

# Visual AI

CPSC 532R/533R – 2019/2020 Term 2

## Lecture 3. Network architectures for image processing (and their optimization)

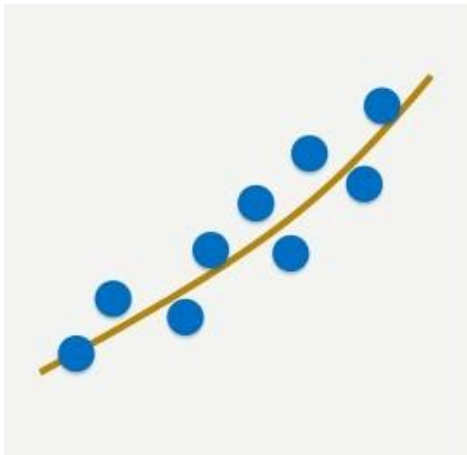
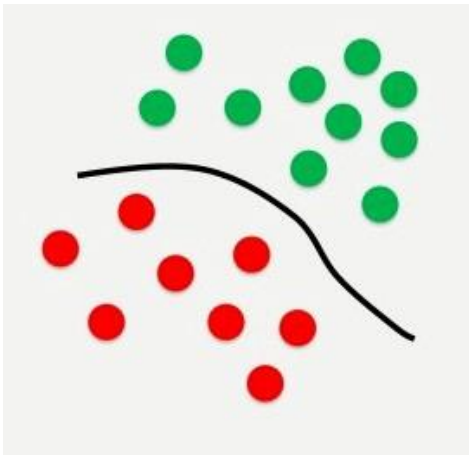
Helge Rhodin



# Classification vs. regression

Classification

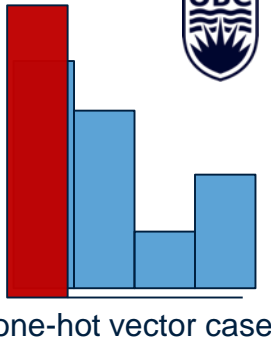
Regression



# Cross-entropy loss / Cross-entropy criterion

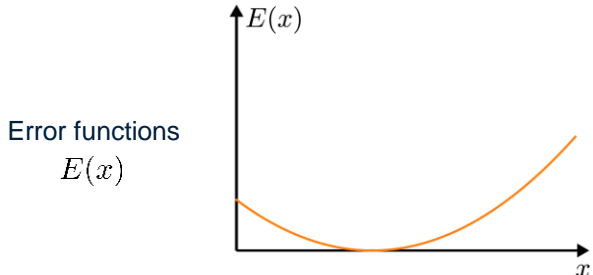
Negative Log Likelihood (NLL) formulation for a *one-hot vector*, target class c

$$l_{\text{NLL}}(x, c) = -\log(f_{[c]}(x))$$



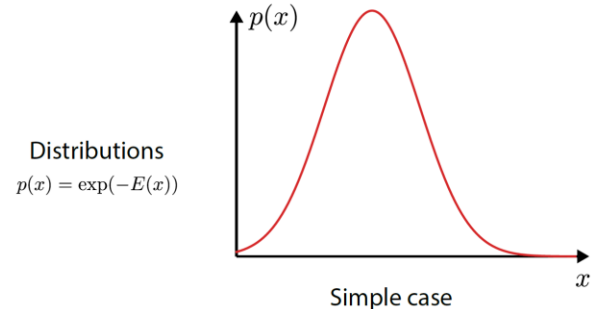
# Regression revisited

Many loss functions are  $-\log$  of probability distributions



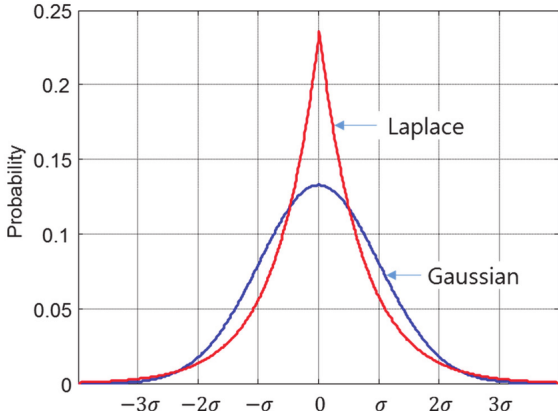
$x^2$   
Mean squared error (MSE)

$|x|$   
Mean absolute error (MAE)



$\exp(-x^2)$   
Gaussian distribution

$\exp(-|x|)$   
Laplace distribution



# Recap done



# Density networks

**Assumption:** Data is distributed according to a random process

$$y \sim N(F(x), \sigma), \text{ with mean } F(x) \text{ and standard deviation } \sigma$$

**Given N samples**  $(x_i, y_i)_{i=1}^N$

**Goal:** Find that function that maximizes the likelihood of the samples  $(x_i, y_i)$

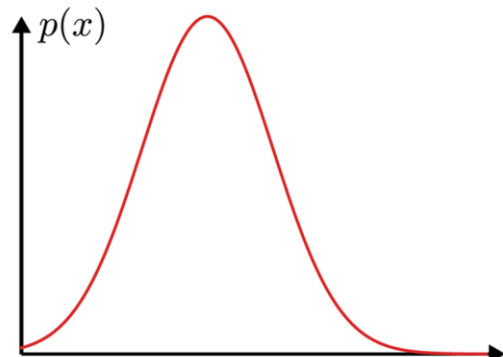
$$L = N(y|f_\theta(x), \sigma), \text{ with mean } f_\theta(x) \text{ a neural network with parameters } \theta$$

