify QCD jet images:



Ice cream classification :)

Lisa Benato

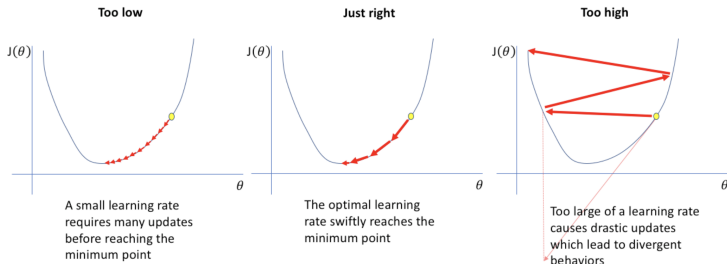


### ③ Optimizing the training procedure

# Optimizing the training procedure

## Reminder

### Convergence depends on learning rate



<https://www.jeremyjordan.me/nn-learning-rate/>

→ Experiment with different orders of magnitude eg.  $10^{-1} \dots 10^{-6}$

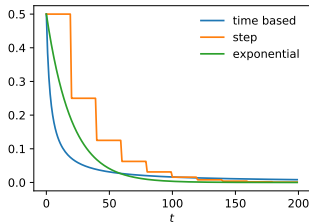
# Learn rate decay

Reduce learning rate over time to improve convergence

Time-Based Decay  $l(t) = \frac{l_0}{1 + k * t}$

Step Decay  $l(t) = l_0 * \lambda^{\text{int}(t/\tau)}$  with  $0 < \lambda < 1$

Exponential Decay  $l(t) = l_0 * e^{-t/\tau}$



# Momentum

Problem: One dimension much steeper than the other

gradient descent

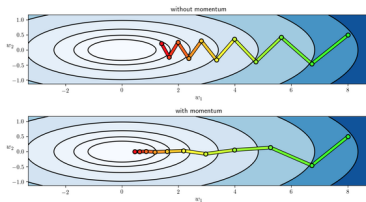
$$\mathbf{W}_t \rightarrow \mathbf{W}_{t+1} = \mathbf{W}_t - \alpha \nabla_{\mathbf{W}_t} \mathcal{L}$$

GD + momentum

$$\mathbf{W}_t \rightarrow \mathbf{W}_{t+1} = \mathbf{W}_t - \alpha v_{dw}$$

$$v_{dw} = \beta v_{dw} + (1 - \beta) \nabla_{\mathbf{W}_t} \mathcal{L}$$

Intuition: ball picks up momentum



[jermwatt.github.io/machine\\_learning\\_refined](https://jermwatt.github.io/machine_learning_refined)

enforces dimensions where gradient points in same direction  
+ reduces oscillation



# Adagrad/RMSprop

Adapt updates to individual parameters

$$\mathbf{W}_t \rightarrow \mathbf{W}_{t+1} = \mathbf{W}_t - \alpha \frac{1}{\sqrt{G} + \epsilon} \nabla_{\mathbf{W}_t} \mathcal{L}$$

→ Different learning rate for each parameter

Adagrad:  $G_{ii,t} = \sum_{t'=0}^t dw_{i,t'}^2$

sum over vector of all past gradients

→ monotonically decreasing learning rate

RMSprop:  $G_{ii,t} = \beta G_{ii,t-1} + (1 - \beta) dw_{i,t}^2$

→  $\beta = 0.9$  → decaying average

# Adam

Adaptive moment estimation

Standard go to option, stable & fast

Combines first moment (momentum) and second moment (RMSprop)

$$\begin{aligned}\mathbf{W}_t &\rightarrow \mathbf{W}_{t+1} = \mathbf{W}_t - \alpha \frac{1}{\sqrt{G} + \epsilon} v_{dw} \\ v_{dw,t} &= \frac{1}{1 - \beta_1} (\beta_1 v_{dw,t-1} + (1 - \beta_1) dw_{i,t}) \\ G_{ii,t} &= \frac{1}{1 - \beta_2} (\beta_2 G_{ii,t-1} + (1 - \beta_2) dw_{i,t}^2)\end{aligned}$$

Others worth exploring!

Might fit your problem better?

