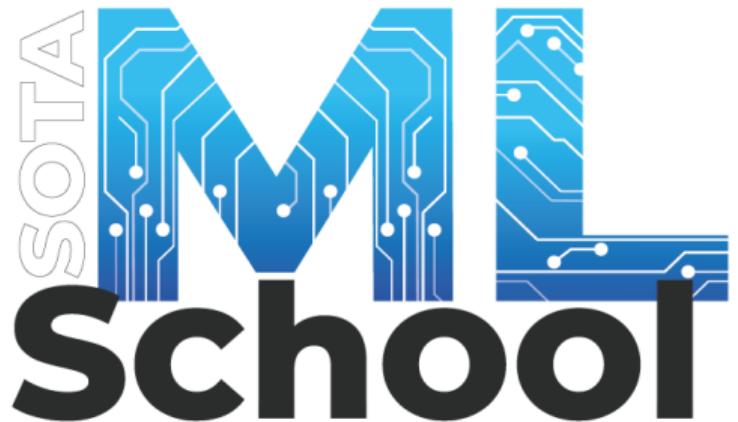


**ITMO.Hack**



**BOTAN**  
INVESTMENTS



Deep Learning Theory Part II

Semyon Polyakov

# Outline

Bias and Variance tradeoff

- Problem setup

- Regularization

Who was after Adam?

- Optimizers

Multi-tasking

- Multi-tasking

- Transfer learning

Compression

- Compression

Last but not least

- Adversarial attacks

- Neural tangent kernels

- Flows

- Final

## Bias and Variance tradeoff

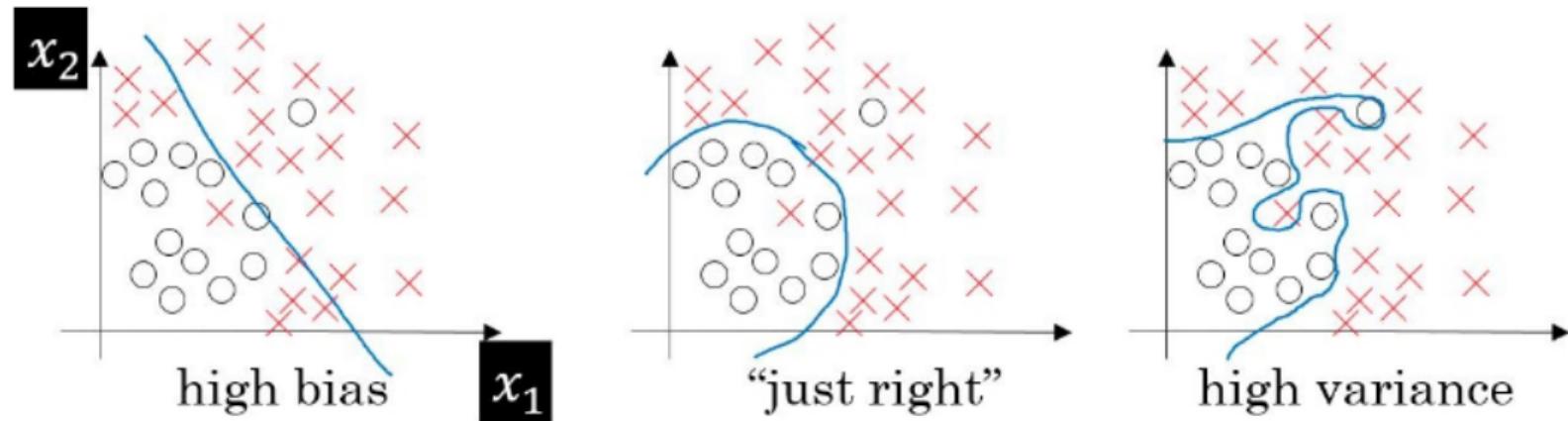


Figure: Bias and variance tradeoff corner cases illustration

# Bias and Variance tradeoff

$$E[(y - \tilde{f}(x))^2] = (\text{Bias}[\tilde{f}(x)])^2 + \text{Var}[\tilde{f}(x)] + \sigma^2$$

$$\text{Bias}[\tilde{f}(x)] = E[\tilde{f}(x) - f(x)]$$

$$\text{Var}[\tilde{f}(x)] = E[(\tilde{f}(x))^2] - E[\tilde{f}(x)]^2$$

## Intuition

The bias error is an error from erroneous assumptions in the learning algorithm (underfitting)

The variance is an error from sensitivity to small fluctuations in the training set (overfitting)

