



# Artificial Neural Networks

[2500WETANN]

José Oramas



# Learning & Optimization

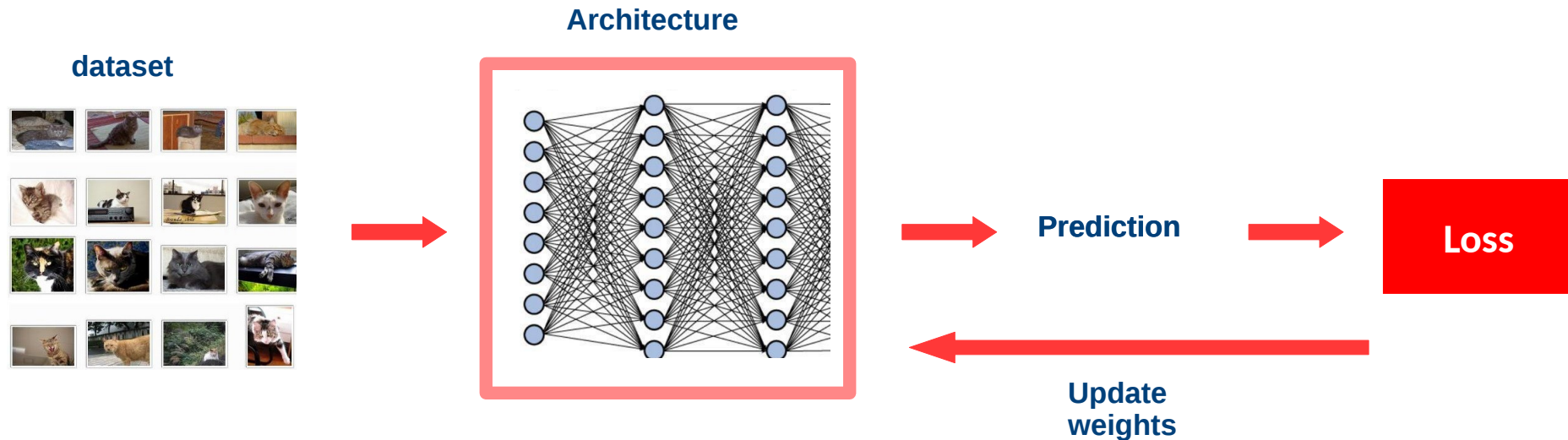
[ ... for Deep Neural Networks ]

José Oramas

# Learning and Optimization

## Training Pipeline

- Techniques Applicable at Different Stages



# Prior-Training

[ The Calm Before the Storm ]



# Data Augmentation

## What?

- Apply a set of operations on a given data sample to produce additional samples

## Benefits

- Increase training data
- Introduce variability



Original Image

Cropped samples



Mirrored samples

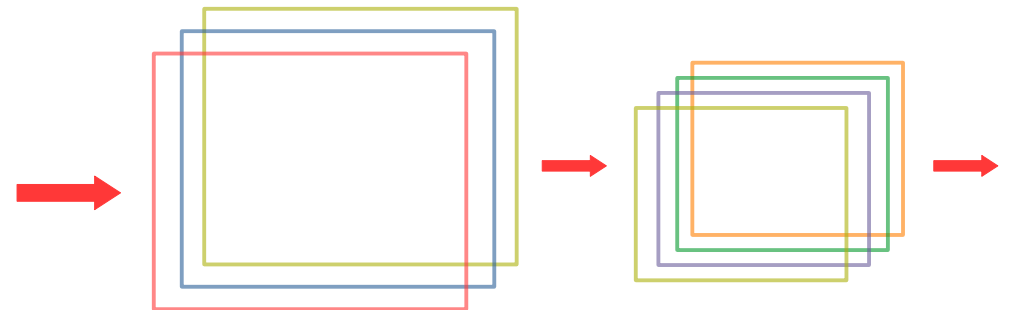
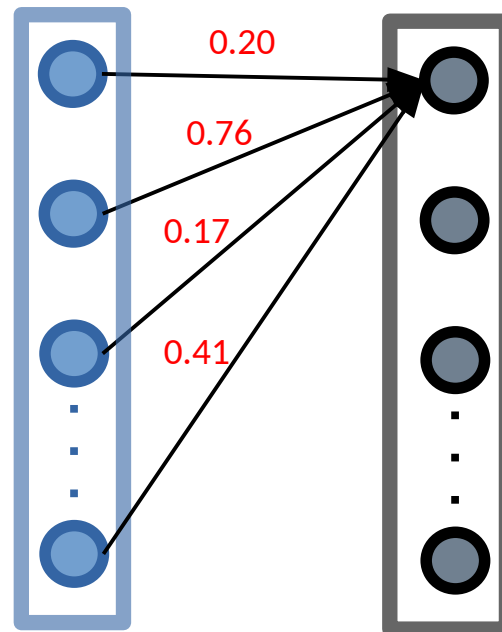


**Q: What happens at test time?**

# Input Normalization

## Common Practice

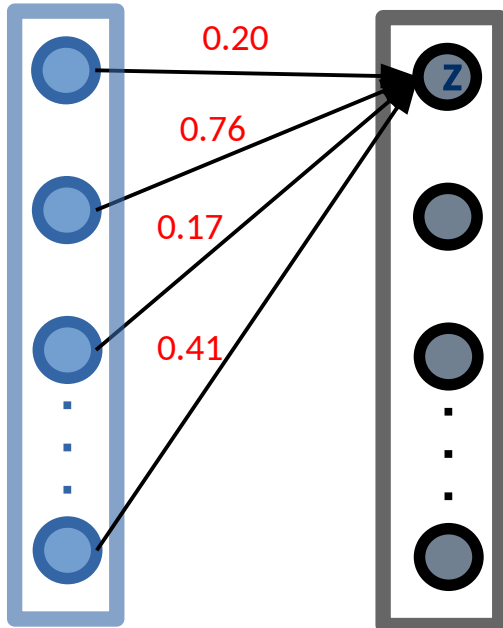
- Remove the “mean image”
- Standardize the inputs



# Initializing Weights of the Network

## Common Practice

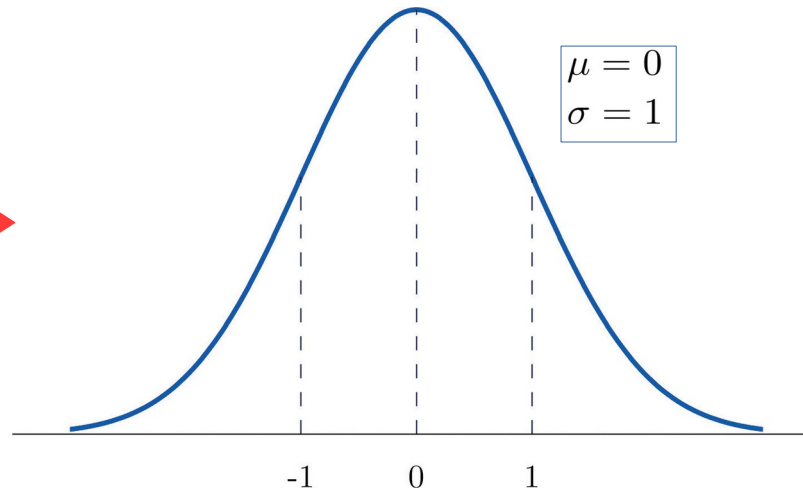
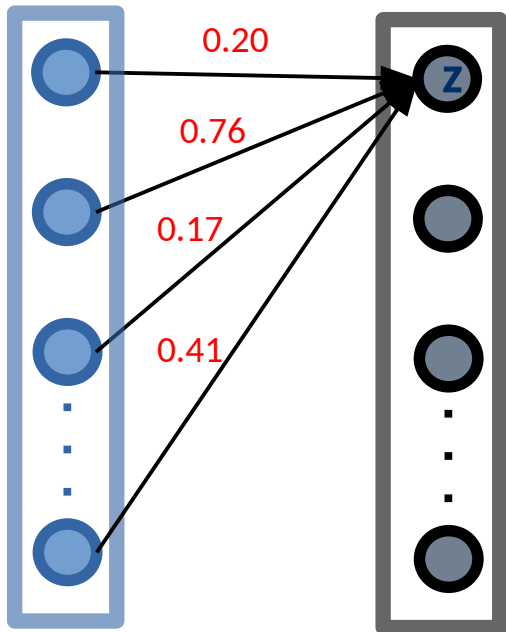
- Random Initialization
- Ensure the weights have a known mean and variance



# Initializing Weights of the Network

## Common Practice

- Random Initialization
- Ensure the weights have a known mean and variance

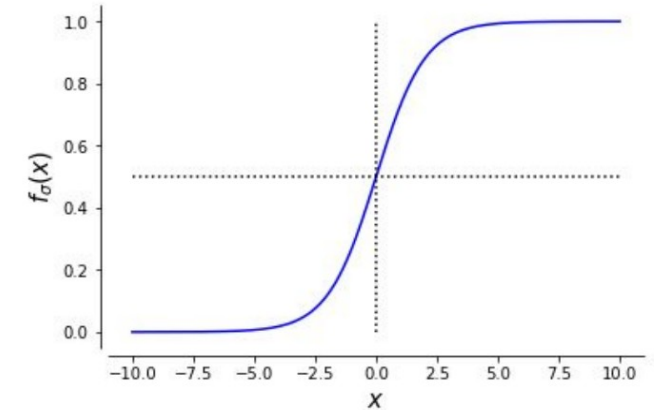
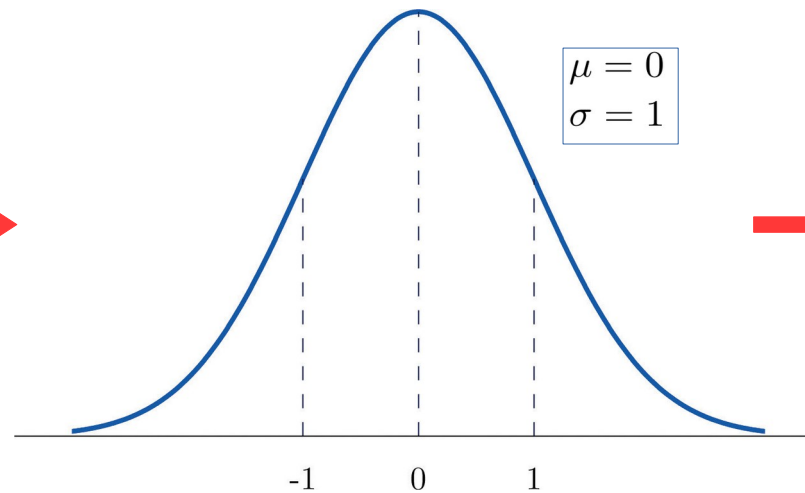
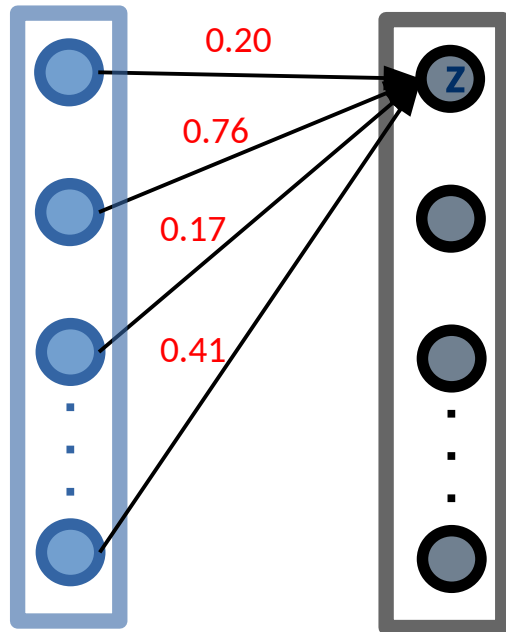




