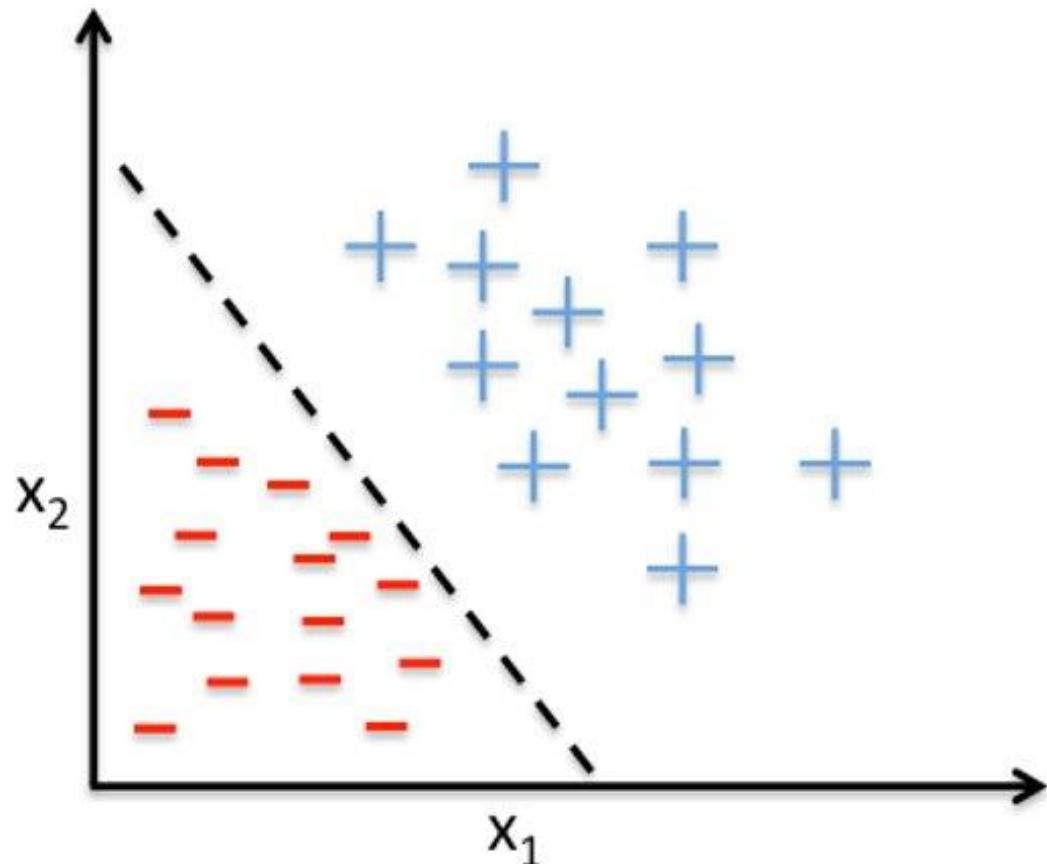


Neural networks

Class 1

- From kernels to Anns: perceptron, svm and kernels
- New hypothesis space (example with invariant functions)
- Shallow and deep networks, terminology;
- Invariant representations, sample complexity, intuition with ReLU.

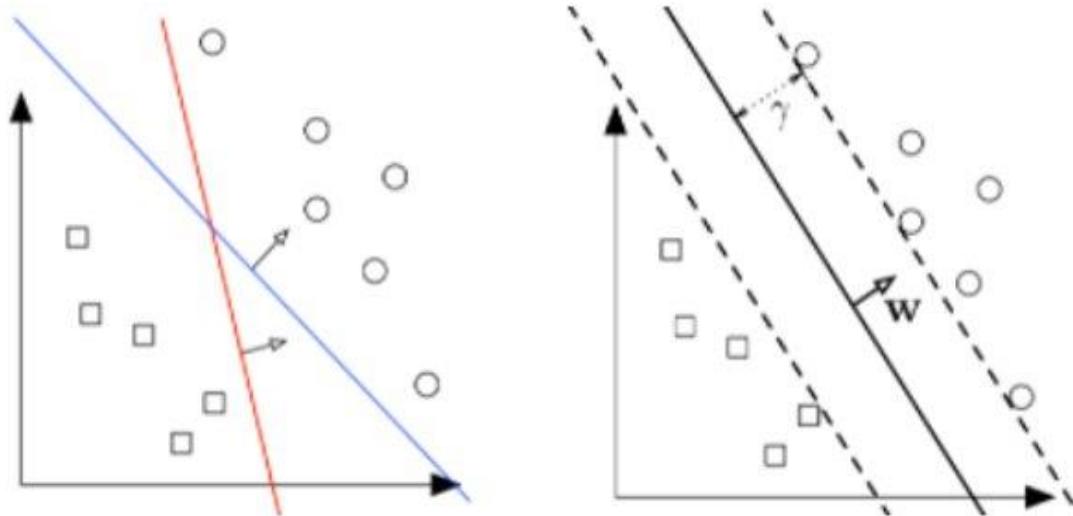
Perceptron: hyperplane separation



$$h(x) = \text{sgn}(w^T x + b)$$

$$L(w) = \frac{1}{N} \sum_i \max(h(x_i)y_i, 0)$$

SVM: maximal margin hyperplane separation

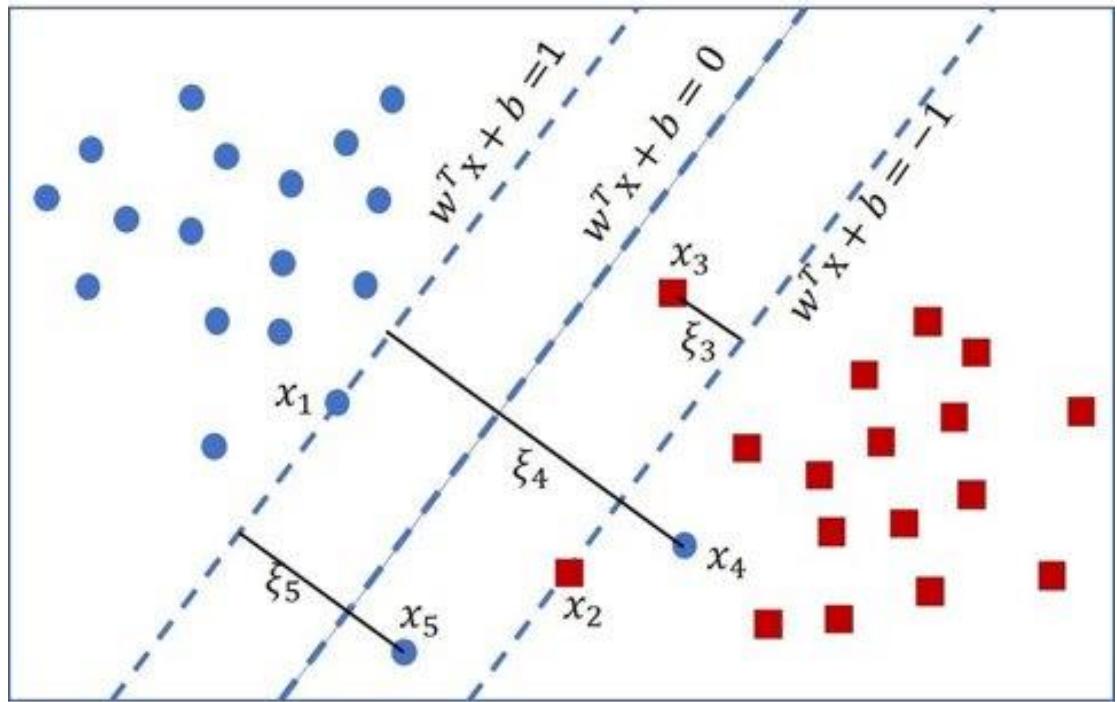


$$\begin{aligned} & \min_w \|w\|_2^2 \\ \text{s.t. } & \forall i \quad y_i h(x_i) \geq 1 \end{aligned}$$

$$\mathcal{L}(w) = \frac{1}{N} \sum_{i=1}^N \max[y_i h(x_i), 0] + \lambda \|w\|_2^2$$

What's 'wrong' in the Loss?

Soft margin SVM and Kernels



$$\begin{aligned} & \min_{\mathbf{w}, b} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^n \xi_i \\ \text{s.t. } & \forall i y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \\ & \forall i \xi_i \geq 0 \end{aligned}$$

Soft margin SVM and Kernels

$$\begin{aligned} & \min_{\mathbf{w}, b} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^n \xi_i \\ \text{s.t. } & \forall i y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \\ & \forall i \xi_i \geq 0 \end{aligned}$$



Dual SVM

$$\begin{aligned} & \min_{\alpha} \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j - \sum_i \alpha_i \\ \text{s.t. } & \sum_i \alpha_i y_i = 0, \quad \alpha_i \in [0, C] \end{aligned}$$

$$\begin{aligned} \mathbf{w} &= \sum_i \alpha_i y_i \mathbf{x}_i & b &= y_i - \sum_j \alpha_j y_j \mathbf{x}_j^\top \mathbf{x}_i \\ f(\mathbf{x}) &= \mathbf{w}^\top \mathbf{x} + b = \sum_i \alpha_i y_i \mathbf{x}_i^\top \mathbf{x} + b & \text{for examples } i \text{ where } 0 < \alpha_i < C \end{aligned}$$

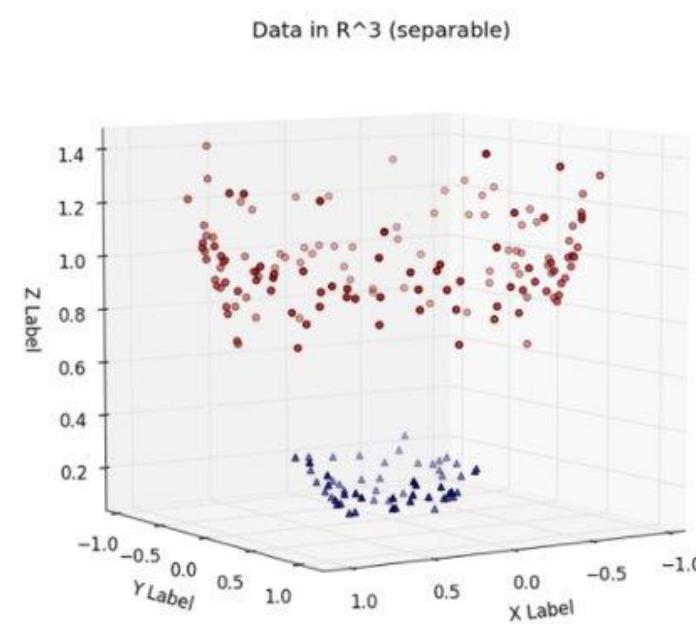
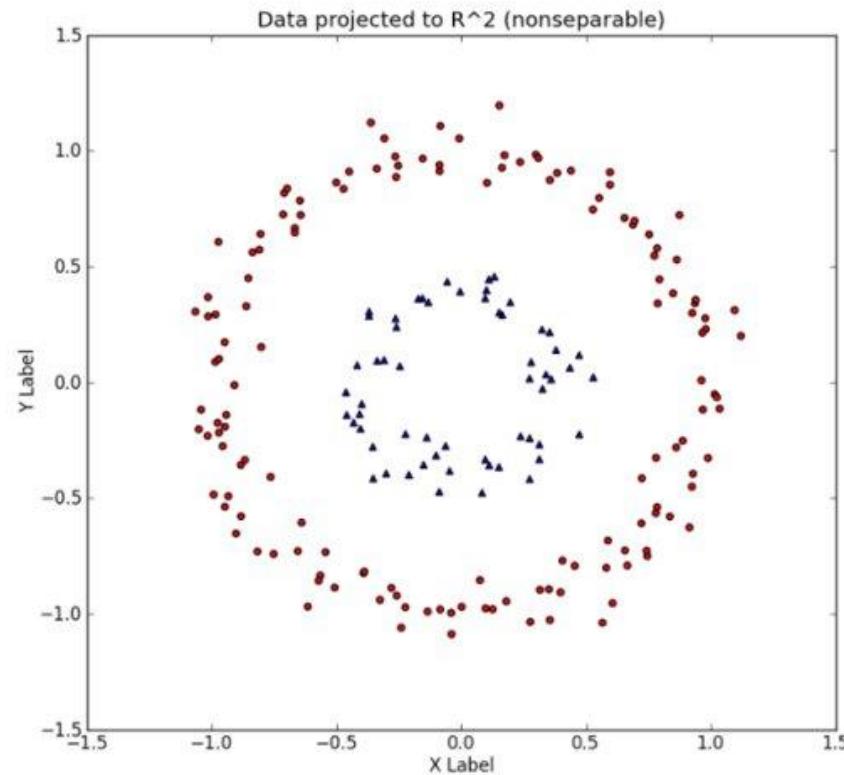
We can express the minimization using dot products of the input

$$K_{i,j} = K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^\top \mathbf{x}_j$$

Kernels

$$x \rightarrow \phi(x)$$

$$K_{i,j} = K(x_i, x_j) = \phi(x_i)^T \phi(x_j)$$



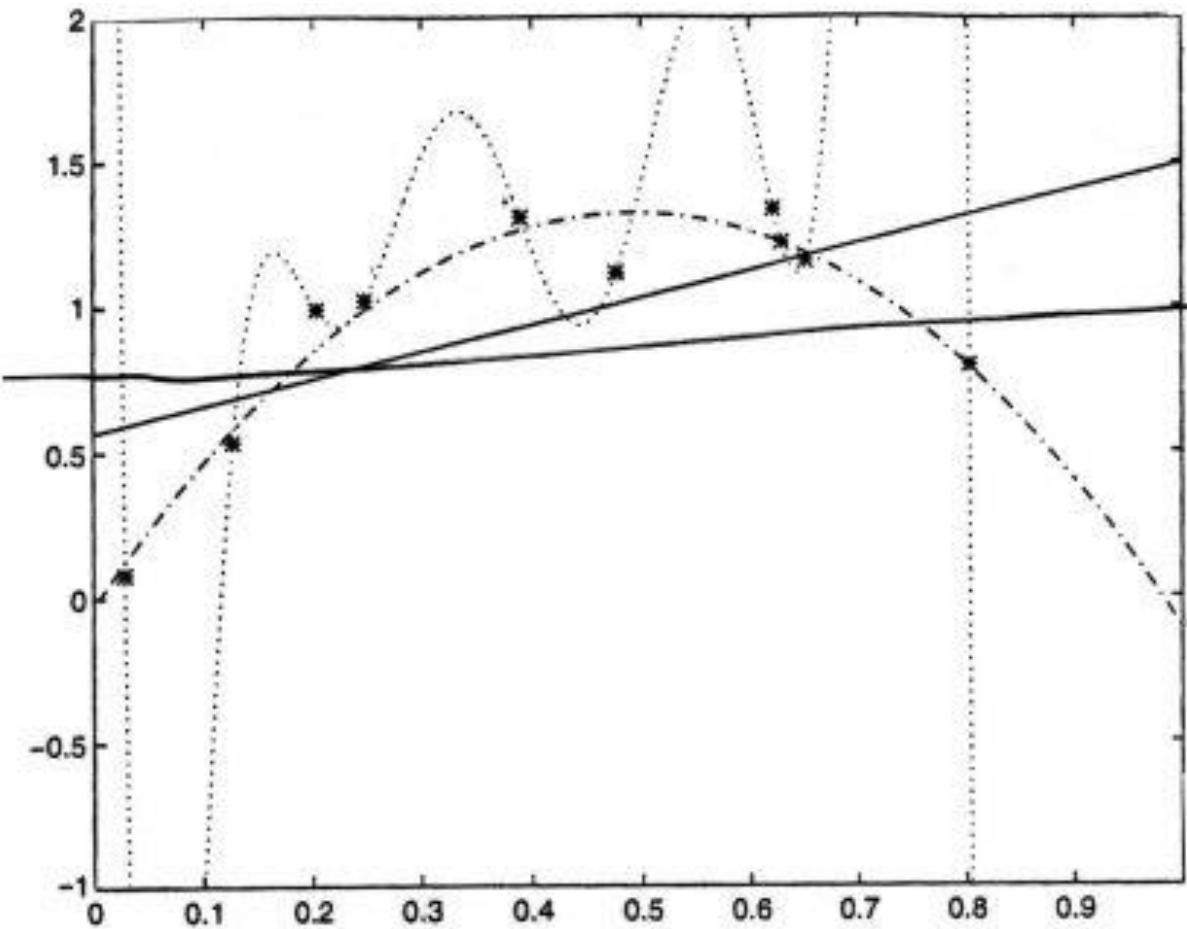
$$(x, y) \rightarrow (x, y, x^2 + y^2)$$

What do we change transforming the coordinates into?

$$(x, y) \rightarrow (x, y, x^2 + y^2)$$

- Through kernels we express our function using a very rich set of features (e.g. we include correlations among the original space input)
- We increase the set of possible learnable functions.

Hypothesis spaces



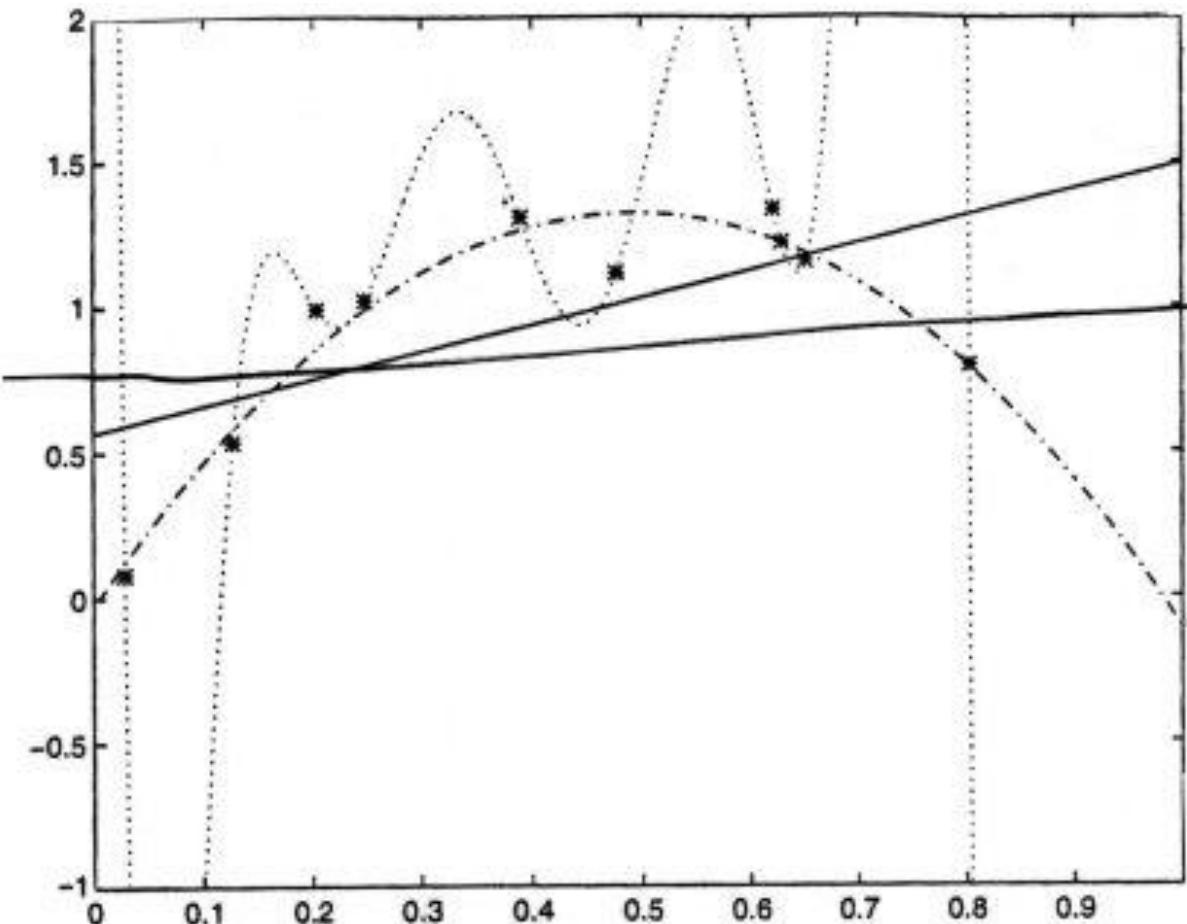
$$H_0 : f(x) = c$$

$$H_1 : f(x) = a + bx$$

$$H_2 : f(x) = a + bx + cx^2$$

What is the kernel
Hypothesis space?

Hypothesis spaces



$$H_0 : f(x) = c$$

$$H_1 : f(x) = a + bx$$

$$H_2 : f(x) = a + bx + cx^2$$

$$H_K : f(x) = \sum_i \alpha_i K(x, x_i)$$

Hypothesis spaces

- Linear separable

$$f(x) = w^T x$$

We learn the separator on data

- Non-linear separable

We learn the separator on a priori hand-crafted features

$$f(x) = w^T \phi(x) = \sum_i \alpha_i K(x, x_i)$$

What is the limitation for kernels?
Hint: how do you choose K ?

- Can we learn the features together with the separator?
- Which features?

$$f(x) = w^T \sigma(W^T x)$$

$$w \in \mathbb{R}^k, \quad W \in \mathbb{R}^{d \times k}, \quad \sigma : \mathbb{R} \rightarrow \mathbb{R}$$

How big is this Hypothesis Space?

Universal approximation theorem

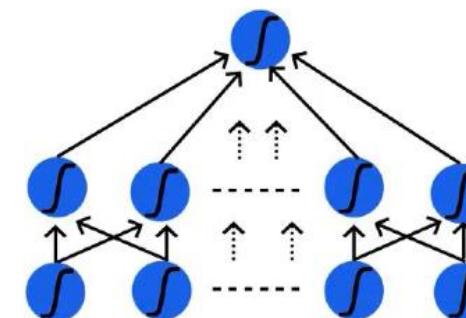
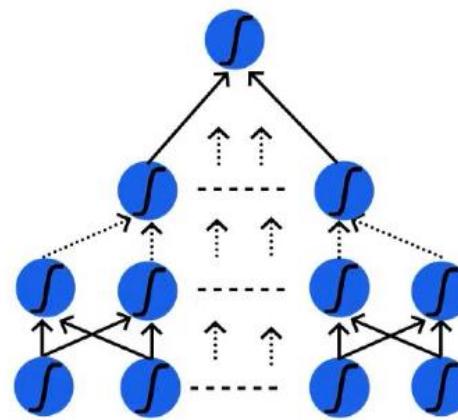
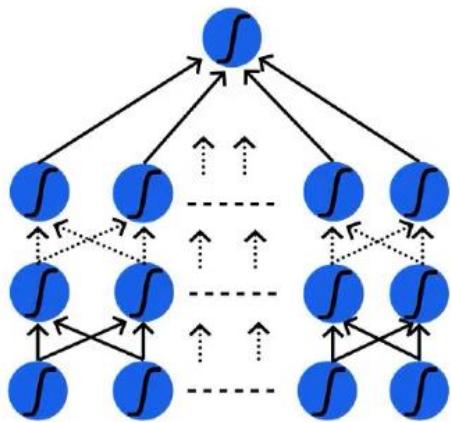
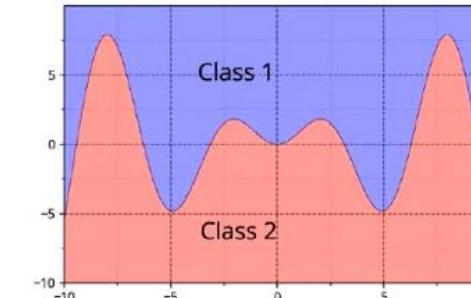
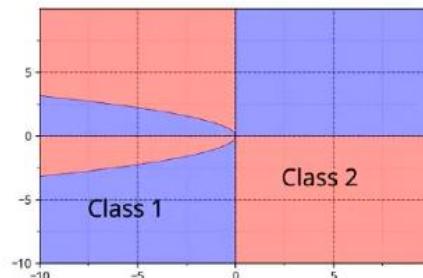
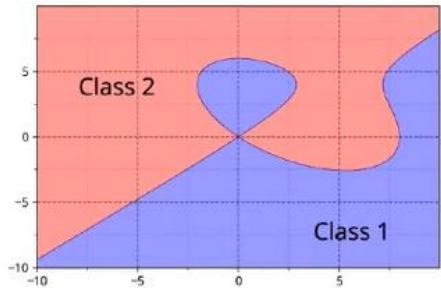
Any **continuous function** defined in the **n-dimensional unit hypercube** may be approximated by a **finite sum** of the type:

$$\sum_{j=1}^N v_j \varphi \left(\vec{\omega}^{(j)} \cdot \vec{x} + b_j \right),$$

wherein $v_j, b_j \in \mathbb{R}$, $\vec{\omega}^{(j)} \in \mathbb{R}^n$, and φ is a **continuous discriminatory function**.

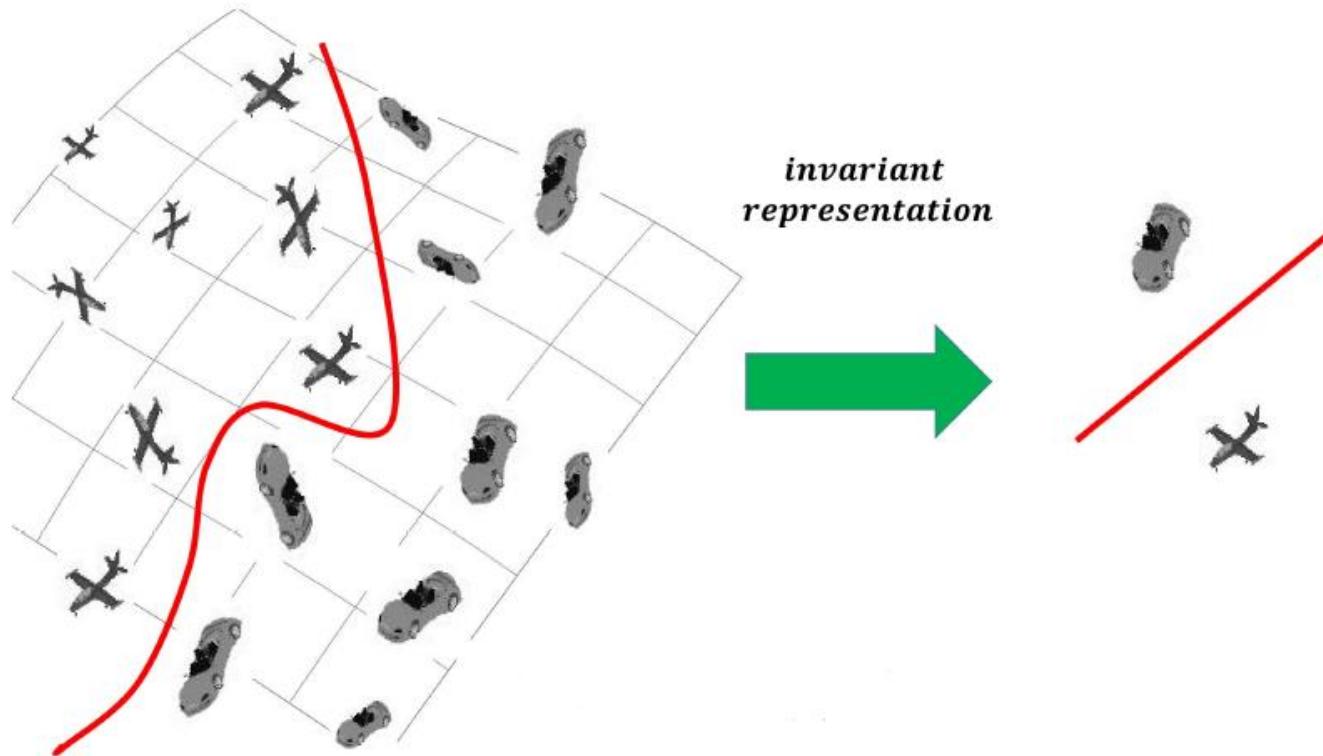
How big is N?

Universal approximation theorem

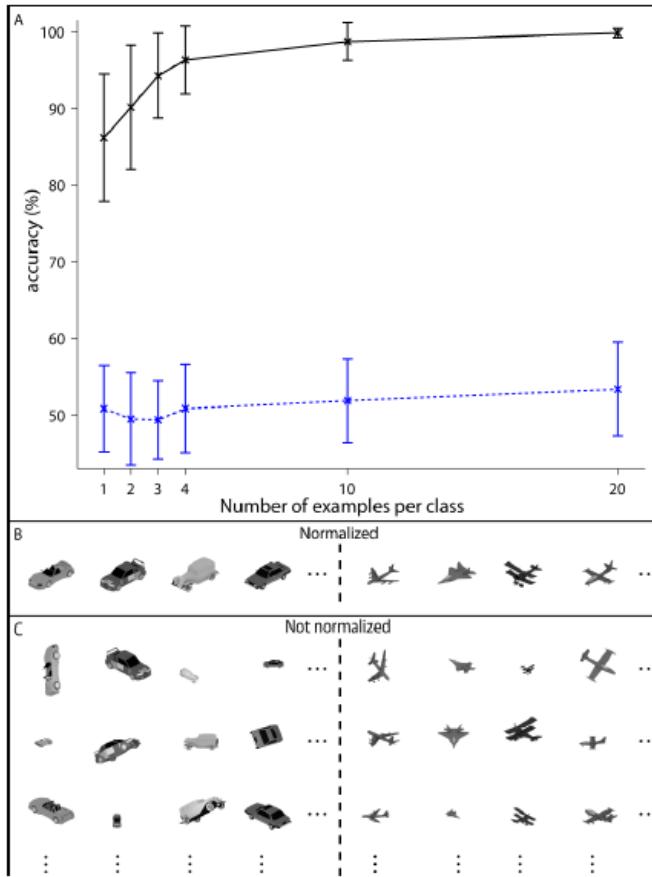


So: choose an ANN and that's it?