**Solution:** minimize the negative log-likelihood (here mean squared error)

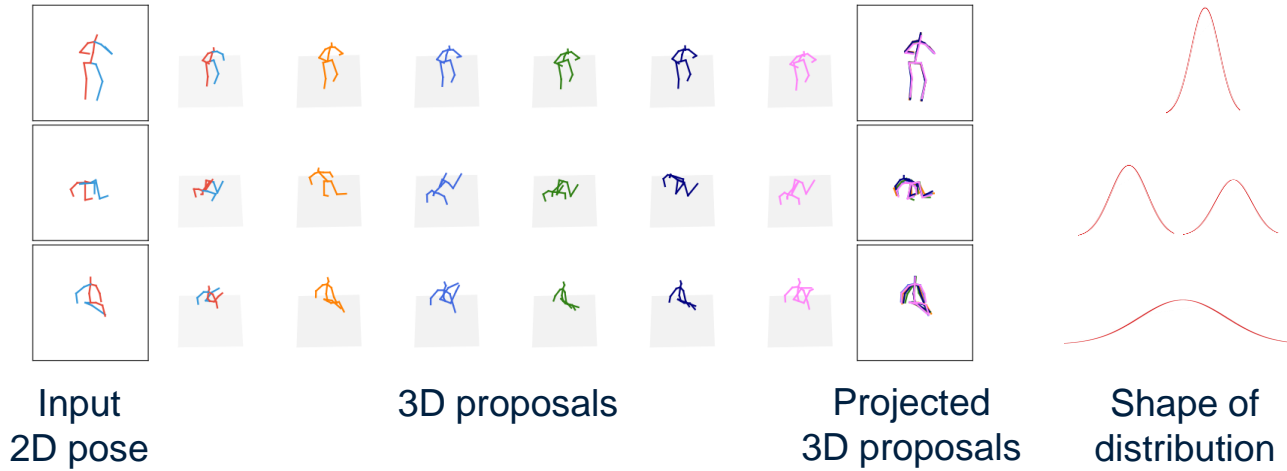
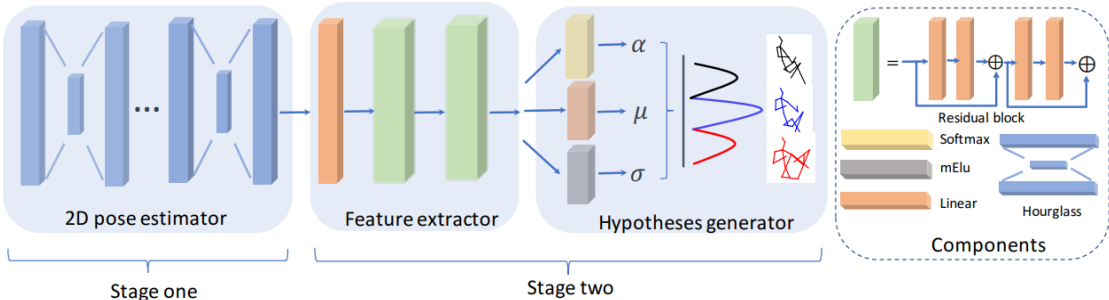
$$E = \frac{1}{2\sigma^2} (y - f_\theta(x))^2$$

**Density Networks:** Predict the mean  $\mu$  and standard deviation  $\sigma$

$$\sigma, \mu = f_\theta(x) \quad L = N(y|\mu, \sigma)$$



# Generating Multiple Hypotheses for 3D Human Pose Estimation with Mixture Density Networks

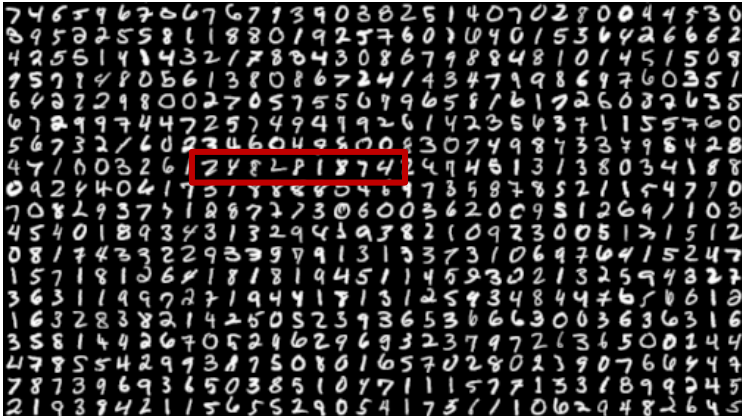


# Separable objective and mini batches

Separable objective over independent samples  $x, y$

$$E(D, \theta) = \sum_{(\mathbf{x}^{(i)}, y^{(i)}) \in D} l(\mathbf{x}^{(i)}, \theta), y^{(i)})$$

Evaluated over mini batches of size  $N$



Stored as tensor, e.g.,

- dim 0:  $N$ , number of images in a batch
- dim 1:  $C$ , number of channels
- dim 2:  $H$ , height of the feature map
- dim 3:  $W$ , width of the feature map



# Optimizers

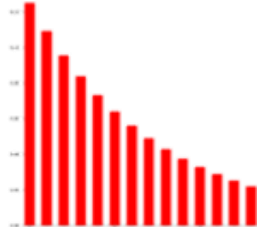
## Stochastic Gradient Descent

- Gradient descent on randomized mini batches (with learning rate alpha)

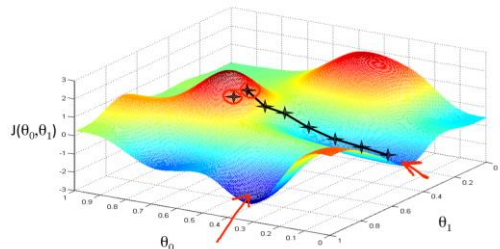
$$\theta_t = \theta_{t-1} - \alpha \sum_{i=1}^n \nabla E_i(\theta) / n,$$

## Adam

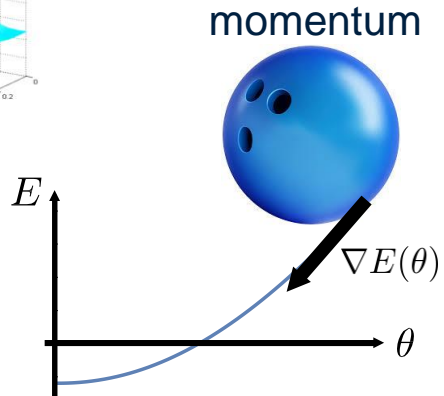
- Momentum-based (continue with larger steps if the previous steps point in the same direction)
- Damp step-length if direction changes often (second moment is high)
- Uses exponential moving average (EMA)



$$\bar{y}_t = \begin{cases} y_1, & t = 1 \\ \beta \cdot \bar{y}_{t-1} + (1 - \beta) \cdot y_t, & t > 1 \end{cases}$$



<https://hackernoon.com/gradient-descent-aynk-7cbe95a778da>



$$\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \hat{m}_t / (\sqrt{\hat{v}_t} + \epsilon)$$