# Initializing Weights of the Network

## Common Practice

- Random Initialization
- Ensure the weights have a known mean and variance

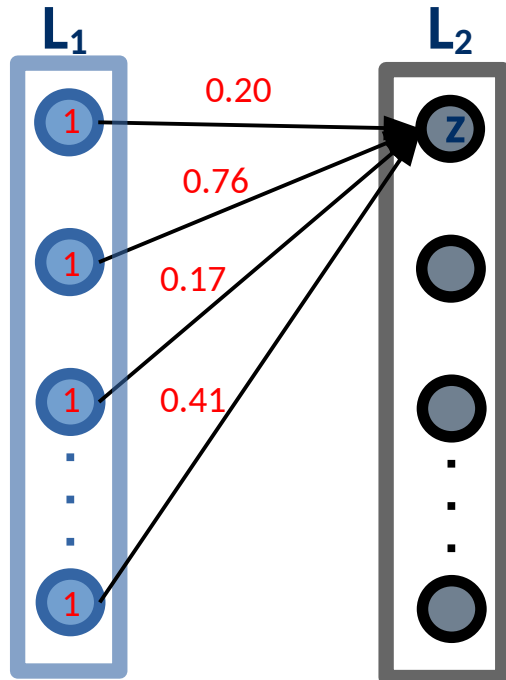


$$f_{\sigma}(x) = \frac{1}{1+e^{-x}} = \frac{e^x}{e^x+1}$$

# Initializing Weights of the Network

## Common Practice – What Happens in the Next Layer?

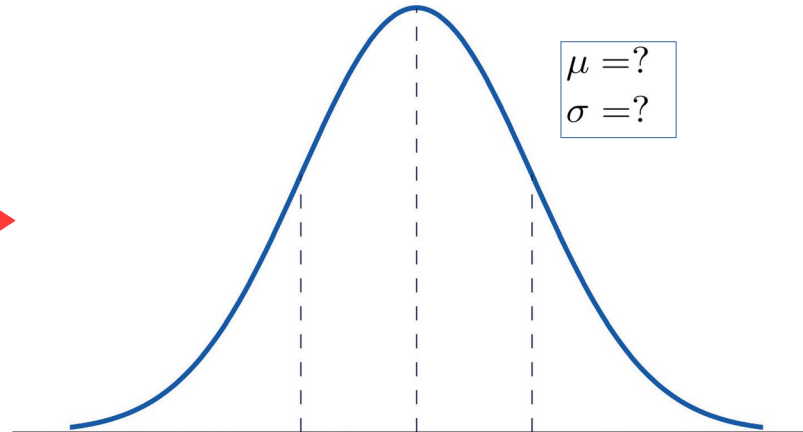
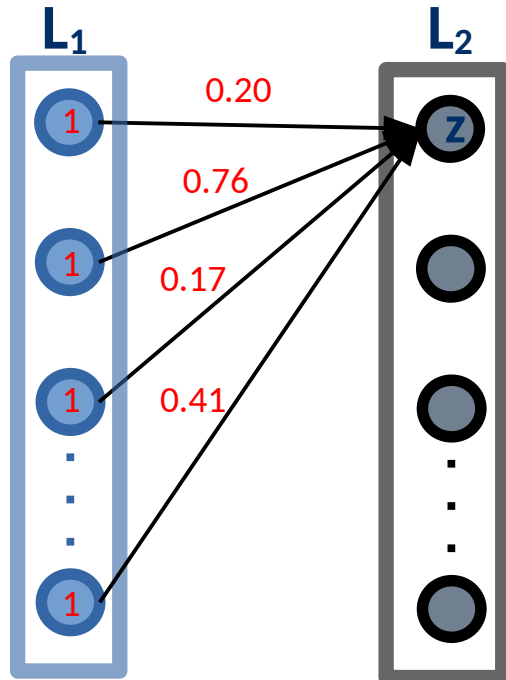
- Random Initialization
- Let's consider a simple setting (n=4 inputs, all set to 1)



# Initializing Weights of the Network

## Common Practice – What Happens in the Next Layer?

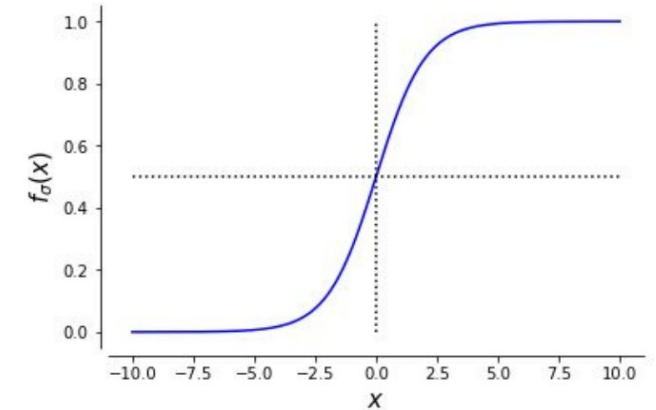
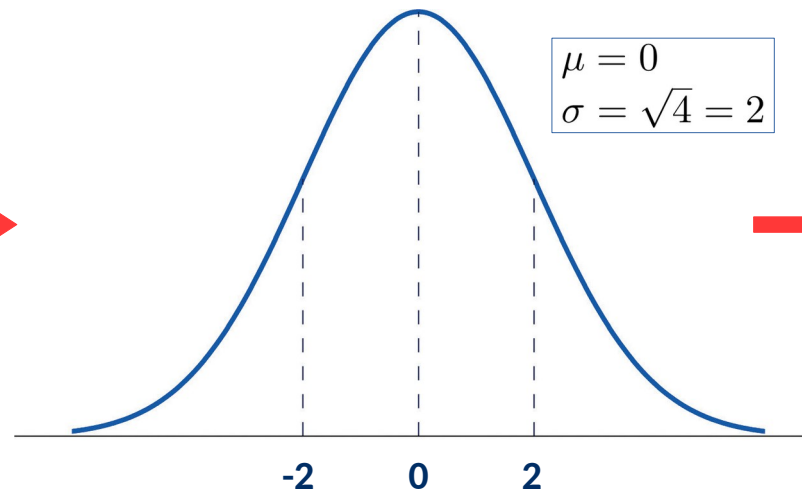
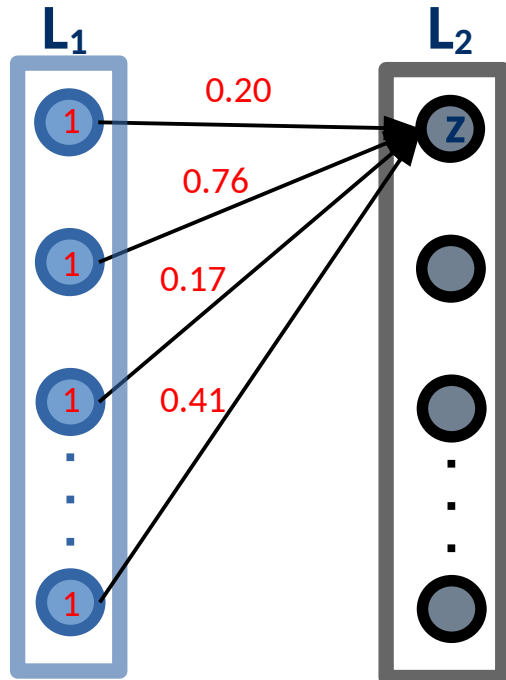
- Random Initialization
- Let's consider a simple setting ( $n=4$  inputs, all set to 1)



# Initializing Weights of the Network

## Common Practice – What Happens in the Next Layer?

- Random Initialization
- Let's consider a simple setting (n=4 inputs, all set to 1)