Importance of choosing the right space of functions: example with invariance



Importance of choosing the right space of functions: example with invariance

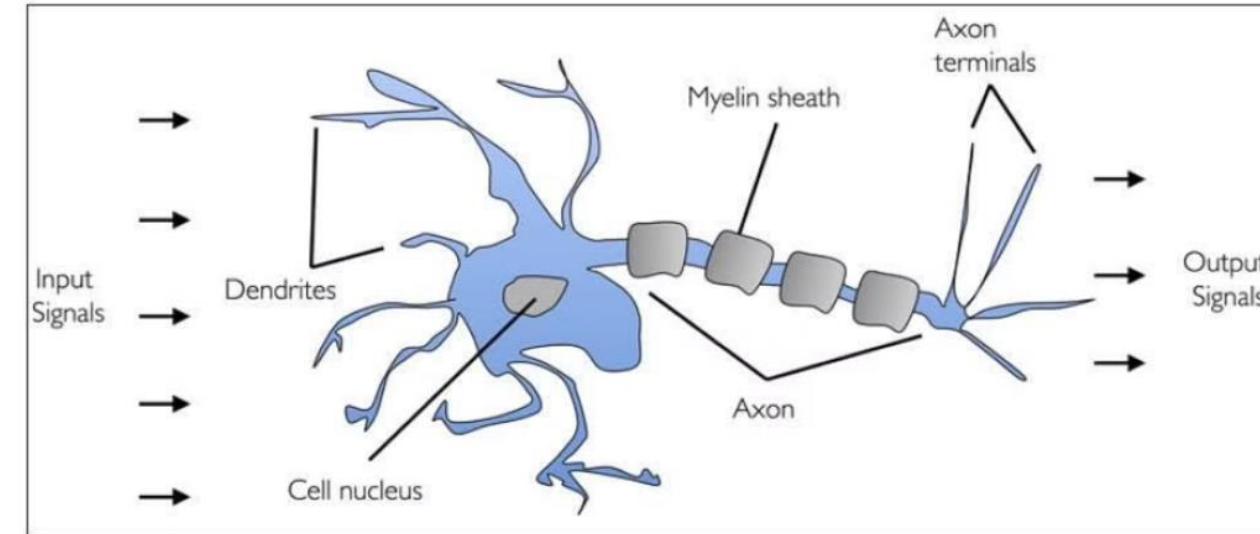
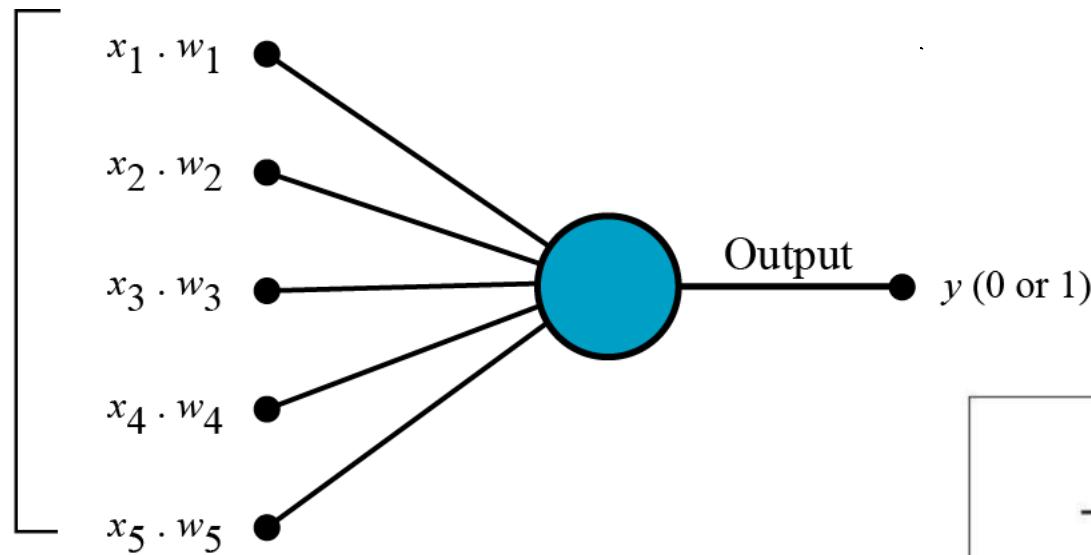


Sample complexity

- The **sample complexity** of a [machine learning](#) algorithm represents the number of training-samples that it needs in order to successfully learn a target function.
- Ex: invariant representations decrease the sample complexity

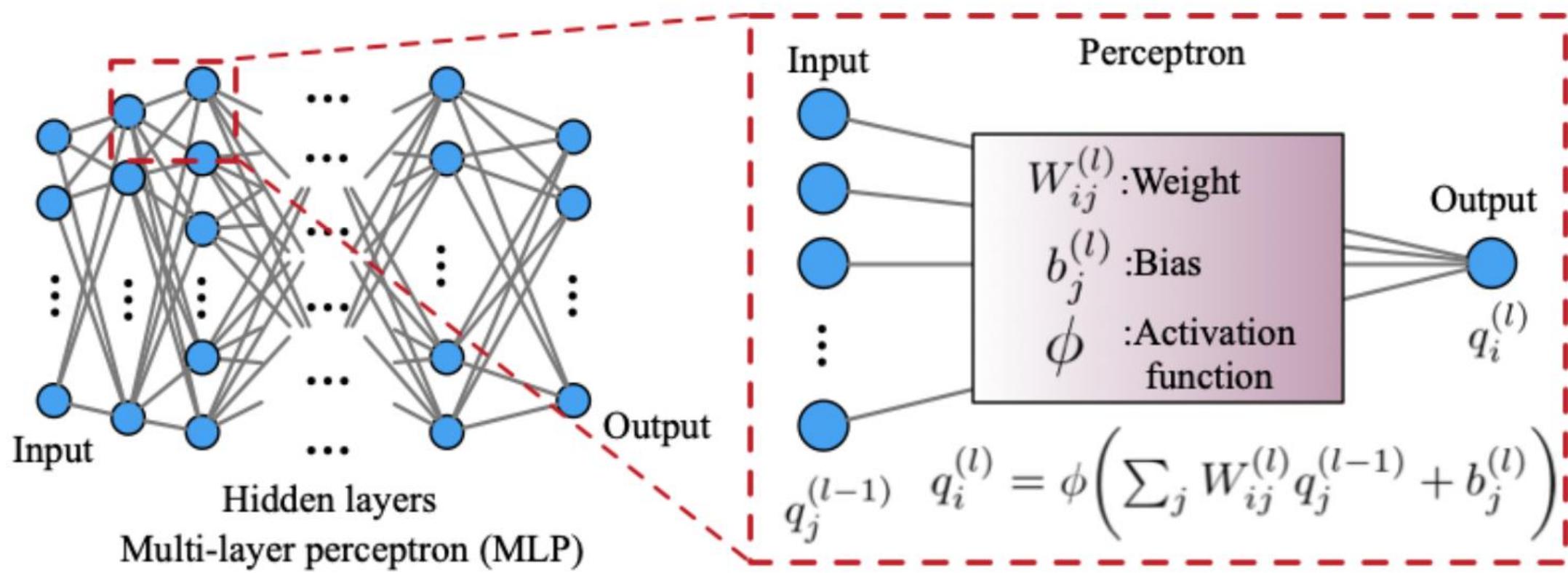
Simplest neural network: perceptron

Inputs

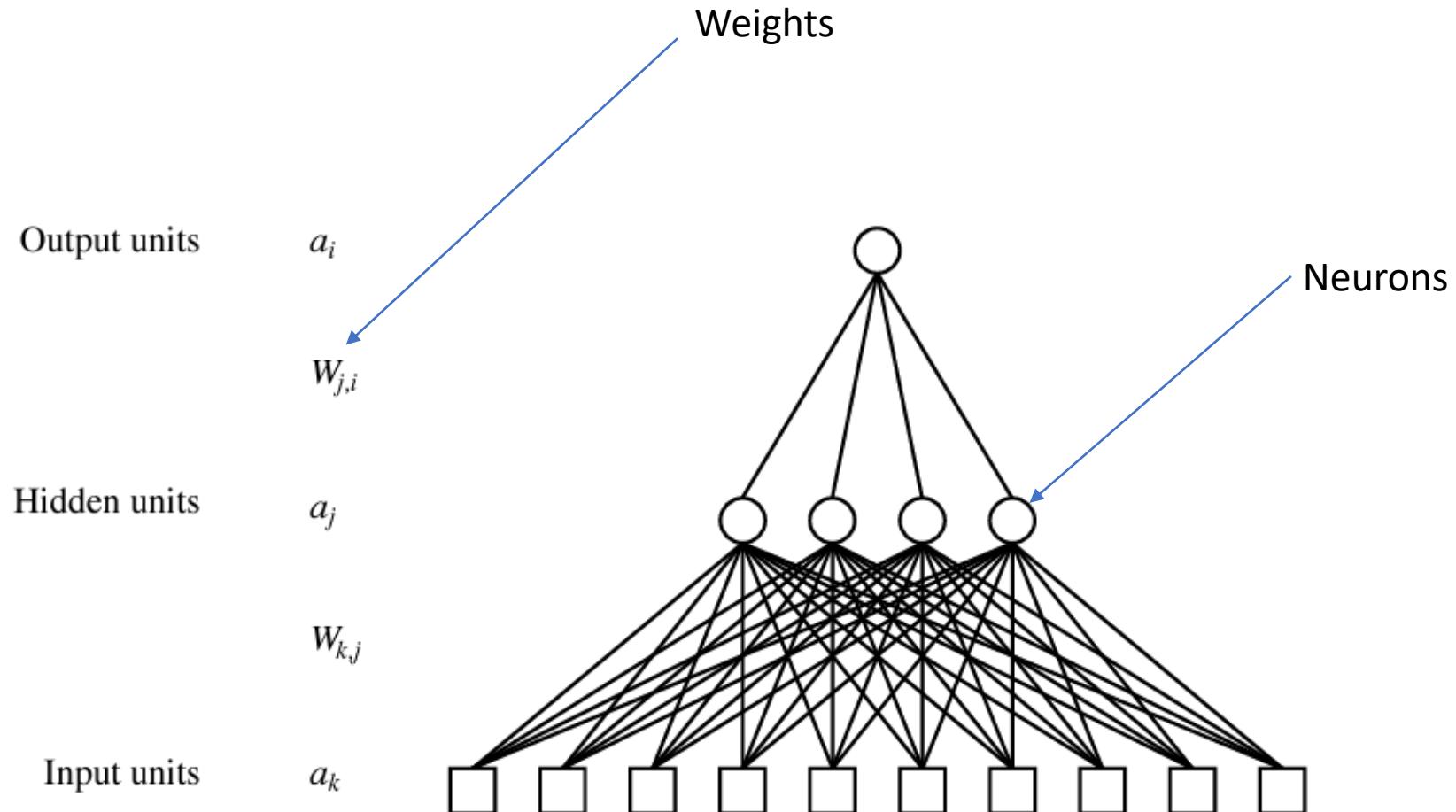


$$f(x) = \text{sgn}(w^T x + b)$$

Generalization: multi layer perceptron



Terminology



(1):Linear filtering

$$y = W^T x + b = \begin{bmatrix} w_1^T x + b_1 \\ w_2^T x + b_2 \\ \vdots \\ w_M^T x + b_M \end{bmatrix}$$

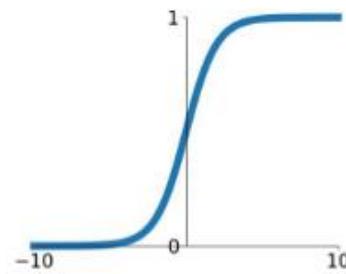
(2):Non-Linear pointwise operation

$$y = \sigma(W^T x + b) = \begin{bmatrix} \sigma(w_1^T x + b_1) \\ \sigma(w_2^T x + b_2) \\ \vdots \\ \sigma(w_M^T x + b_M) \end{bmatrix}$$

A zoo of “activation functions”

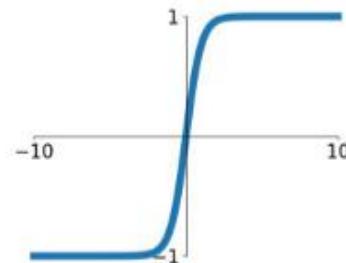
Sigmoid

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



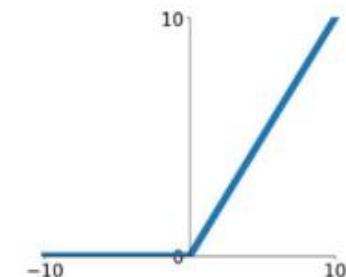
tanh

$$\tanh(x)$$



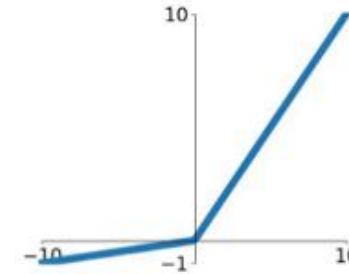
ReLU

$$\max(0, x)$$



Leaky ReLU

$$\max(0.1x, x)$$

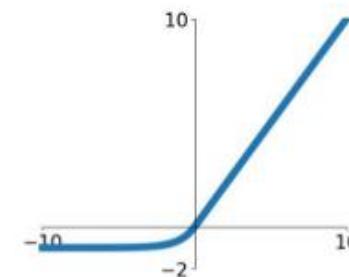


Maxout

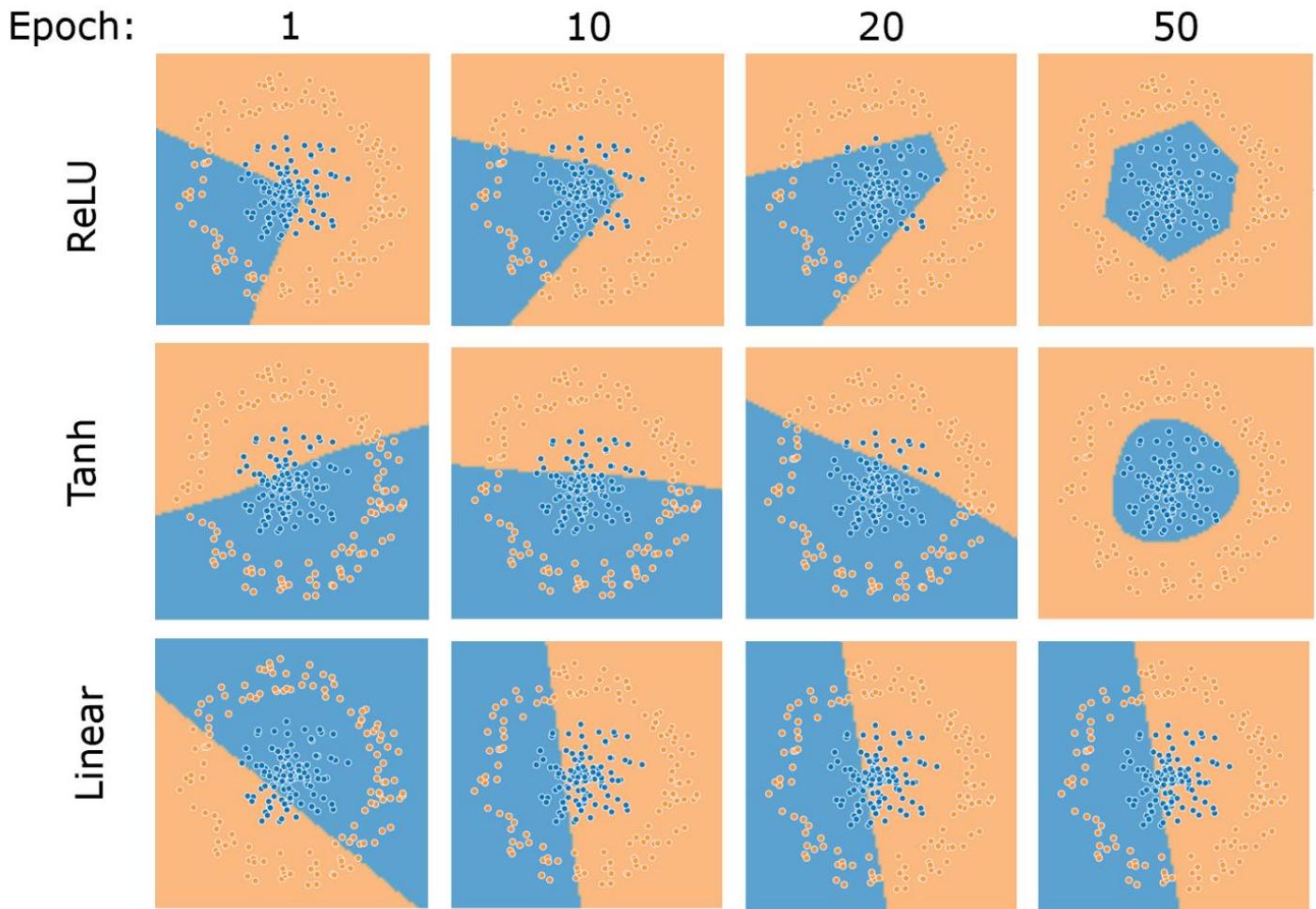
$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$

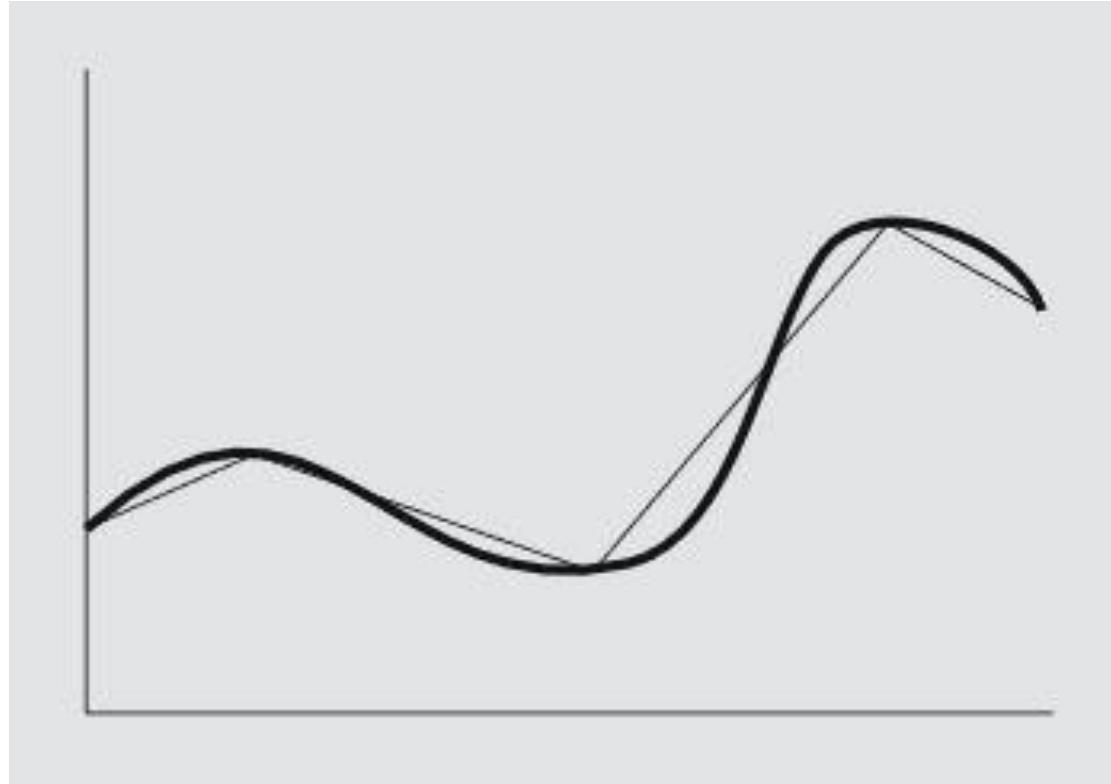


Examples



Intuition
with a ReLu
network

$$f(x) = w^T \sigma(W^T x) = \sum_i w_i \max(\tilde{W}_i^T x + b_i, 0)$$





$$n_1(x) = \text{Relu}(-5x - 7.7)$$

$$n_2(x) = \text{Relu}(-1.2x - 1.3)$$

$$n_3(x) = \text{Relu}(1.2x + 1)$$

$$n_4(x) = \text{Relu}(1.2x - .2)$$

$$n_5(x) = \text{Relu}(2x - 1.1)$$

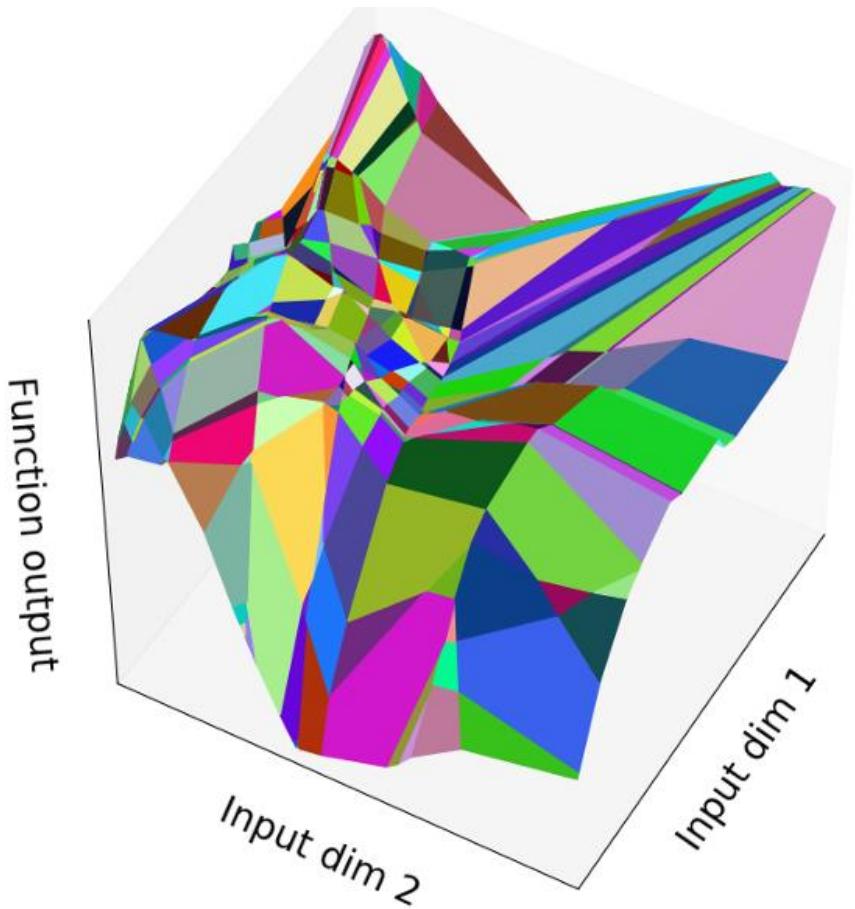
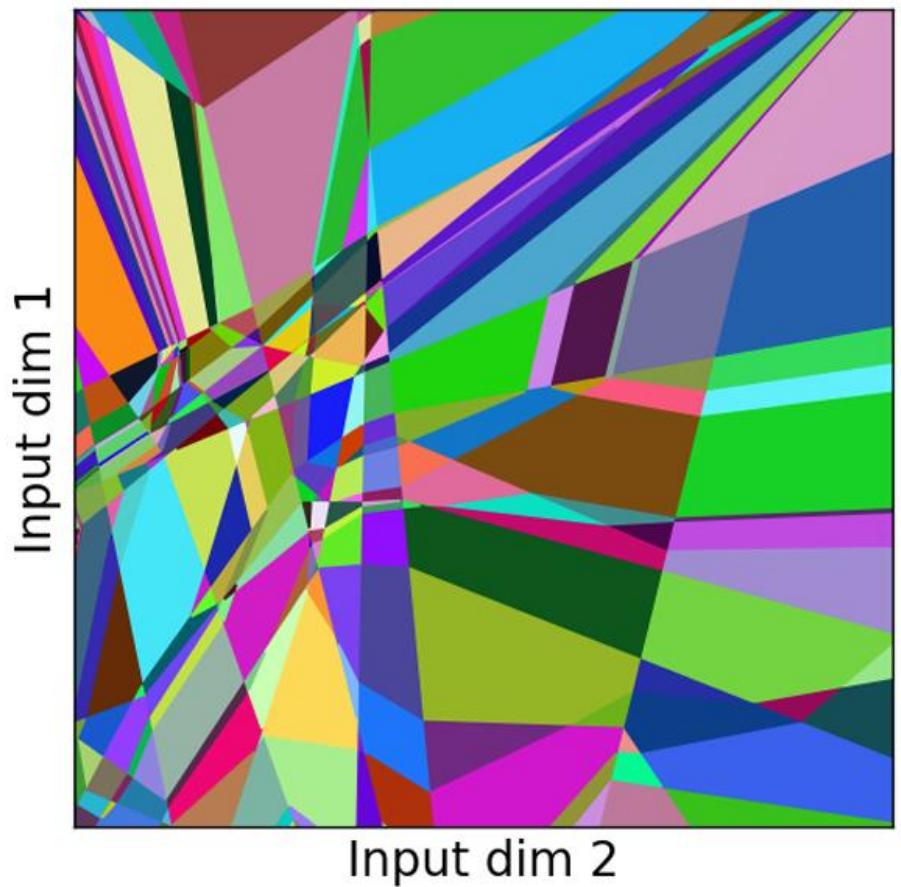
$$n_6(x) = \text{Relu}(5x - 5)$$

$$\begin{aligned}Z(x) = & -n_1(x) - n_2(x) - n_3(x) \\& + n_4(x) + n_5(x) + n_6(x)\end{aligned}$$

2 Dimensional

$$f(x) = w^T \sigma(W^T x) = \sum_i w_i \max(\tilde{W}_i^T x + b_i, 0)$$





Loss functions examples: penalizing errors

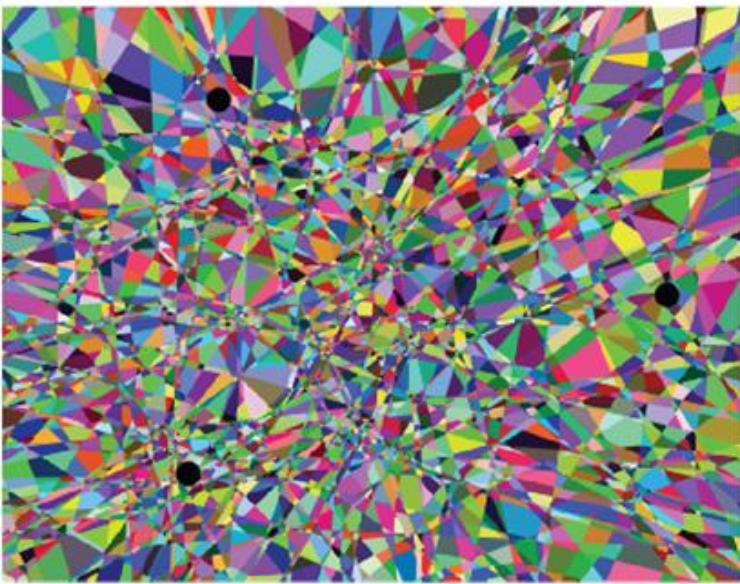
$$\mathcal{L}(w, b) = \frac{1}{N} \sum_i (y_i - (w^T x + b))^2$$

What is the main difference from an optimization point of view?

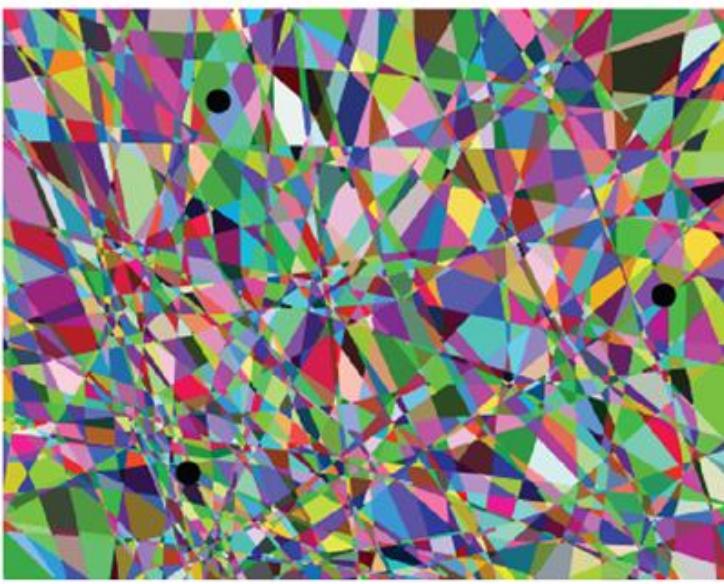
$$\mathcal{L}(w, W) = \frac{1}{N} \sum_{i=1}^N \left(y_i - w^T \sigma(W^T x) \right)^2$$

The Loss is minimized by gradient descent

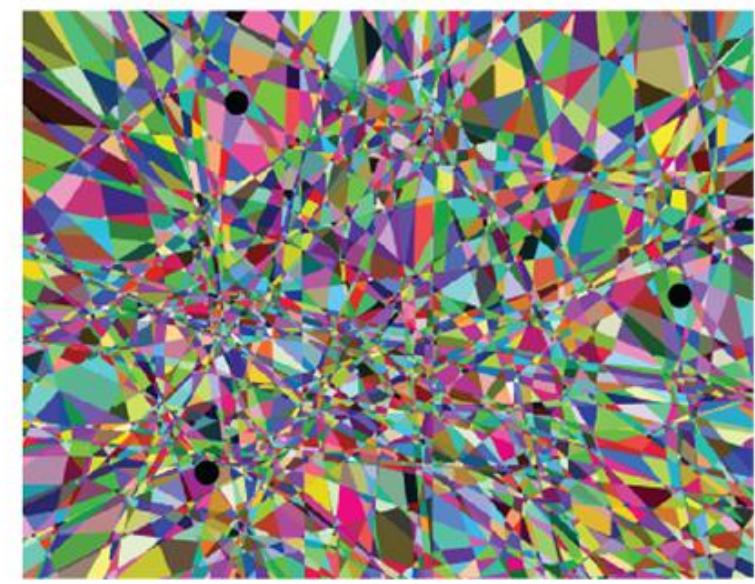
Epoch 0: 9744 regions



Epoch 1: 4196 regions



Epoch 20: 8541 regions

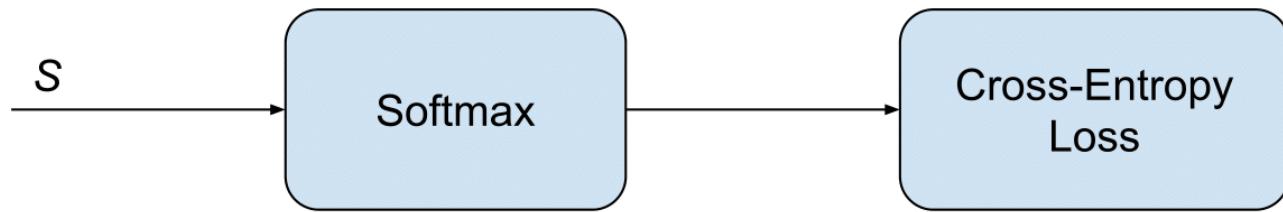


Other examples

\mathcal{L}_1	L ₁ loss	$\ \mathbf{y} - \mathbf{o}\ _1$
\mathcal{L}_2	L ₂ loss	$\ \mathbf{y} - \mathbf{o}\ _2^2$
$\mathcal{L}_1 \circ \sigma$	expectation loss	$\ \mathbf{y} - \sigma(\mathbf{o})\ _1$
$\mathcal{L}_2 \circ \sigma$	regularised expectation loss ¹	$\ \mathbf{y} - \sigma(\mathbf{o})\ _2^2$
$\mathcal{L}_\infty \circ \sigma$	Chebyshev loss	$\max_j \sigma(\mathbf{o})^{(j)} - \mathbf{y}^{(j)} $
hinge	hinge [13] (margin) loss	$\sum_j \max(0, \frac{1}{2} - \hat{\mathbf{y}}^{(j)} \mathbf{o}^{(j)})$
hinge ²	squared hinge (margin) loss	$\sum_j \max(0, \frac{1}{2} - \hat{\mathbf{y}}^{(j)} \mathbf{o}^{(j)})^2$
hinge ³	cubed hinge (margin) loss	$\sum_j \max(0, \frac{1}{2} - \hat{\mathbf{y}}^{(j)} \mathbf{o}^{(j)})^3$
log	log (cross entropy) loss	$-\sum_j \mathbf{y}^{(j)} \log \sigma(\mathbf{o})^{(j)}$
log ²	squared log loss	$-\sum_j [\mathbf{y}^{(j)} \log \sigma(\mathbf{o})^{(j)}]^2$
tan	Tanimoto loss	$-\sum_j \sigma(\mathbf{o})^{(j)} \mathbf{y}^{(j)}$
D _{CS}	Cauchy-Schwarz Divergence [3]	$\frac{\ \sigma(\mathbf{o})\ _2^2 + \ \mathbf{y}\ _2^2 - \sum_j \sigma(\mathbf{o})^{(j)} \mathbf{y}^{(j)}}{\ \sigma(\mathbf{o})\ _2 \ \mathbf{y}\ _2}$
		$- \log \frac{\sum_j \sigma(\mathbf{o})^{(j)} \mathbf{y}^{(j)}}{\ \sigma(\mathbf{o})\ _2 \ \mathbf{y}\ _2}$

- Each loss is weighting errors in different ways and its choice depends on the problem
- Suppose e.g. we want to output a vector of probabilities for different classes...