# Regularization

$$E_D(w) + \lambda E_W(w)$$

$$E_W = \sum_{j=1}^M |w_j|^q$$

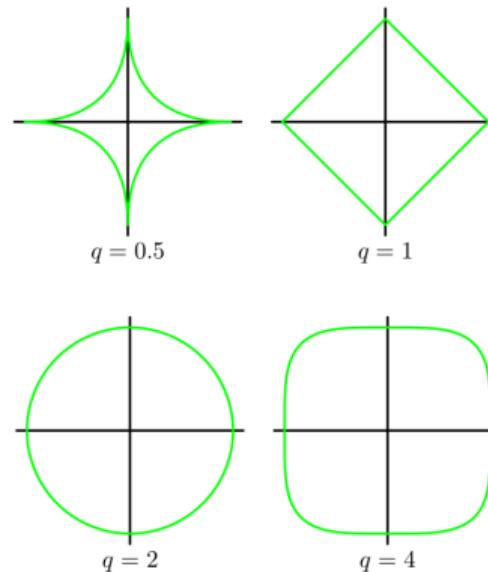


Figure: Polynomial regularizations - L1, L2, L3, L4

# Regularization

$$E_D(w) + \lambda E_W(w)$$

$$E_W = \sum_{j=1}^M |w_j|^q$$

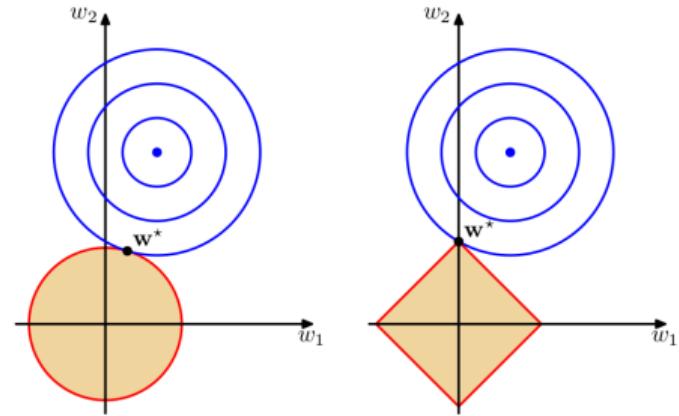


Figure: L1 tends to zeroing least valuable weights

# Regularization

Other ways for regularization:

- ▶ Dropout (optimal is near 0.5)
- ▶ BatchNormalization (+layer, instance and group normalization)

## BatchNorm limitations

- ▶ Works bad for small batches
- ▶ Too computationally expensive for recurrent networks

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_1 \dots m\}$ ;  
Parameters to be learned:  $\gamma, \beta$

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

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

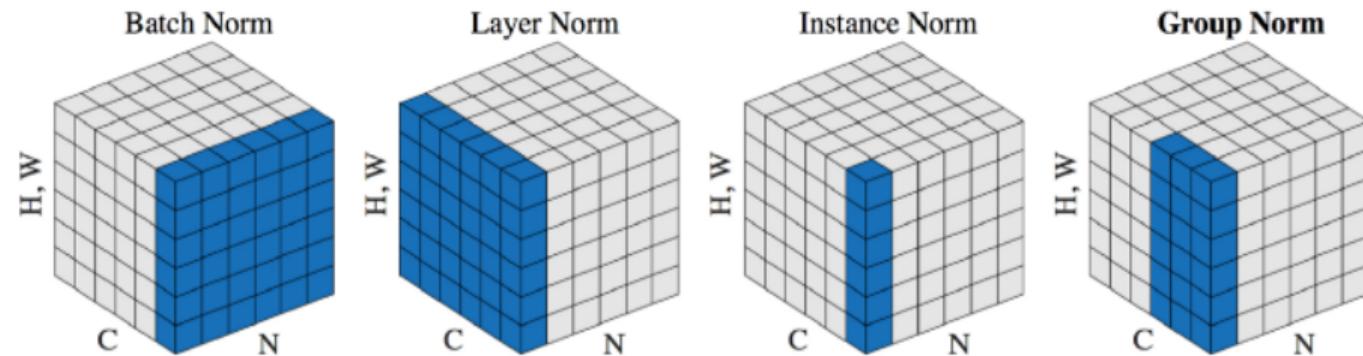
$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{mini-batch variance}$$

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

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

Figure: BatchNorm algorithm

# Regularization



**Figure:** Different types of normalization in NN

# Optimizers

**Optimizers, warmup and LR scheduling**

# Optimizers

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$$
$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$$

## Adam

- ▶ One of the best optimizers, but often loses to SGD with momentum
- ▶ Has 3 hyperparameters
- ▶ There are plenty of improved Adam versions (NAdam, RAdam, AdaMod, Lamb, LARS, ...)

$$\hat{m}_t = \frac{m_t}{1 - \beta_1}$$

$$\hat{v}_t = \frac{v_t}{1 - \beta_2}$$

$$w_t = w_{t-1} - \eta \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon}$$