Adam update rule

*Smaller batch size can be better; it induces more noise!*

# Adam and co.

- Adam is my current favorite
    - Not that sensitive to learning rate
    - No scheduler necessary
    - Intuitive motivation
- Disadvantage: Properly tuned SGD can be more accurate
- Recent alternative
    - Learning with Random Learning Rates [Bluer et al.,]
      - give each neuron a different learning rate
      - those with inappropriate rates will die (constant output for all feasible input values)
      - parameter free, more stable training

---

```

Require:  $\alpha$ : Stepsize
Require:  $\beta_1, \beta_2 \in [0, 1)$ : Exponential decay rates for the moment estimates
Require:  $f(\theta)$ : Stochastic objective function with parameters  $\theta$ 
Require:  $\theta_0$ : Initial parameter vector
 $m_0 \leftarrow 0$  (Initialize 1st moment vector)
 $v_0 \leftarrow 0$  (Initialize 2nd moment vector)
 $t \leftarrow 0$  (Initialize timestep)
while  $\theta_t$  not converged do
   $t \leftarrow t + 1$ 
   $g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$  (Get gradients w.r.t. stochastic objective at timestep  $t$ )
   $m_t \leftarrow \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$  (Update biased first moment estimate)
   $v_t \leftarrow \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2$  (Update biased second raw moment estimate)
   $\hat{m}_t \leftarrow m_t / (1 - \beta_1^t)$  (Compute bias-corrected first moment estimate)
   $\hat{v}_t \leftarrow v_t / (1 - \beta_2^t)$  (Compute bias-corrected second raw moment estimate)
   $\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \hat{m}_t / (\sqrt{\hat{v}_t} + \epsilon)$  (Update parameters)
end while
return  $\theta_t$  (Resulting parameters)
  
```

---

[Kingma, Diederik P., and Jimmy Ba. "Adam: A method for stochastic optimization. ICLR 2015]

# Automatic differentiation and backpropagation

Forward pass

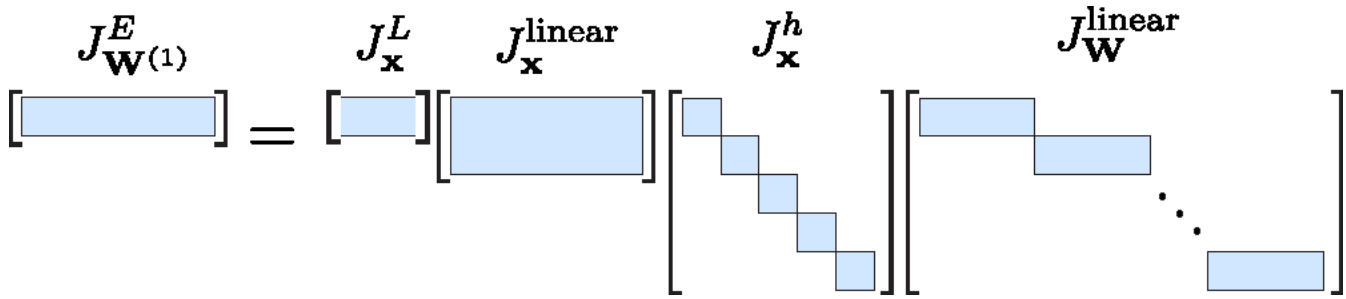
$$L(h(\text{linear}(h(\text{linear}(x, W^{(1)})), W^{(2)})))$$

$$= Lh \left( \begin{array}{|c|c|} \hline W^* & b^* \\ \hline \end{array} h \left( \begin{array}{|c|c|c|} \hline W^* & b^* & x^* \\ \hline \end{array} \right) \right)$$

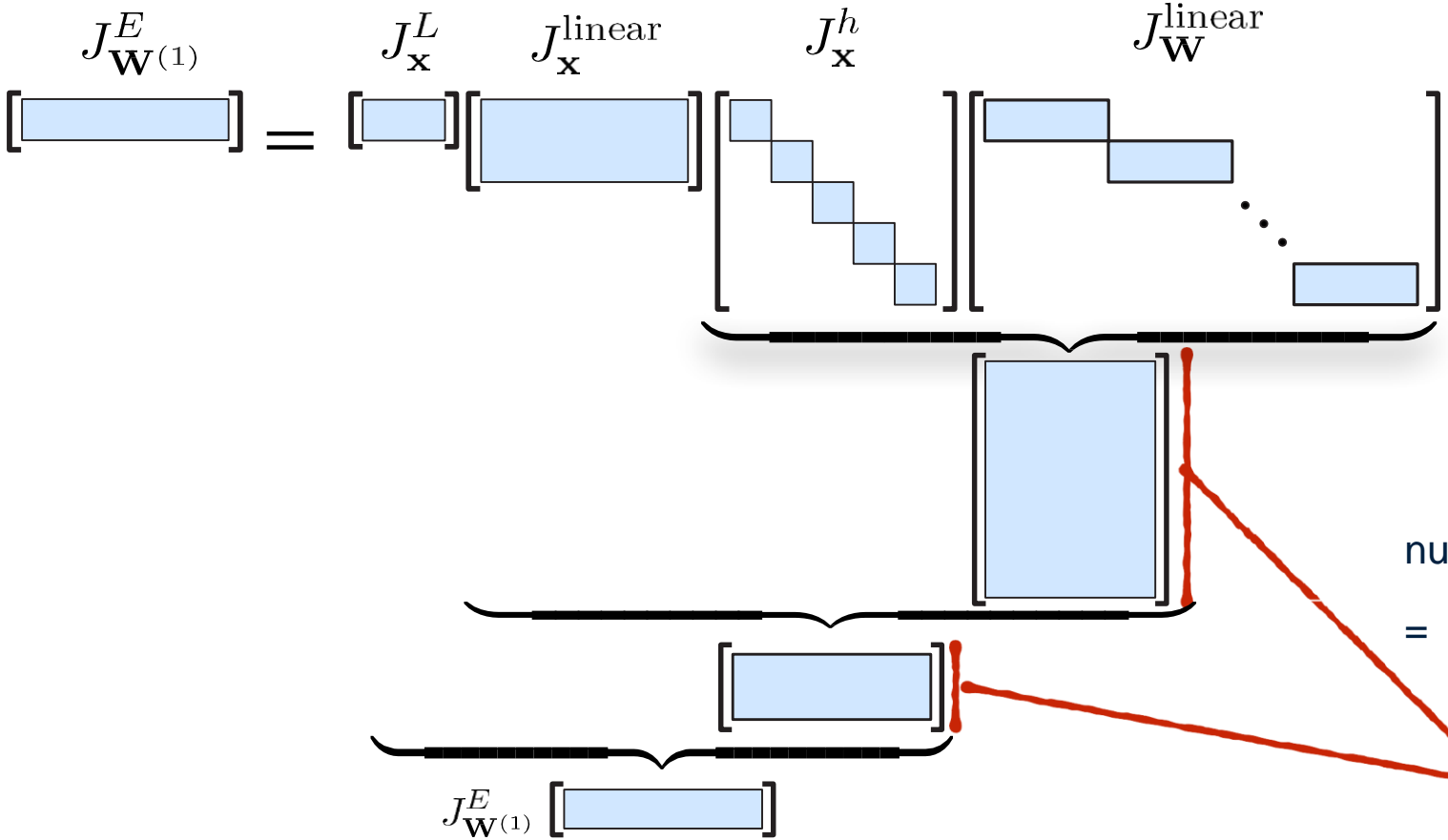
$$J_x^f = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