Example of Classification Loss: cross entropy loss and softmax

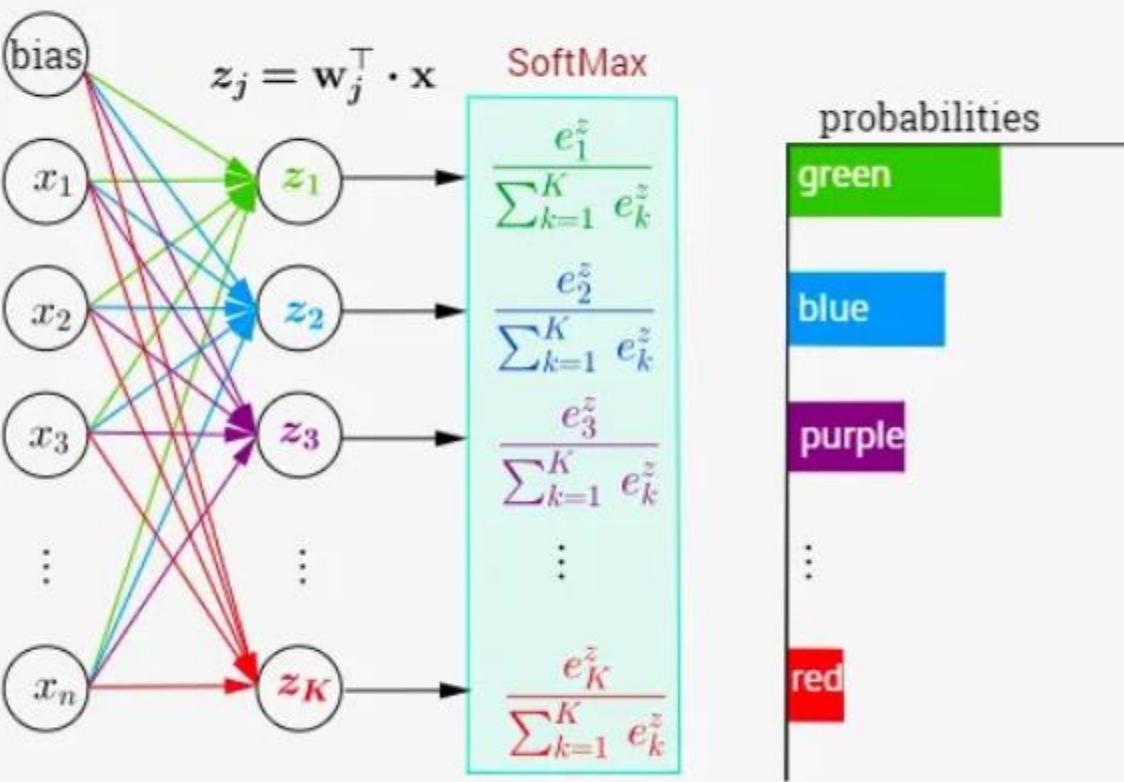


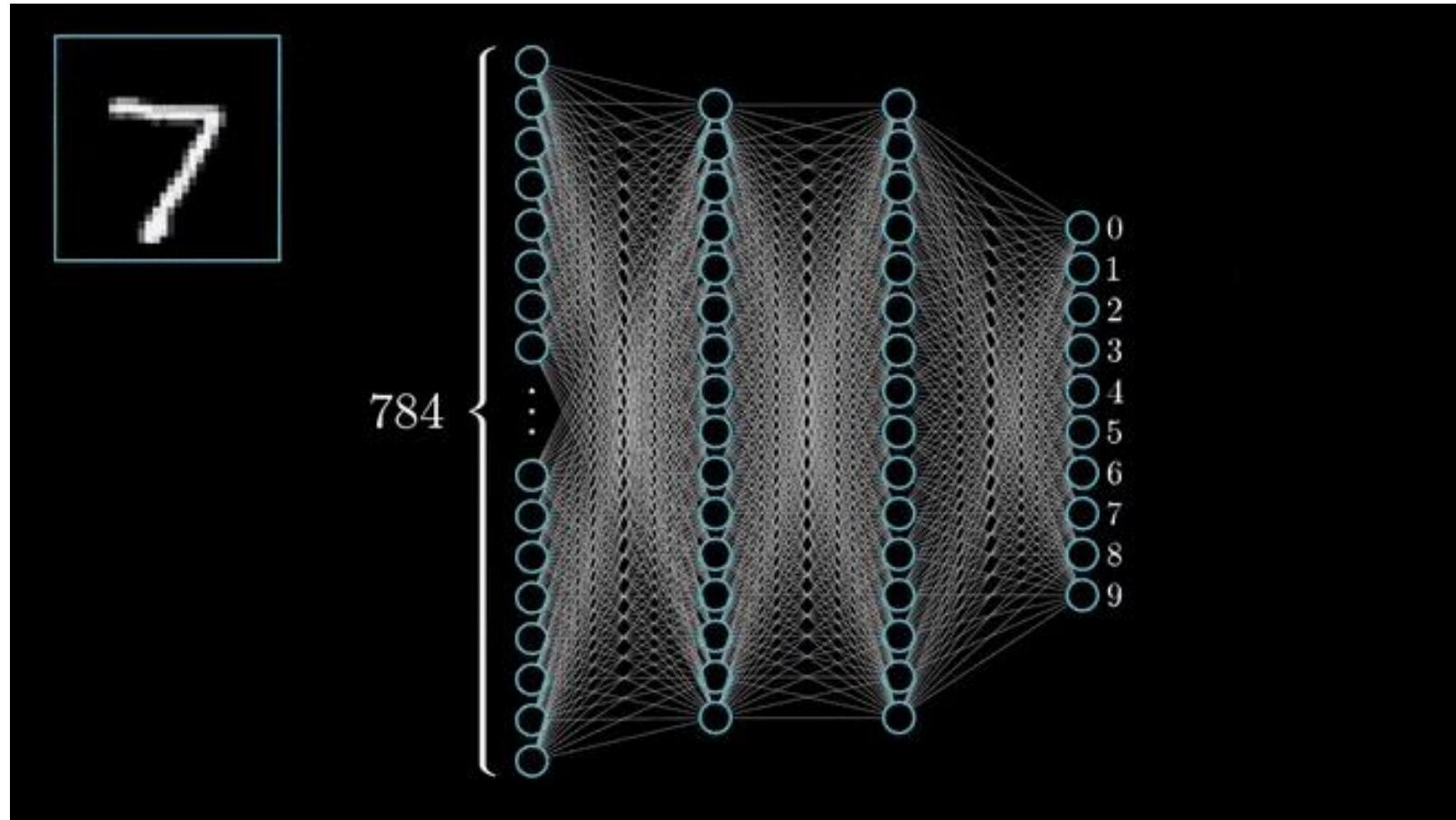
$$\hat{y}_i = \frac{e^{s_i}}{\sum_j e^{s_j}} - \frac{1}{\text{output size}} \sum_{i=1}^{\text{output size}} y_i \cdot \log \hat{y}_i + (1 - y_i) \cdot \log (1 - \hat{y}_i)$$

$$\mathbf{z} = \begin{pmatrix} z_1 \\ z_2 \\ z_3 \\ z_4 \end{pmatrix} \quad \left\{ \begin{array}{c} 1.1 \rightarrow s_i = \frac{e^{z_i}}{\sum_l e^{z_l}} \rightarrow 0.224 \\ 2.2 \rightarrow \qquad \qquad \qquad \rightarrow 0.672 \\ 0.2 \rightarrow \qquad \qquad \qquad \rightarrow 0.091 \\ -1.7 \rightarrow \qquad \qquad \qquad \rightarrow 0.013 \end{array} \right\} \mathbf{s} = \begin{pmatrix} s_1 \\ s_2 \\ s_3 \\ s_4 \end{pmatrix}$$

Multi-Class Classification with NN and SoftMax Function

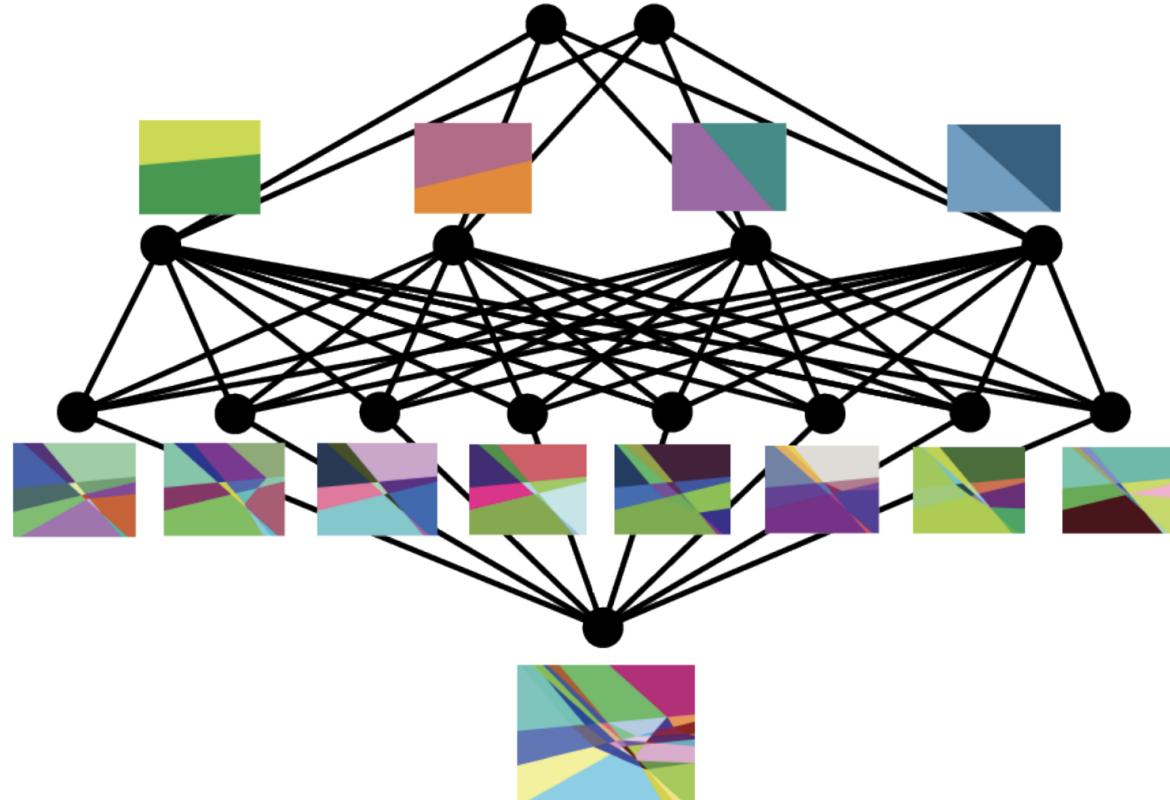
$$\mathbf{z} = \begin{bmatrix} z_1 \\ z_2 \\ z_3 \\ \vdots \\ z_K \end{bmatrix} = \begin{bmatrix} \mathbf{w}_1^\top \\ \mathbf{w}_2^\top \\ \mathbf{w}_3^\top \\ \vdots \\ \mathbf{w}_K^\top \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix}$$





Deep networks

$$f(x) = w^T \sigma(W_2^T \sigma(W_1^T x))$$



Is the number of regions increasing linearly with depth?

Intuition why depth: approximation reusability

From universal approximation theorem you need **exponentially many units** but you increase exponentially the repeated approximations with depth, **reusing** the previous layer approx.

$$f(x) = w^T \sigma(W_L^T (\sigma(W_{L-1}^T \cdots \sigma(W_1^T x) \cdots))$$

Learned features

$$w \in \mathbb{R}^{k_L}, \quad W \in \mathbb{R}^{d \times k_\ell}, \quad \sigma : \mathbb{R} \rightarrow \mathbb{R}$$

How neural networks build up their understanding of images



Edges (layer conv2d0)

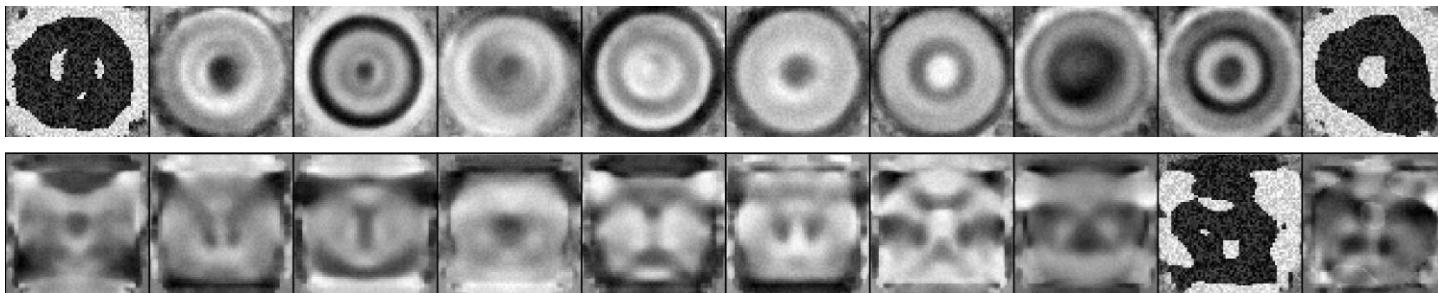
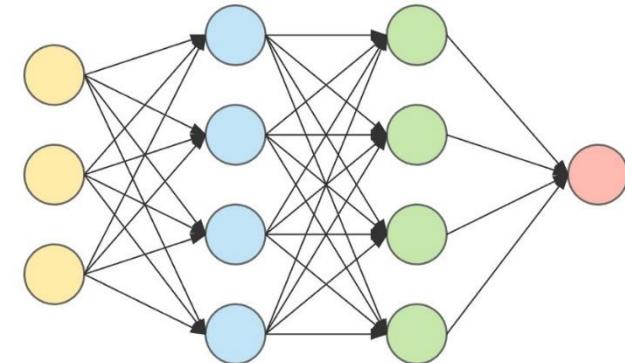
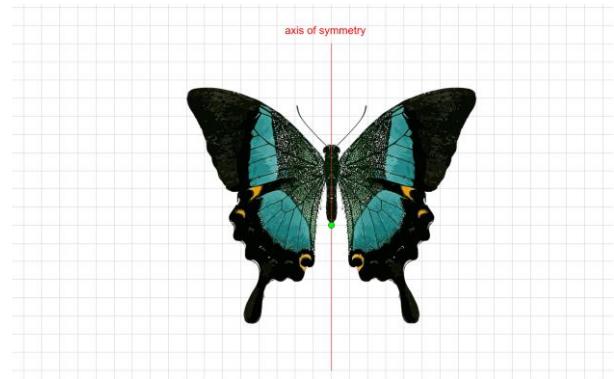
Textures (layer mixed3a)

Patterns (layer mixed4a)

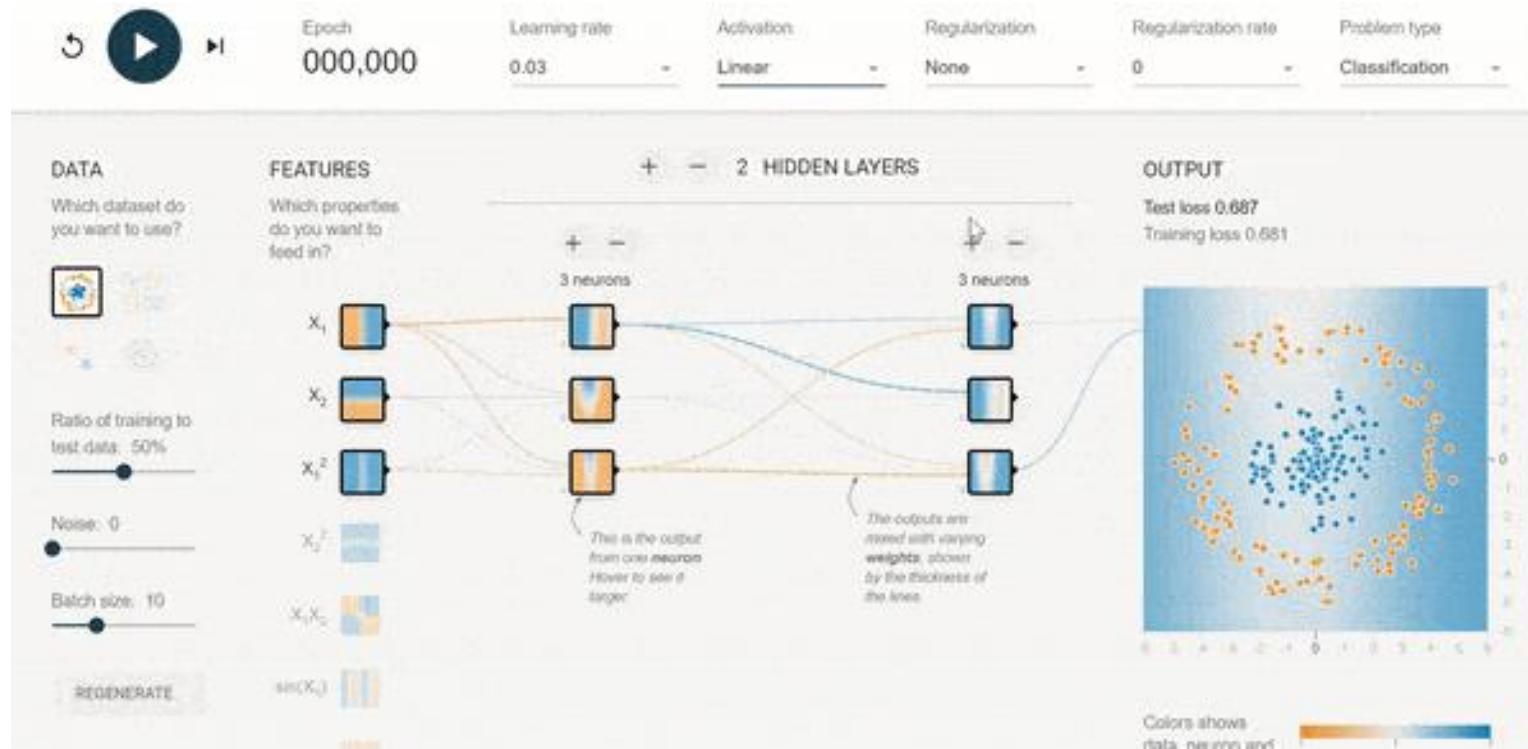
Parts (layers mixed4b & mixed4c)

Objects (layers mixed4d & mixed4e)

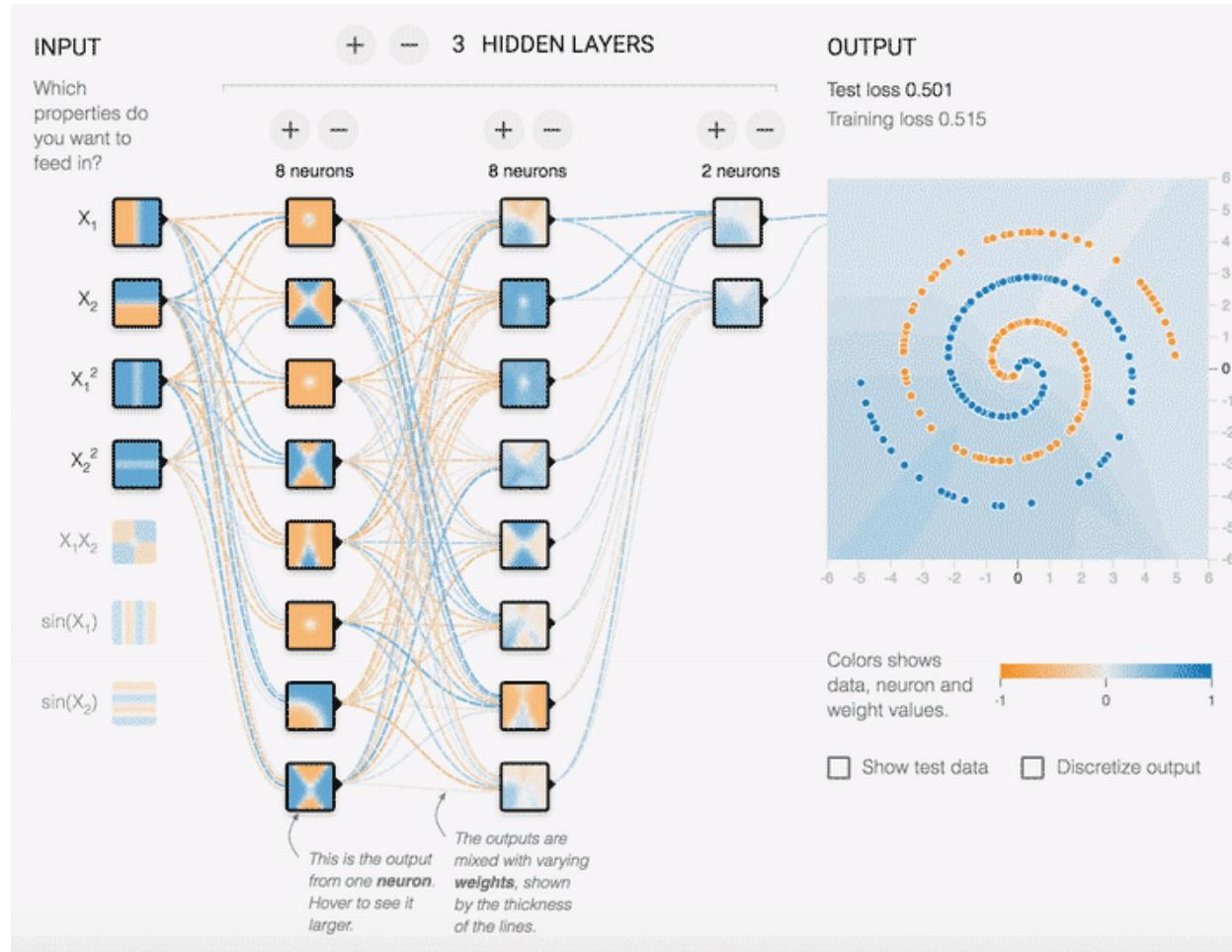
Symmetries/regularities in data/data statistics are reflected in the learned weights of a network.



Compositional functions power



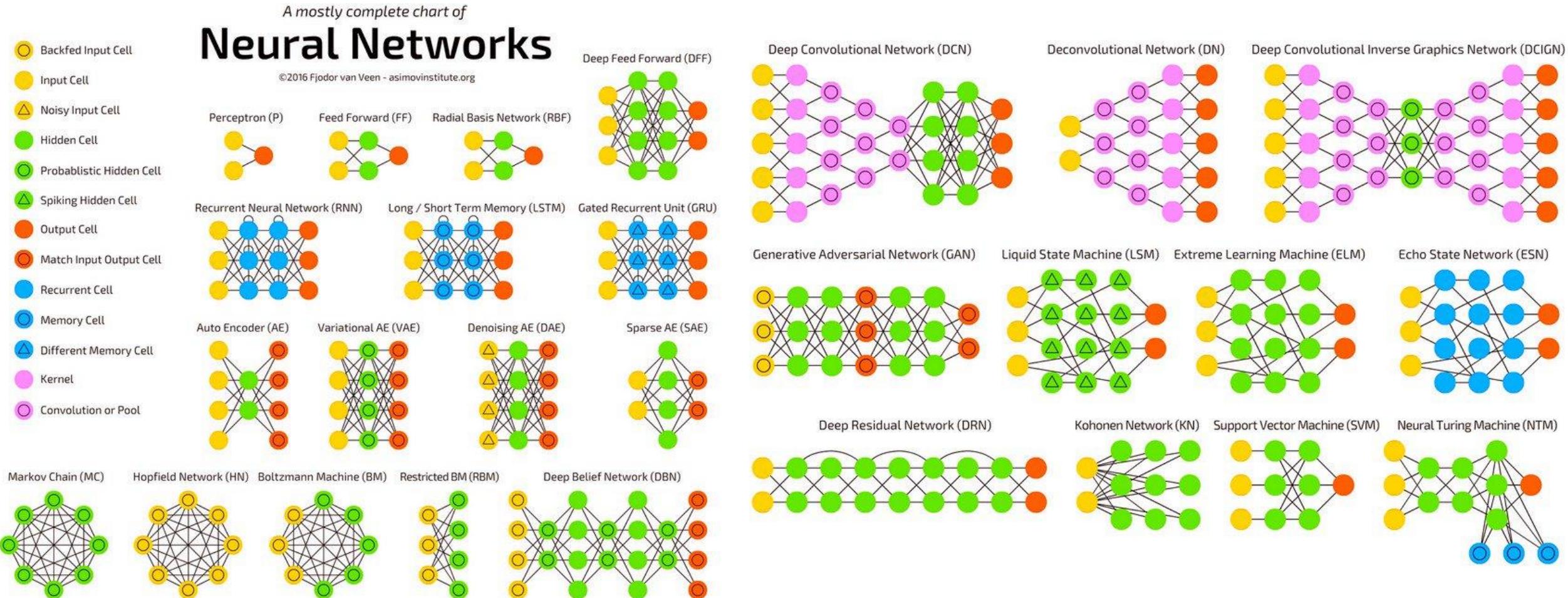
Compositional functions power



Tensor flow playground

- <https://playground.tensorflow.org/#activation=tanh&batchSize=10&dataset=circle®Dataset=reg-plane&learningRate=0.03®ularizationRate=0&noise=0&networkShape=4,2&seed=0.31138&showTestData=false&discretize=false&percTrainData=50&x=true&y=true&xTimesY=false&xSquared=false&ySquared=false&cosX=false&sinX=false&cosY=false&sinY=false&collectStats=false&problem=classification&initZero=false&hideText=false>

A zoo of networks



<https://modelzoo.co/>

Class 2: training neural networks

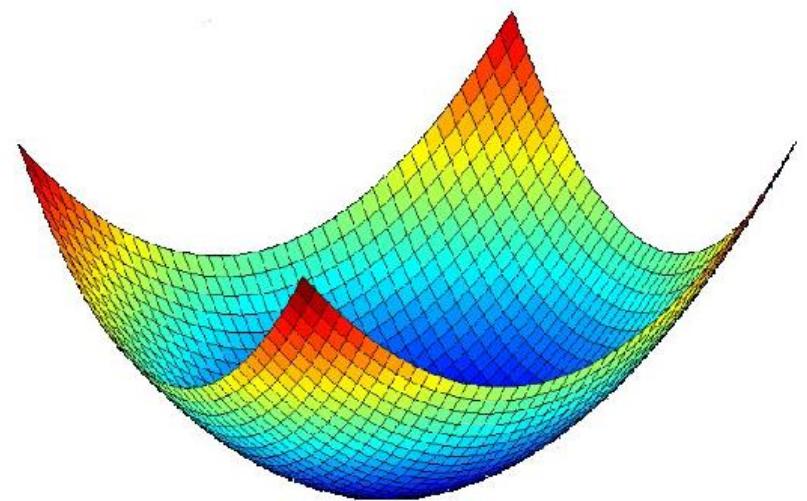
- Learning: forward and backward pass
- Gradient descent and stochastic gradient descent

How do we learn w and W?

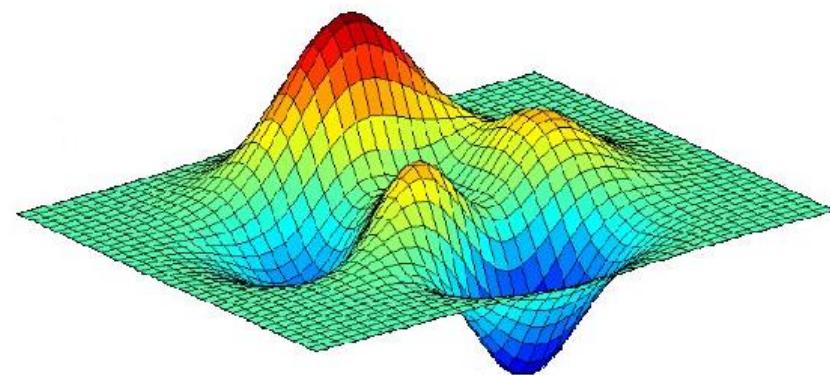
$$\mathcal{L}(w, W) = \frac{1}{N} \sum_{i=1}^N \left(y_i - w^T \sigma(W^T x) \right)^2$$

What's new w.r.t. the LSM?

$$\mathcal{L}(w, W) = \frac{1}{N} \sum_{i=1}^N \left(y_i - w^T \sigma(W^T x) \right)^2$$

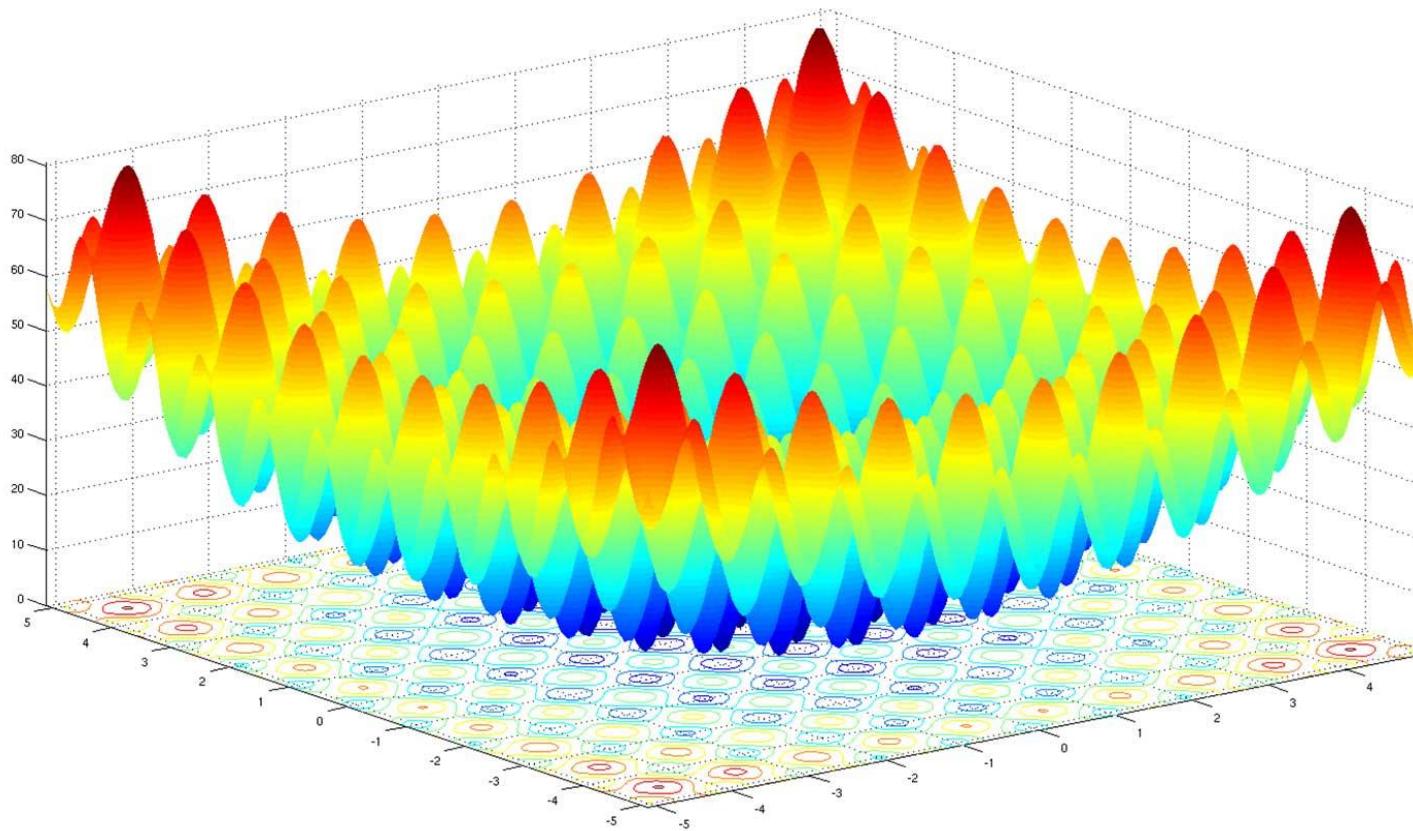


convex function



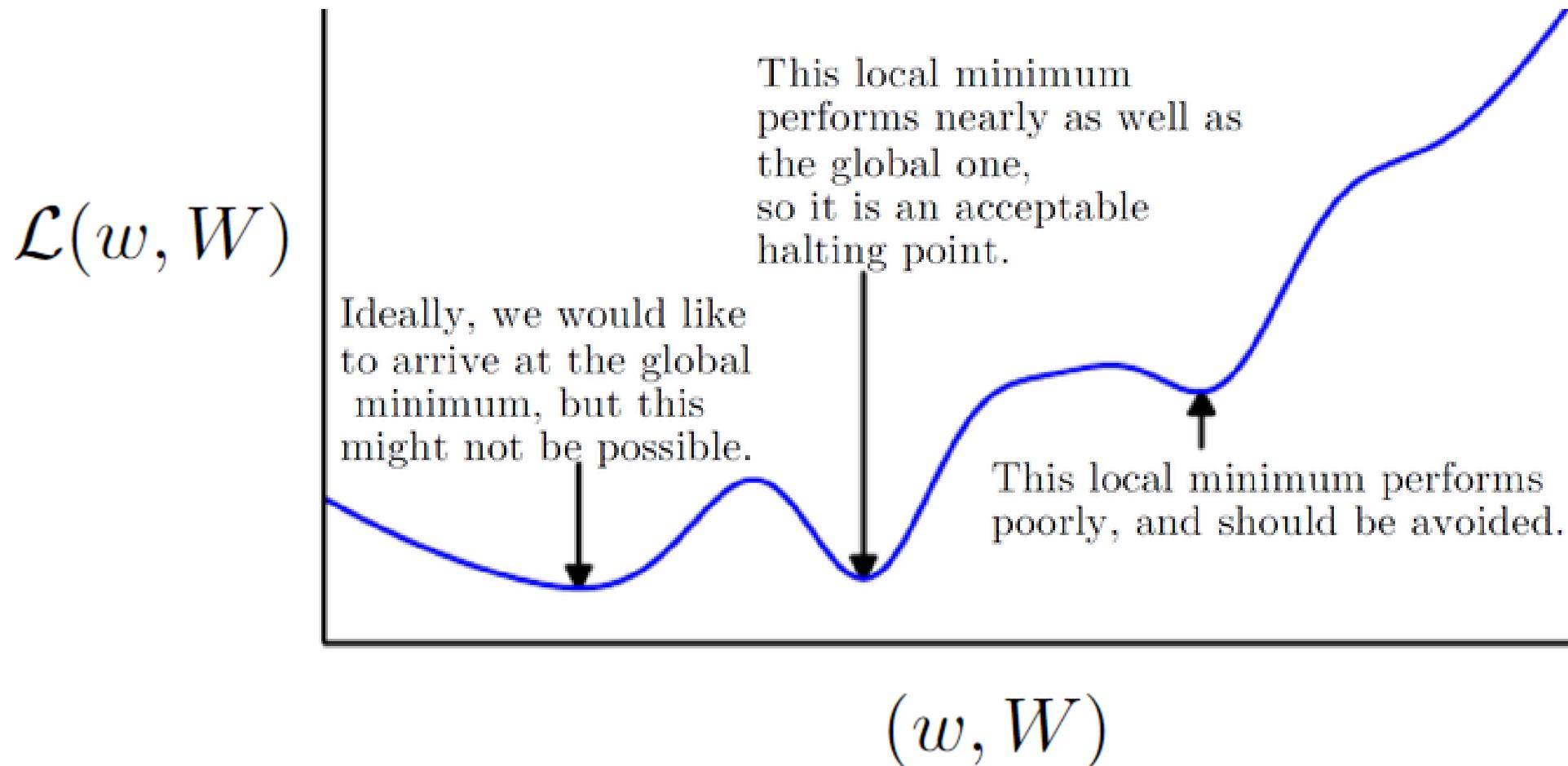
non-convex function

$$\mathcal{L}(w, W) = \frac{1}{N} \sum_{i=1}^N \left(y_i - w^T \sigma(W^T x) \right)^2$$