Figure: Adam optimizer

# Optimizers

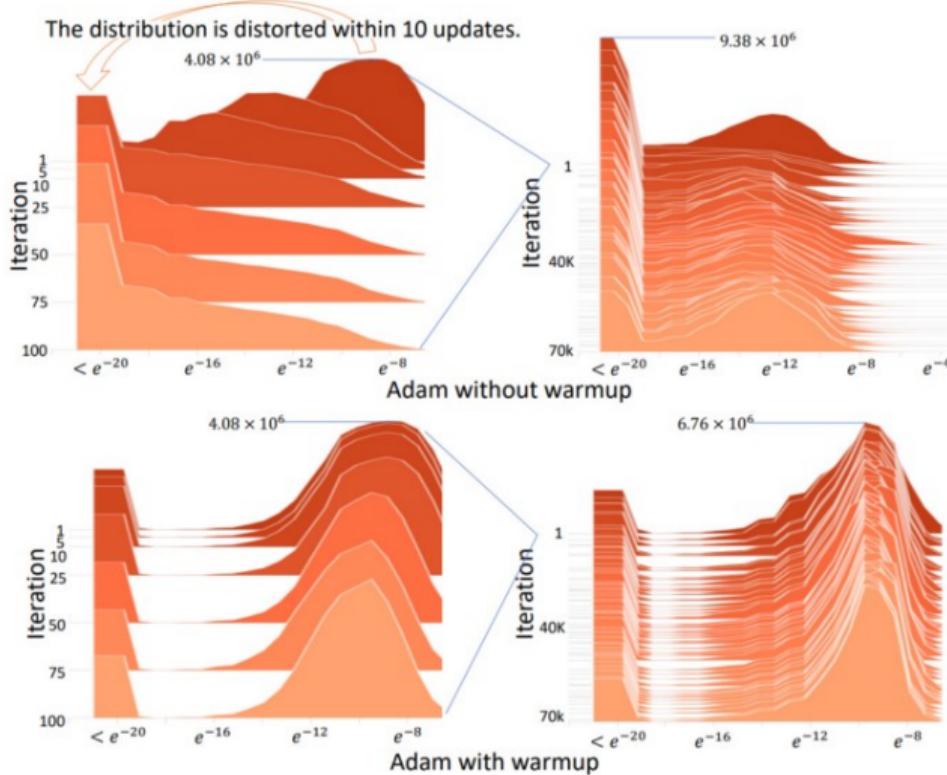


Figure: Adam warmup influence

# Optimizers

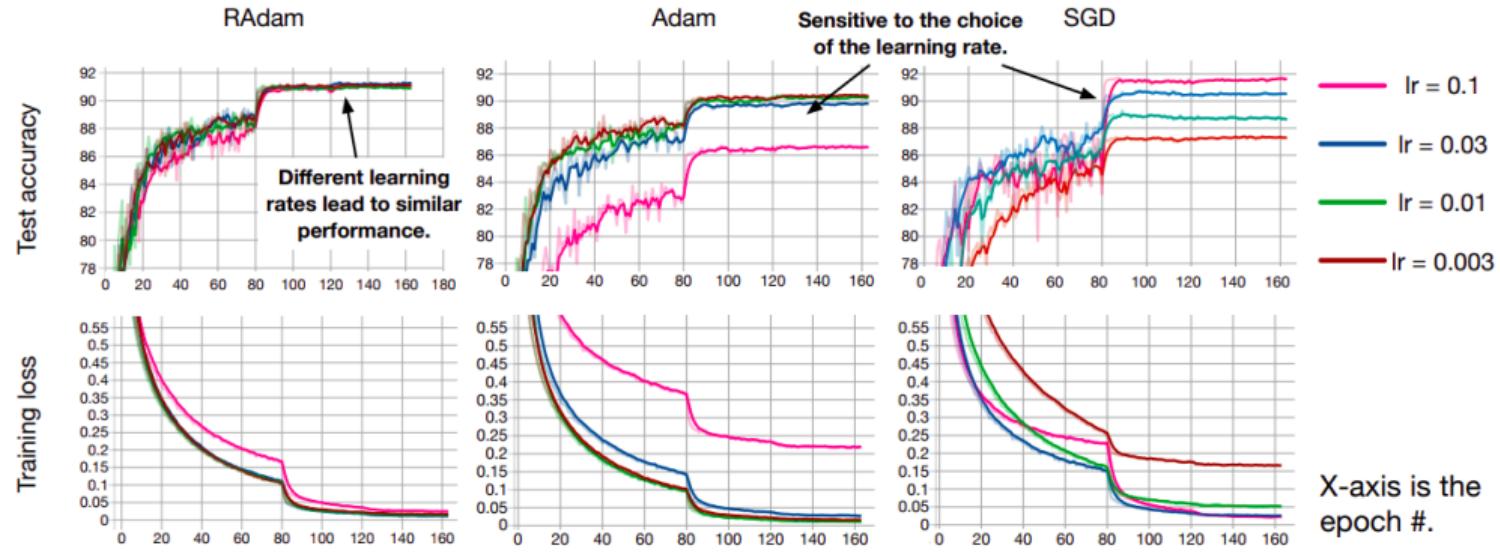


Figure: RAdam performance for CIFAR-10 dataset

# Optimizers

---

**Algorithm 2:** Rectified Adam. All operations are element-wise.

---

**Input:**  $\{\alpha_t\}_{t=1}^T$ : step size,  $\{\beta_1, \beta_2\}$ : decay rate to calculate moving average and moving 2nd moment,  $\theta_0$ : initial parameter,  $f_t(\theta)$ : stochastic objective function.