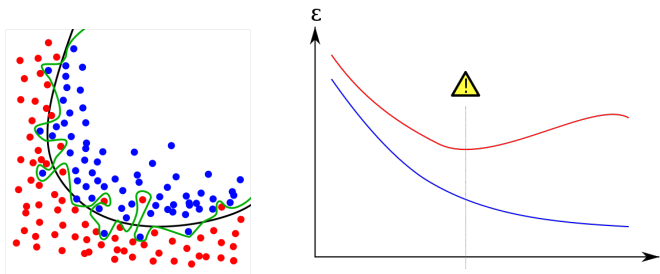
- Nesterov accelerated gradient
- Adadelta
- AMSGrad



## 4 Regularization

## Reminder: Overtraining

Overtraining: networks picks up irrelevant features



Control with validation/test data

We refer to this **regularization** technique as **early stopping**  
not always applicable  $\rightarrow$  consider alternatives

# Regularization

## Modify network

- Dropout
- Batch normalization

## Modify loss

- $l_1$  regularization
- $l_2$  regularization
- gradient penalty

# Modifying the loss

The network is constrained by punishing large weight values

$$\mathcal{L} = \mathcal{L}(\mathbf{y}, \tilde{\mathbf{y}}) + \alpha \Omega(\mathbf{W})$$

$l_1$

$$\Omega(\mathbf{W}) = \|\mathbf{W}\|_1 = \sum_{ij} |W_{ij}|$$

$l_2$

$$\Omega(\mathbf{W}) = \|\mathbf{W}\|_2^2 = \sum_{ij} W_{ij}^2$$

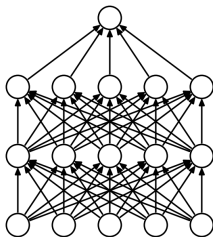
gradient penalty

$$\Omega(\mathbf{W}) \sim \|\nabla_{\mathbf{x}} \tilde{\mathbf{y}}\|_2^2$$

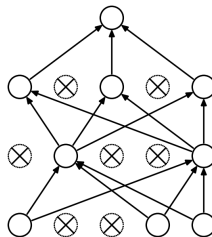
# Dropout

Randomly switching off nodes during training

Intuition: Train many different models, then average for evaluation



(a) Standard Neural Net



(b) After applying dropout.

2014, N. Srivastava et. al

# Batch normalization

Idea: fix mean and variance of layer output

$$\mathbf{x} \rightarrow \mathbf{x}_{norm} = \frac{\mathbf{x} - \bar{\mathbf{x}}}{\sqrt{\sigma(\mathbf{x}) + \epsilon}}$$
$$\mathbf{y} = \gamma \mathbf{x}' + \beta$$

trainable parameters  $\gamma, \beta$

During training: normalization per batch

For inference: normalization from full dataset

Why it works is subject of current research!

Smoothness of optimization landscape? Length-Direction decoupling?



## 5 Hyperparameter tuning

# Hyperparameter tuning

How can we find the best settings for the training?

Problem: We can not compute a gradient!

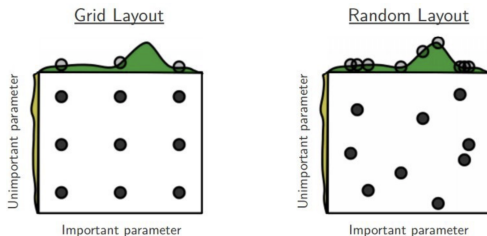
# Hyperparameter tuning

How can we find the best settings for the training?

Problem: We can not compute a gradient!

- ① by hand → underrated, helps to build experience
- ② Grid search
- ③ Random (blind)
- ④ Bayesian optimization (educated guess, advanced)

# Advantage of random vs grid search



Advantages: easy to code, run parallel

Disadvantage: no use of information from previous iterations, curse of dimensionality

# All the things you can do to your ML setup

- ① Data preprocessing
  - Rescaling, PCA
- ② Network initialization
  - Glorot/HE, Normal/uniform
- ③ Optimizing the training procedure
  - Learning rate scheduling, momentum, Adagrad, Adam
- ④ Regularization
  - Via early stopping, additional loss, dropout, or Batchnorm
- ⑤ Hyperparameter tuning
  - get a feeling for the network, random search, Bayesian optimization

## Ready to try it out?

- colab
- gitHub
- dhrou
- HEPMLtutorials
- HEPML\_HandsOn\_NN.ipynb

Big thank you to David Rousseau for sharing this tutorial!

## Corrections

```
D = Model(inputs=inputs, outputs=Dx)
class_weight = {
0: class_weights[0],
1: class_weights[1],
}
D.fit(
X_train,
y_train.values,
epochs=10,
verbose=0,
class_weight=class_weight
)
```