Jacobian matrix

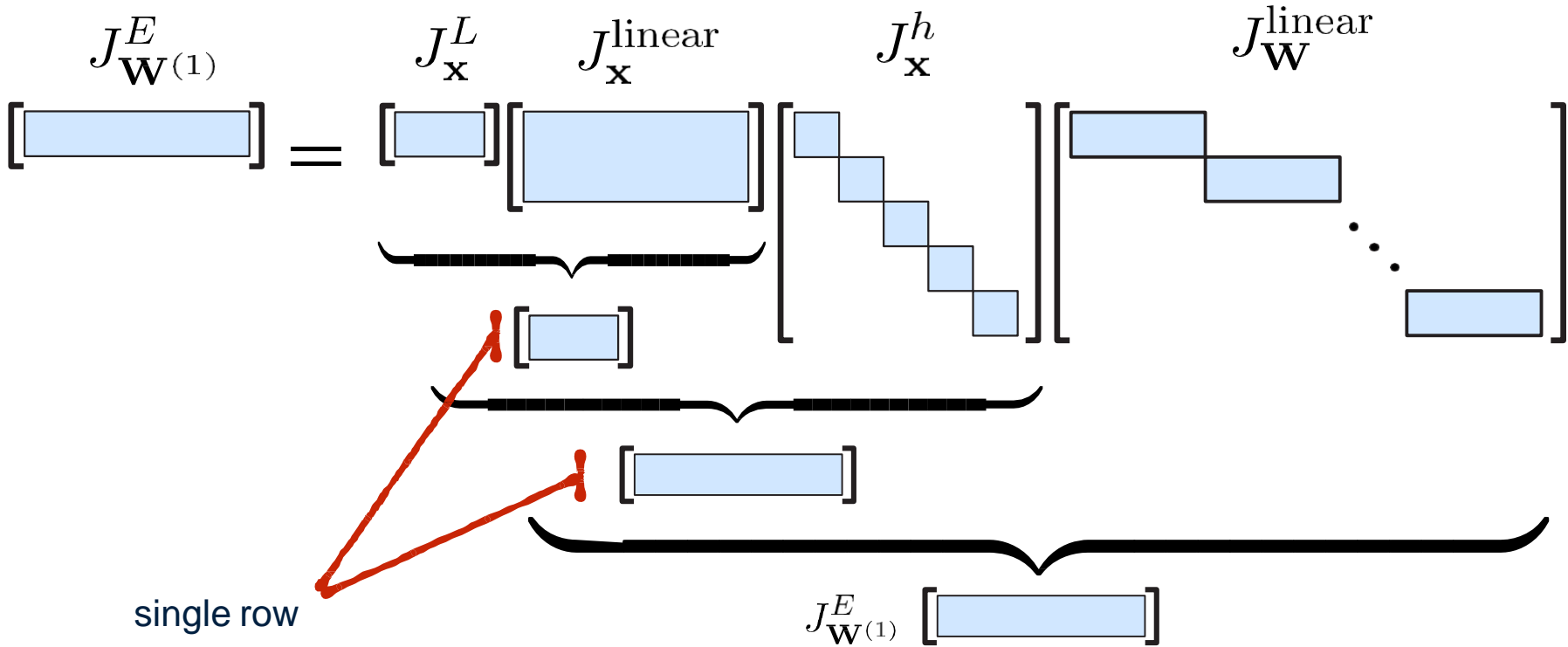
Backwards pass to  $W^{(1)}$



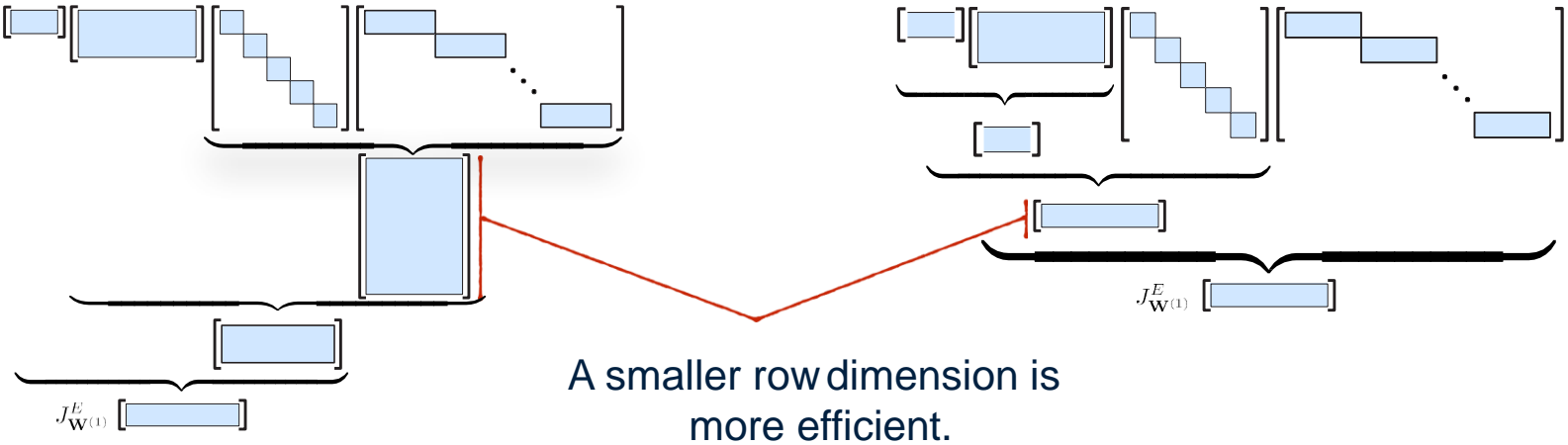
# Forward propagation



# Reverse mode - backpropagation



# Forward vs. reverse mode



Forward accumulation is more efficient for functions that have more outputs than inputs.