**Output:**  $\theta_T$ : resulting parameters

```
1  $m_0, v_0 \leftarrow 0, 0$  (Initialize moving 1st and 2nd moment)
2  $\rho_\infty \leftarrow 2/(1 - \beta_2) - 1$  (Compute the maximum length of the approximated SMA)
3 while  $t = \{1, \dots, T\}$  do
4    $g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$  (Calculate gradients w.r.t. stochastic objective at timestep t)
5    $v_t \leftarrow 1/\beta_2 v_{t-1} + (1 - \beta_2)g_t^2$  (Update exponential moving 2nd moment)
6    $m_t \leftarrow \beta_1 m_{t-1} + (1 - \beta_1)g_t$  (Update exponential moving 1st moment)
7    $\widehat{m}_t \leftarrow m_t / (1 - \beta_1^t)$  (Compute bias-corrected moving average)
8    $\rho_t \leftarrow \rho_\infty - 2t\beta_2^t / (1 - \beta_2^t)$  (Compute the length of the approximated SMA)
9   if the variance is tractable, i.e.,  $\rho_t > 4$  then
10     $l_t \leftarrow \sqrt{(1 - \beta_2^t)/v_t}$  (Compute adaptive learning rate)
11     $r_t \leftarrow \sqrt{\frac{(\rho_t - 4)(\rho_t - 2)\rho_\infty}{(\rho_\infty - 4)(\rho_\infty - 2)\rho_t}}$  (Compute the variance rectification term)
12     $\theta_t \leftarrow \theta_{t-1} - \alpha_t r_t \widehat{m}_t l_t$  (Update parameters with adaptive momentum)
13  else
14     $\theta_t \leftarrow \theta_{t-1} - \alpha_t \widehat{m}_t$  (Update parameters with un-adapted momentum)
15 return  $\theta_T$ 
```

---

# Optimizers

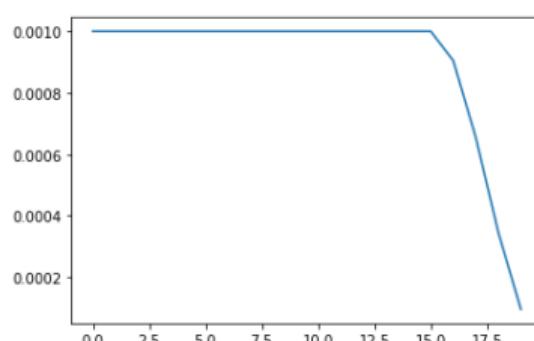
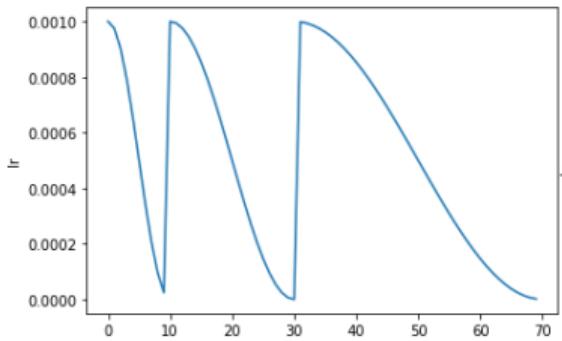
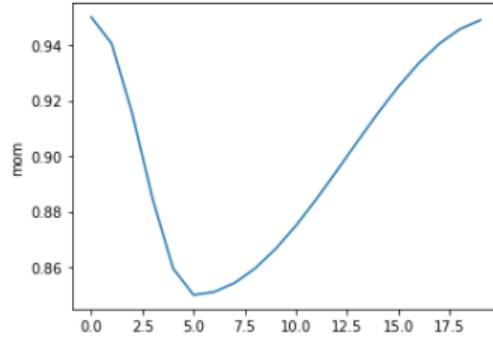
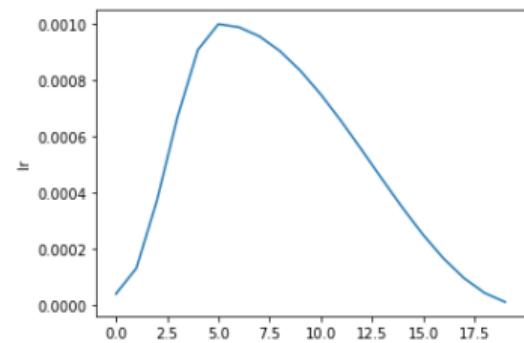


Figure: Learning rate scheduling examples from fastai

# Optimizers

## Notes:

- ▶ No free lunch
- ▶ Combinations of ideas can help (Nesterov momentum, exponential weighting, ..)
- ▶ Layer-wise learning rate wins global learning rate
- ▶ Slow, but effective: SGD with momentum + warmup/LR scheduling + tuning