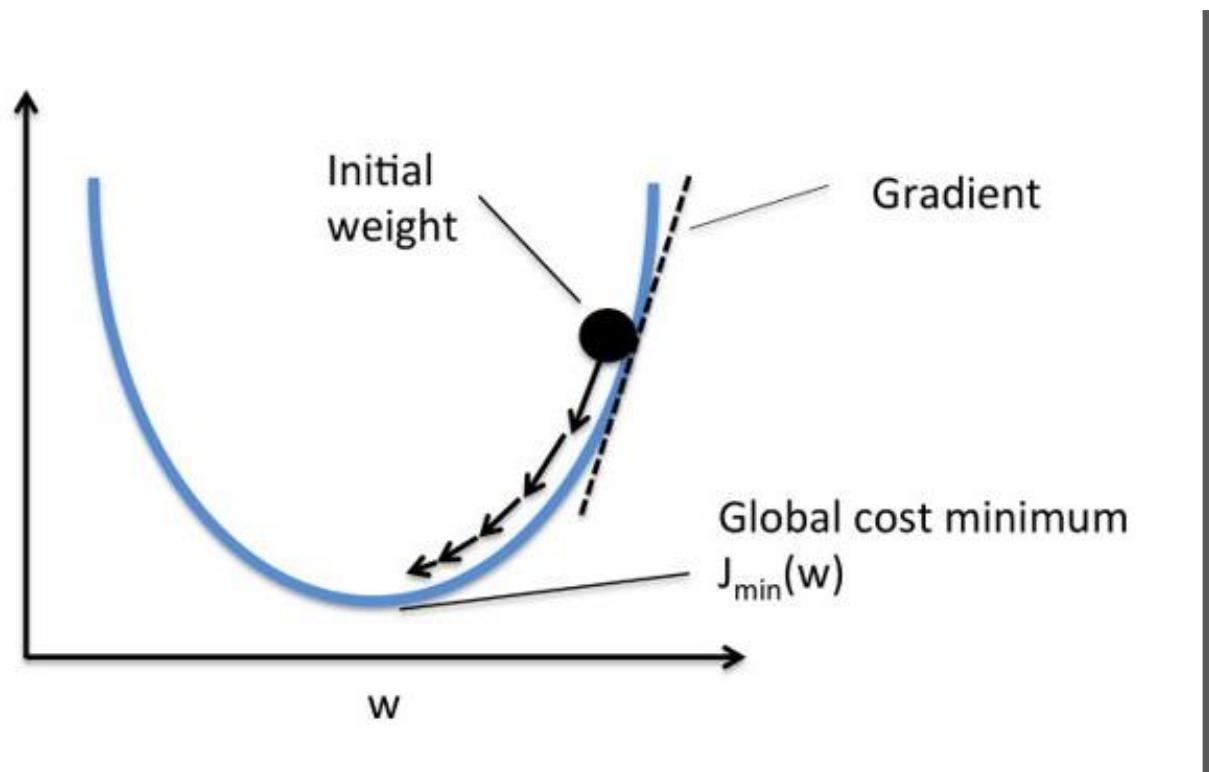


Local
minima!!!!

Dependence on initialization



Gradient descent



$$w_{t+1} = w_t - \gamma \nabla_w \mathcal{L}(w, W)$$

$$W_{t+1} = w_t - \gamma \nabla_W \mathcal{L}(w, W)$$

For neural networks the procedure is called **backpropagation**

GD variations

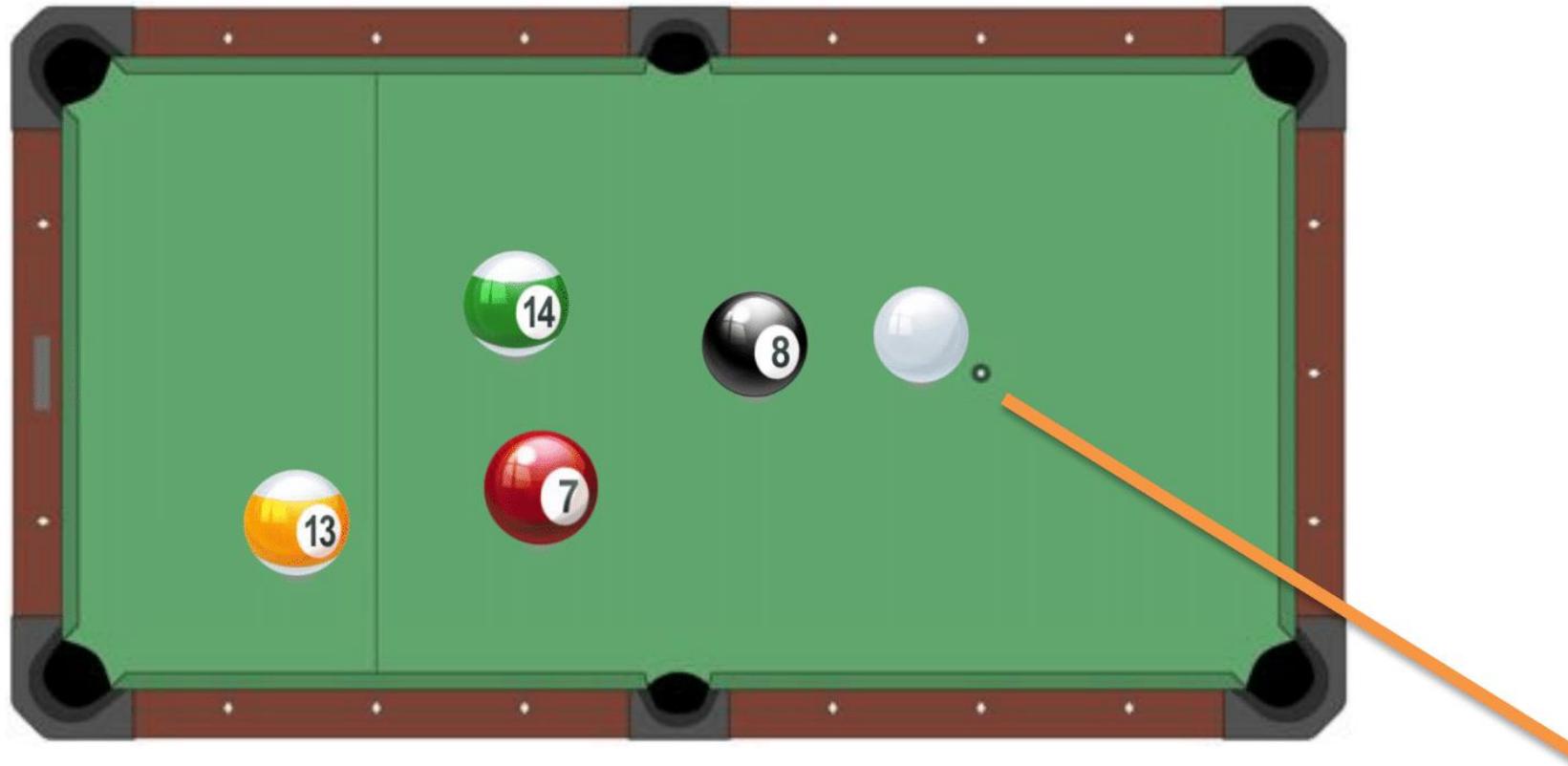
$$\mathcal{L}(w, X) = \sum_i l_i(w, x_i)$$

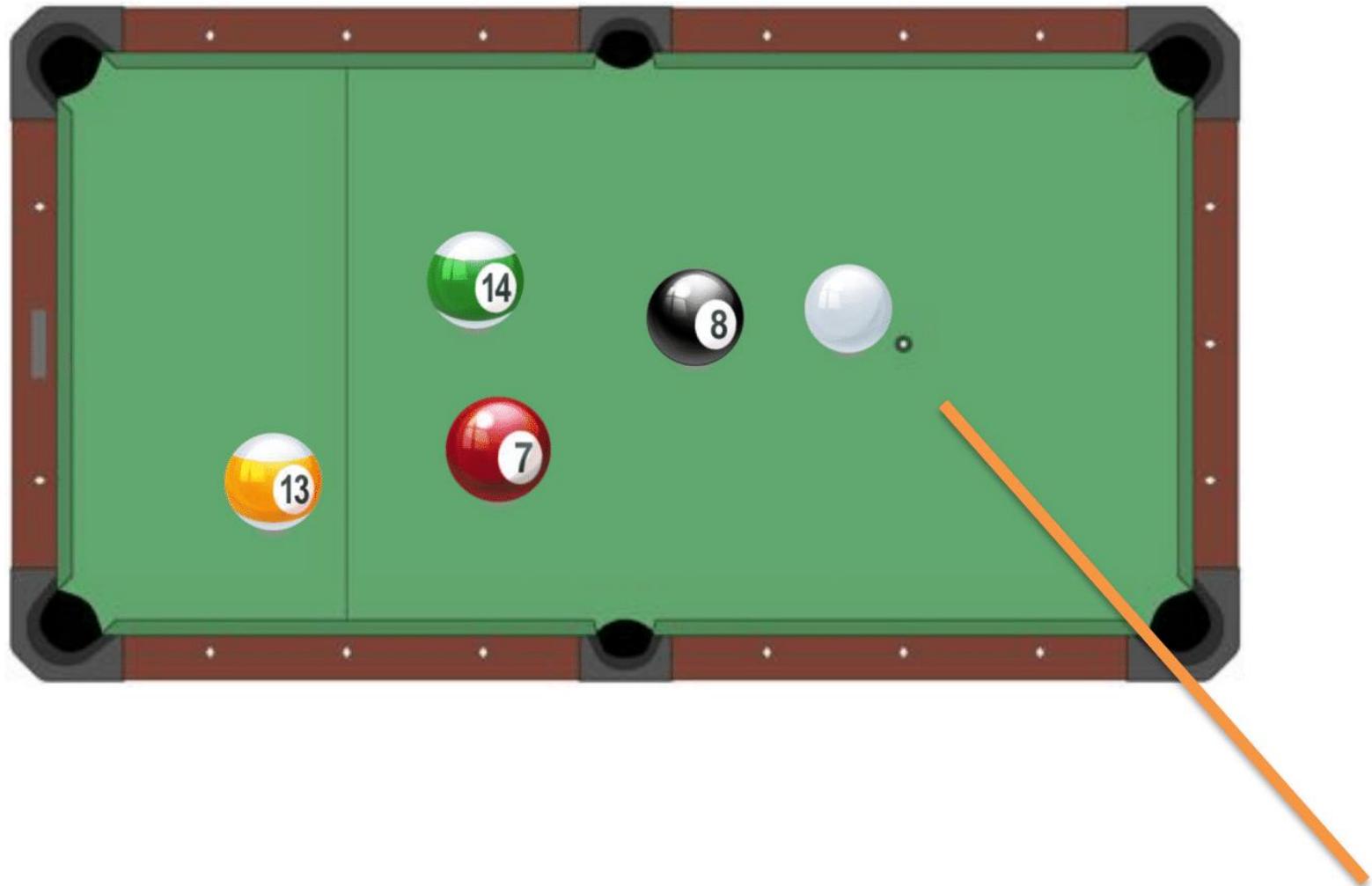
$$w_{t+1} = w_t - \eta \nabla_w \ell_i \quad \text{Stochastic GD}$$

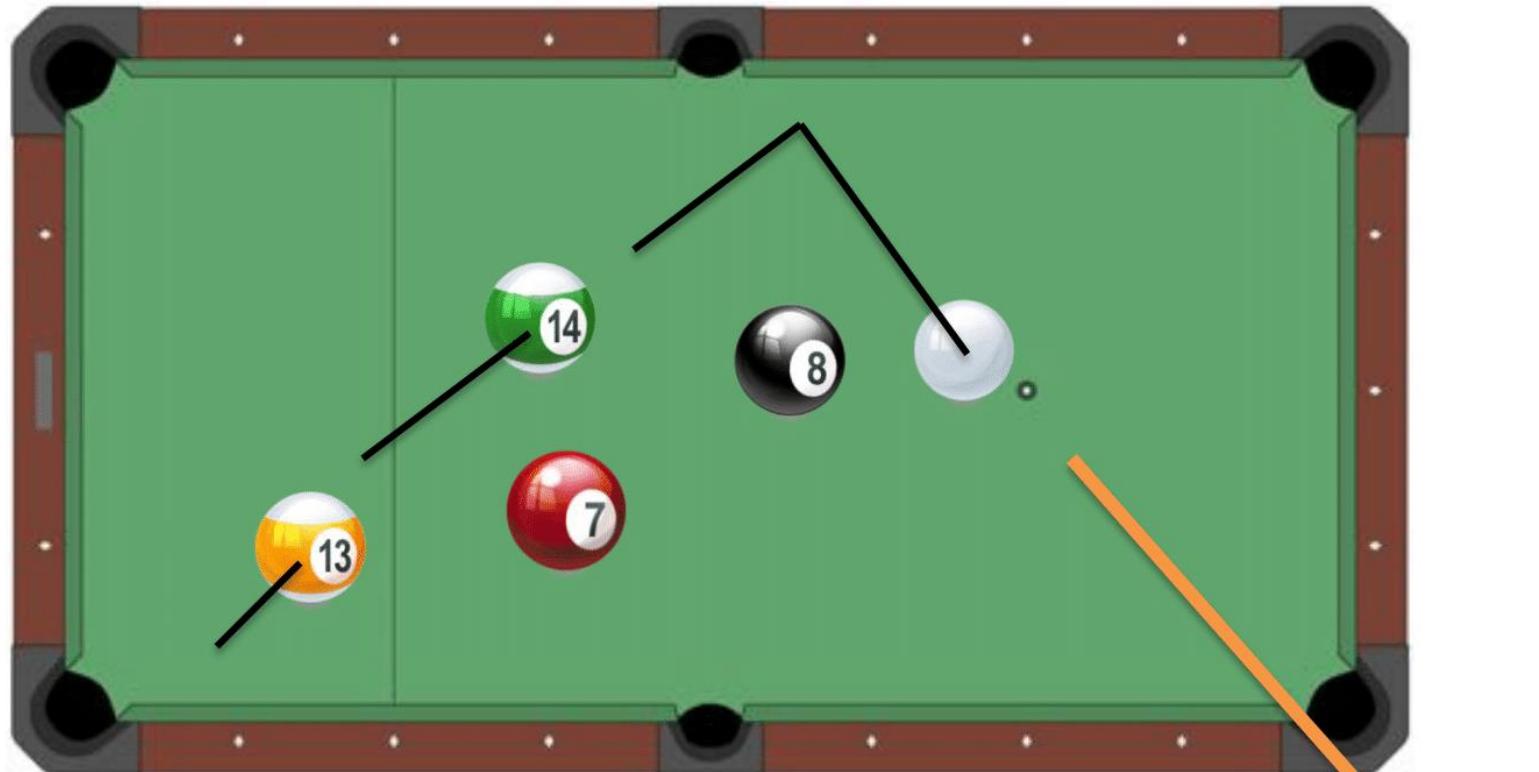
$$w_{t+1} = w_t - \eta \sum_{i \in I} \nabla_w \ell_i \quad \text{Mini-batch GD}$$

$$w_{t+1} = w_t - \eta \sum_i \nabla_w \ell_i \quad \text{Full-batch GD}$$

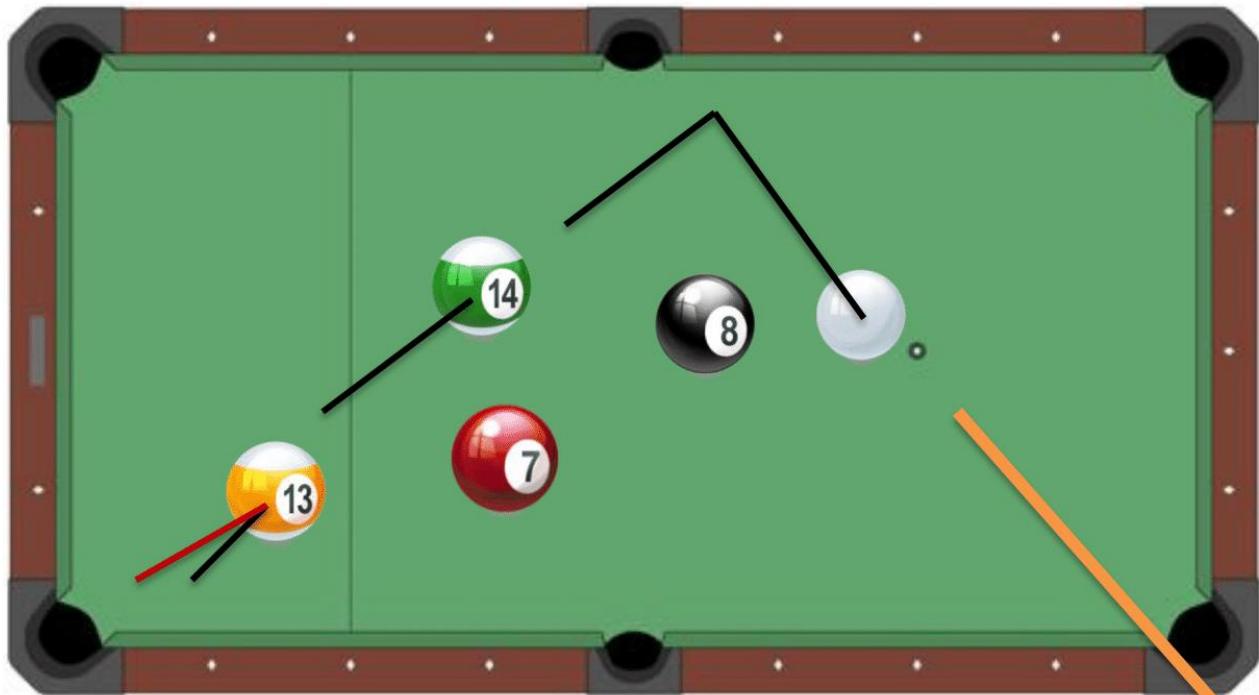
Intuition



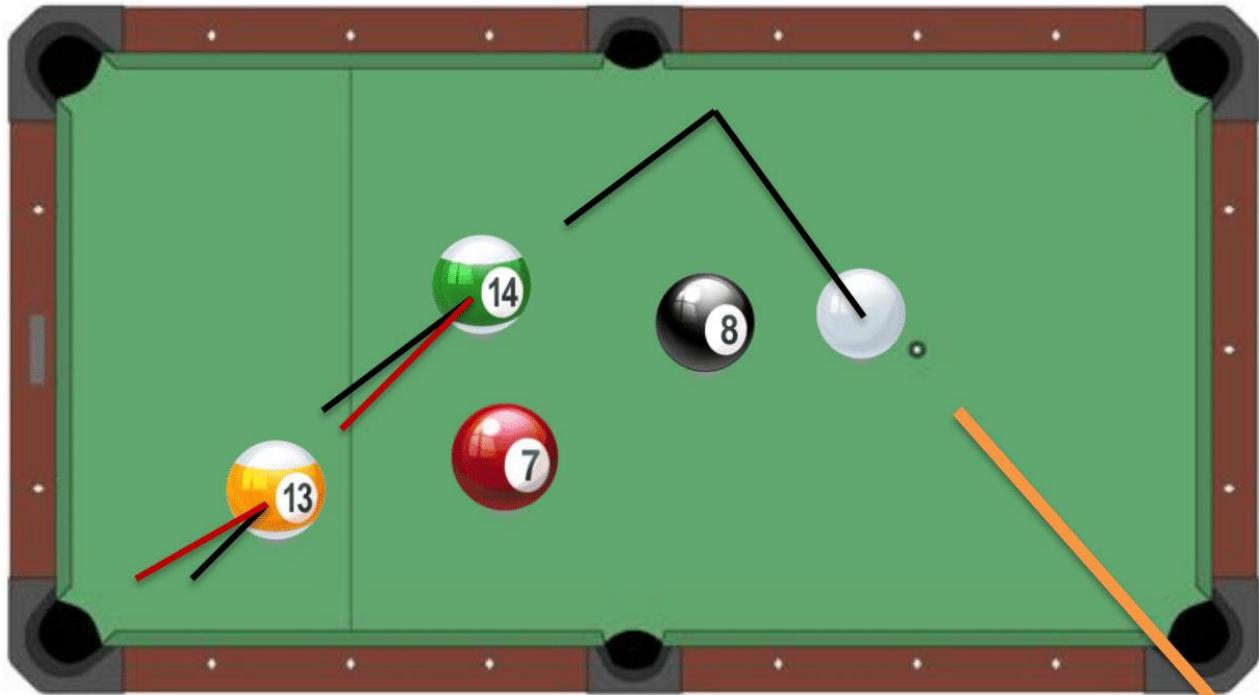




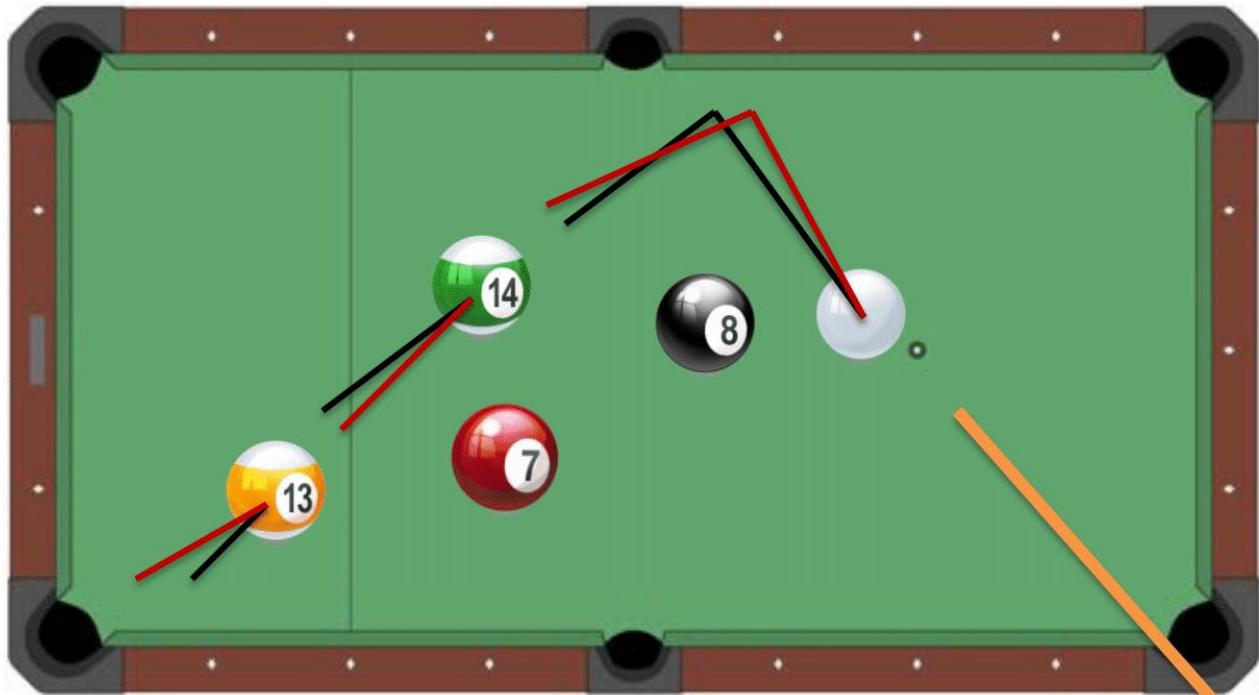
Forward pass



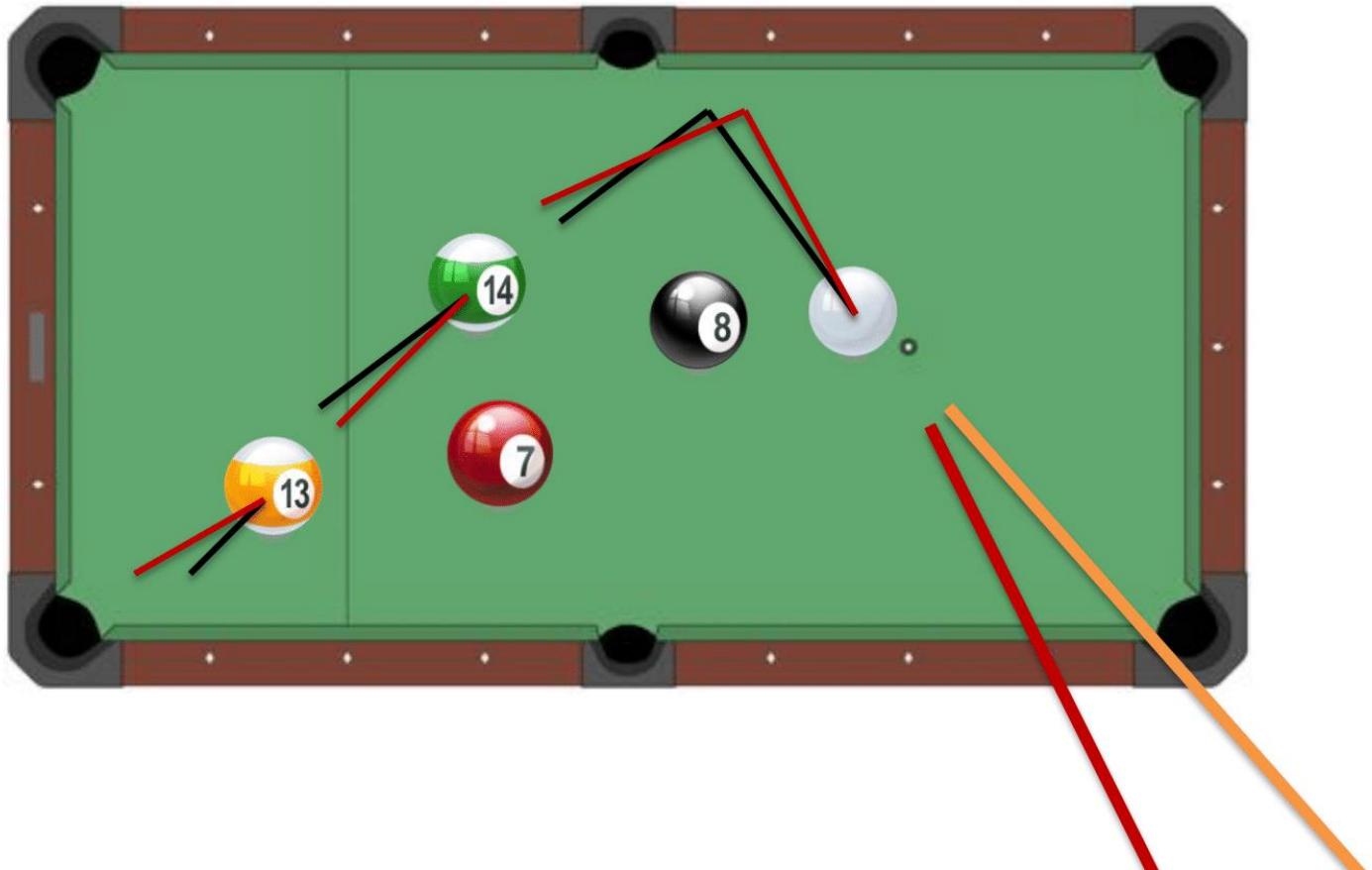
Backward pass



Backward pass

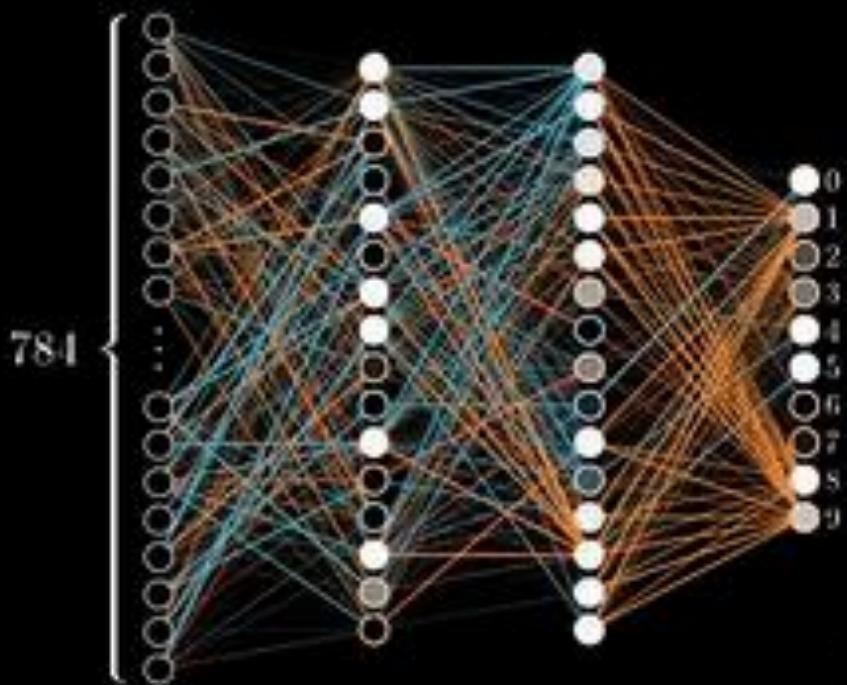


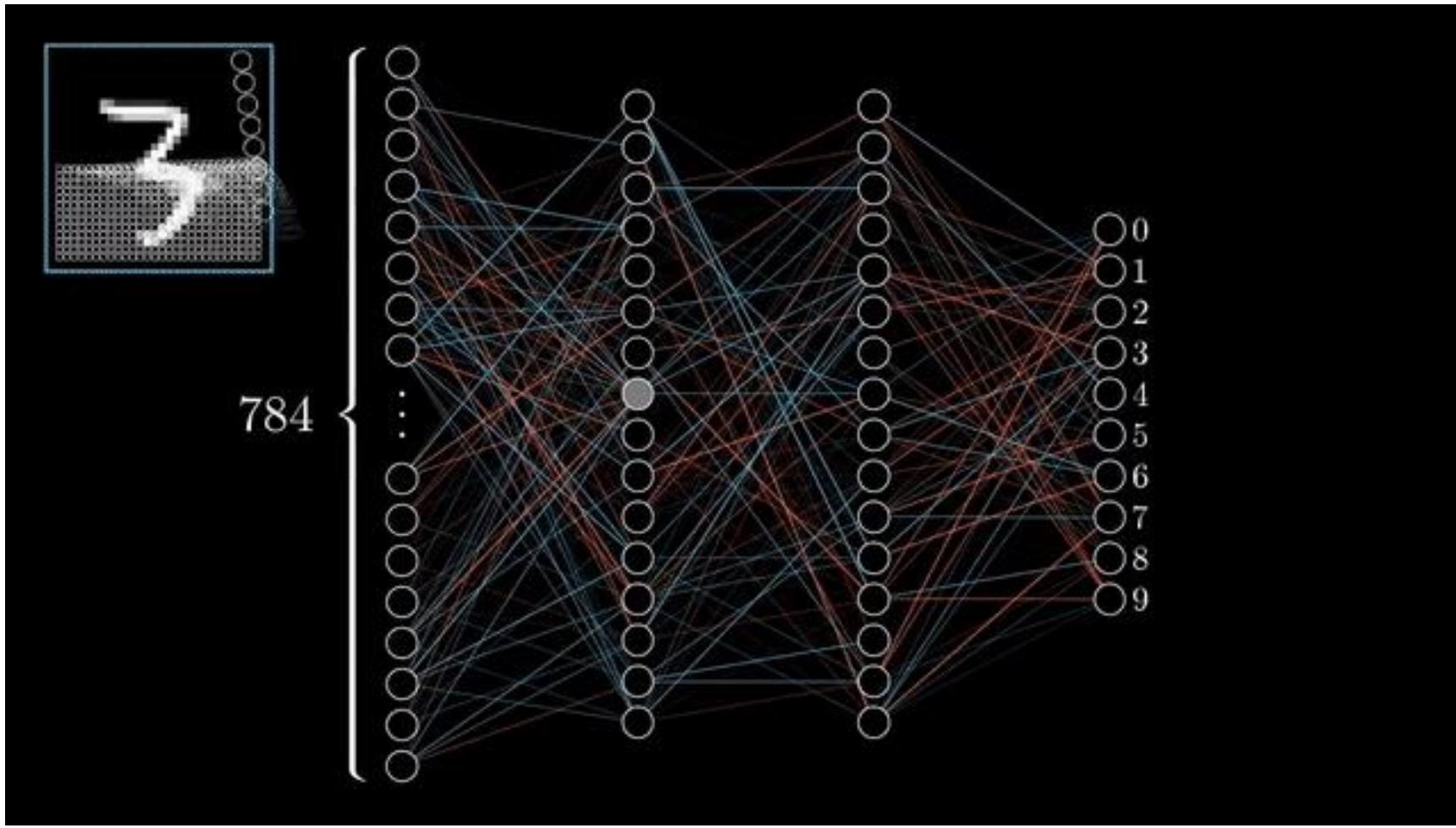
Backward pass



Backward pass

Training in
progress. . .