Reverse accumulation is more efficient for functions that have more inputs than outputs.



# More optimizations

Creating the Jacobian matrices is expensive. Instead, matrix products can be simplified.

## Backpropagation through activationfunction

$$\begin{bmatrix} J_x^E \\ h' \\ h' \\ h' \\ h' \\ h' \\ J_x^{h(1)} \end{bmatrix} = \begin{bmatrix} J_x^E \end{bmatrix} * \begin{bmatrix} h' & h' & h' & h' & h' & h' \end{bmatrix}$$

elementwise multiplication

## Backpropagation through linearlayer

$$\begin{bmatrix} J_x^E \\ \mathbf{x} \\ \mathbf{x} \\ \dots \\ \mathbf{x} \\ \mathbf{x} \end{bmatrix} J_W^{\text{linear}(\mathbf{x}, \mathbf{W}, \mathbf{b})} = \begin{bmatrix} J_x^{E\top} \end{bmatrix} \begin{bmatrix} \mathbf{x} \end{bmatrix}$$

needs to be flattened



# Advantage of backpropagation

Backpropagation is a form of reverse automatic differentiation, where the Jacobi matrix is not explicitly computed. The gradient is propagated by simpler equivalent operations.

## Jacobian formulation

$$\begin{aligned}
 J_{\mathbf{W}^{(1)}}^E &= J_{\mathbf{x}}^L J_{\mathbf{x}}^{\text{linear}} J_{\mathbf{x}}^h J_{\mathbf{W}}^{\text{linear}} \\
 \left[ \text{---} \right] &= \left[ \text{---} \right] \left[ \text{---} \right] \left[ \begin{array}{c} \square \\ \square \\ \square \\ \square \\ \square \end{array} \right] \left[ \begin{array}{c} \text{---} \\ \text{---} \\ \dots \\ \text{---} \end{array} \right]
 \end{aligned}$$

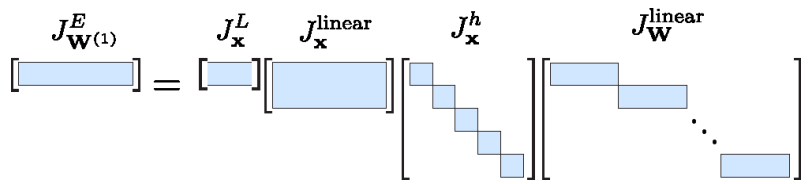
## Compact backpropagation

$$\begin{aligned}
 J_{\mathbf{W}^{(1)}}^E &= J_{\mathbf{x}}^L * \left[ \mathbf{W}^{(2)} \right] * \left[ h' h' h' h' h' h' \right]^T \left[ \mathbf{x} \right] \\
 \left[ \text{---} \right] &= \left[ \text{---} \right] * \left[ \text{---} \right] * \left[ \text{---} \right]^T \left[ \text{---} \right]
 \end{aligned}$$

# Vanishing gradients problem

The objective function

$$O(x, y) = L(h(1(h(1(x, W^{(1)}))), W^{(2)})), y)$$

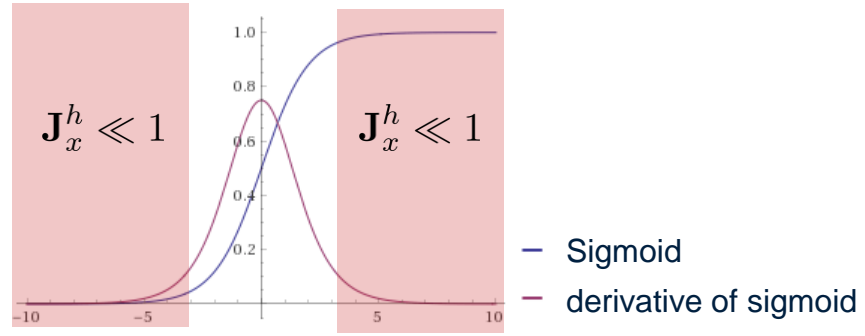


The gradient of O with respect to  $W^{(2)}$

$$O(x, y) = \mathbf{J}_x^L \mathbf{J}_x^h \mathbf{J}_{W^{(2)}}^1$$

The gradient of O with respect to  $W^{(1)}$

$$O(x, y) = \mathbf{J}_x^L \mathbf{J}_x^h \mathbf{J}_x^1 \mathbf{J}_x^h \mathbf{J}_{W^{(1)}}^1$$



*Use ReLU rather than sigmoid in deep neural networks!*

The gradient vanishes exponentially with respect to the number of layers if  $\mathbf{J}_x^h < 1$

# Mini break

Find a team partner for your course project



# Input and output normalization

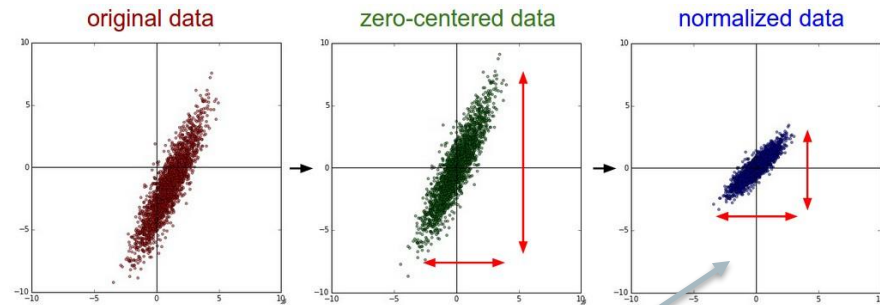
Goal: Normalize input and output variables to have  $\mu=0$  and  $\sigma=1$

$$\tilde{\mathbf{x}} = \frac{\mathbf{x} - \mu}{\sigma}$$

- For an image, normalize each pixel by the std and mean color (averaged over the **training set**)

Related to data whitening

- whitening transforms a random vector to have zero mean and unit diagonal covariance
- by contrast, the default normalization for deep learning is element wise, neglecting dependency
  - the resulting covariance is not diagonal!