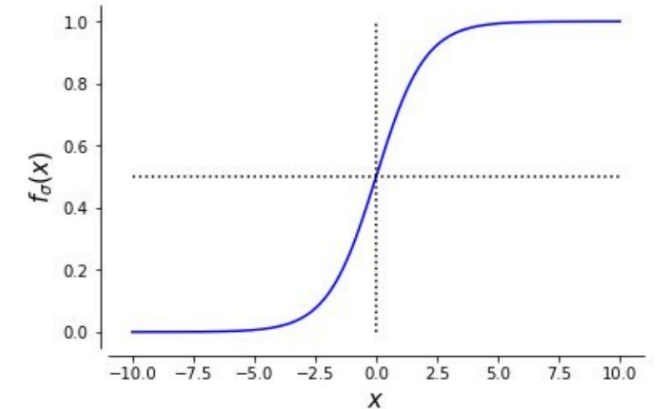
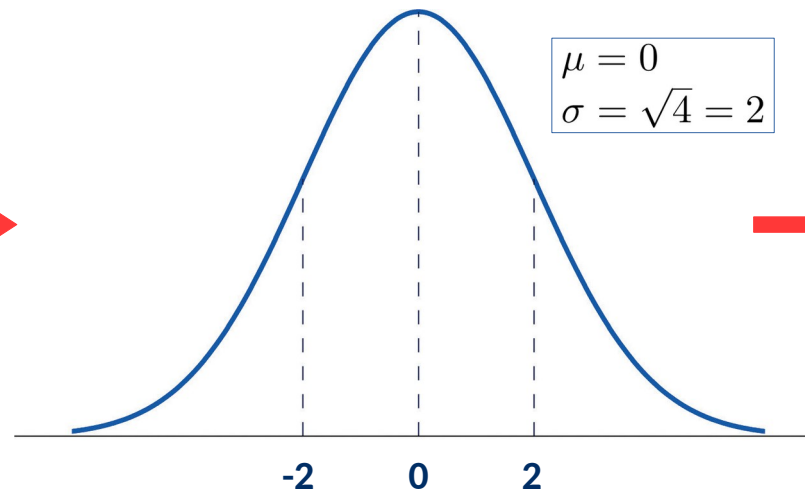
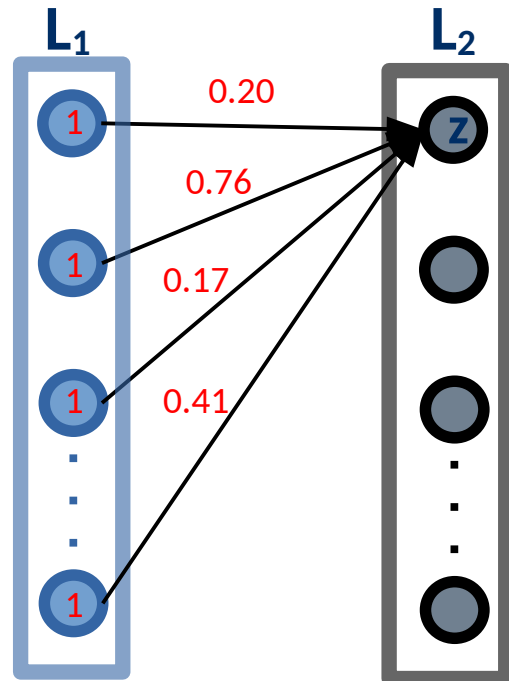


$$f_\sigma(x) = \frac{1}{1+e^{-x}} = \frac{e^x}{e^x+1}$$

# Initializing Weights of the Network

## Common Practice – What Happens in the Next Layer?

- Random Initialization
- Let's consider a simple setting (n=4 inputs, all set to 1)



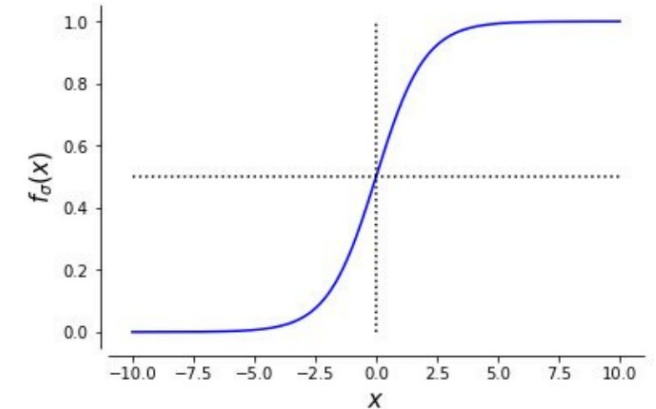
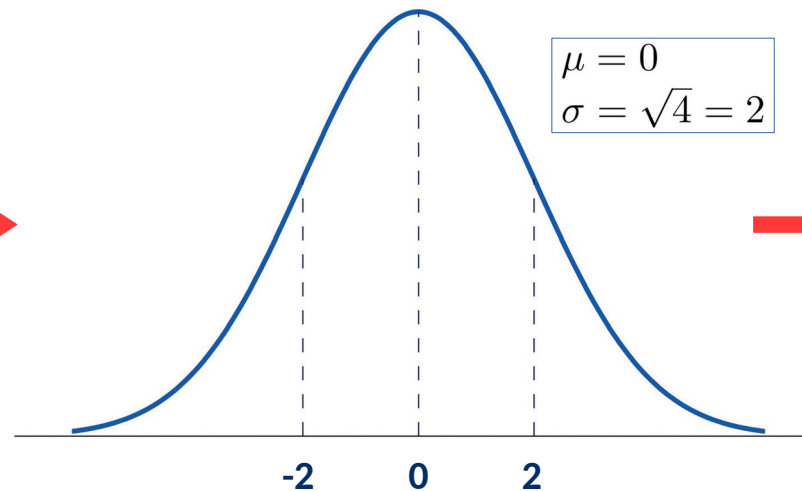
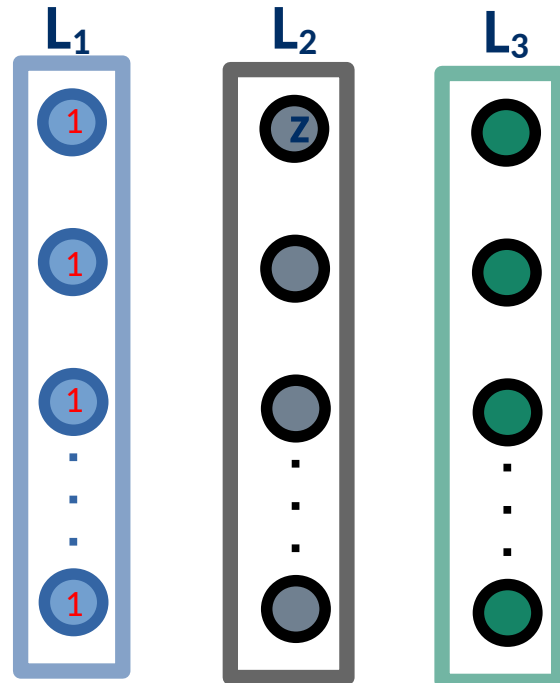
$$f_\sigma(x) = \frac{1}{1+e^{-x}} = \frac{e^x}{e^x+1}$$

Q1: What would happen if my architecture is wider?  
Q2: What would happen if my architecture is deeper?

# Initializing Weights of the Network

## Common Practice – What Happens in the Next Layer?

- Random Initialization
- Let's consider a simple setting (n=4 inputs, all set to 1)



$$f_{\sigma}(x) = \frac{1}{1+e^{-x}} = \frac{e^x}{e^x+1}$$

Q1: What would happen if my architecture is wider?  
Q2: What would happen if my architecture is deeper?



# Initializing Weights of the Network

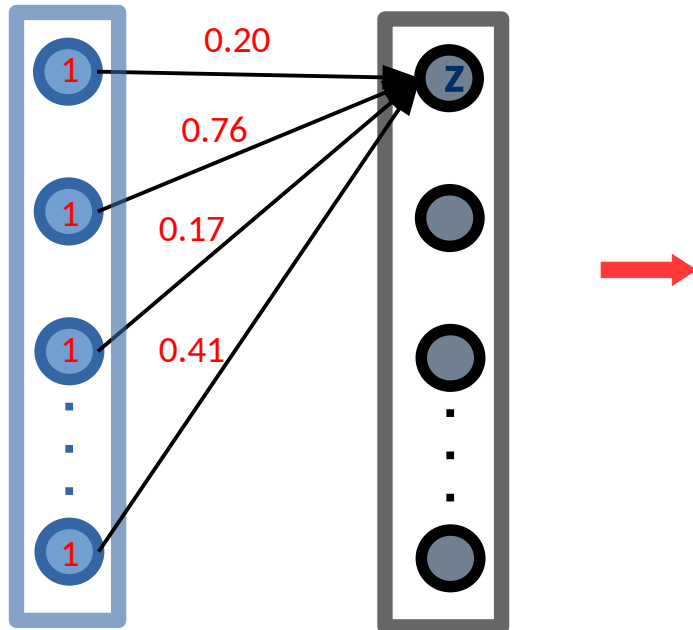
## Solution – Xavier/Glorot Initialization

- **Problem:** large variance –  $\text{var}(z)$
- **Solution:** let's make it smaller  $\rightarrow \text{var}(z) = 1/n$

# Initializing Weights of the Network

## Solution – Xavier/Glorot Initialization

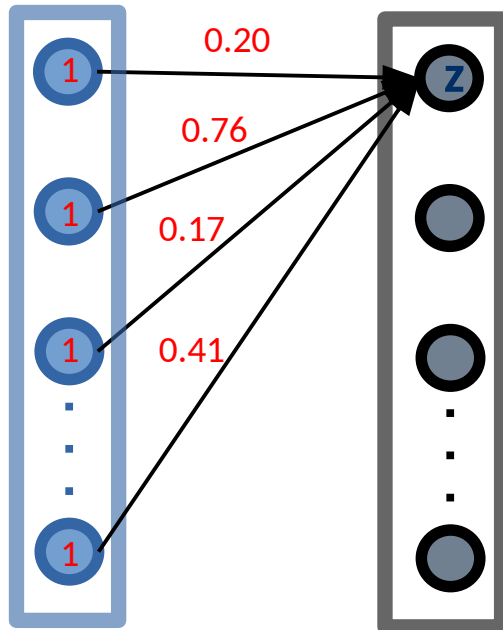
- **Problem:** large variance –  $\text{var}(z)$
- **Solution:** let's make it smaller  $\rightarrow \text{var}(z) = 1/n$



# Initializing Weights of the Network

## Solution – Xavier/Glorot Initialization

- **Problem:** large variance –  $\text{var}(z)$
- **Solution:** let's make it smaller  $\rightarrow \text{var}(z) = 1/n$



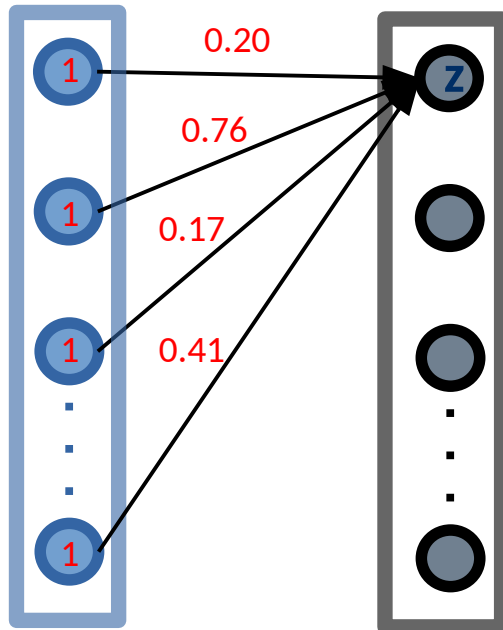
How to do it?