How do we change the weights of the network to improve classification?

Chain rule

$$\frac{d[f \circ g]}{dx} = \frac{df(g)}{dg} \frac{dg}{dx}$$

$$\frac{d}{dx} \cos(\sin(x)) = \frac{d}{da} \cos(a) \frac{da}{dx} = \cos'(\sin(x)) \sin'(x), \quad a = \sin(x)$$

Training a 1 hidden layer neural network for classification task

Consider a one-layer NN:

$$\hat{y} = \text{softmax}(U \text{ReLU}(Wx + b_1) + b_2)$$

- Input: x
- Linear transformation: $z = Wx + b_1$
- Activation function: $h = \text{ReLU}(z)$
- Second linear transformation: $\theta = Uh + b_2$
- Output prediction: $\hat{y} = \text{softmax}(\theta)$
- Loss function: $J = CE(y, \hat{y})$

Dimensions: $x \in \mathbb{R}^{D_x \times 1}$, $b_1 \in \mathbb{R}^{D_h \times 1}$, $W \in \mathbb{R}^{D_h \times D_x}$, $b_2 \in \mathbb{R}^{N_c \times 1}$, $U \in \mathbb{R}^{N_c \times D_h}$

We want to calculate

- $\frac{\partial J}{\partial U}$
- $\frac{\partial J}{\partial b_2}$
- $\frac{\partial J}{\partial W}$
- $\frac{\partial J}{\partial b_1}$
- $\frac{\partial J}{\partial x}$

To start with, recall that $\text{ReLU}(x) = \max(x, 0)$. This means:

$$\text{ReLU}'(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases} = \text{sgn}(\text{ReLU}(x))$$

where sgn is the signum function.

Now, let's write out the chain rule for $\frac{\partial J}{\partial U}$ and $\frac{\partial J}{\partial b_2}$:

$$\frac{\partial J}{\partial U} = \frac{\partial J}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \theta} \frac{\partial \theta}{\partial U}$$

$$\frac{\partial J}{\partial b_2} = \frac{\partial J}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \theta} \frac{\partial \theta}{\partial b_2}$$

Notice that $\frac{\partial J}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \theta} = \frac{\partial J}{\partial \theta}$ is present in both gradients.

To make computations more efficient, we define intermediate variables to represent the derivatives:

$$\delta_1 = \frac{\partial J}{\partial \theta} = (\hat{y} - y)^T$$

$$\delta_2 = \frac{\partial J}{\partial z} = \frac{\partial J}{\partial \theta} \frac{\partial \theta}{\partial h} \frac{\partial h}{\partial z}$$

Using the chain rule and substituting δ_1 , we find:

$$\delta_2 = \delta_1 U \circ \text{ReLU}'(z) = \delta_1 U \circ \text{sgn}(h)$$

A good way of checking our work is by looking at the dimensions of the Jacobians:

$$\frac{\partial J}{\partial z} = \delta_1 U \circ \text{sgn}(h)$$

Here, δ_1 , U , and $\text{sgn}(h)$ should align in dimensions to ensure the correct computation of gradients.

Gradient with respect to U :

$$\frac{\partial J}{\partial U} = \frac{\partial J}{\partial \theta} \frac{\partial \theta}{\partial U} = \delta_1 \frac{\partial \theta}{\partial U} = \delta_1^T h^T$$

Gradient with respect to b_2 :

$$\frac{\partial J}{\partial b_2} = \frac{\partial J}{\partial \theta} \frac{\partial \theta}{\partial b_2} = \delta_1 \frac{\partial \theta}{\partial b_2} = \delta_1^T$$

Gradient with respect to W :

$$\frac{\partial J}{\partial W} = \frac{\partial J}{\partial \theta} \frac{\partial z}{\partial W} = \delta_2 \frac{\partial z}{\partial W} = \delta_2^T x^T$$

Gradient with respect to b_1 :

$$\frac{\partial J}{\partial b_1} = \frac{\partial J}{\partial \theta} \frac{\partial z}{\partial b_1} = \delta_2 \frac{\partial z}{\partial b_1} = \delta_2^T$$

More in general:

$$\phi_L(x) = \sigma(W_L^T \sigma(W_{L-1}^T \cdots (\sigma(W_1^T x) \cdots))$$

$$act_l = W_l^T \sigma(act_{l-1})$$

$$\frac{\partial \sigma(act_L)}{\partial W_i} = \frac{\partial \sigma(act_L)}{\partial act_L} \frac{\partial act_L}{\partial \sigma(act_{L-1})} \cdots \frac{\partial \sigma(act_i)}{\partial act_i} \frac{\partial act_i}{\partial W_i}$$

$$\frac{\partial \sigma(act_l)}{\partial act_l} = \sigma'(act_l)$$

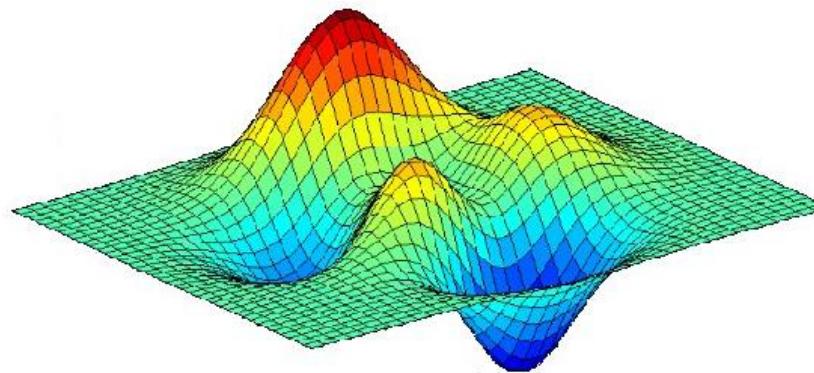
$$\frac{\partial act_L}{\partial \sigma(act_{l-1})} = W_l^T$$

Training neural networks: backward pass, W

$$\begin{aligned}\frac{\partial \sigma(\text{act}_L)}{\partial W_i} &= \sigma'(\text{act}_l)W_l^T \dots \frac{\partial \sigma(\text{act}_i)}{\partial \text{act}_i} \frac{\partial \text{act}_i}{\partial W_i} \\ &= \sigma'(\text{act}_l)W_l^T \dots \frac{\partial \sigma(W_i^T \sigma(\text{act}_{i-1}))}{\partial W_i^T \sigma(\text{act}_{i-1})} \frac{\partial \text{act}_i}{\partial W_i} \\ &= \sigma'(\text{act}_l)W_l^T \dots \sigma'(W_i^T \sigma(\text{act}_{i-1})) \sigma(\text{act}_{i-1})\end{aligned}$$

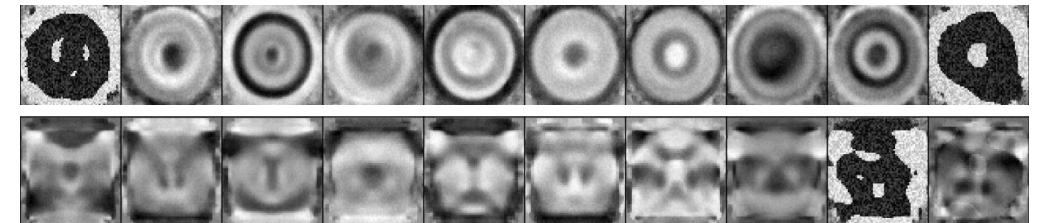
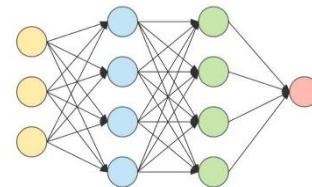
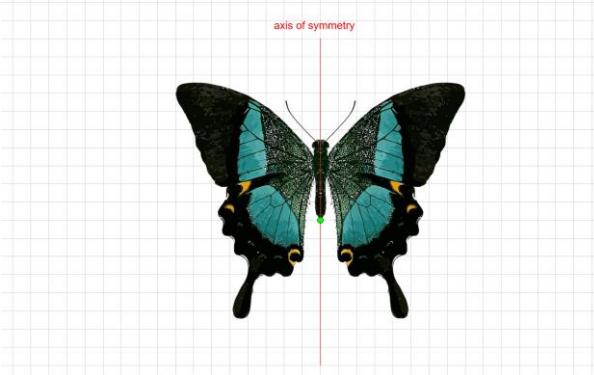
- We calculate the effect of the network on the Loss J with the forward pass
- We update the variables using the derivative at each stage of the calculation
- The update is done using the chain rule and the updates at the previous derivative

Why do we find good minima with backprop?



non-convex function

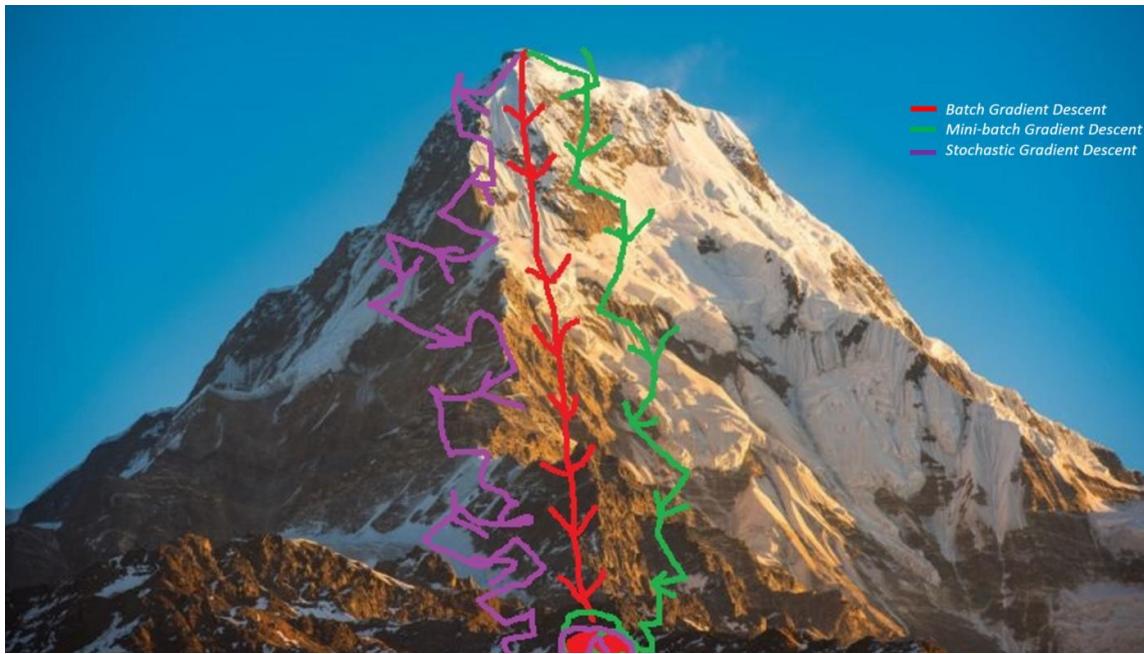
GD and data Symmetries



$$\begin{aligned} W_{t+1} - W_t &\propto \nabla_W \mathcal{L}(W_t; X, y) = \\ &= \frac{1}{R} \sum_{i=1}^R \ell'(v^T \sigma \langle W_t, x_i \rangle), y_i) v^T \sigma' \langle W_t, x_i \rangle x_i. \end{aligned}$$

Class 3: Stochastic gradient descent, implicit bias and good solutions

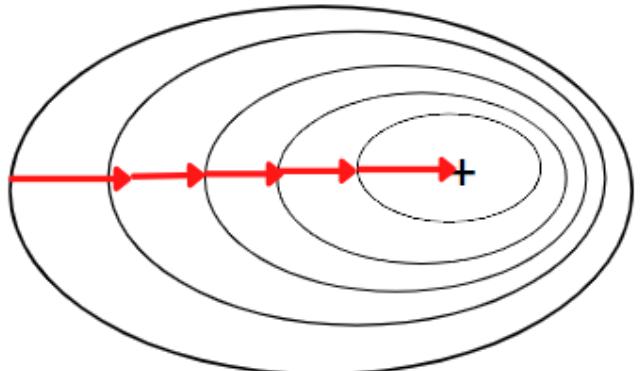
Mini-batch and Stochastic gradient decent



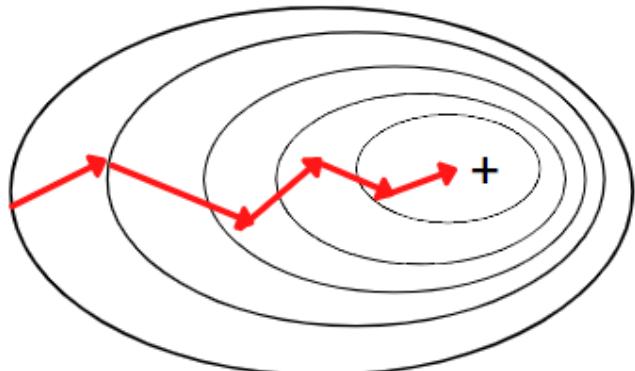
$$W_{t+1} = w_t - \gamma \nabla_W \frac{1}{B} \sum_{i=1}^B \left(y_i - w^T \sigma(W^T x_i) \right)^2$$

$$W_{t+1} = w_t - \gamma \nabla_W \left(y_i - w^T \sigma(W^T x_i) \right)^2$$

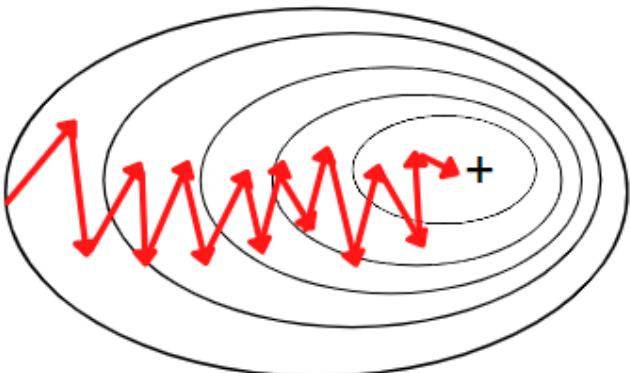
Batch Gradient Descent



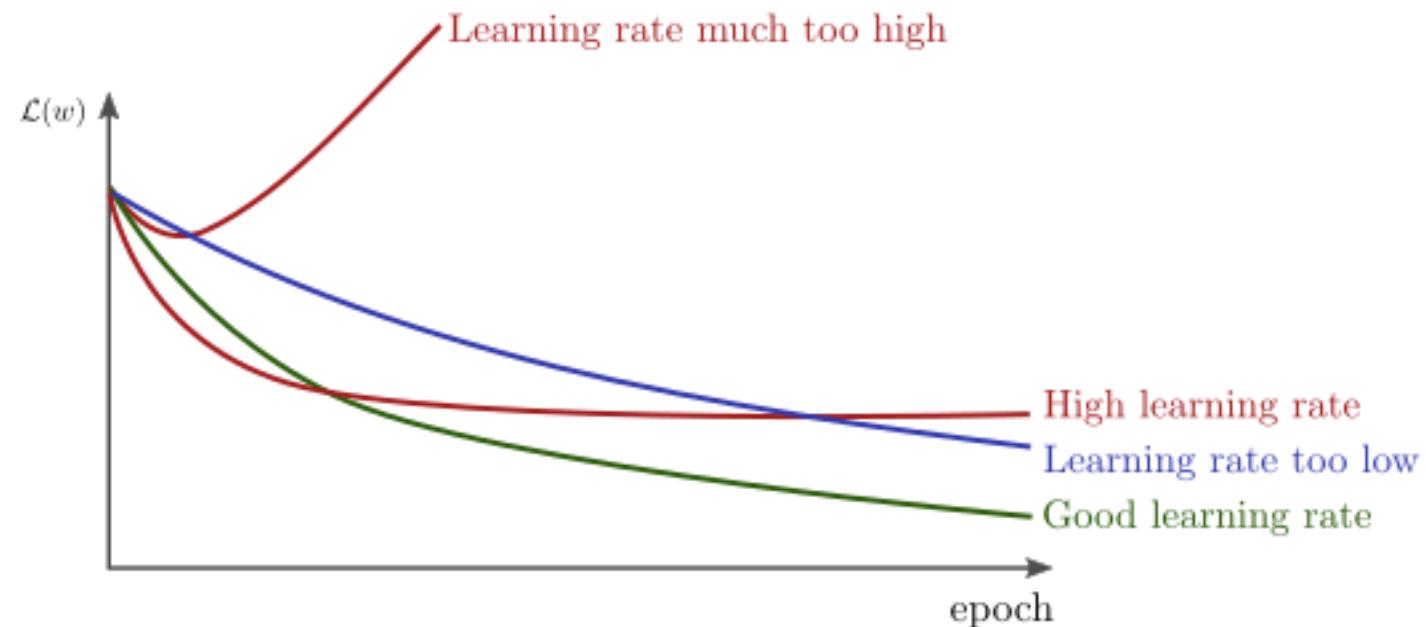
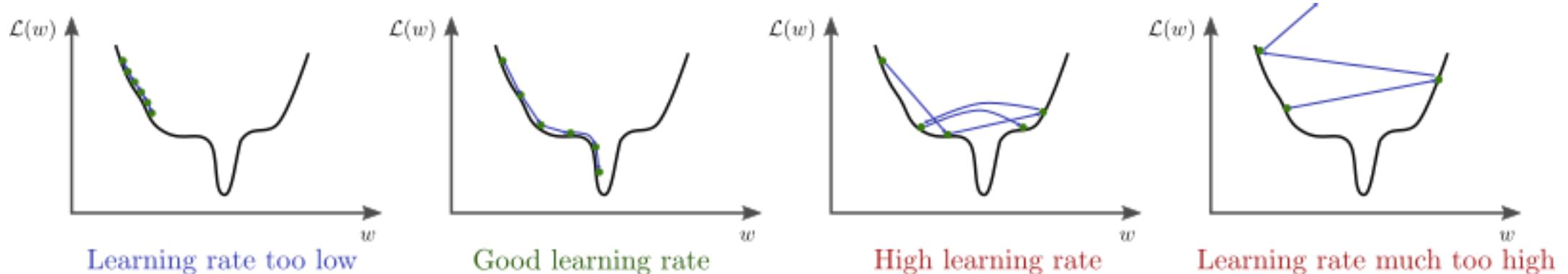
Mini-Batch Gradient Descent



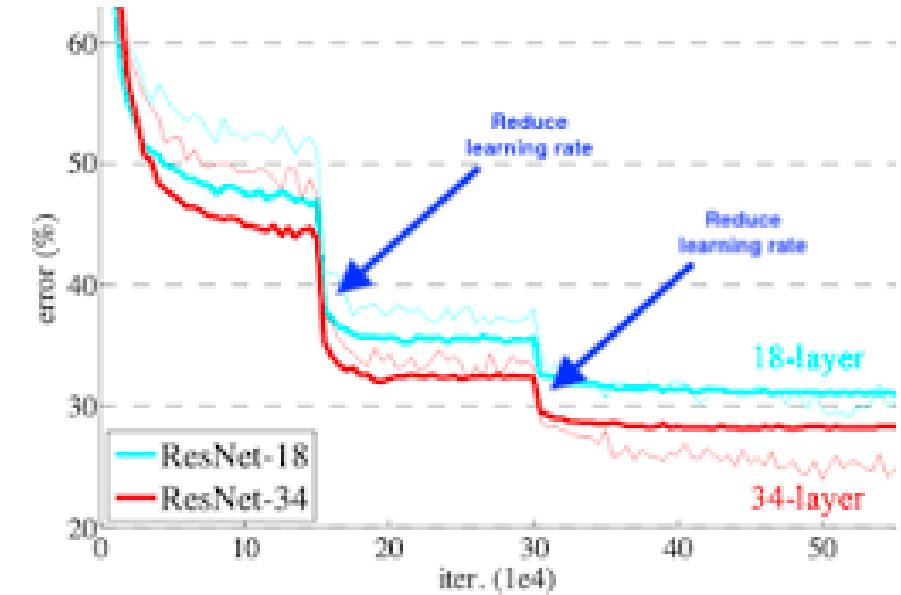
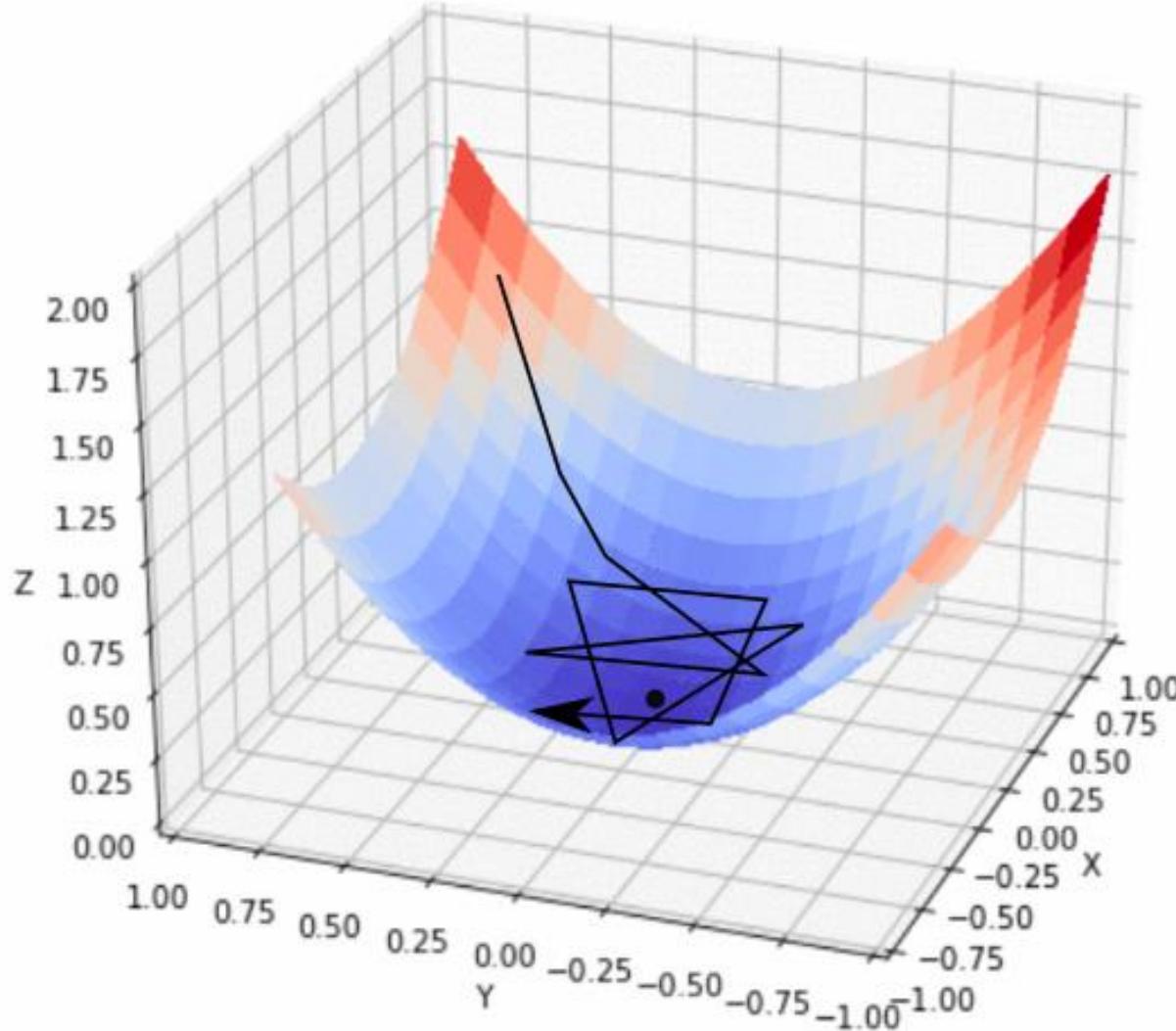
Stochastic Gradient Descent



Learning rate

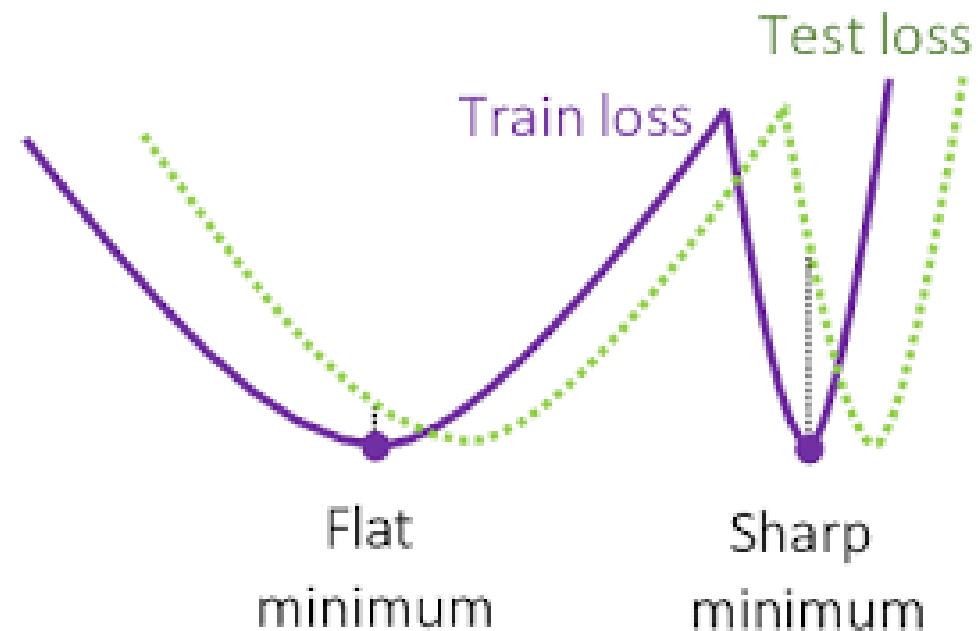


Learning rate for SGD: scheduling

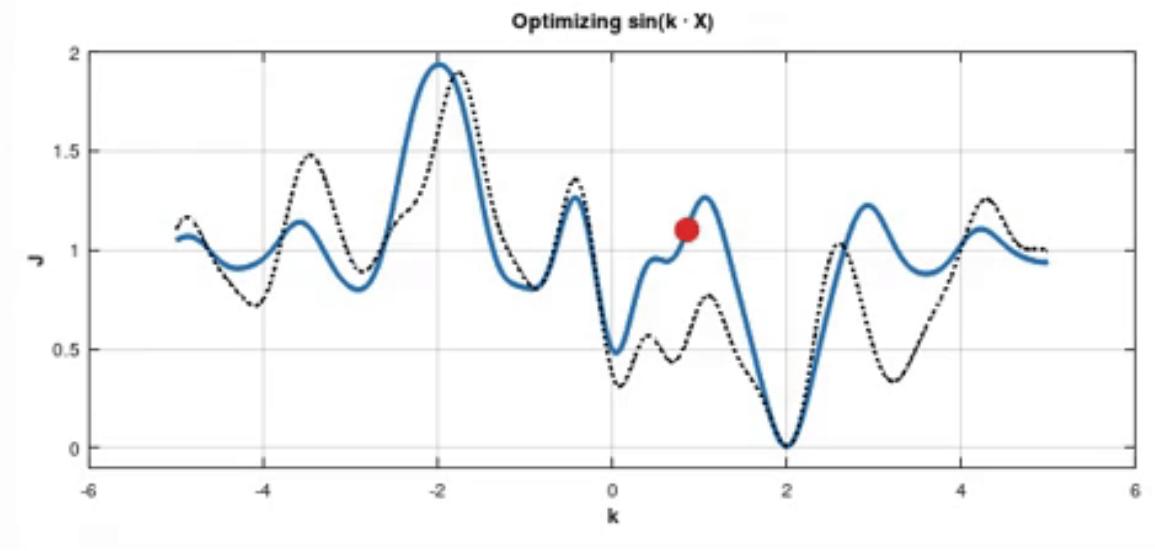


Advantages of SGD

- Skip bad local minima, converge to flat minima (basins of attraction). Robust to data perturbations. Good for **generalization, robustness**.