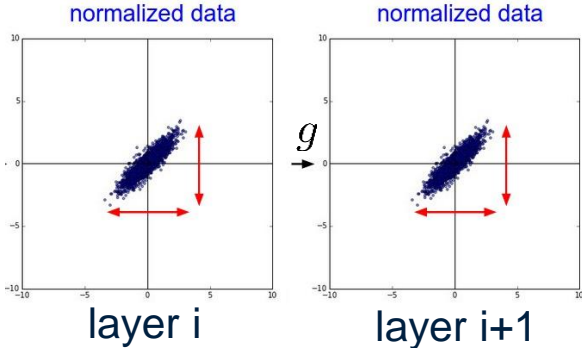
<http://cs231n.github.io/neural-networks-2/>

# Neural network initialization

Goal: preserve mean and variance through the network

- Assume that the input is a random variable with  $\text{var}(x) = 1$  and  $\text{mean}(x) = 0$
- Derive the function  $g$  that describes the change of variance and mean between layers

$$\begin{pmatrix} \tilde{\mu} \\ \tilde{\nu} \end{pmatrix} = g \begin{pmatrix} \mu \\ \nu \end{pmatrix}$$



Initialize the neural network weights (weights of linear layers) such that  $g$  is the identity function

- For the linear neuron with  $K$  incoming neurons

$$\begin{aligned} \text{Var}(\mathbf{w} \cdot \mathbf{x}) &= \sum_{i=1}^K \text{Var}(\mathbf{w}_i) \text{Var}(\mathbf{x}) \\ &= K \text{Var}(\mathbf{w}_i) \text{Var}(\mathbf{x}) \\ &= K \text{Var}(\mathbf{w}_i) \end{aligned} \quad \Rightarrow \quad \text{Var}(\mathbf{w}) = \frac{1}{K}$$

$$\text{Var}(x + y) = \text{Var}(x) + \text{Var}(y)$$

$$\text{Var}(xy) = \text{Var}(x) \text{Var}(y)$$

for variables with zero mean

Xavier Initialization: Initialize with samples from a (Gaussian) distribution with  $\text{std} = \sqrt{1/K}$

# Neural network initialization II

The activation function changes the distribution

- the mean of ReLU(x) is nonzero
  - hence, the *variance of product* equation does not apply
  - instead, it holds

$$\text{Var}(xy) = \text{Var}(x)E(y^2)$$

for  $x$  zero mean and  $y$  arbitrary

- and, assuming that  $y$  is from a symmetric distribution,

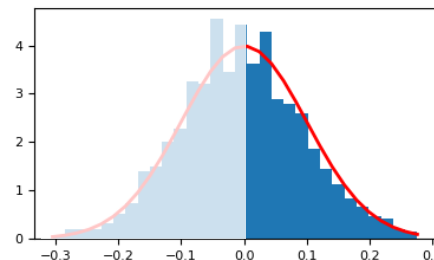
$$E(\text{ReLU}(\mathbf{y})^2) = \frac{1}{2}\text{Var}(\mathbf{y})$$

- the variance transformation of linear layer + activation becomes

$$\text{Var}(\mathbf{w} \cdot \text{ReLU}(\mathbf{x})) = \frac{K}{2}\text{Var}(\mathbf{w}_i) \quad \Rightarrow \quad \text{Var}(\mathbf{w}) = \frac{2}{K}$$

He et al. Initialization: Initialize with samples from a (Gaussian) distribution with  $\text{std} = \sqrt{2/K}$

~~$\text{Var}(xy) = \text{Var}(x)\text{Var}(y)$   
for variables with zero mean~~



# Batch normalization

[Batch Normalization: Accelerating Deep Network Training ~~by Reducing Internal Covariate Shift~~]

- Normalize after each linear + activation function
  - normalize across minibatch, to have  $\mu=0$  and  $\sigma=1$

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

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2$$

- Strict normalization reduces performance, hence, add back a learnable offset and scale

$$y_i \leftarrow \gamma \hat{x}_i + \beta$$

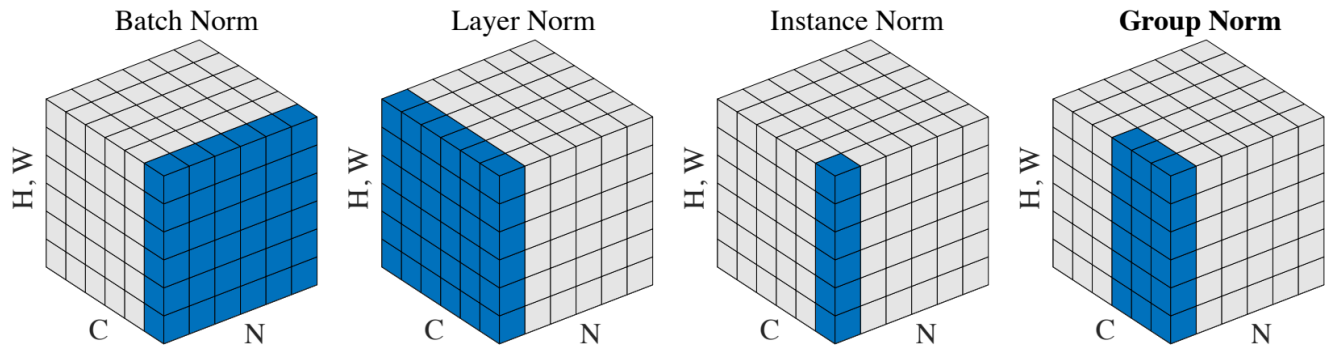
- What if we only have a single image at inference time?
  - Re-apply mean and variance recorded during training (using exponential moving average)

# Batch normalization effect and variants

What is the benefit of first normalizing and then ‘denormalizing’?