# Multi-tasking

**Transfer learning, Multi-tasking, one-shot learning**

# Preliminaries

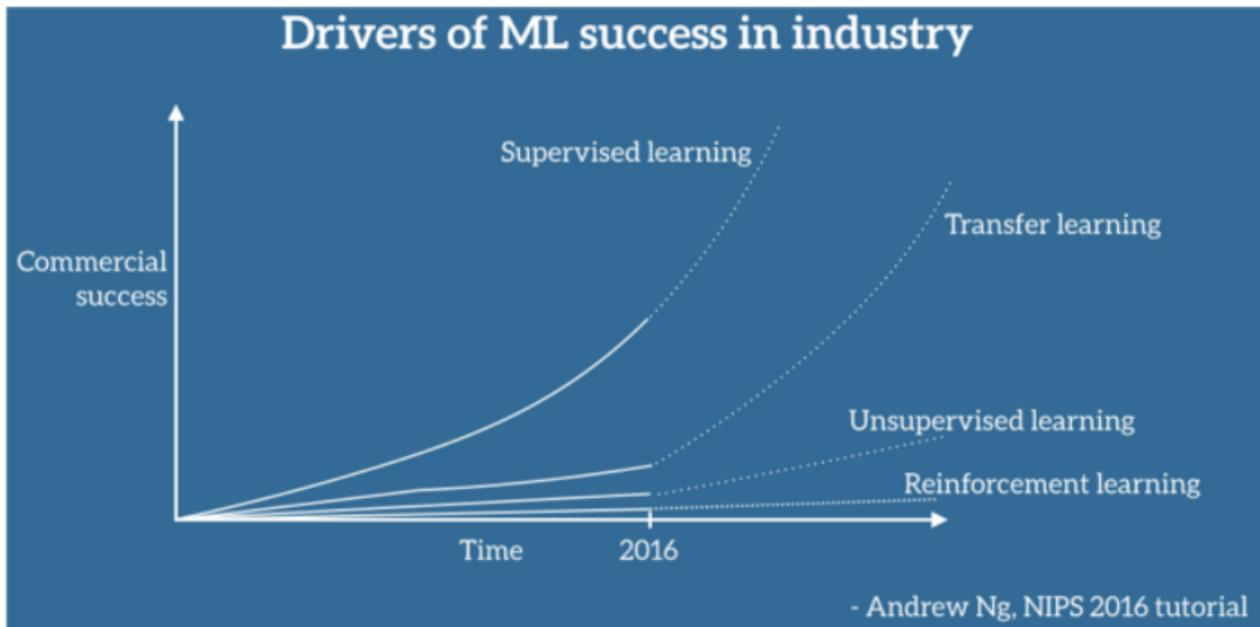


Figure: ML areas success plot by Andrew NG

# Transfer learning

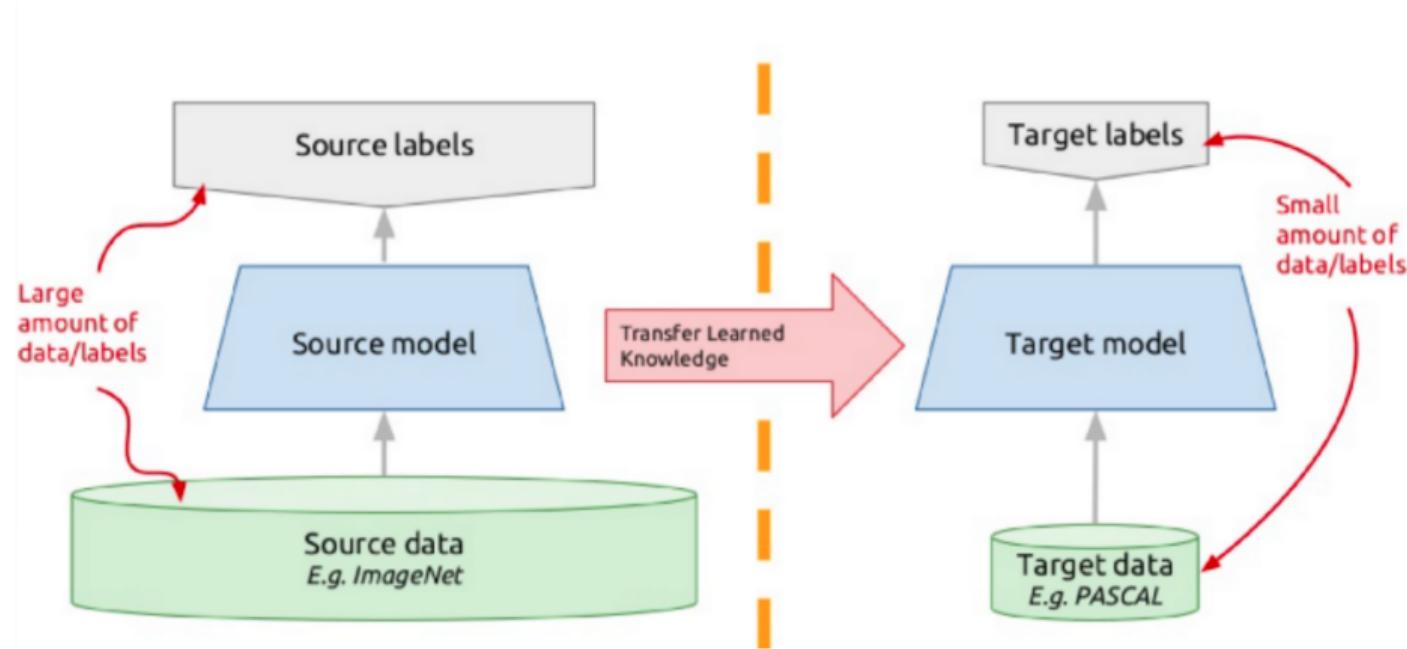


Figure: Transfer learning principal scheme

# Multi-tasking

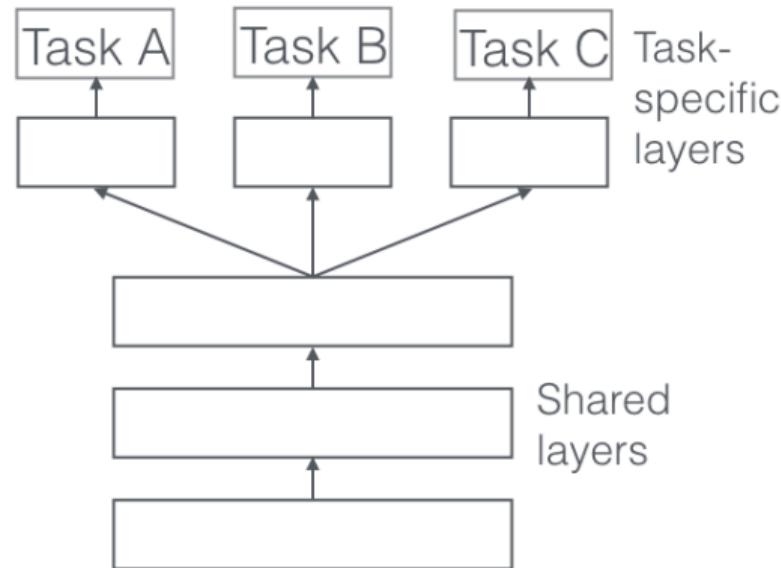


Figure: Shared layers for multi-task learning