$$w_i = w_i \times \sqrt{1/n}$$

# Initializing Weights of the Network

## Solution – Xavier/Glorot Initialization

- **Problem:** large variance –  $\text{var}(z)$
- **Solution:** let's make it smaller  $\rightarrow \text{var}(z) = 1/n$



How to do it?

$$w_i = w_i \times \sqrt{1/n}$$

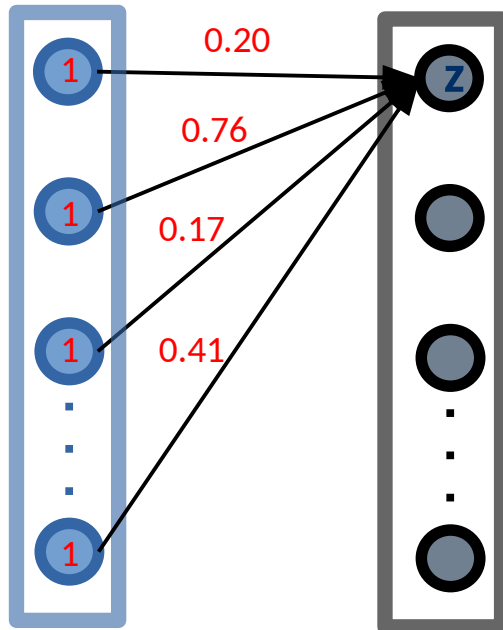
For ReLU layers

$$w_i = w_i \times \sqrt{2/n}$$

# Initializing Weights of the Network

## Solution – Xavier/Glorot Initialization

- **Problem:** large variance –  $\text{var}(z)$
- **Solution:** let's make it smaller  $\rightarrow \text{var}(z) = 1/n$



How to do it?

$$w_i = w_i \times \sqrt{1/n}$$

For ReLU layers

$$w_i = w_i \times \sqrt{2/n}$$

Previously (for reference)

$$w_i = w_i \times \sqrt{\frac{2}{n_{in} + n_{out}}}$$

# During Training

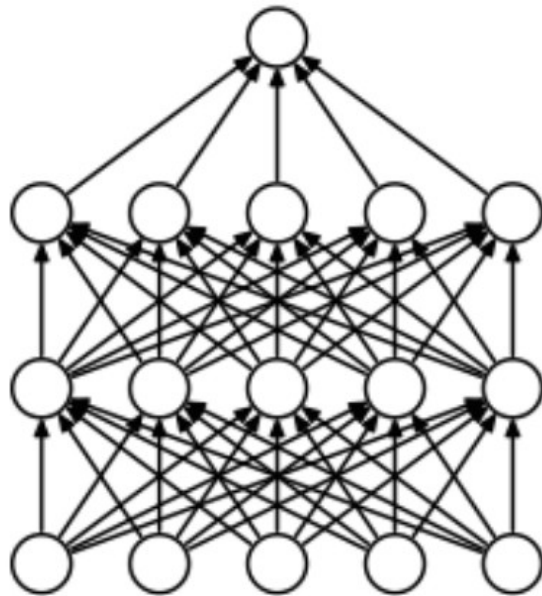
[ While the Computer is Hard at Work ]



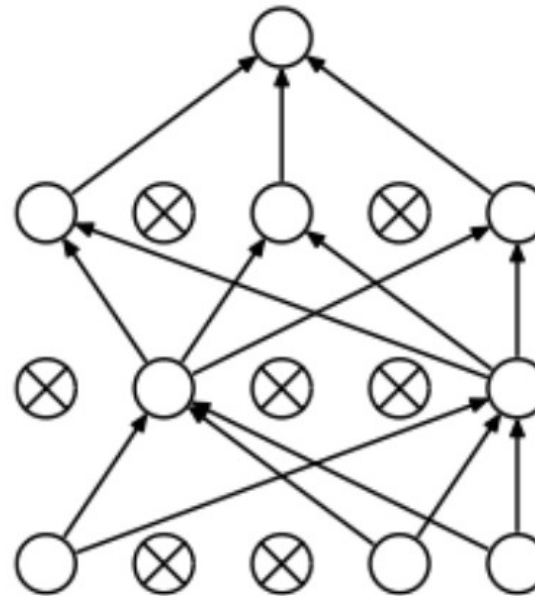
# Updating Weights of the Network

## Dropout

- **Problem:** decrease dependence of a given feature
- **Solution:** randomly deactivate neurons



Standard Neural Net



After applying dropout.

## Benefits

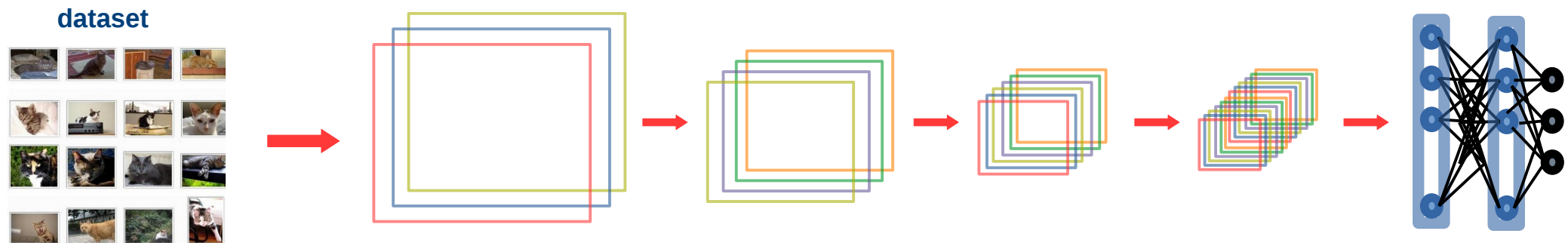
Avoid over-fitting  
Promote ensemble learning

**Q: What happens at test time?**

# Updating Weights of the Network

## During Training

- **Problem:** Updates on weights at a later layers should take into account changes at earlier layers (covariance shift)
  - Introduces changes in the distribution of internal activations
  - Requires careful initialization and a small learning rate



**Desiderata: Distribution of inputs ( $x$ ) should be fixed over time**

# Updating Weights of the Network

## Batch Normalization

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

Parameters to be learned:  $\gamma, \beta$

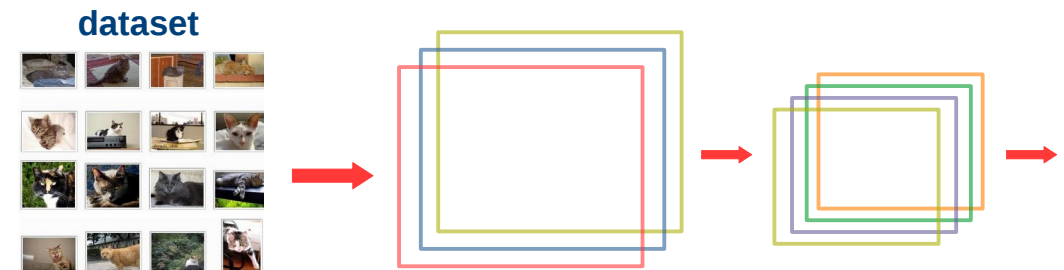
**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

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

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

- **Solution:** Normalize internal activations by considering dataset statistics

- Stochastic optimization  
→ batch-level statistics



# Updating Weights of the Network

## Batch Normalization

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

Parameters to be learned:  $\gamma, \beta$

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

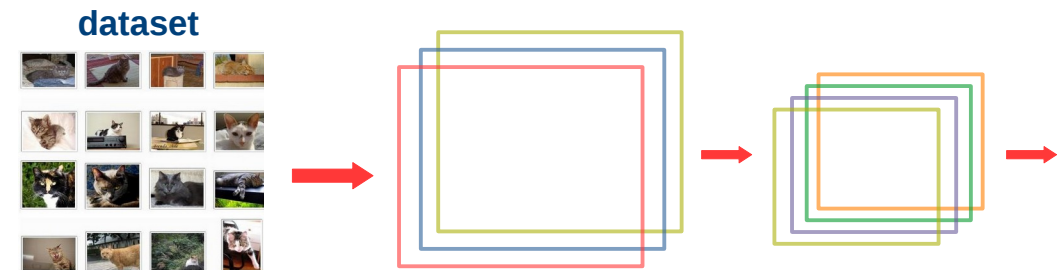
$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

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

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

- **Solution:** Normalize internal activations by considering dataset statistics

- Stochastic optimization  
→ batch-level statistics



# Updating Weights of the Network

## Batch Normalization

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

Parameters to be learned:  $\gamma, \beta$

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

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

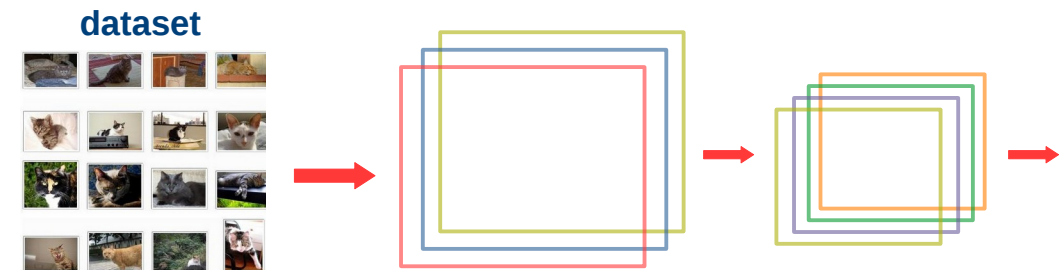
$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

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

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

- **Solution:** Normalize internal activations by considering dataset statistics

- Stochastic optimization  
→ batch-level statistics



# Updating Weights of the Network

## Batch Normalization

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

Parameters to be learned:  $\gamma, \beta$

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

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

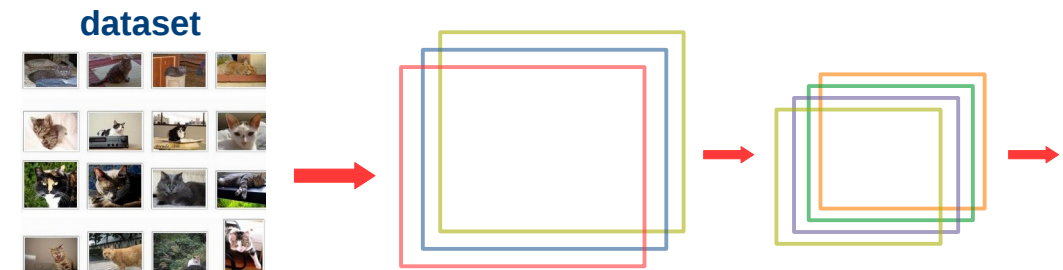
$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

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

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

## Benefits

- Less sensitivity to initialization
- Allows using larger learning rates ( faster training )





# Computing the Loss

[ While Checking How Well It Works ]

# Break

[ See you in 15 mins. ]

# Computing the Loss

[ While Checking How Well It Works ]

# Training with Weighted Examples

Let's revisit the computation of the gradient of the loss

- Gradient Descend

$$\theta_{t+1} = \theta_t - \alpha_t \nabla_{\theta} L(\theta_t)$$

where,

$$\begin{aligned} \nabla_{\theta} L(\theta_t) &= \nabla_{\theta} \sum_i l(f(x^{(i)}, \theta_t), y^{(i)}) \\ &= \sum_i \nabla_{\theta} l(f(x^{(i)}, \theta_t), y^{(i)}) \end{aligned}$$

# Training with Weighted Examples

Let's revisit the computation of the gradient of the loss

- Gradient Descend

$$\theta_{t+1} = \theta_t - \alpha_t \nabla_{\theta} L(\theta_t)$$

where,

$$\begin{aligned} \nabla_{\theta} L(\theta_t) &= \nabla_{\theta} \sum_i l(f(x^{(i)}, \theta_t), y^{(i)}) \\ &= \sum_i \nabla_{\theta} l(f(x^{(i)}, \theta_t), y^{(i)}) \end{aligned}$$

Q: Any potential problem/weakness?