Intuition



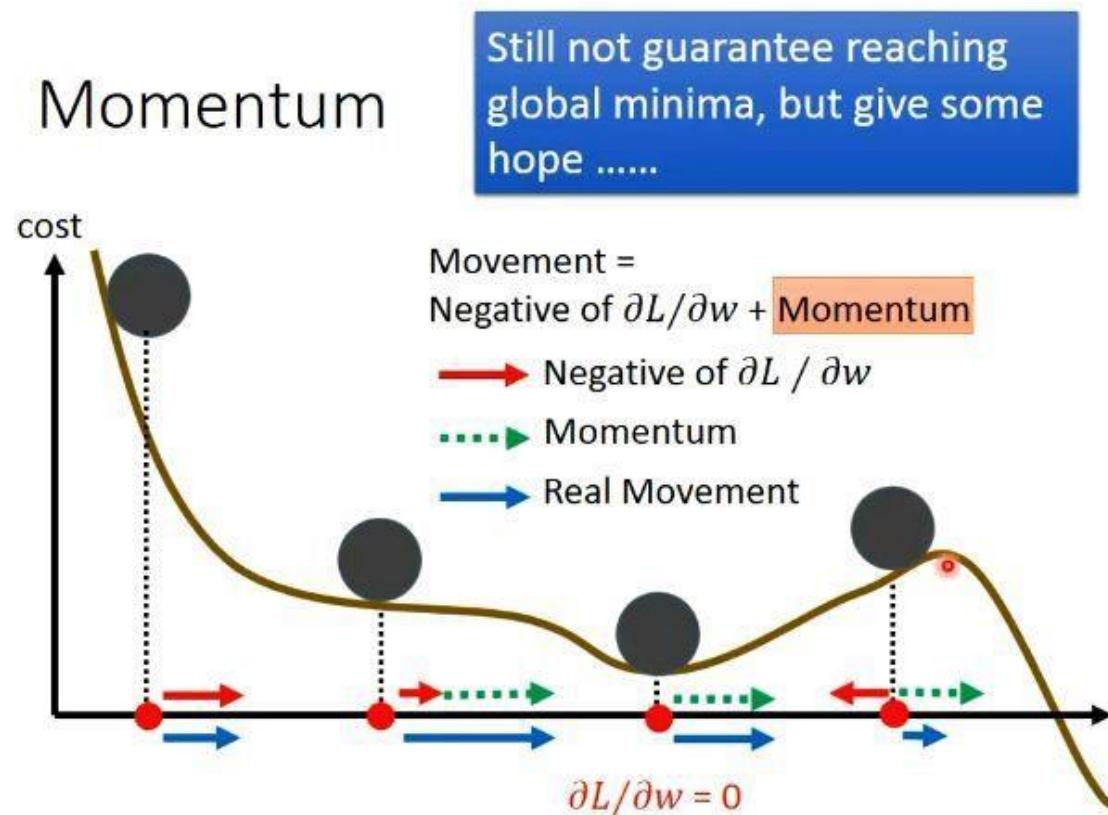
At each step the SGD is looking to a “different, sampled” function

On average the main effect is dominated by big “basins of attraction”

Advantages of SGD

- We need way less computations (proportional to 1 or Batch size)
- Converges to good minima (invariant to data perturbations)

Other strategies to reach good minima: add momentum

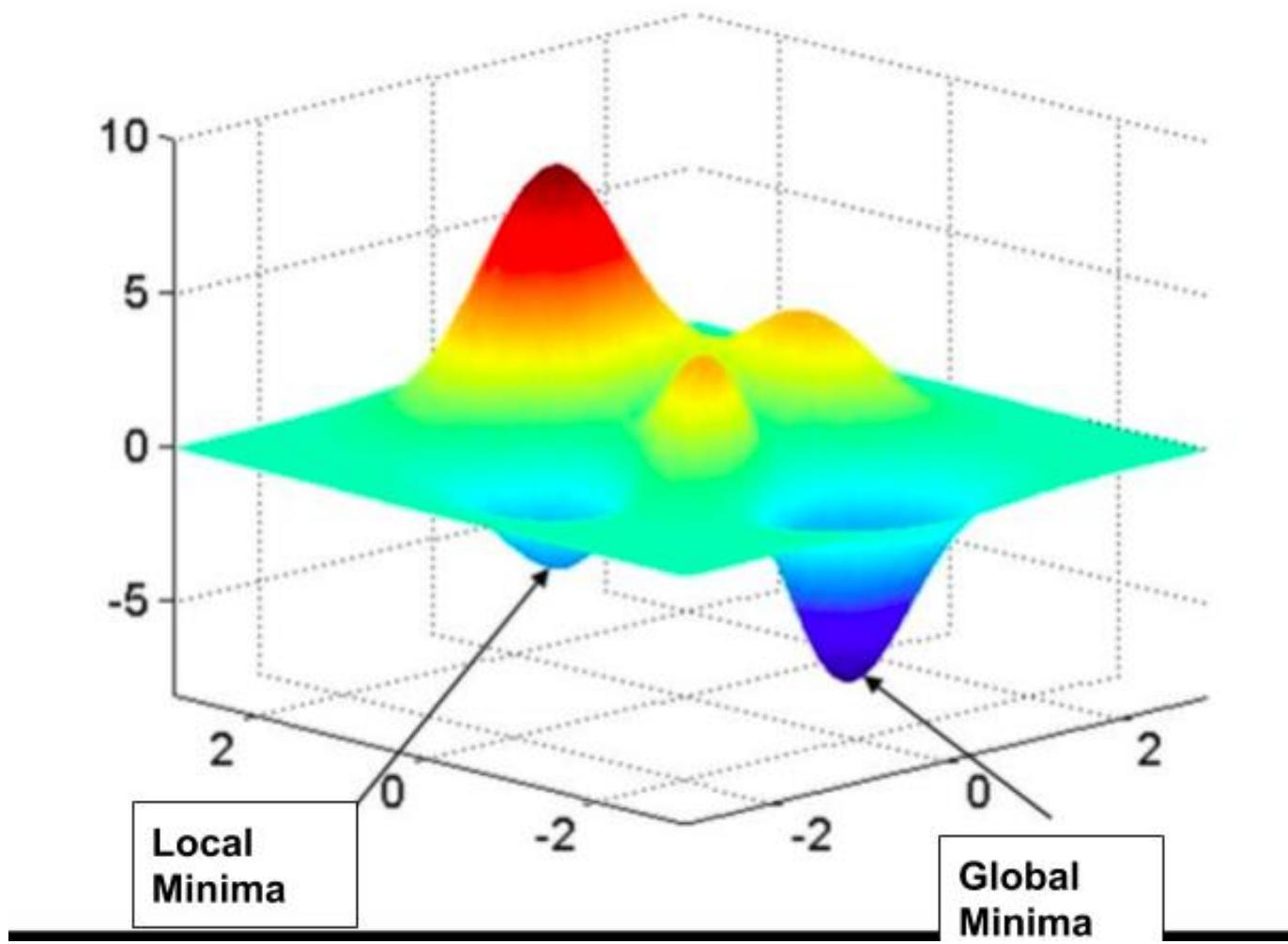


$$v_{t+1} = \rho v_t + \nabla_w \mathcal{L}(w_t)$$

$$w_{t+1} = w_t - \gamma v_{t+1}$$

Going out from flat bad minima

Importance of initialization



Which other factors influence the convergence to good minima?

Implicit regularization

- **Explicit regularization** imposes an explicit penalty on the parameters of the model (typically some measure of complexity or sensitivity)
- **Implicit regularization** occurs when the dynamics of training lead to certain minima rather than others

- Algorithm
- Initialization
- Architecture



Example with regression: minimum norm solution

one particular solution is

$$x_{\text{ln}} = A^T(AA^T)^{-1}y$$

(AA^T is invertible since A full rank)

in fact, x_{ln} is the solution of $y = Ax$ that minimizes $\|x\|$

i.e., x_{ln} is solution of optimization problem

$$\begin{aligned} & \text{minimize} && \|x\| \\ & \text{subject to} && Ax = y \end{aligned}$$

(with variable $x \in \mathbf{R}^n$)

Example with regression: minimum norm solution

suppose $Ax = y$, so $A(x - x_{\text{ln}}) = 0$ and

$$\begin{aligned}(x - x_{\text{ln}})^T x_{\text{ln}} &= (x - x_{\text{ln}})^T A^T (AA^T)^{-1} y \\&= (A(x - x_{\text{ln}}))^T (AA^T)^{-1} y \\&= 0\end{aligned}$$

i.e., $(x - x_{\text{ln}}) \perp x_{\text{ln}}$, so

$$\|x\|^2 = \|x_{\text{ln}} + x - x_{\text{ln}}\|^2 = \|x_{\text{ln}}\|^2 + \|x - x_{\text{ln}}\|^2 \geq \|x_{\text{ln}}\|^2$$

i.e., x_{ln} has smallest norm of any solution



Example with regression: minimum norm solution with GD

Let cost function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be defined by

$$f(\mathbf{x}) := \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2$$

whose gradient is

$$\nabla f(\mathbf{x}) = \mathbf{A}^\top (\mathbf{Ax} - \mathbf{b})$$

Using gradient descent with step $\mu > 0$,

$$\begin{aligned} \mathbf{x}_{k+1} &= \mathbf{x}_k - \mu \nabla f(\mathbf{x}_k) \\ &= (\mathbf{I} - \mu \mathbf{A}^\top \mathbf{A}) \mathbf{x}_k + \mu \mathbf{A}^\top \mathbf{b} \end{aligned}$$

Hence,

$$\mathbf{x}_k = (\mathbf{I} - \mu \mathbf{A}^\top \mathbf{A})^k \mathbf{x}_0 + \mu \sum_{\ell=0}^{k-1} (\mathbf{I} - \mu \mathbf{A}^\top \mathbf{A})^\ell \mathbf{A}^\top \mathbf{b}$$

$$\mathbf{A} = \mathbf{U} \Sigma \mathbf{V}^\top = \mathbf{U} \begin{bmatrix} \Sigma_1 & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{V}_1^\top \\ \mathbf{V}_2^\top \end{bmatrix} = \mathbf{U} \Sigma_1 \mathbf{V}_1^\top$$

Letting $\mathbf{y} := \mathbf{V}^\top \mathbf{x}$, we rewrite

$$\begin{aligned} \mathbf{y}_k &= (\mathbf{I} - \mu \Sigma^\top \Sigma)^k \mathbf{y}_0 + \mu \sum_{\ell=0}^{k-1} (\mathbf{I} - \mu \Sigma^\top \Sigma)^\ell \Sigma^\top \mathbf{U}^\top \mathbf{b} \\ &= \begin{bmatrix} (\mathbf{I} - \mu \Sigma_1^2)^k & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \mathbf{y}_0 + \mu \sum_{\ell=0}^{k-1} \begin{bmatrix} (\mathbf{I} - \mu \Sigma_1^2)^\ell & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \Sigma_1 \\ \mathbf{0} \end{bmatrix} \mathbf{U}^\top \mathbf{b} \\ &= \begin{bmatrix} (\mathbf{I} - \mu \Sigma_1^2)^k & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \mathbf{y}_0 + \mu \sum_{\ell=0}^{k-1} \begin{bmatrix} (\mathbf{I} - \mu \Sigma_1^2)^\ell \Sigma_1 \\ \mathbf{0} \end{bmatrix} \mathbf{U}^\top \mathbf{b} \end{aligned}$$



Example with regression: minimum norm solution with GD

Choosing $\mu > 0$ such that all eigenvalues of $\mathbf{I} - \mu \Sigma_1^2$ are strictly inside the unit circle, then $\mathbf{y}_k \rightarrow \mathbf{y}_\infty$, where

$$\mathbf{y}_\infty = \begin{bmatrix} \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{I} \end{bmatrix} \mathbf{y}_0 + \mu \sum_{\ell=0}^{\infty} \begin{bmatrix} (\mathbf{I} - \mu \Sigma_1^2)^\ell \Sigma_1 \\ \mathbf{O} \end{bmatrix} \mathbf{U}^\top \mathbf{b}$$

where

$$\mu \sum_{\ell=0}^{\infty} (\mathbf{I} - \mu \Sigma_1^2)^\ell \Sigma_1 = \mu (\mathbf{I} - \mathbf{I} + \mu \Sigma_1^2)^{-1} \Sigma_1 = \Sigma_1^{-1}$$

and, thus,

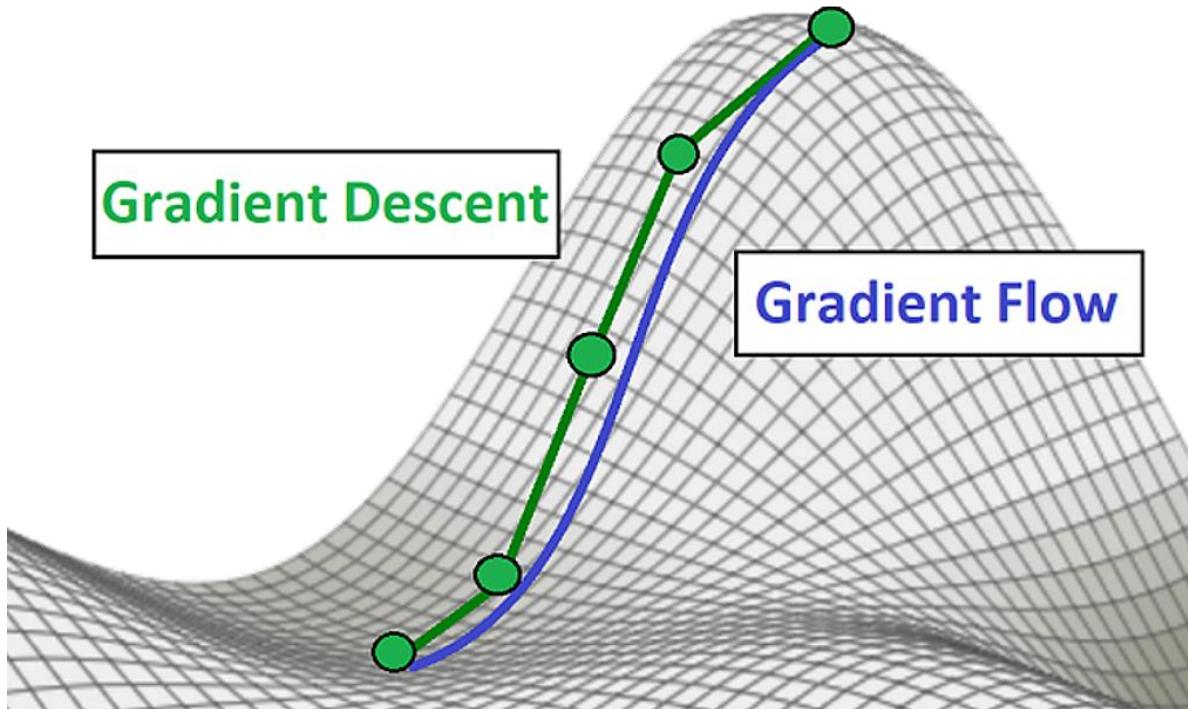
$$\mathbf{y}_\infty = \begin{bmatrix} \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{I} \end{bmatrix} \mathbf{y}_0 + \begin{bmatrix} \Sigma_1^{-1} \\ \mathbf{O} \end{bmatrix} \mathbf{U}^\top \mathbf{b}$$

Since $\mathbf{x} := \mathbf{V}\mathbf{y}$,

$$\mathbf{x}_\infty = \mathbf{V}_2 \mathbf{V}_2^\top \mathbf{x}_0 + \underbrace{\mathbf{V}_1 \Sigma_1^{-1} \mathbf{U}^\top \mathbf{b}}_{=\mathbf{x}_{LN}}$$

Therefore, we conclude that if \mathbf{x}_0 is orthogonal to the null space of \mathbf{A} , then gradient descent will converge to the least-norm solution.

Implicit regularization for neural networks with discrete GD

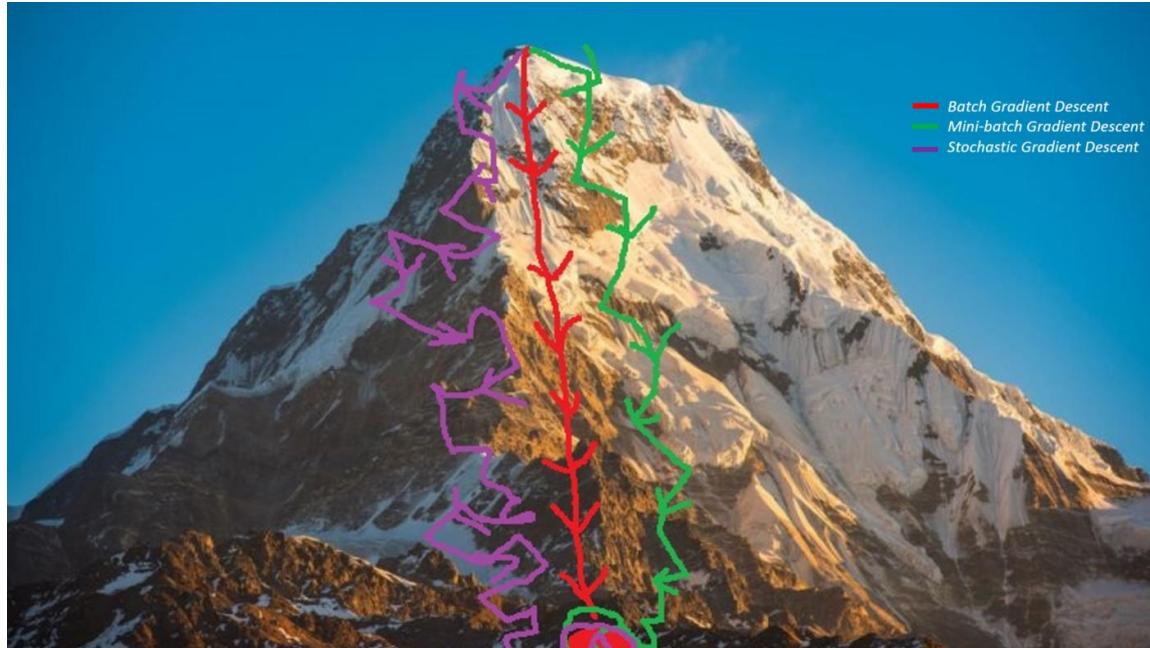


Effective Loss due to GD discretization

$$\tilde{\mathcal{L}}_{GD}(\mathbf{w}) = \mathcal{L}(\mathbf{w}) + \frac{\epsilon}{4} \|\nabla \mathcal{L}(\mathbf{w})\|^2$$

$$\dot{\mathbf{w}} = -\nabla_{\mathbf{w}} \tilde{\mathcal{L}}_{GD}(\mathbf{w}),$$

Implicit regularization for neural networks due to SGD



Effective Loss for SGD + discretization

$$\tilde{\mathcal{L}}_{SGD}(\mathbf{w}) = \tilde{\mathcal{L}}_{GD}(\mathbf{w}) + \frac{\epsilon}{4} \sum_{k=1}^m \|\nabla \mathcal{L}_k(\mathbf{w}) - \mathcal{L}(\mathbf{w})\|^2$$

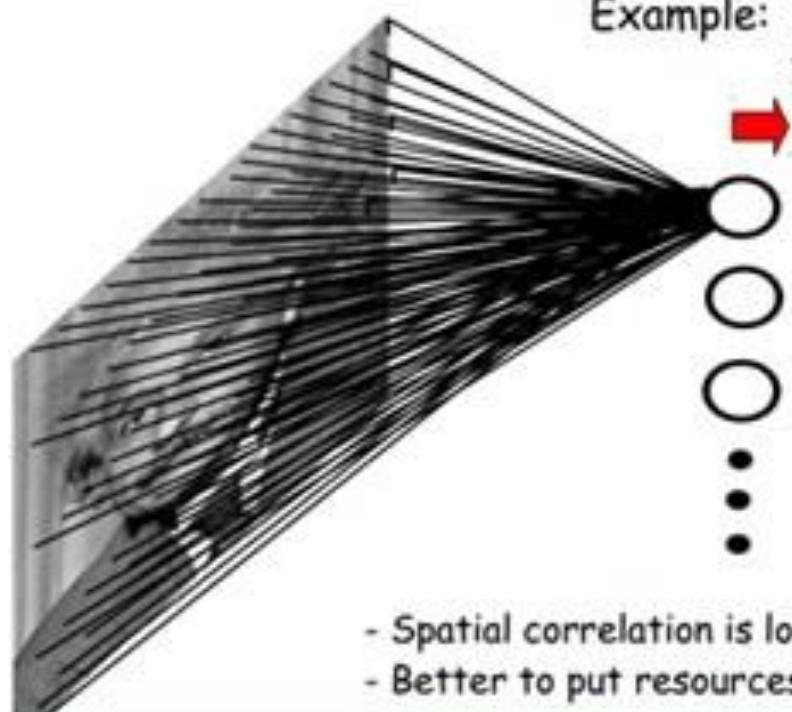
The second term estimates the variance of the minibatch gradients, which is a measure of stability, and hence of generalization ability.

Class 4

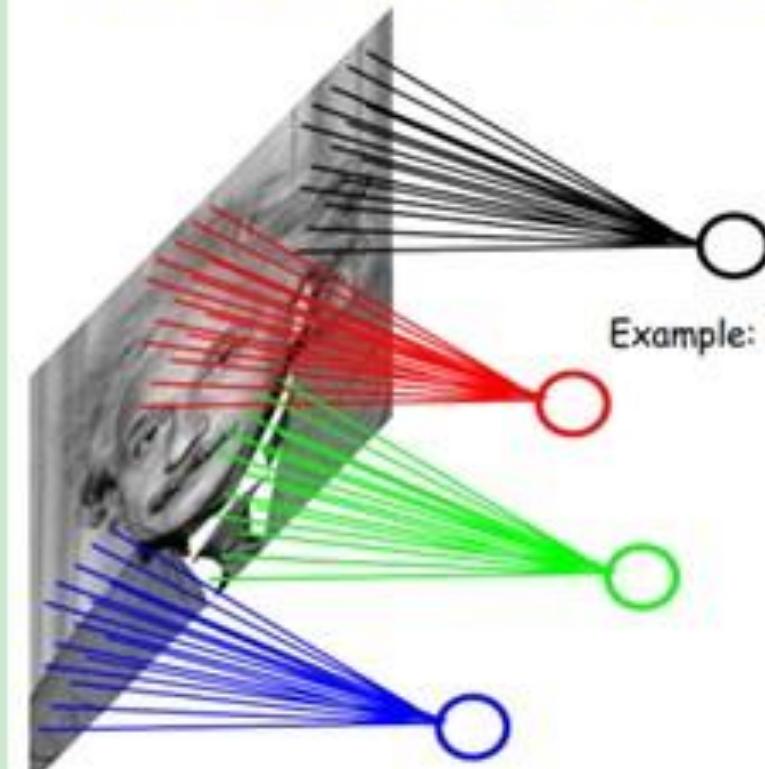
- Convolution neural networks: invariance, equivariance, locality.
- Implementing other types of invariances: data augmentations

Convolutional vs fully connected networks

FULLY CONNECTED NEURAL NET



LOCALLY CONNECTED NEURAL NET

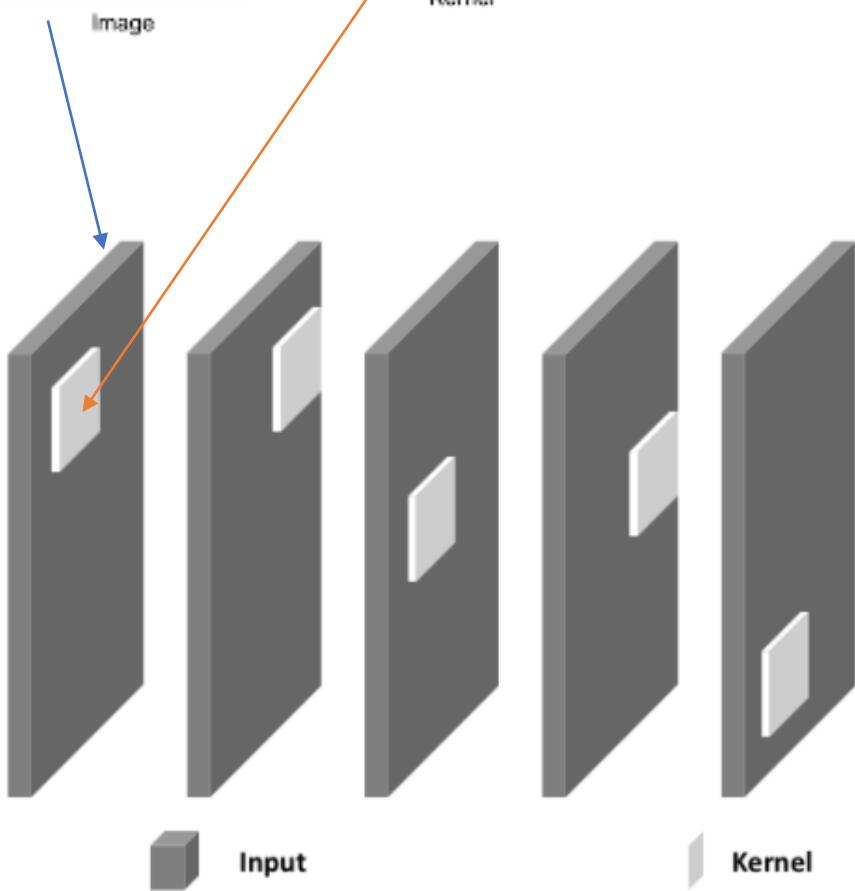


Convolution operation

$$\begin{matrix} 2 & 4 & 9 & 1 & 4 \\ 2 & 1 & 4 & 4 & 6 \\ 1 & 1 & 2 & 9 & 2 \\ 7 & 3 & 5 & 1 & 3 \\ 2 & 3 & 4 & 8 & 5 \end{matrix} \times \begin{matrix} 1 & 2 & 3 \\ -4 & 7 & 4 \\ 2 & -5 & 1 \end{matrix} = \begin{matrix} 51 & & \\ & & \\ & & \end{matrix}$$

Image Filter / Kernel Feature

$w * x$



Input image

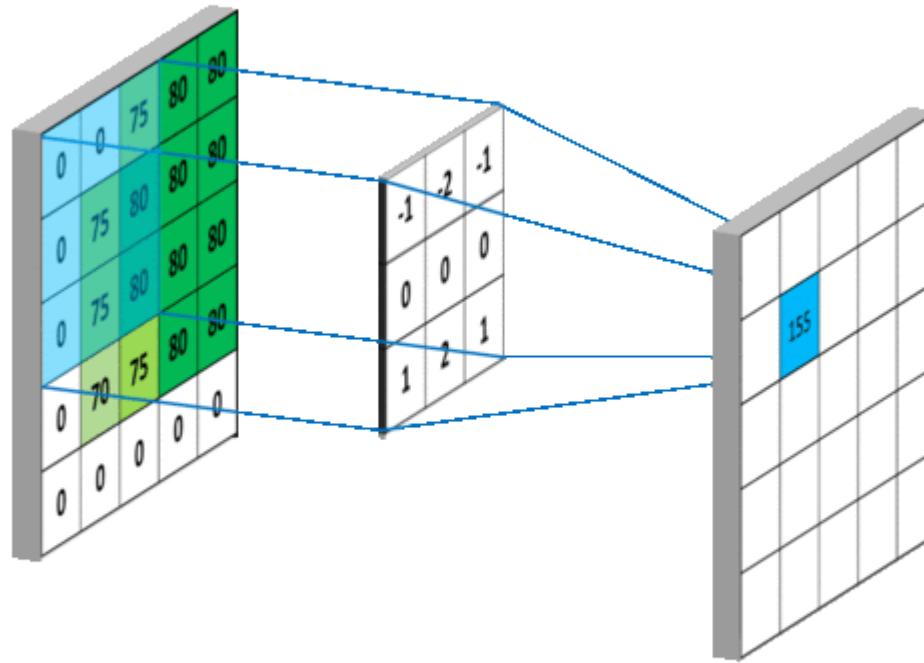


Convolution Kernel

$$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$$

Feature map

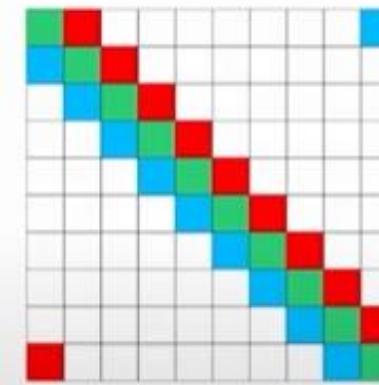




$$(f * g)(t) \stackrel{\text{def}}{=} \int_{-\infty}^{\infty} f(\tau)g(t - \tau) \, d\tau = \int_{\infty}^{-\infty} f(\tau)T_tg(-\tau)d\tau = \langle f, T_tg \rangle$$

Discrete version of the convolution

$$C(w) = \begin{bmatrix} c_1 & c_3 & c_2 \\ c_2 & c_1 & c_3 \\ c_3 & c_2 & c_1 \end{bmatrix}_w$$



Why convolutions? Invariance and equivariance

$$(w * x)_i = (W^T x)_i = \langle w_i, x \rangle = \langle g_i w, x \rangle$$

$$f(x) = \sigma(W^T x) = \begin{bmatrix} \sigma \langle g_1 w, x \rangle \\ \vdots \\ \sigma \langle g_N w, x \rangle \end{bmatrix}$$

Why convolutions? Invariance and equivariance

Group: set \mathcal{G} with composition/multiplication operation satisfying:

- ① closure: $gg' \in \mathcal{G}, \forall g, g' \in \mathcal{G}$
- ② associativity: $(gg')g'' = g(g'g'') \in \mathcal{G}, \forall g, g', g'' \in \mathcal{G}$
- ③ identity: there exists $Id \in \mathcal{G}$ such that $Idg = gId = g \quad \forall g \in \mathcal{G}$
- ④ invertibility: $\forall g \in \mathcal{G}$ there exists $g^{-1} \in \mathcal{G} : gg^{-1} = Id.$

Rotation



Why convolutions? Invariance and equivariance

$$f(gx) = \begin{bmatrix} \sigma \langle g_1 w, gx \rangle \\ \vdots \\ \sigma \langle g_N w, gx \rangle \end{bmatrix} = \begin{bmatrix} \sigma \langle w, g_1^T gx \rangle \\ \vdots \\ \sigma \langle w, g_N^T gx \rangle \end{bmatrix} = P_g \begin{bmatrix} \sigma \langle w, gx \rangle \\ \vdots \\ \sigma \langle w, gx \rangle \end{bmatrix} = P_g f(x)$$



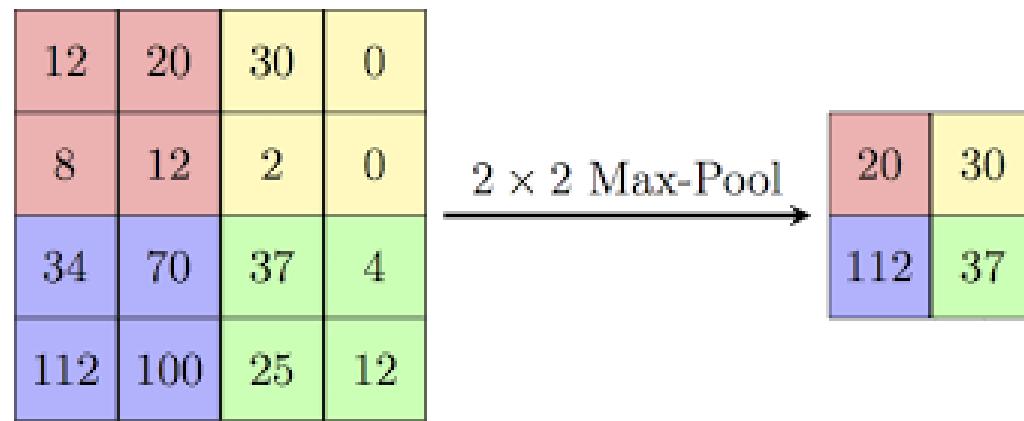
All learnable functions

*All learnable
equivariant
functions*

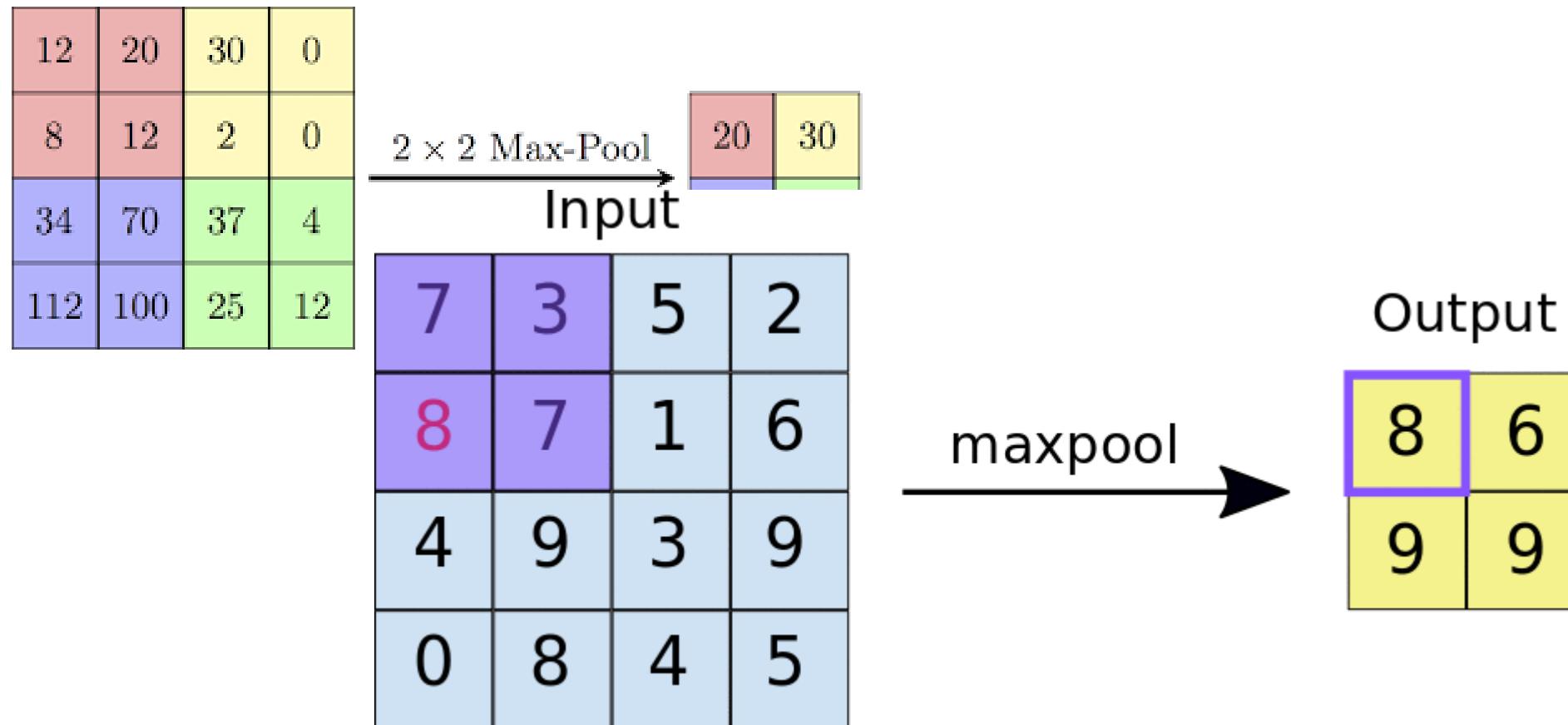
**Functions you
actually wanted
to learn.**

*All learnable
functions
constrained
by your data.*

Pooling operation



Pooling operation



Can you think of other pooling operations?