- noise from other images regularizes
- it separates learning of the variance (scale) and bias (offset) from the values itself
- Empirical: training deeper networks, with sigmoid activation, higher learning rate, and faster convergence

Variants normalize over different slices of the feature tensor:



[Wu and He. Group Normalization]



# Regularization

## Dropout

- randomly zero out activations
- re-weight the non-zero ones to maintain the distribution of the unmodified activations
  - induced noise reduces overfitting

## Weight decay

$$\tilde{\mathbf{w}} = (1 - \tau)\mathbf{w} \text{ with } \tau \text{ small}$$

*Weight decay and square prior are equivalent under certain conditions (vanilla SGD without momentum)*

## Prior on neural network weights

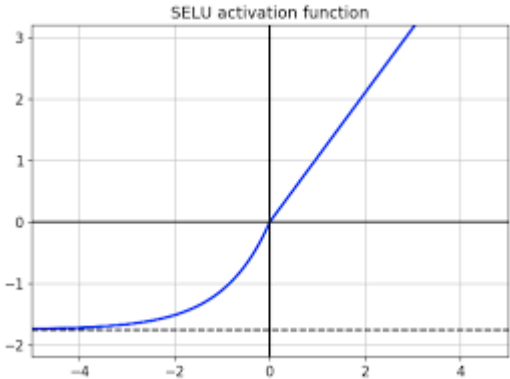
$$\tilde{E}(\mathbf{w}) = E(\mathbf{w}) + \frac{\lambda}{2} \mathbf{w}^2$$

# Self-normalizing neural networks

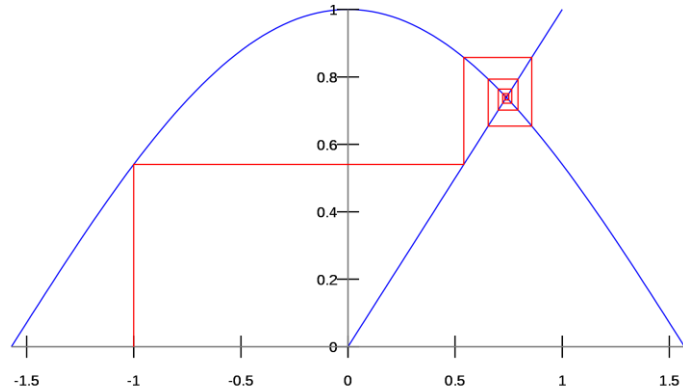
## Self-normalizing Neural Networks

[Klambauer et al.]

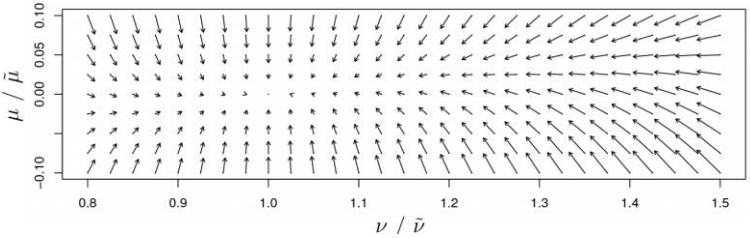
- fixed point enforced by choice of activation function (SELUs)
- stable and attracting fixed point for the function  $g$  that maps mean and variance from one layer to the next
- Possibility to train deep fully connected NNs



$$\text{selu}(x) = \lambda \begin{cases} x & \text{if } x > 0 \\ \alpha e^x - \alpha & \text{if } x \leq 0 \end{cases}$$



Fixed point iterations for  $\cos(x)$



Mapping of the function  $g$  towards  $\mu=0$  and  $v=1$

# Residual networks and skip connections

- Deep networks are hard to train

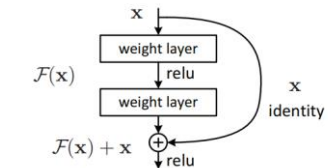
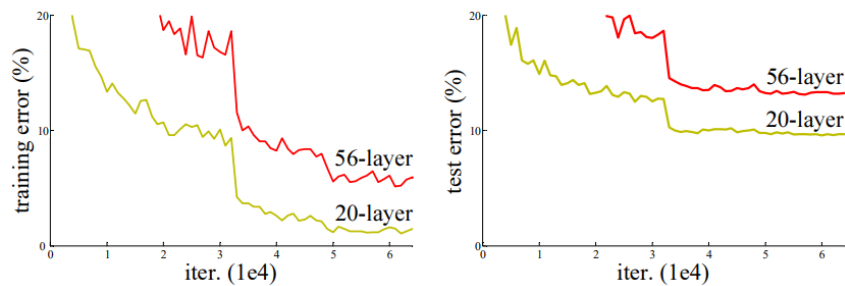


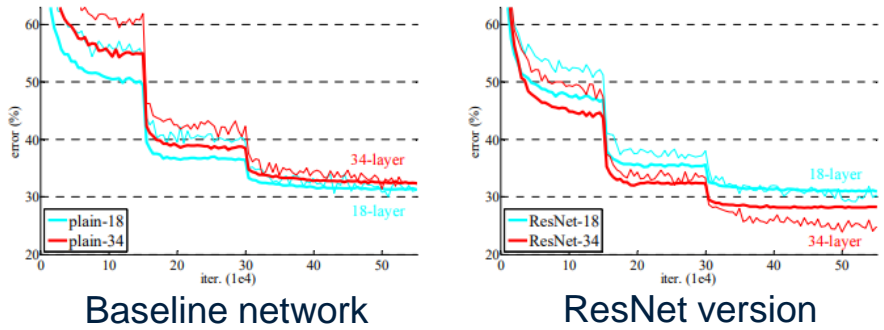
Figure 2. Residual learning: a building block.