# Training with Weighted Examples

Let's revisit the computation of the gradient of the loss

- Gradient Descend

$$\theta_{t+1} = \theta_t - \alpha_t \nabla_{\theta} L(\theta_t)$$

where,

$$\begin{aligned} \nabla_{\theta} L(\theta_t) &= \nabla_{\theta} \sum_i l(f(x^{(i)}, \theta_t), y^{(i)}) \\ &= \sum_i \nabla_{\theta} l(f(x^{(i)}, \theta_t), y^{(i)}) \end{aligned}$$

Q: Any potential problem/weakness?

Q: What would happen if the  $x_i$  examples are all from the same class ?



# Training with Weighted Examples

Let's revisit the computation of the gradient of the loss

- Gradient Descend

$$\theta_{t+1} = \theta_t - \alpha_t \nabla_{\theta} L(\theta_t)$$

where,

$$\begin{aligned} \nabla_{\theta} L(\theta_t) &= \nabla_{\theta} \sum_i l(f(x^{(i)}, \theta_t), y^{(i)}) \\ &= \sum_i \nabla_{\theta} l(f(x^{(i)}, \theta_t), y^{(i)}) \end{aligned}$$

re-weighting,

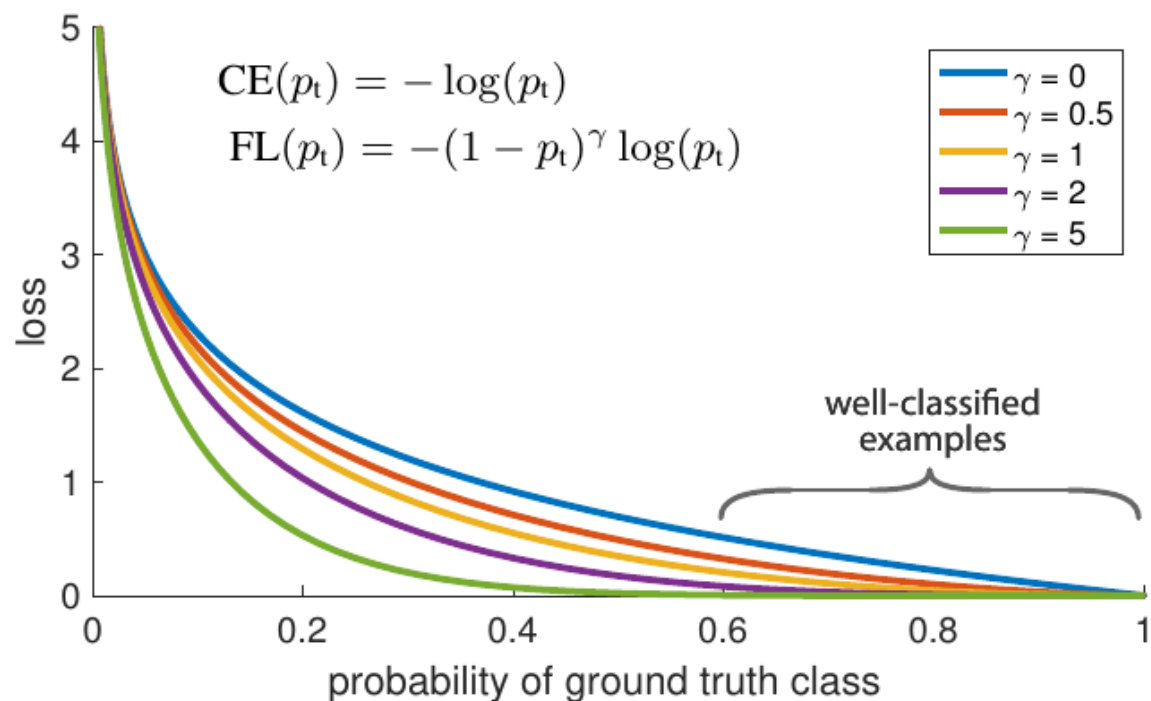
$$= \sum_i \frac{1}{Z(y^{(i)})} \nabla_{\theta} l(f(x^{(i)}, \theta_t), y^{(i)})$$

Q: Any potential problem/weakness?

Q: What would happen if the  $x_i$  examples are all from the same class ?

# Training with Examples of Different Complexity

## Focal Loss



$$FL(p_t) = -(1 - p_t)^\gamma \log(p_t)$$

## What it does?

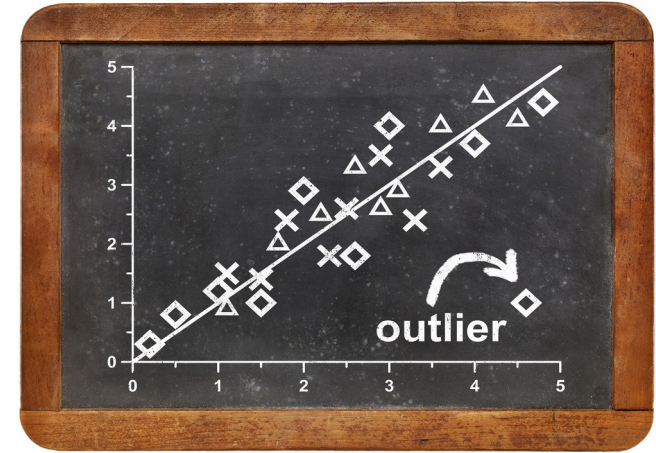
- Down-weights the loss from well-classified examples
- Focusses training on sparse set of hard examples

# Training with Examples of Different Complexity

## Focal Loss

Where it could be useful?

- Dense predictions tasks
- In the presence of outliers



# Learning Representations by Comparison

## Triplet Loss

- Given three examples (Anchor, Positive, Negative)
- Learn a representation that  $\text{distance}(\text{positive}, \text{anchor}) < \text{distance}(\text{negative}, \text{anchor})$



# Learning Representations by Comparison

## Triplet Loss

- Given three examples (Anchor, Positive, Negative)
- Learn a representation that  $\text{distance}(\text{positive}, \text{anchor}) < \text{distance}(\text{negative}, \text{anchor})$



Definition

$$L(A, P, N) = \max( d(A, P) - d(A, N) + \alpha, 0 )$$

# Learning Representations by Comparison

## Triplet Loss

- Given three examples (Anchor, Positive, Negative)
- Learn a representation that  $\text{distance}(\text{positive}, \text{anchor}) < \text{distance}(\text{negative}, \text{anchor})$



### Definition

$$L(A, P, N) = \max( d(A, P) - d(A, N) + \alpha, 0 )$$

### Using the L2-distance

$$L(A, P, N) = \max( \|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \alpha, 0 )$$

# Learning Representations by Comparison

## Triplet Loss

- Given three examples (Anchor, Positive, Negative)
- Learn a representation that  $\text{distance}(\text{positive}, \text{anchor}) < \text{distance}(\text{negative}, \text{anchor})$



### Definition

$$L(A, P, N) = \max( d(A, P) - d(A, N) + \alpha, 0 )$$

### Using the L2-distance

$$L(A, P, N) = \max( \|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \alpha, 0 )$$