From equivariance to invariance

$$\sum_i (f(gx))_i = \sum_i (P_g f(x))_{\textcolor{red}{i}}$$

Why is this expression invariant to g transformations?

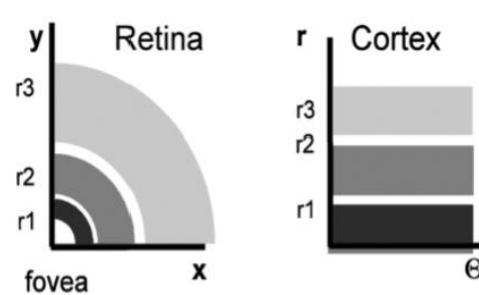
Why do we care? Mimicking the brain

Equivariance and invariance are important properties not only
in deep nets but also in visual cortex signal processing

Invariant object recognition

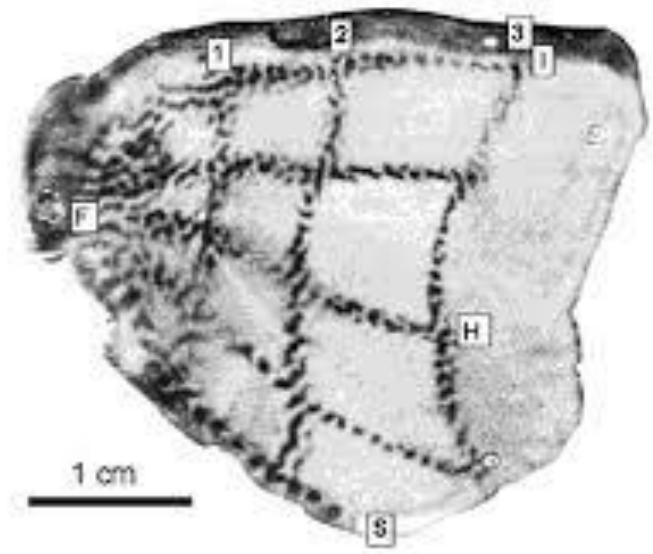
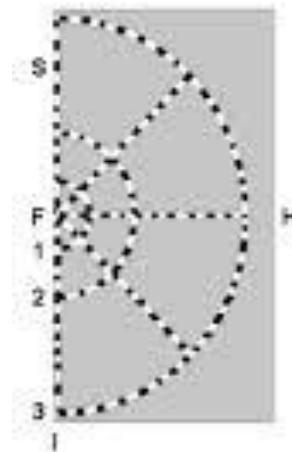


Retinotopic mapping



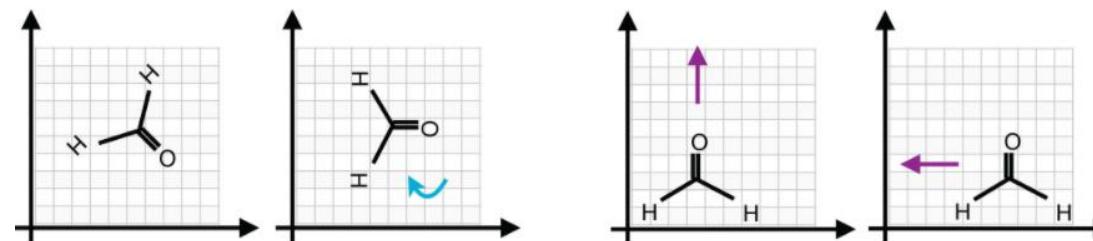
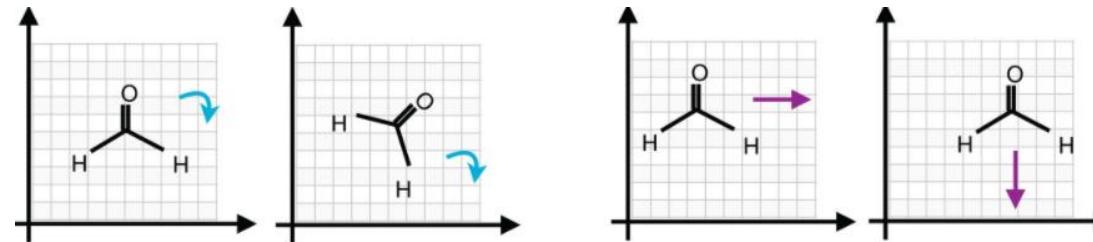
Invariance: $f(gx) = f(x)$.

Equivariance: $f(gx) = gf(x)$



While in CNNs invariance/equivariance are natural (at least for translations) how these properties could possibly be implemented in visual cortex is not clear.

Why do we care? Molecule symmetries



Rotations

Translations

$$\begin{bmatrix} \text{H} & 0 & 1 & 0 \\ 0 & \text{H} & 1 & 0 \\ 1 & 1 & \text{C} & 2 \\ 0 & 0 & 2 & \text{O} \end{bmatrix} \quad \begin{array}{c} 4 \\ \leftarrow \\ 1 \end{array} \quad \begin{array}{c} 3 \\ \leftarrow \\ 2 \end{array}$$

$$\begin{array}{c} 3 \\ \rightarrow \\ 1 \end{array} \quad \begin{bmatrix} \text{H} & 0 & 0 & 1 \\ 0 & \text{H} & 0 & 1 \\ 0 & 0 & \text{O} & 2 \\ 1 & 1 & 2 & \text{C} \end{bmatrix}$$

Permutations

Convolutional networks and convolution kernels

Convolution layer representation (with average pooling),

$$\Phi(x) = \left(\sum_g s(\langle x, gt_1 \rangle), \dots, \sum_g s(\langle x, gt_T \rangle) \right), \quad x \in \mathcal{X}.$$

The **associated kernel** can be written as

$$K(x, x') = \langle \Phi(x), \Phi(x') \rangle = \sum_t \left(\sum_g s(\langle x, gt_1 \rangle) \cdot \sum_g s(\langle x', gt_1 \rangle) \right)$$

Convolutional networks and convolution kernels (cont.)

Using linearity

$$K(x, x') = \sum_g \sum_{g'} \underbrace{\left(\sum_t s(\langle x, g t_1 \rangle) \cdot s(\langle x', g' t_1 \rangle) \right)}_{K'(x, x')}.$$

Convolution kernels

$$K(x, x') = \sum_g \sum_{g'} K'(gx, g'x')$$

- ▶ Kernel of the above form are a special case of **convolution kernels**.
- ▶ If g_1, \dots, g_G form a **group** then the kernel is invariant
 $K(x, x') = K(x, gx')$, i.e. $\Phi(x) = \Phi(gx)$.

What if we do not have a group transformation?

Group: set \mathcal{G} with composition/multiplication operation satisfying:

- ① closure: $gg' \in \mathcal{G}, \forall g, g' \in \mathcal{G}$
- ② associativity: $(gg')g'' = g(g'g'') \in \mathcal{G}, \forall g, g', g'' \in \mathcal{G}$
- ③ identity: there exists $Id \in \mathcal{G}$ such that $Idg = gId = g \quad \forall g \in \mathcal{G}$
- ④ invertibility: $\forall g \in \mathcal{G}$ there exists $g^{-1} \in \mathcal{G} : gg^{-1} = Id$.

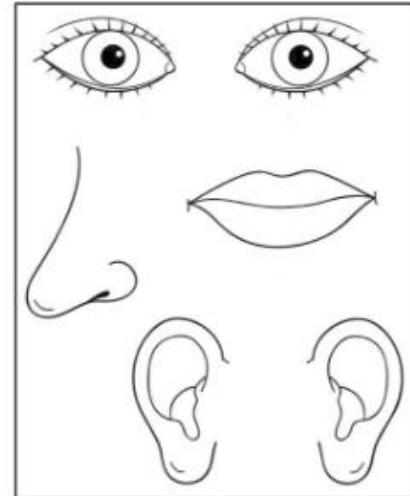
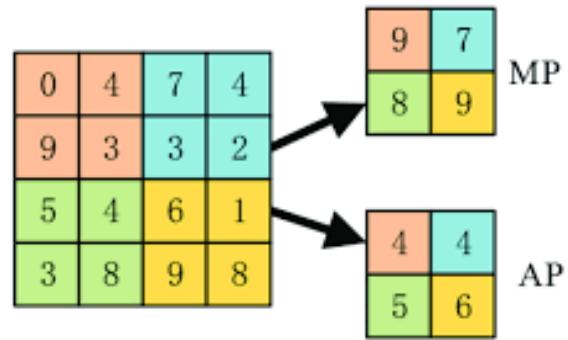
Rotation



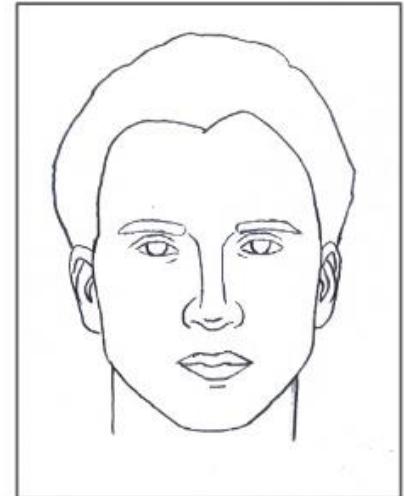
Why local filters?

- Robustness to (small) deformations
- Parts reusability in a hierarchical model

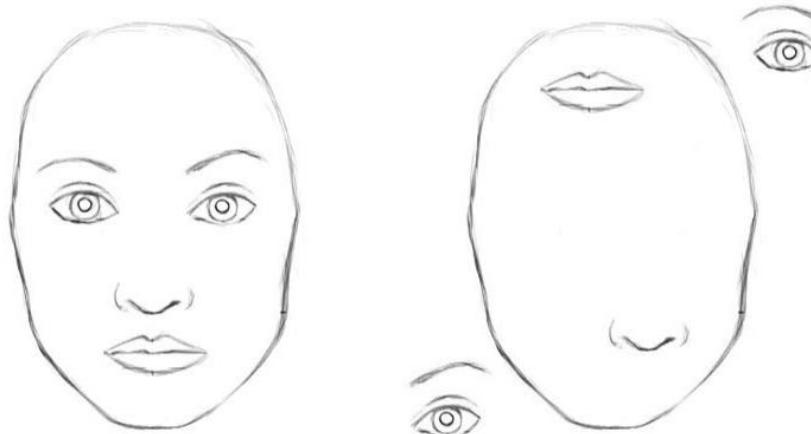
Robustness to deformations: pooling+ locality



Face

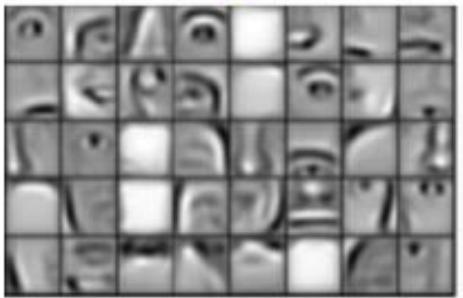


Face

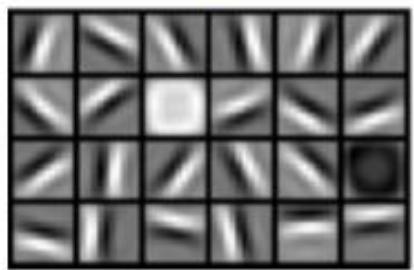




(4) object models



(3) object parts
(combination
of edges)



(2) edges

(1) pixels

Locality, parts,
compositionality:
reusability

Faces



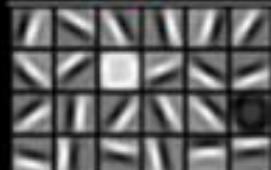
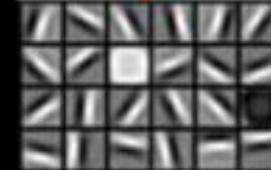
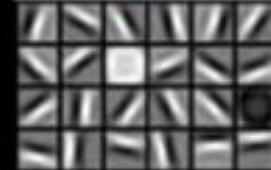
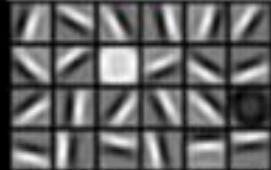
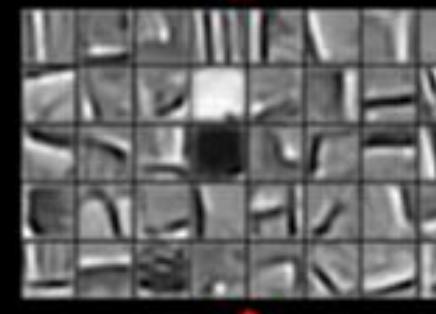
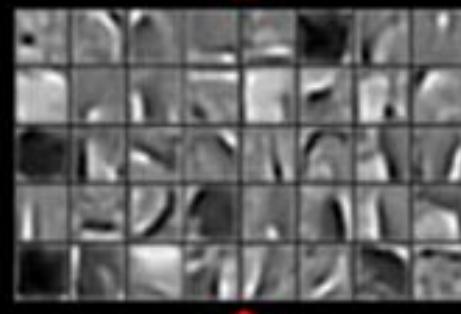
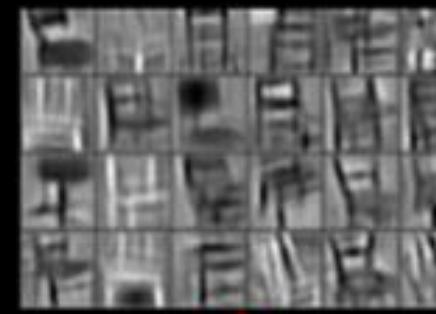
Cars



Elephants



Chairs

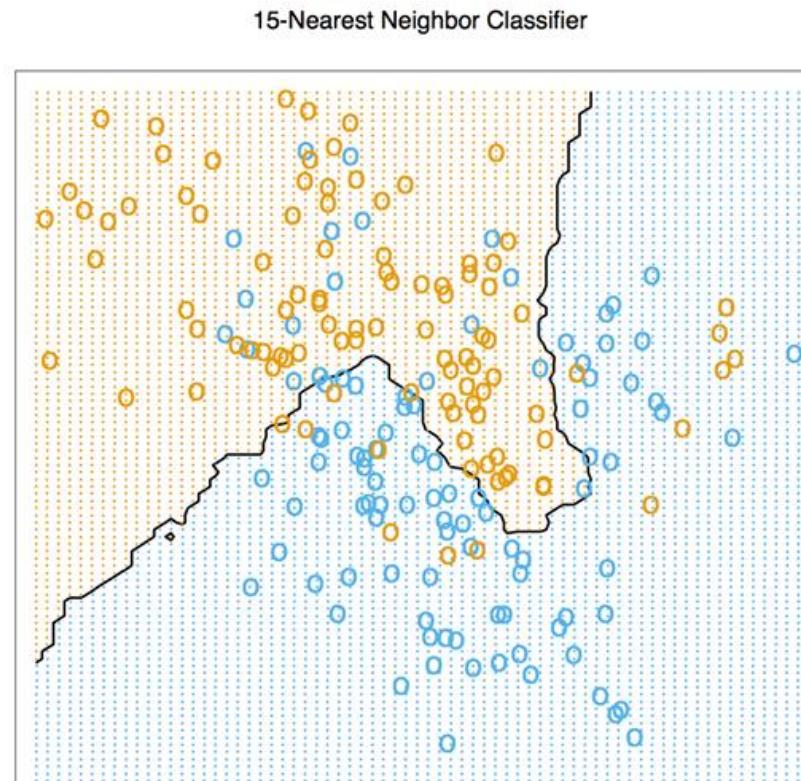


Why convolutionality, locality and compositionality help to have a "better" learning?

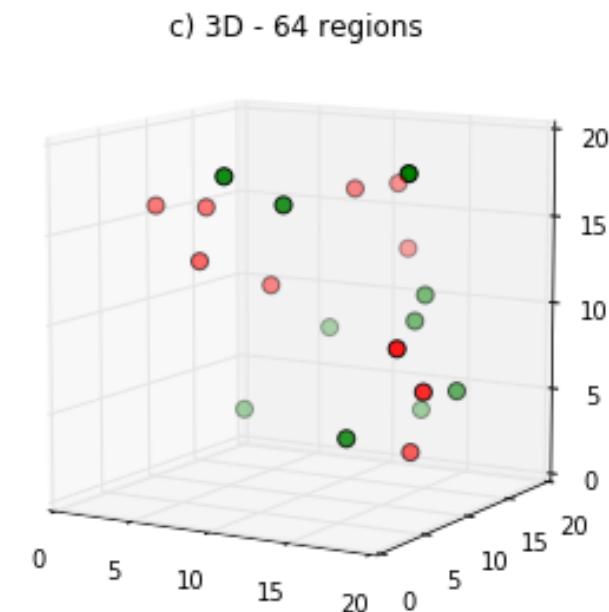
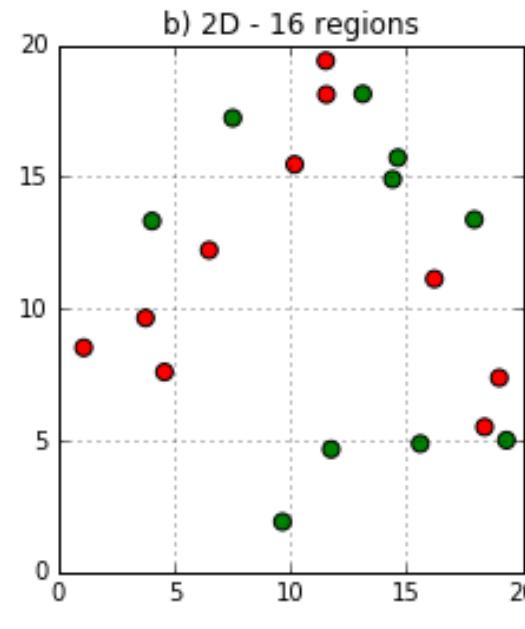
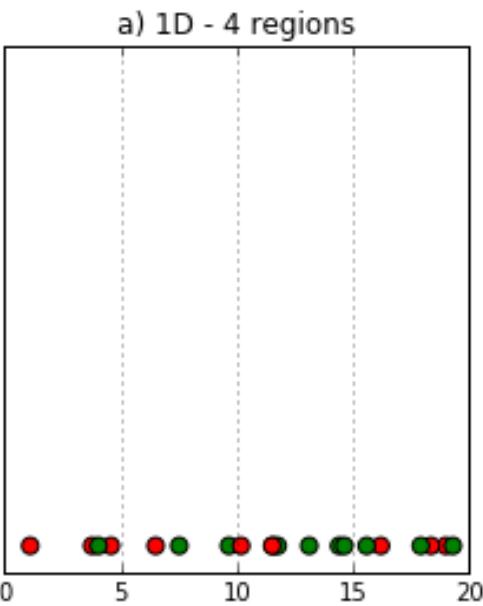
How do you estimate a function? Example with NNC

Supervised classification or regression often rely on local averages:

- **Classification** : you know the classes of n points from your learning database, you can classify a new point x by computing the most represented class in the neighborhood of x .



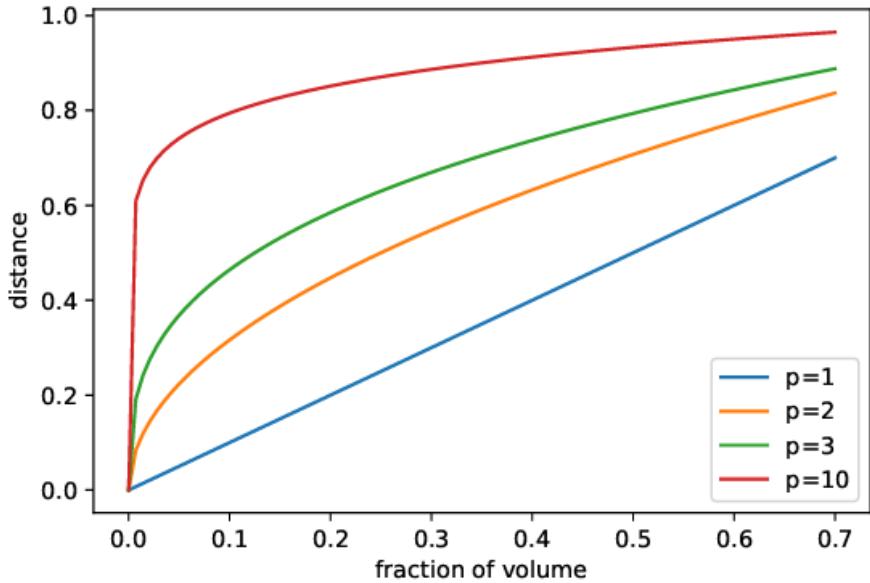
However...



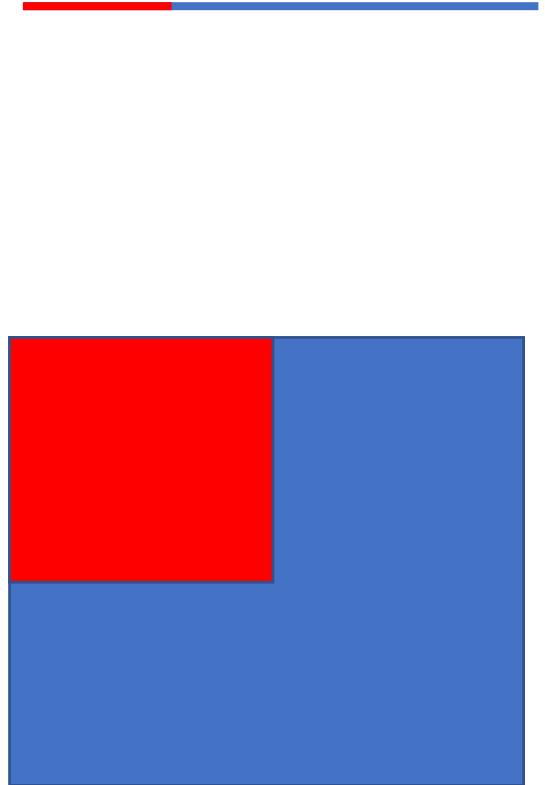
High dimensional spaces are empty

Assume your data lives in $[0, 1]^p$. To capture a neighborhood which represents a fraction s of the hypercube volume, you need the edge length to be $s^{1/p}$

- $s = 0.1, p = 10, s^{1/p} = 0.63$
- $s = 0.01, p = 10, s^{1/p} = 0.8$



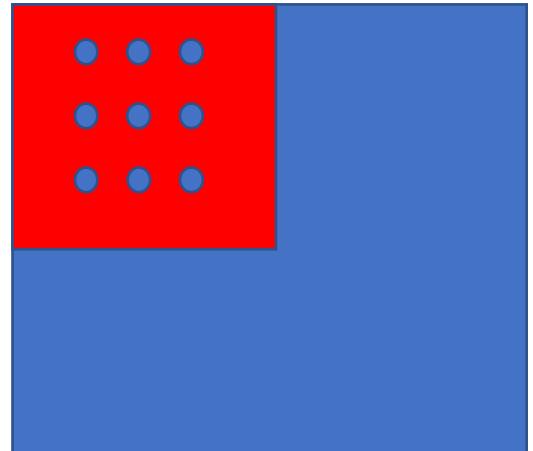
Neighbors are no longer local



Curse of dimensionality



To approximate a (Lipschitz) continuous function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ with ϵ accuracy one needs $O(\epsilon^{-d})$ samples

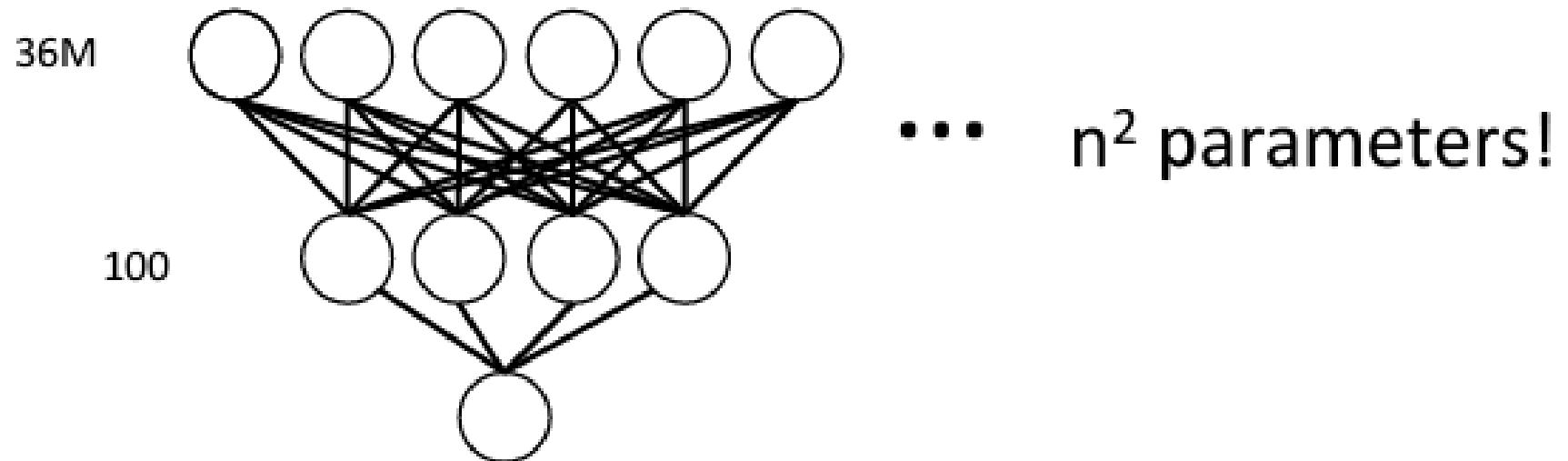


Input image resolution = 12 Mpixel * 3 channels = 36M elements

With $\epsilon \sim 0.1$, we need $10^{36000000}$ samples (10^{78} to 10^{82} atoms in the known, observable universe)

Same sample density (same precision): we need P^d samples

Curse of dimensionality: networks

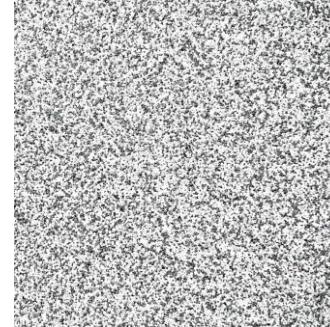


- Input image resolution = 12 Mpixel * 3 channels = 36M elements
- MLP with one hidden layer of width=100 has **3.6B parameters**

Do we need all parameters? This is true if the image space has all possible images



but also



But the statistics of the images is translation and scale invariant

So why weight sharing helps is helping in reducing the curse of dimensionality?