- Residual blocks with shortcut/skip connections

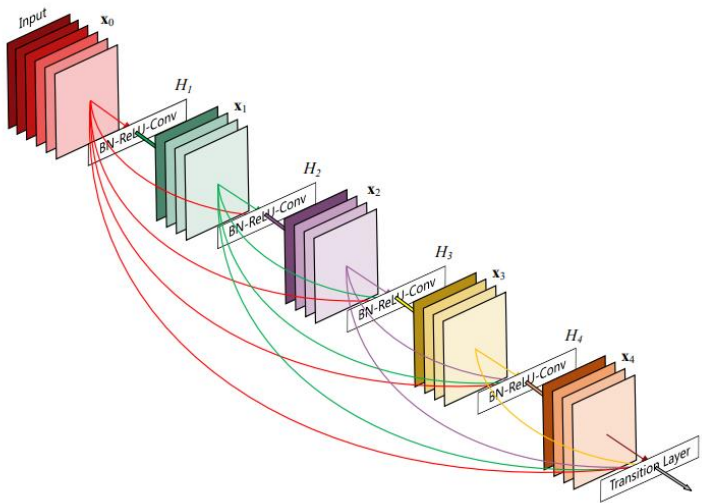
$$y = F(x) + x$$

- no extra parameters
- enables training of deep neural networks

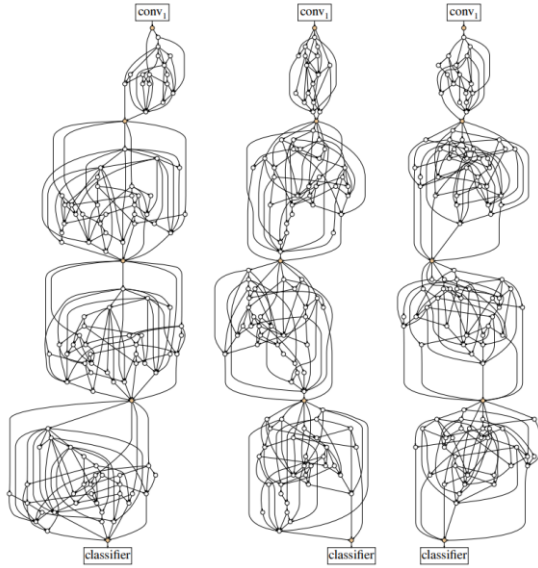
## Image net training



# Other network architectures



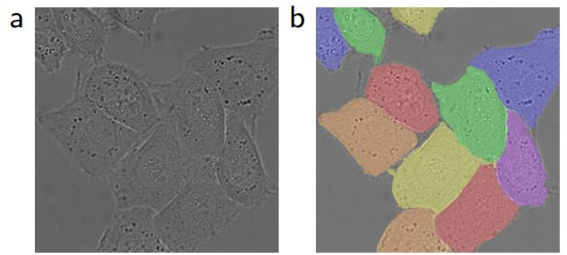
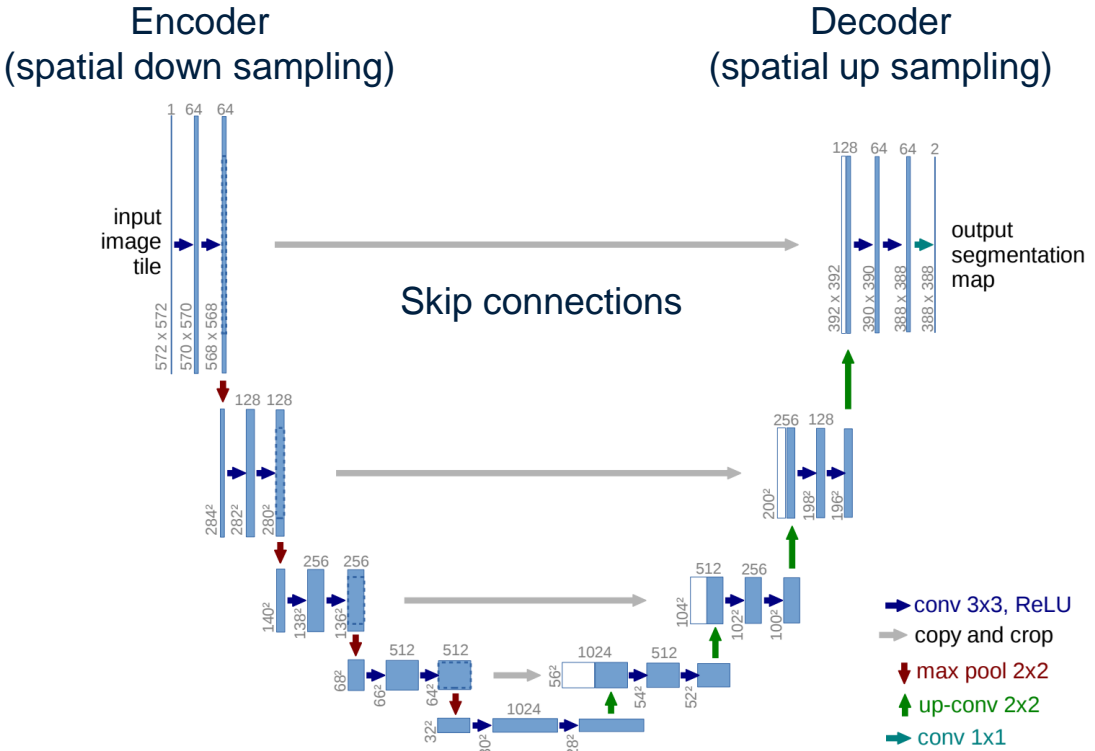
**DenseNet**  
(skip connection to all future layers)



**Randomly wired networks**  
(search for best wiring among candidates)

# U-Net architecture

- Similar input and output resolution
- A global encoding is learned by down sampling (to 32 x 32 px)
- Progressive increase of channels maintains throughput / capacity
- Skip connections preserve details

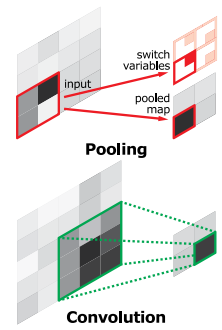


[U-Net: Convolutional Networks for Biomedical Image Segmentation]

# Spatial down and upsampling

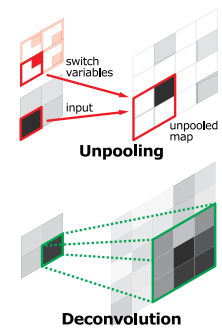
## Spatial downsampling

- max-pooling
- average pooling
- convolution with stride



## Spatial upsampling

- max-unpooling
- (bilinear) interpolation
- deconvolution



[Learning Deconvolution Network for Semantic Segmentation]

# Hidden questions

1. What is the difference between a strong and a weak hypothesis?

2. What is the difference between a strong and a weak hypothesis?

3. What is the difference between a strong and a weak hypothesis?