Using in a cost function  $= \sum_{i=1}^M L( A^{(i)}, P^{(i)}, N^{(i)} )$

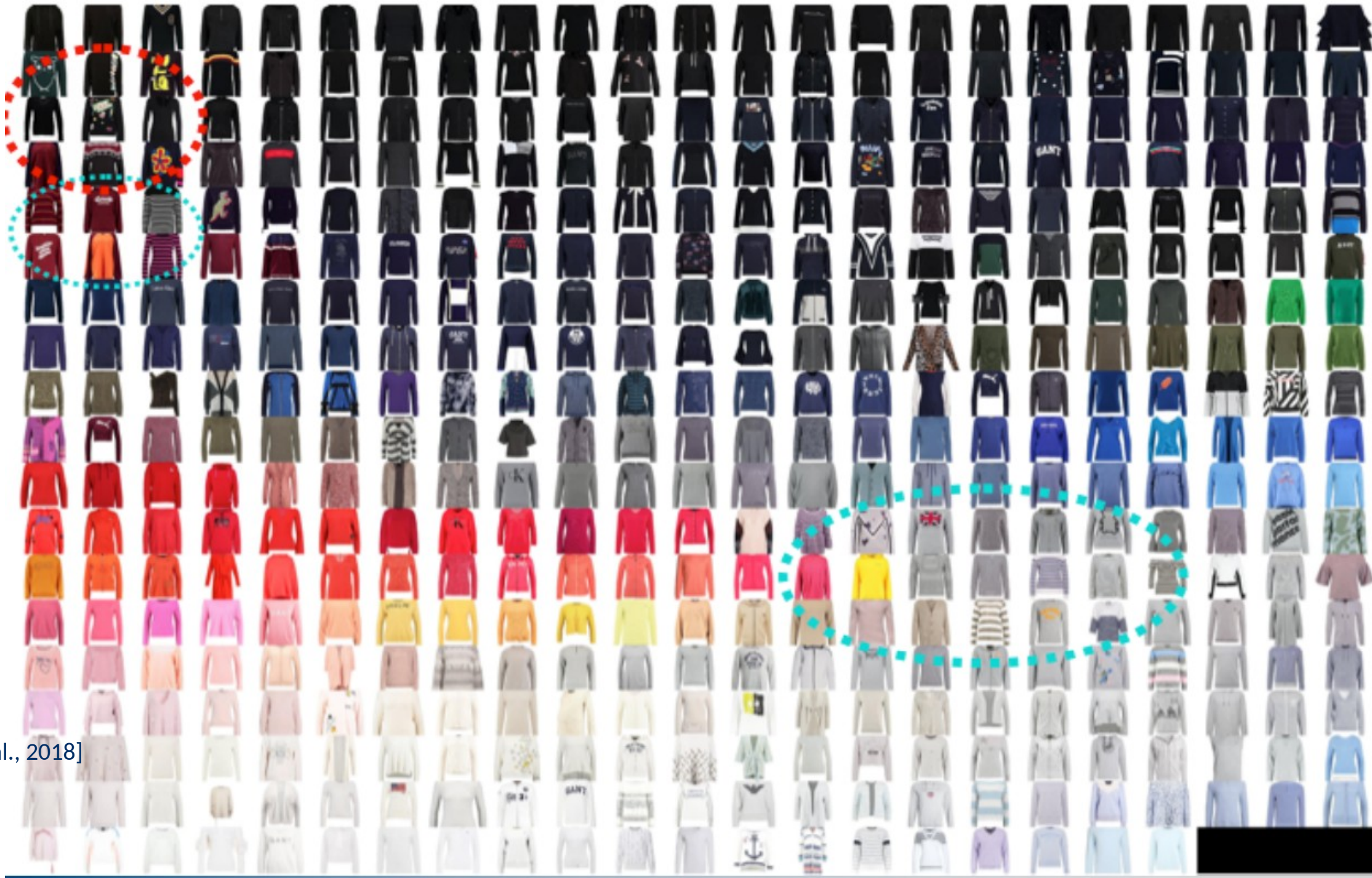


# Triplet Loss



[Wang et al., 2018]



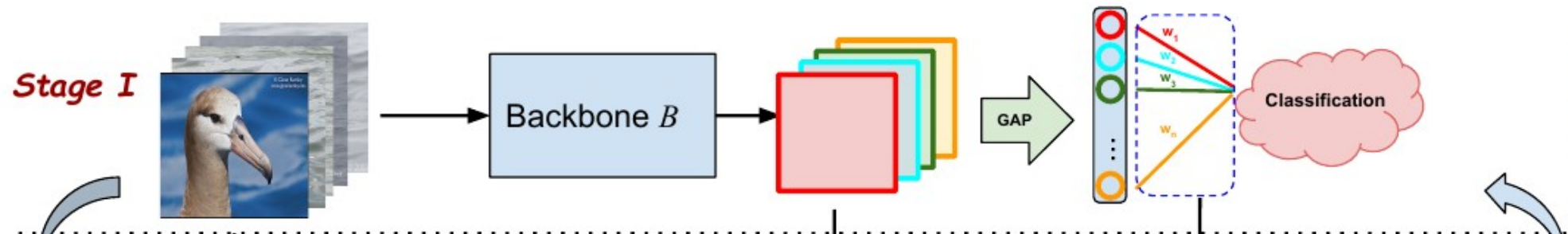


[Wang et al., 2018]

# Using Multiple Loss Functions

## Object Localization – MinMaxCAM [Wang et al., 2021]

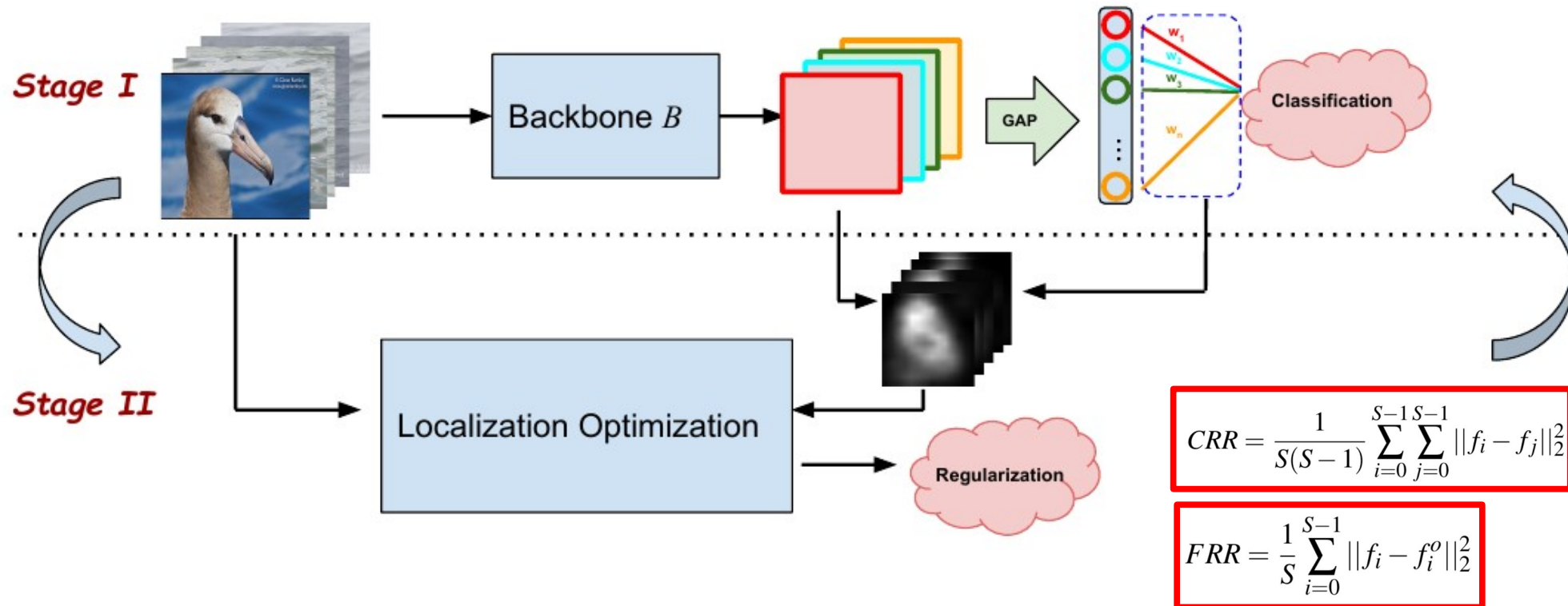
Idea: Regularize a high-performing classifier to enable localization



# Using Multiple Loss Functions

## Object Localization – MinMaxCAM [Wang et al., 2021]

Idea: Regularize a high-performing classifier to enable localization

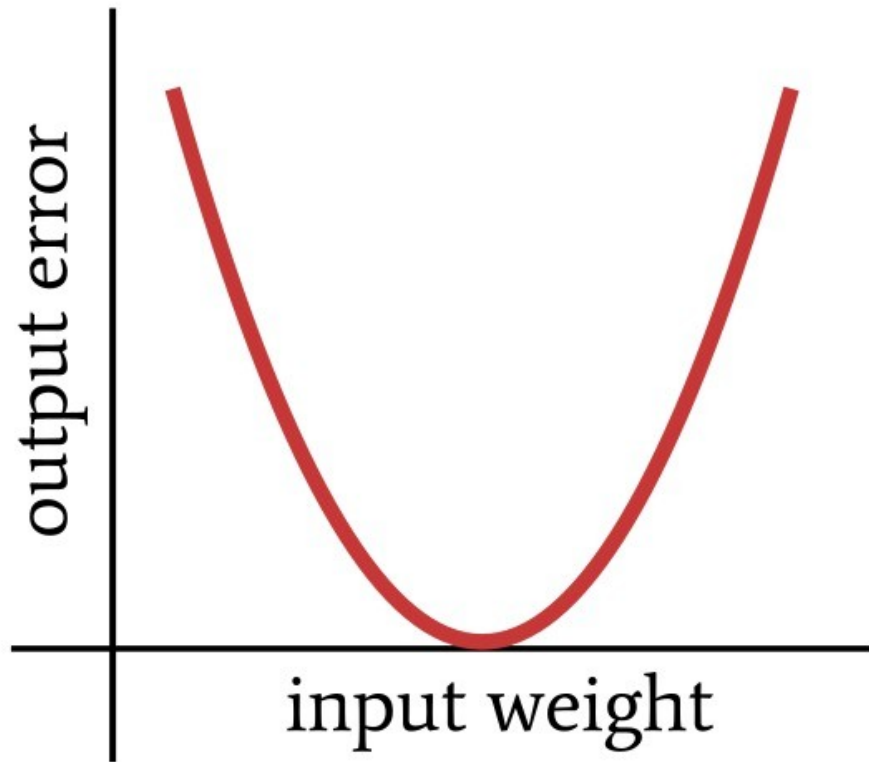


# Optimization

[ While Searching for the Best Solution ]