# Multi-tasking

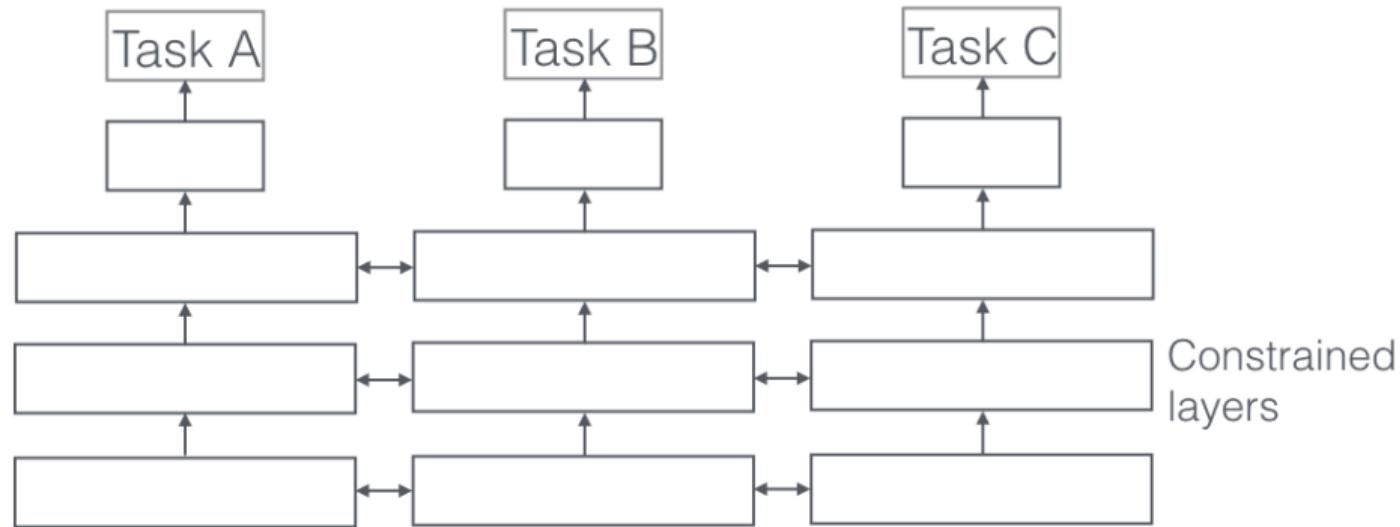


Figure: Weakly shared layers for multi-task learning

# Transfer learning

## Notes:

- ▶ One-shot learning = Metric learning + transfer learning
- ▶ Augmentation helps even with 1 example
- ▶ Zero-shot learning = Metric learning + (external data)?
- ▶ Center loss for better regularization