# Optimizing the Training Procedure

## Fixed VS Variable Learning Rate



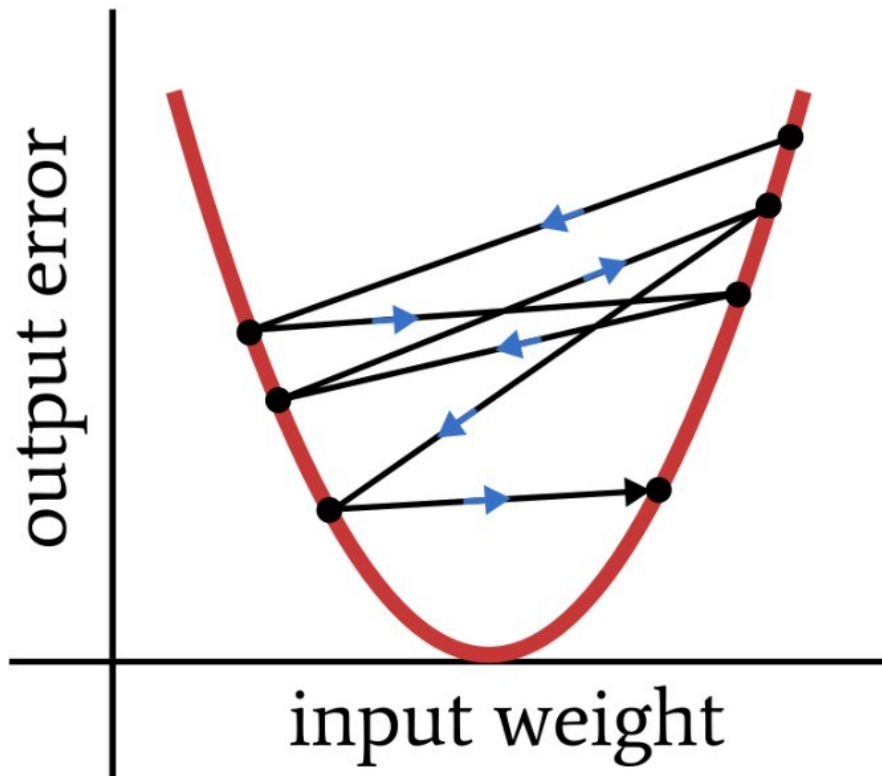
- **Problem:** As training progresses taken steps might be too large to reach the optimum

$$\theta_{t+1} = \theta_t - \alpha_t \nabla_{\theta} L(\theta_t)$$



# Optimizing the Training Procedure

## Fixed VS Variable Learning Rate

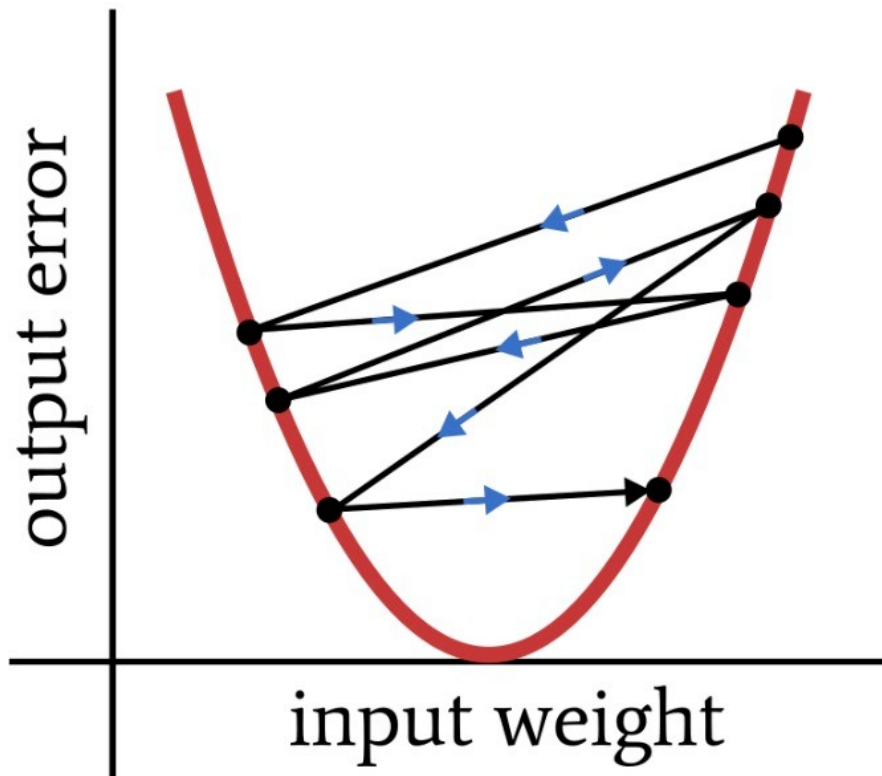


- **Problem:** As training progresses taken steps might be too large to reach the optimum

$$\theta_{t+1} = \theta_t - \alpha_t \nabla_{\theta} L(\theta_t)$$

# Optimizing the Training Procedure

## Fixed VS Variable Learning Rate



$$\theta_{t+1} = \theta_t - \alpha_t \nabla_{\theta} L(\theta_t)$$

- **Problem:** As training progresses taken steps might be too large to reach the optimum
- **Solution:** decrease the learning rate as training progresses.  
(Annealing)

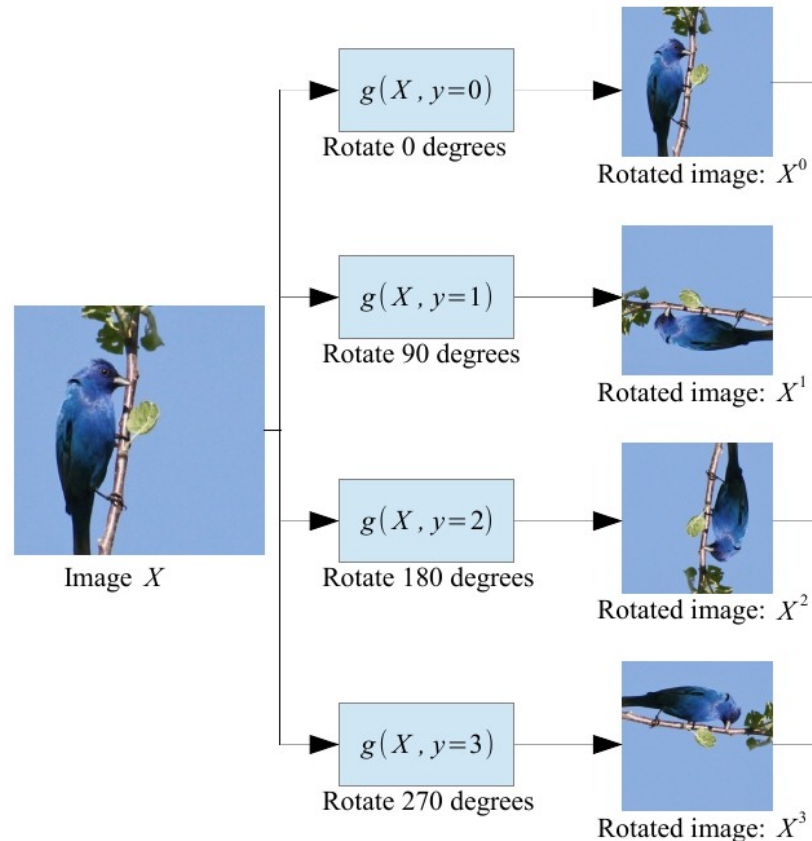
# Combinations

[ A bit of everything ]



# Self-Supervised Learning

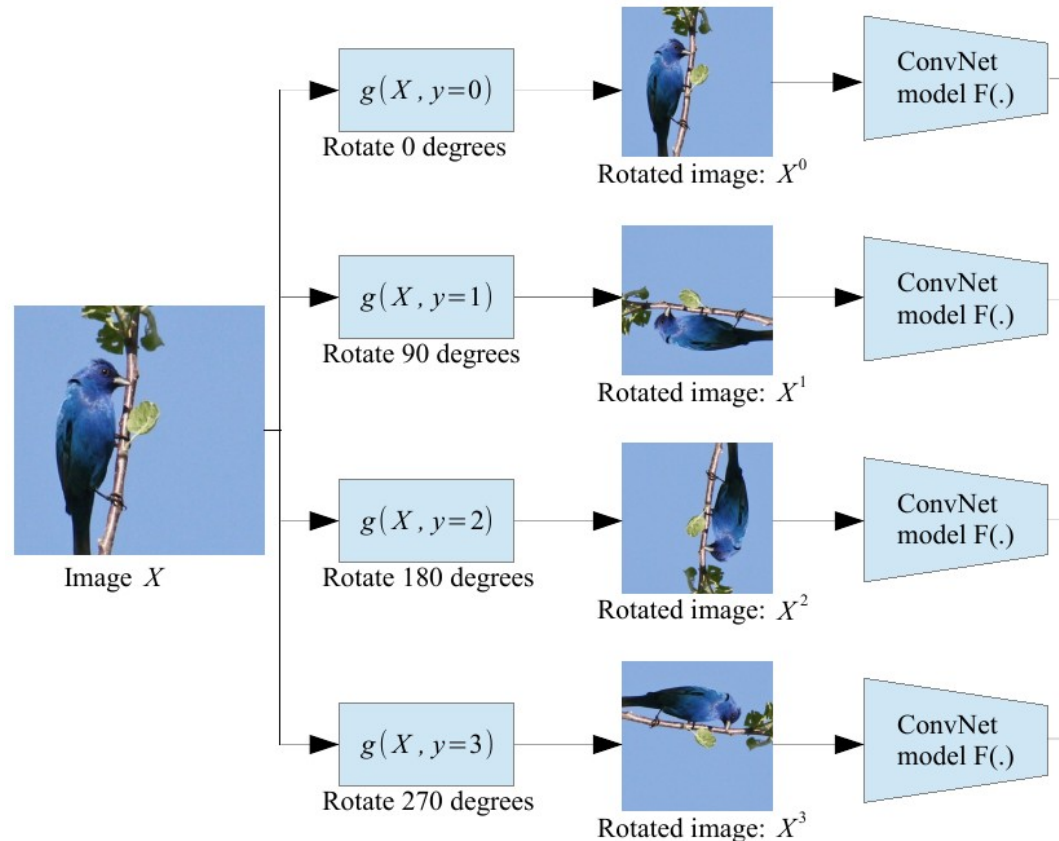
- Problem: data annotation is expensive
- Solution: supervise using labels generated from data (without manual annotation)



[Gidaris et al., 2018]