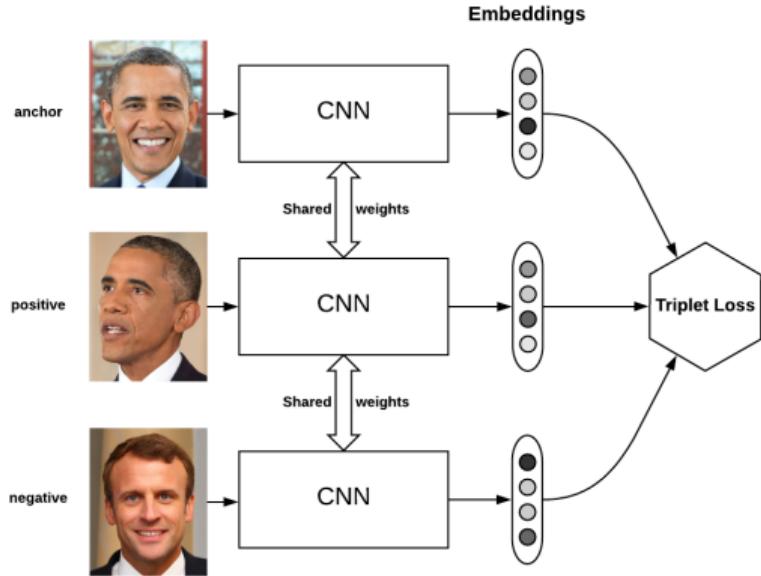


Figure: Triplet loss example

# Compression

**Neural networks compression and speed up**

# Compression and speed up

## Notes:

- ▶ Quantization (16bit, 8bit, ..)
- ▶ Factorization
- ▶ Pruning (weights, neurons, filters)
- ▶ Knowledge distillation

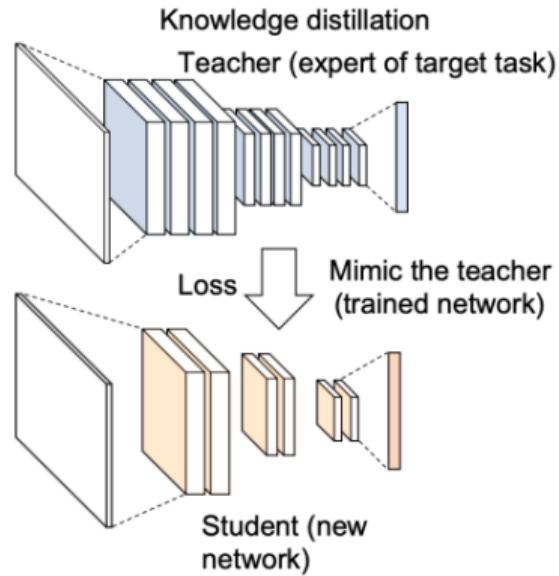


Figure: Triplet loss example

## Recent neural networks research

# Adversarial attacks

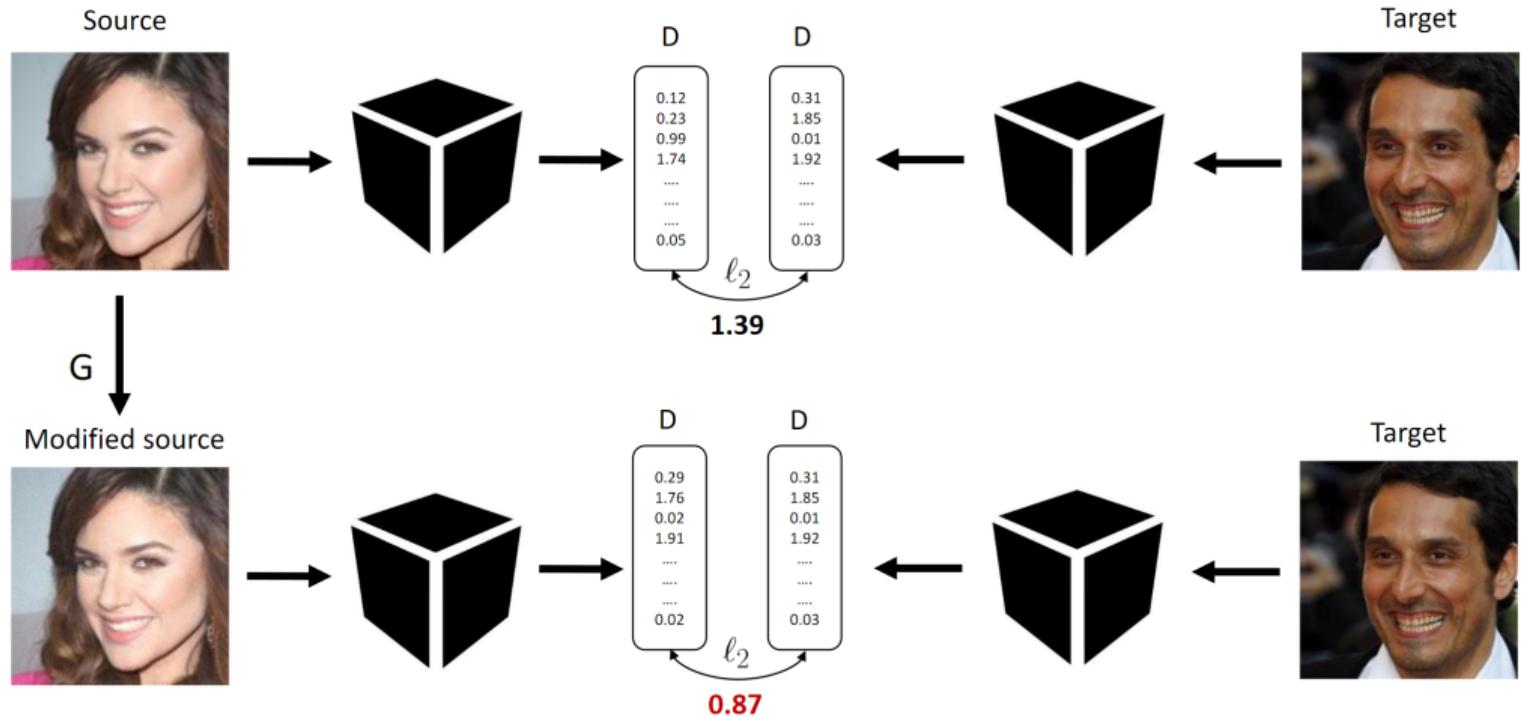


Figure: Adversarial attack example

# Neural tangent kernels

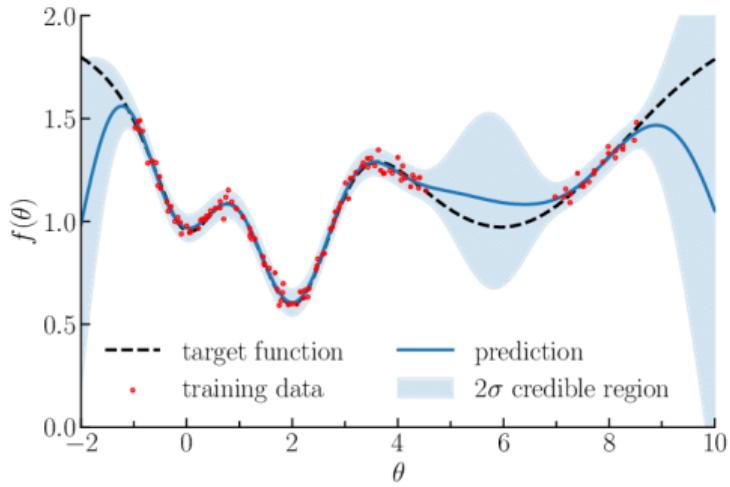


Figure: Gaussian process

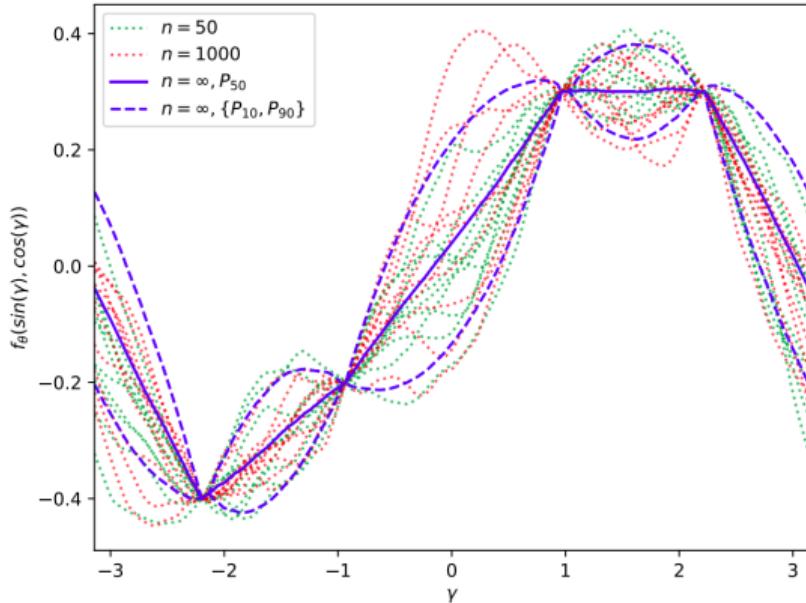


Figure: Neural network asymptotically approximates GP

# Flows

$$x \sim P(X)$$

$$z \sim p_\theta(z)$$

$$x = g_\theta(z)$$

$$z = f_\theta(x) = g^{-1}_\theta(x)$$

$$f = f_1 \circ f_2 \circ \dots \circ f_k$$

Figure: Flow-based networks



Figure: Glow: Flow-based face generation network results

**During the research, not a single neural network was harmed**

# Thank you for attention!