# Self-Supervised Learning

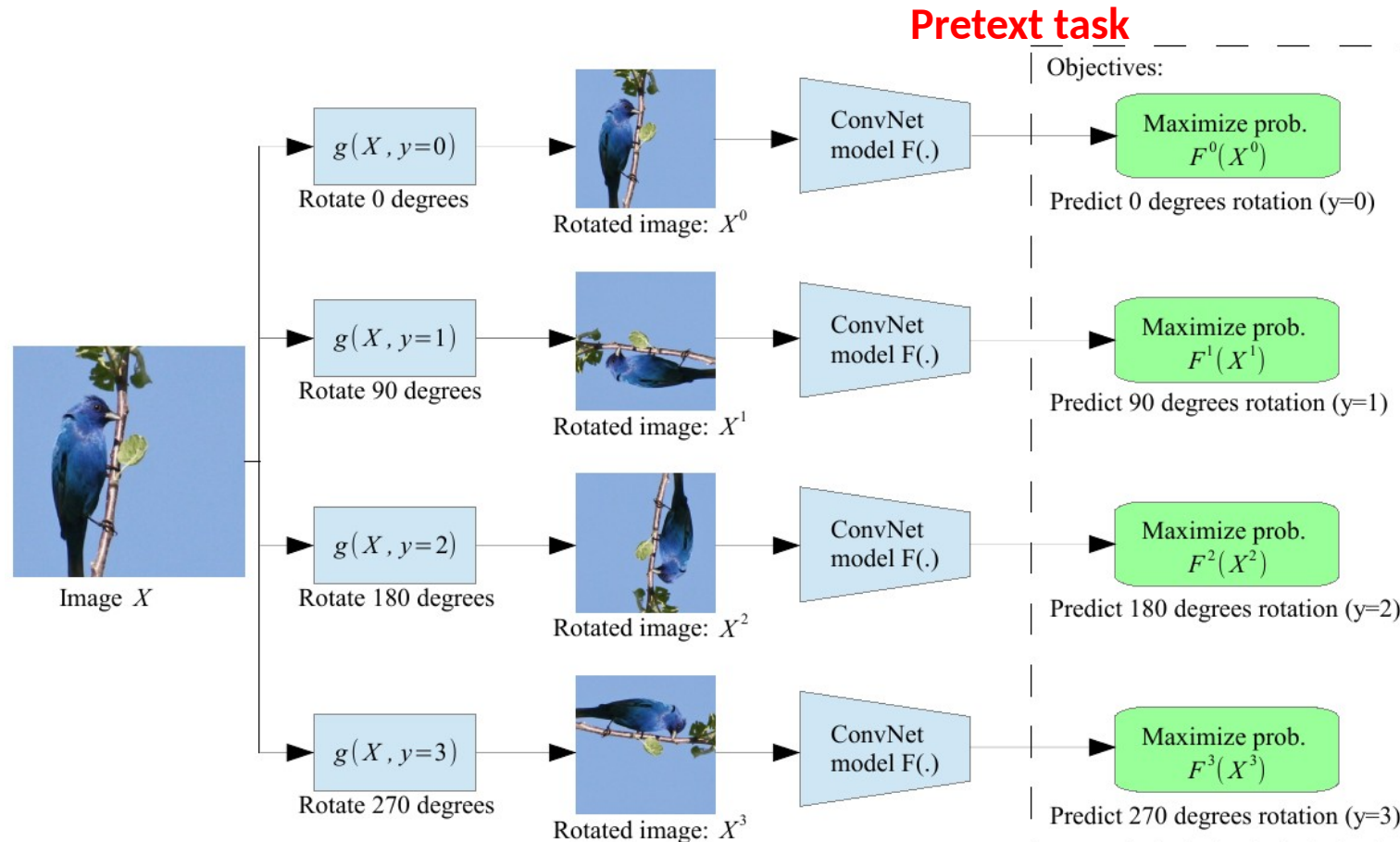
- Problem: data annotation is expensive
- Solution: supervise using labels generated from data (without manual annotation)



[Gidaris et al., 2018]

# Self-Supervised Learning

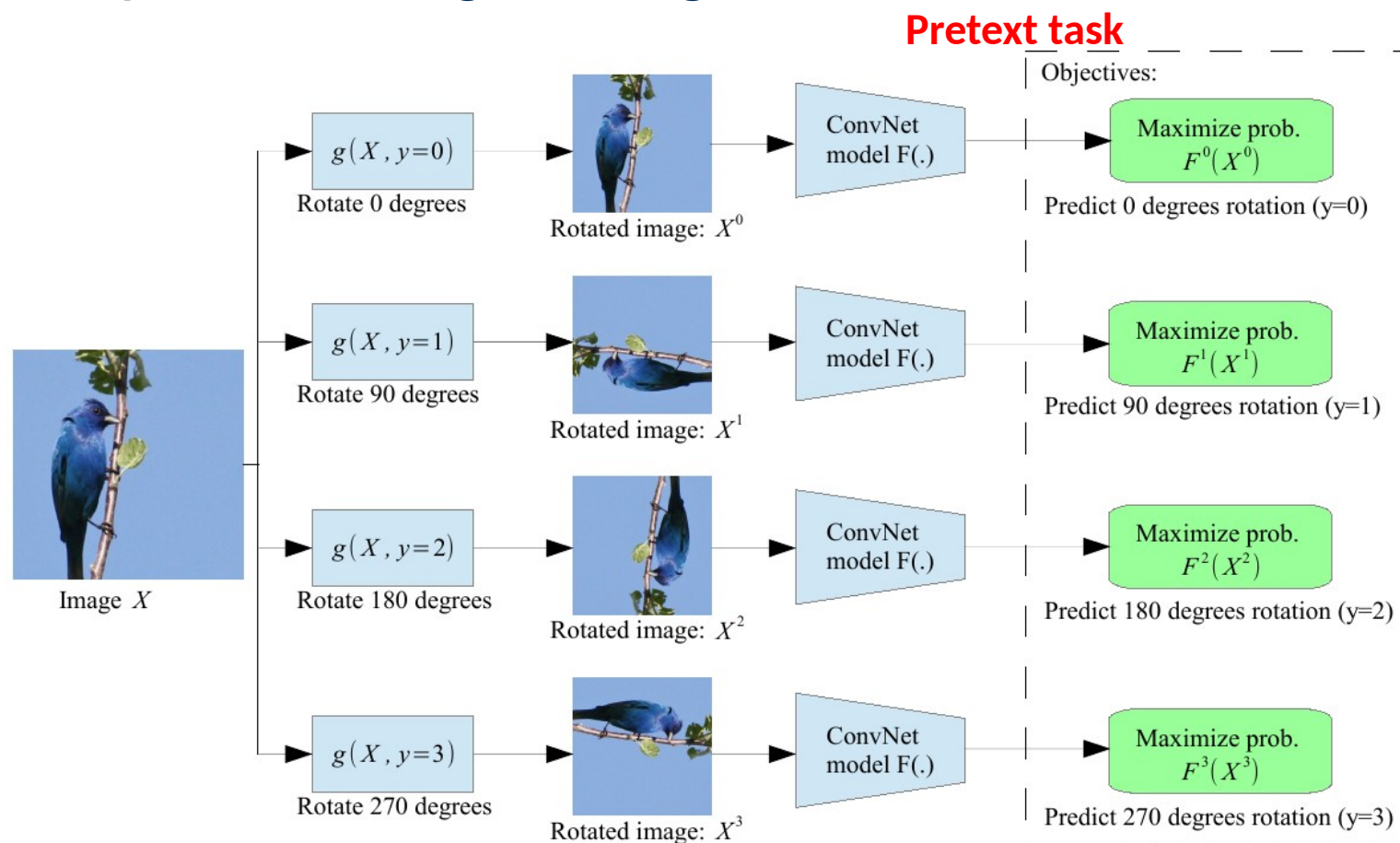
- Problem: data annotation is expensive
- Solution: supervise using labels generated from data (without manual annotation)



[Gidaris et al., 2018]

# Self-Supervised Learning

- Problem: data annotation is expensive
- Solution: supervise using labels generated from data (without manual annotation)



[Gidaris et al., 2018]

# Summarizing

[ Finally :D ]

# Summarizing

- **Different techniques possible**

Additional adaptation to the problem at hand might be required

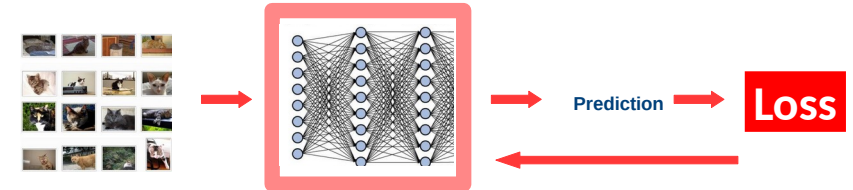
# Summarizing

- **Different techniques possible**

Additional adaptation to the problem at hand might be required

- **Can be applied at different stages of training**

pre/during training | computing the loss | updating the weights



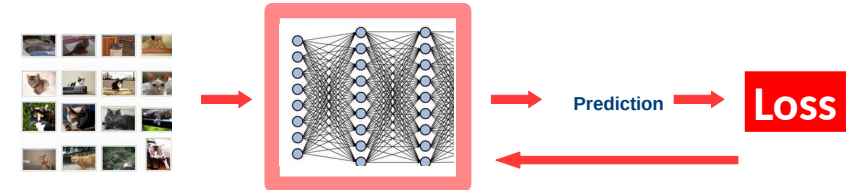
# Summarizing

- **Different techniques possible**

Additional adaptation to the problem at hand might be required

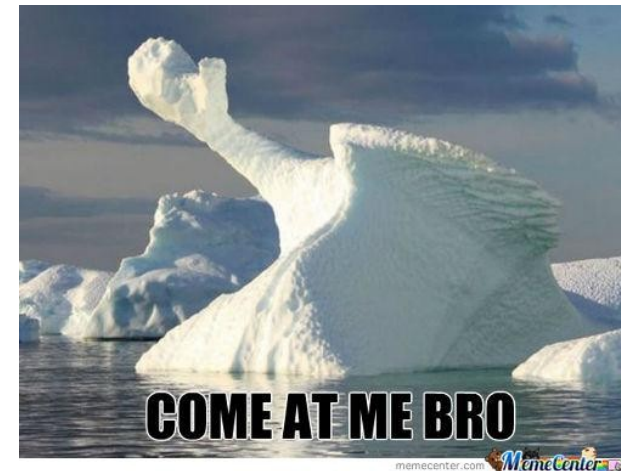
- **Can be applied at different stages of training**

pre/during training | computing the loss | updating the weights



- **Nothing but the tip of very huge iceberg**

lots of other techniques available





# Pay Attention to...

- When the techniques are applied?
  - Additional actions to be taken at different times
- What is the problem they address?
- How these techniques operate and what is their effect
  - How to adapt them to specific architectures?

# Questions?

# References

- S. Ioffe and C. Szegedy  
**Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift.** ICML 2015.
- A. Krizhevsky, I. Sutskever, G. E. Hinton.  
**ImageNet Classification with Deep Convolutional Neural Networks.** NeurIPS 2012
- T. Lin, P. Goyal, R. Girshick, K. He, P. Dollár  
**Focal Loss for Dense Object Detection.** ICCV 2017
- K. Wang, L. Ma, J. Oramas, L. Van Gool and T. Tuytelaars  
**Unsupervised shape transformer for image translation and cross-domain retrieval.** ArXiv:1812.02134, 2018
- S. Wei, V. Ramakrishna, T. Kanade, Y. Sheikh,  
**Convolutional Pose Machines,** CVPR 2016.
- Spyros Gidaris, Praveer Singh, Nikos Komodakis.  
**Unsupervised Representation Learning by Predicting Image Rotations.** ICLR 2018



# Artificial Neural Networks

[2500WETANN]

José Oramas