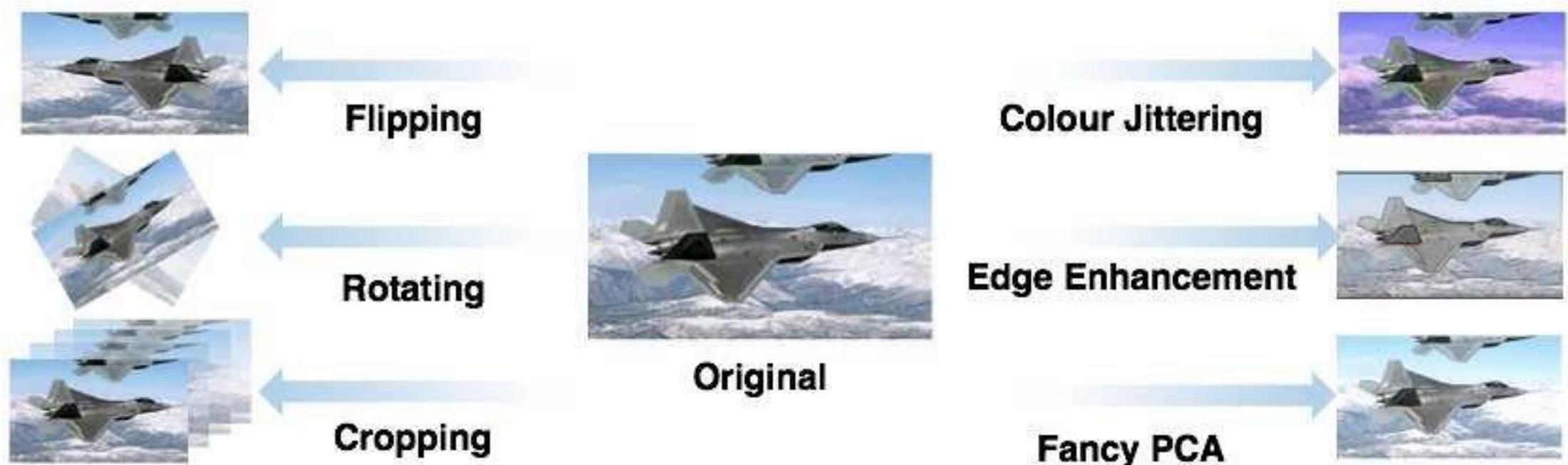
Data augmentation and other types of invariances: building an invariant representation

$$f(x) = \sum_i \sigma \langle g_i w, x \rangle = f(gx)$$

$$\{x_i\}_i \rightarrow \{g_j x_i\}_{ij}$$

$$f(x) = \sum \sigma \langle w, g_i x \rangle = f(gx)$$

Data augmentation and other types of invariances

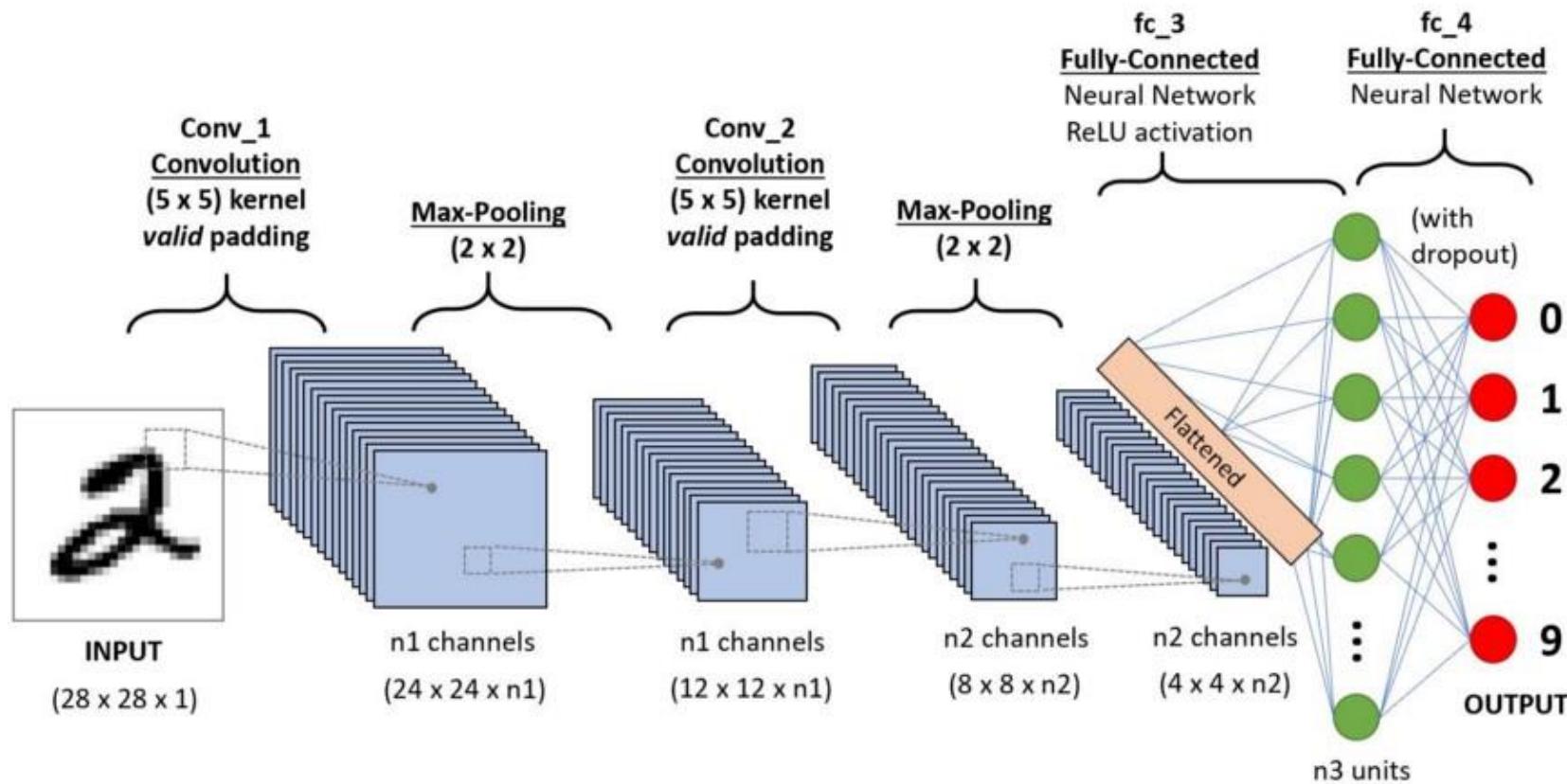


An example of a convolutional architecture: Lenet

MNIST dataset



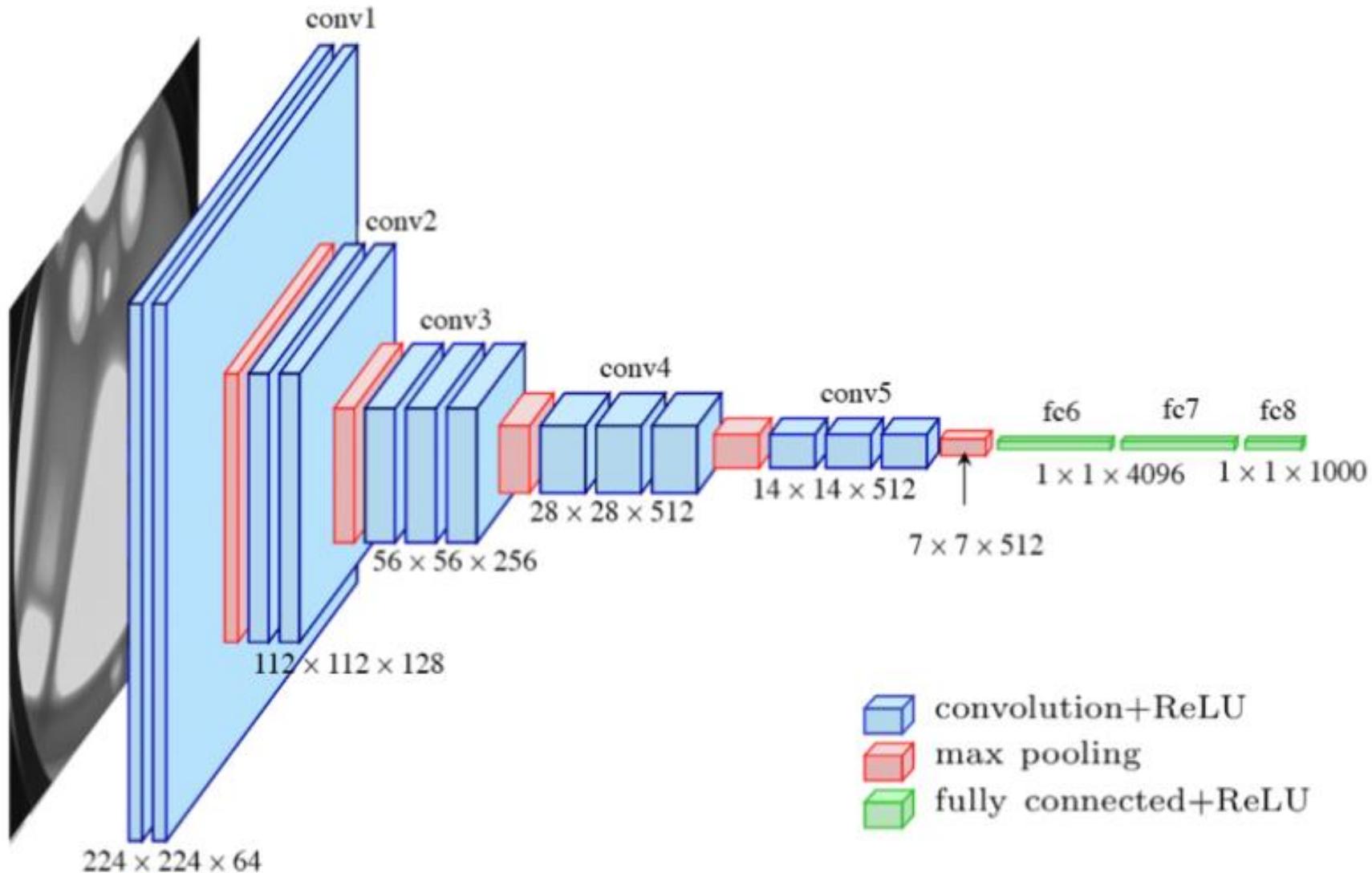
CNNs started the ML revolution



Many CNNs and image datasets



VGG architecture



Beyond convolution: Resnet

[He et al., 2015]

Hypothesis: the problem is an *optimization* problem, deeper models are harder to optimize

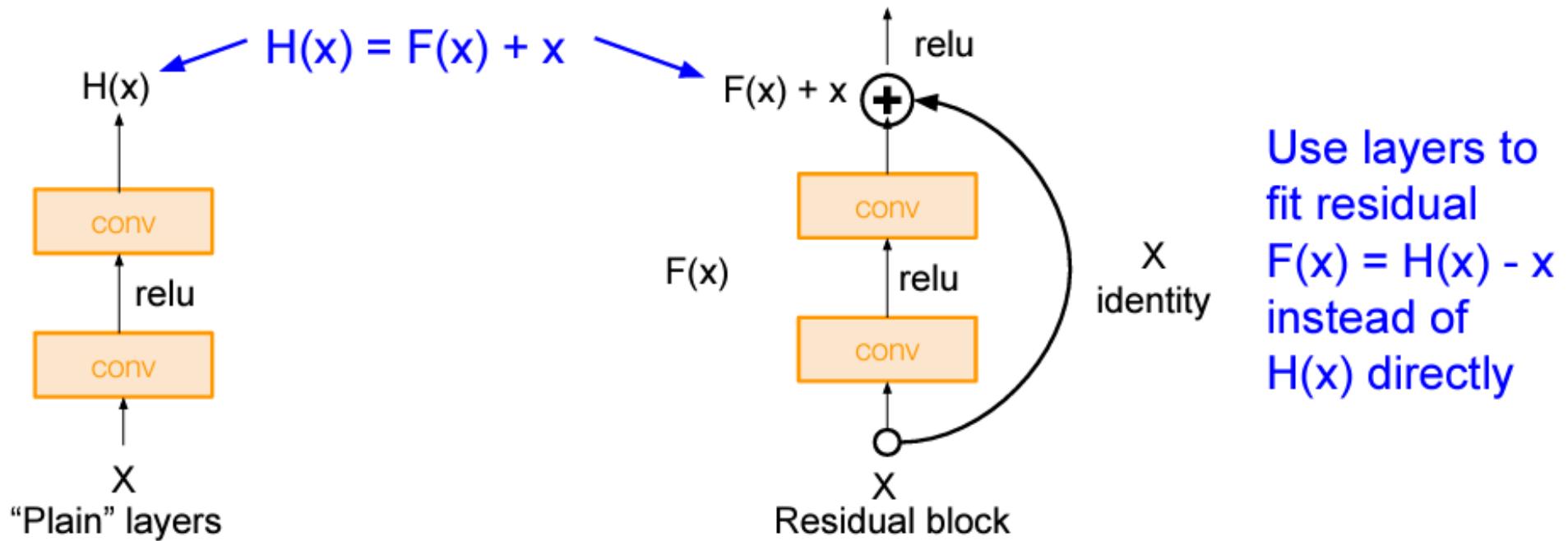
The deeper model should be able to perform at least as well as the shallower model.

A solution by construction is copying the learned layers from the shallower model and setting additional layers to identity mapping.

Resnet: intuition

[He et al., 2015]

Solution: Use network layers to fit a residual mapping instead of directly trying to fit a desired underlying mapping



ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners

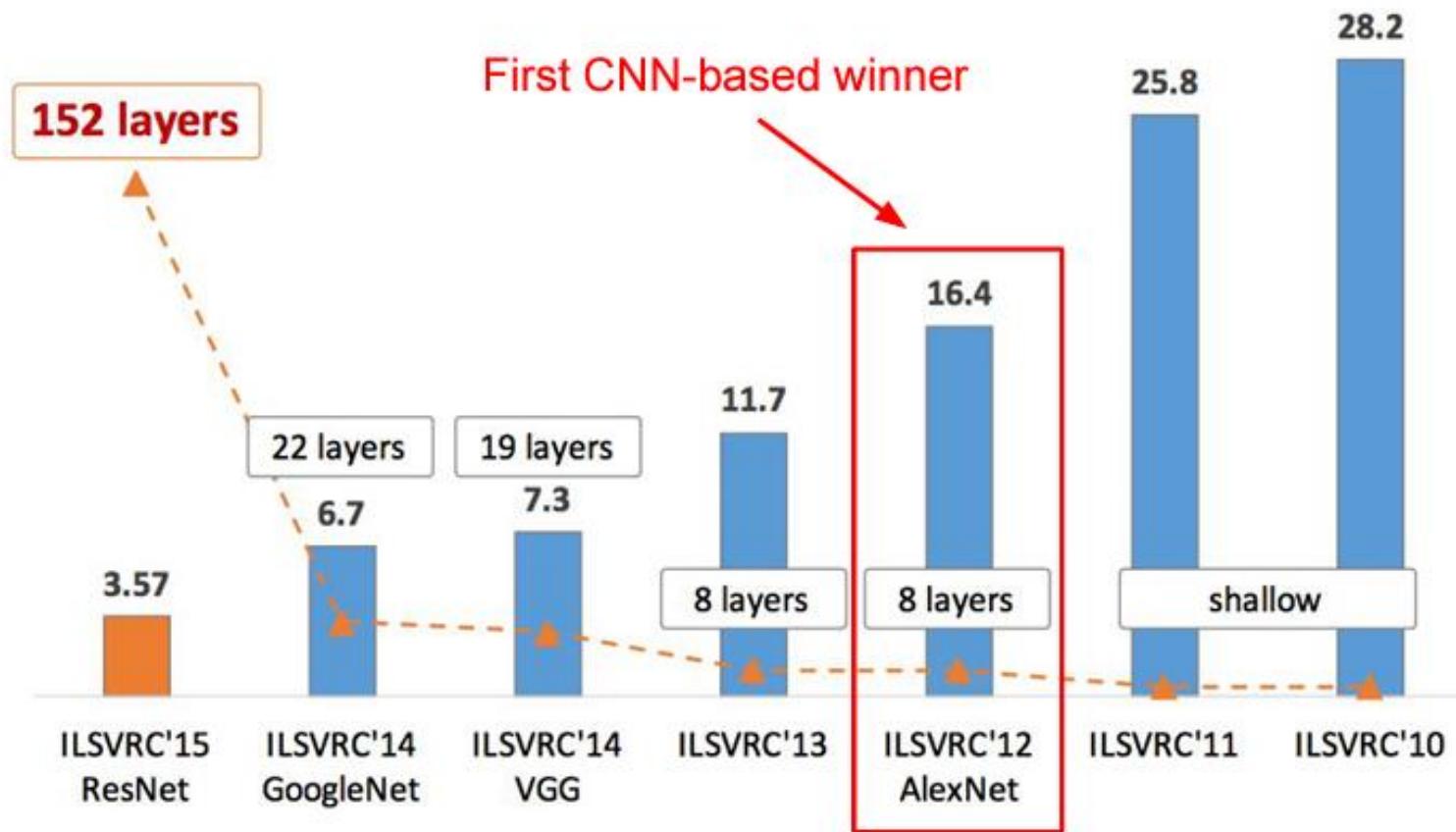


Figure copyright Kaiming He, 2016. Reproduced with permission.

Class 5

Adversarial attacks

What is an adversarial example?



Sheepdog or Mop?

What is an adversarial example?



Sloth or Pain au chocolat?

What is an adversarial example?



Puppy or Bagel?

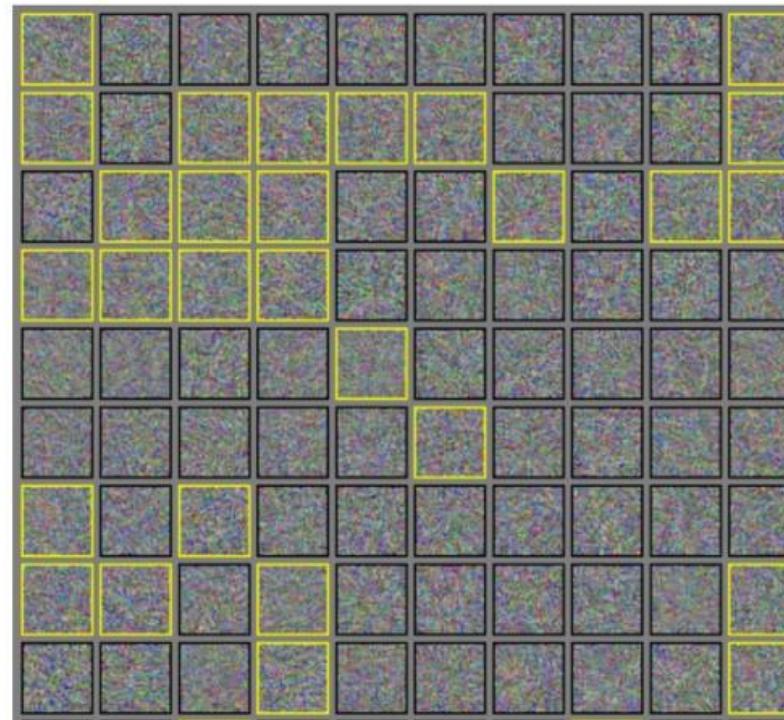
What is an adversarial example?



Chihuahua or Muffin?

Suppose you train a net to classify airplanes vs cars. Then...

- If you start with random noise and take one gradient step, you can often produce a confident classification as some category.
- The images highlighted in yellow are classified as “airplane” with > 50% probability.



- A variant: search for the image closest to the original one which is misclassified as a particular category (e.g. ostrich).
- This is called a **targeted adversarial example**, since it targets a particular category.
- The following adversarial examples are misclassified as ostriches. (Middle = perturbation $\times 10$.)

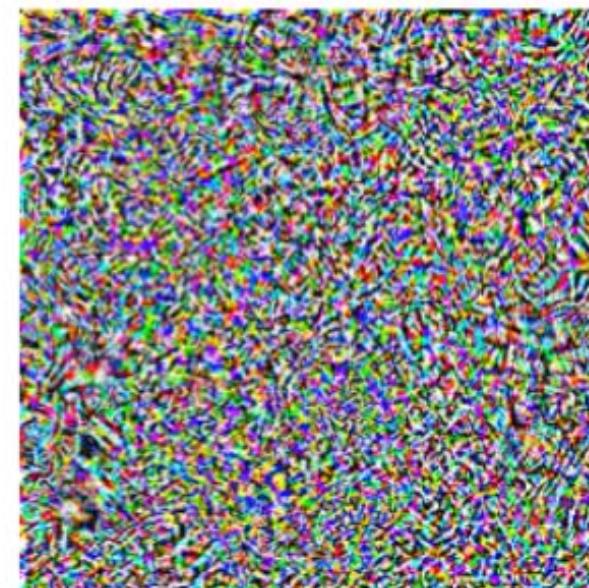


Adversarial attacks

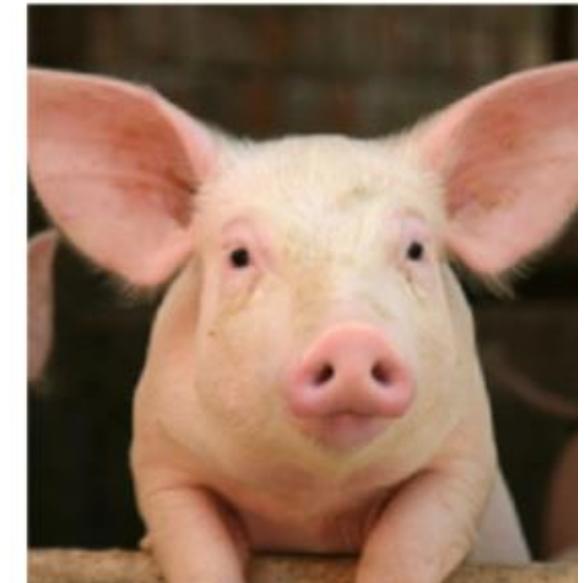
“pig”



+ 0.005 x



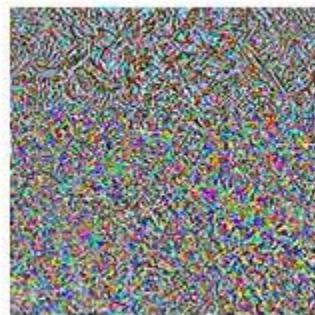
“airliner”



Adversarial examples for CNNs



Alps: 94%



Dog: 100%

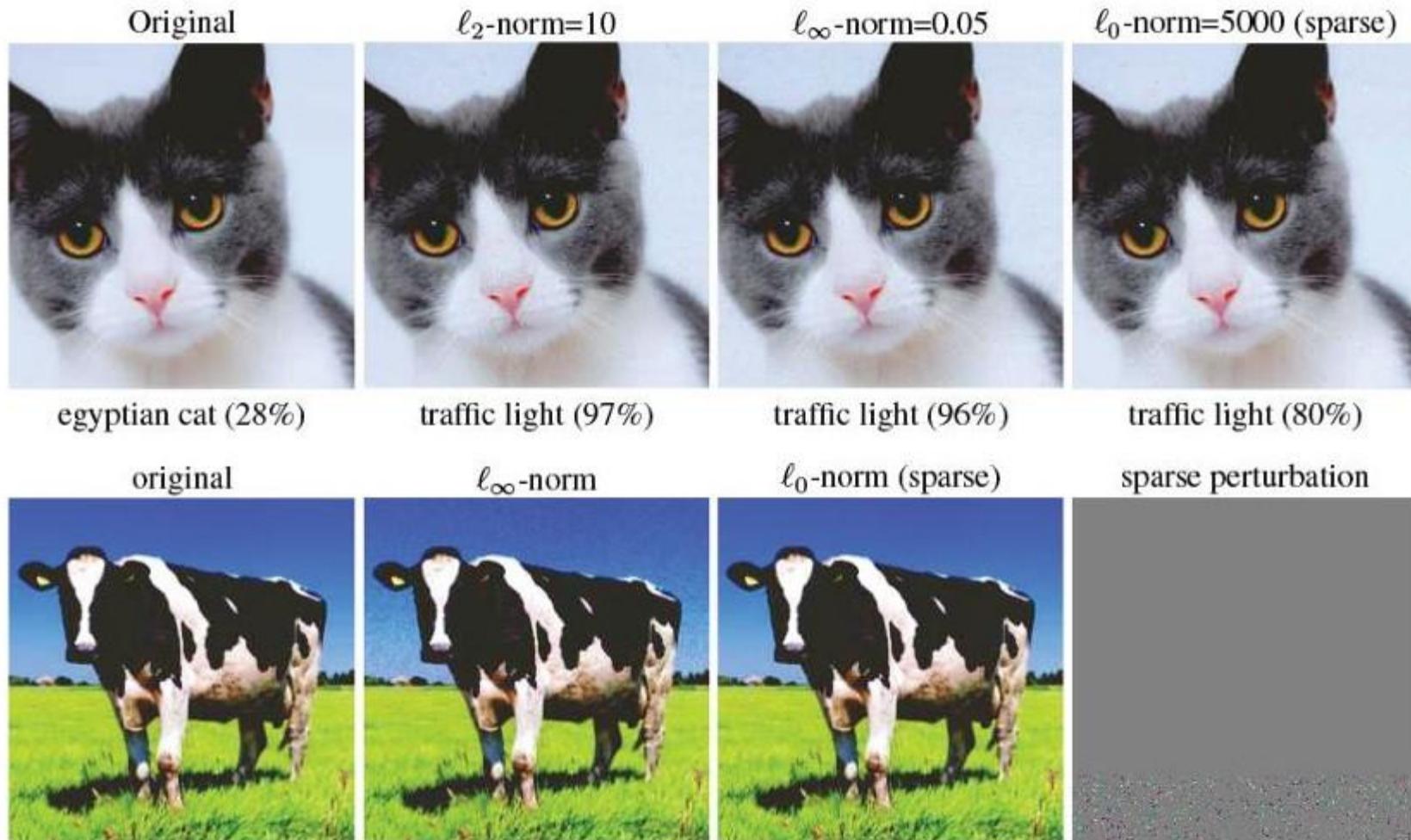


Puffer: 98%



Crab: 100%

Adversarial Examples for different norms



Shafahi et al. (2019). Are adversarial examples inevitable? ICLR (to appear)

- The paper which introduced adversarial examples (in 2013) was titled “Intriguing Properties of Neural Networks.”
- Now they’re regarded as a serious security threat.
 - Nobody has found a reliable method yet to defend against them.
 - Adversarial examples transfer to different networks trained on a disjoint subset of the training set!
 - You don’t need access to the original network; you can train up a new network to match its predictions, and then construct adversarial examples for that.
 - Attack carried out against proprietary classification networks accessed using prediction APIs (MetaMind, Amazon, Google)

Adversarial examples are transferable!!!!

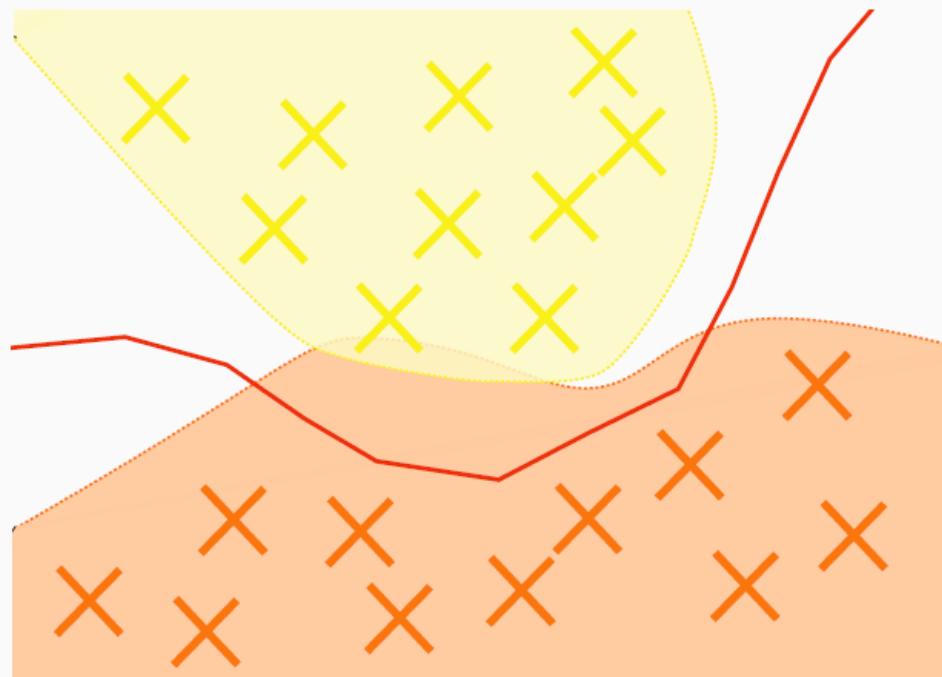
- adversarial examples are highly transferable
- it is very likely that an adversarial example of one network can fool another network
- transferability depends on the type of attack

OOD generalization problems



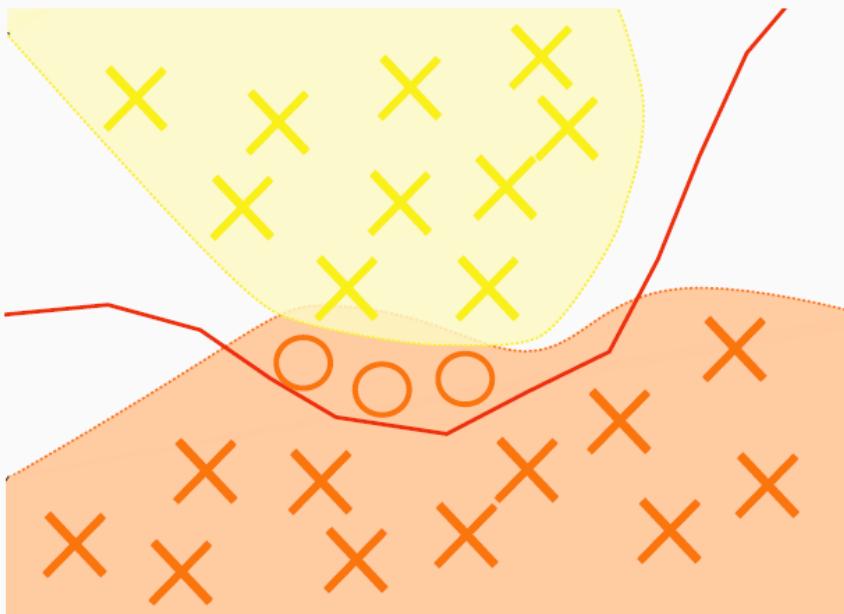
What is the intuition for adversarial attacks?

Decision Boundary of the Model



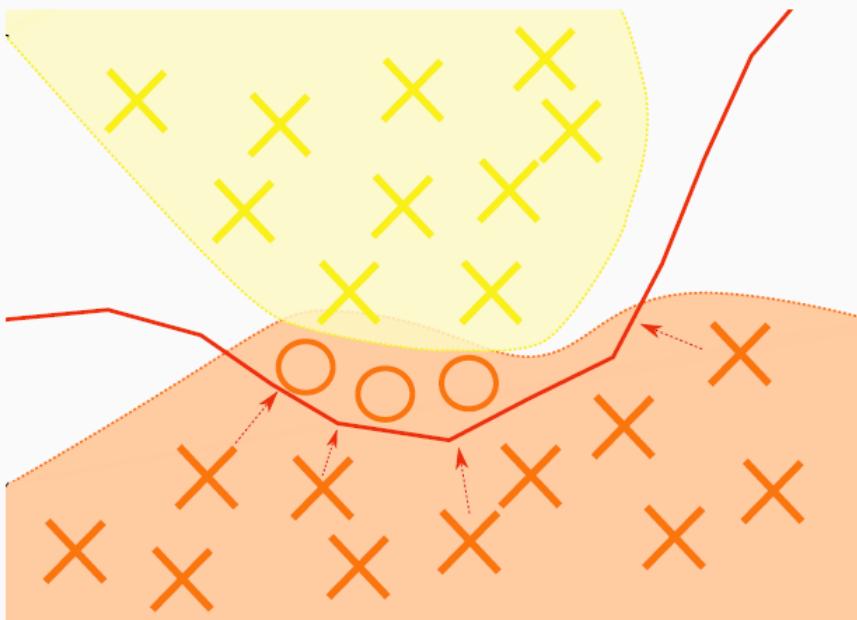
The **learned model** slightly differs from the **true** data distribution...

The Space of Adversarial Examples



... which makes room for **adversarial examples**.

Attack: Use the Adversarial Directions



- Most attacks try to move inputs across the boundary.
- Attacking with a random distortion doesn't work well in practice.

Definition 1 (Adversarial Attack). Let $\mathbf{x}_0 \in \mathbb{R}^d$ be a data point belong to class \mathcal{C}_i . Define a target class \mathcal{C}_t . An **adversarial attack** is a mapping $\mathcal{A} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ such that the perturbed data

$$\mathbf{x} = \mathcal{A}(\mathbf{x}_0)$$

is misclassified as \mathcal{C}_t .

Definition 2 (Additive Adversarial Attack). Let $\mathbf{x}_0 \in \mathbb{R}^d$ be a data point belong to class \mathcal{C}_i . Define a target class \mathcal{C}_t . An **additive** adversarial attack is an addition of a perturbation $\mathbf{r} \in \mathbb{R}^d$ such that the perturbed data

$$\mathbf{x} = \mathbf{x}_0 + \mathbf{r}$$

is misclassified as \mathcal{C}_t .

Targeted vs untargeted attacks

Non-targeted adversarial examples: mislead the model to provide **any wrong** prediction

$$\begin{aligned} & \max_{x^*} \ell(f_\theta(x^*), y) \\ \text{s.t. } & d(x, x^*) \leq B \end{aligned}$$

Max distance from the true prediction

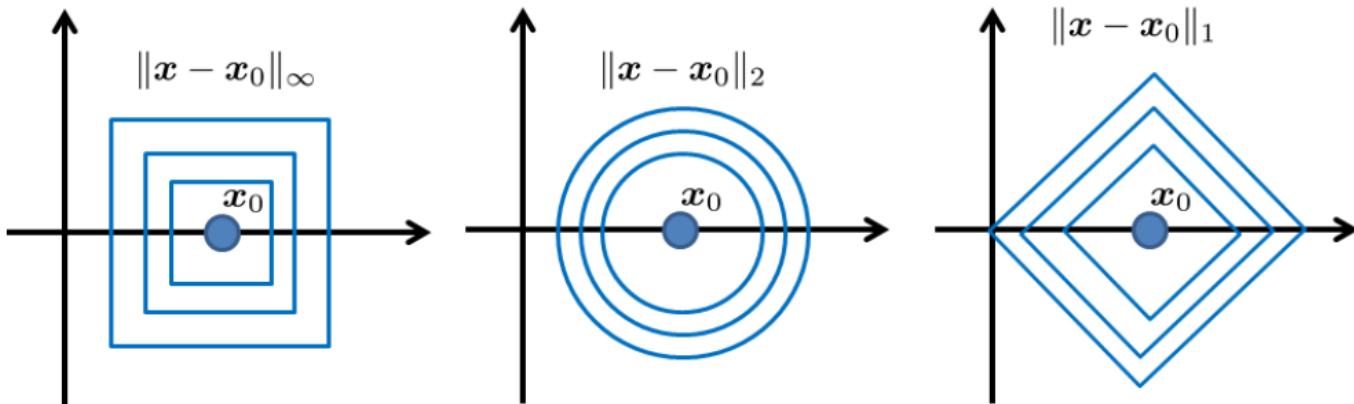
Targeted adversarial examples: mislead the model to provide the **target prediction $y^* \neq y$** specified by the adversary

$$\begin{aligned} & \min_{x^*} \ell(f_\theta(x^*), y^*) \\ \text{s.t. } & d(x, x^*) \leq B \end{aligned}$$

Min distance from the wrong prediction

$d(x, x^*)$ is an ℓ_n norm in most existing work

Minimum norm attack



Theorem 2 (Minimum ℓ_2 Norm Attack for Two-Class Linear Classifier). *The adversarial attack to a two-class linear classifier is the solution of*

$$\underset{x}{\text{minimize}} \quad \|x - x_0\|^2 \quad \text{subject to} \quad w^T x + w_0 = 0, \quad (3.25)$$

which is given by

$$x^* = x_0 - \left(\frac{w^T x_0 + w_0}{\|w\|_2} \right) \frac{w}{\|w\|_2}. \quad (3.26)$$

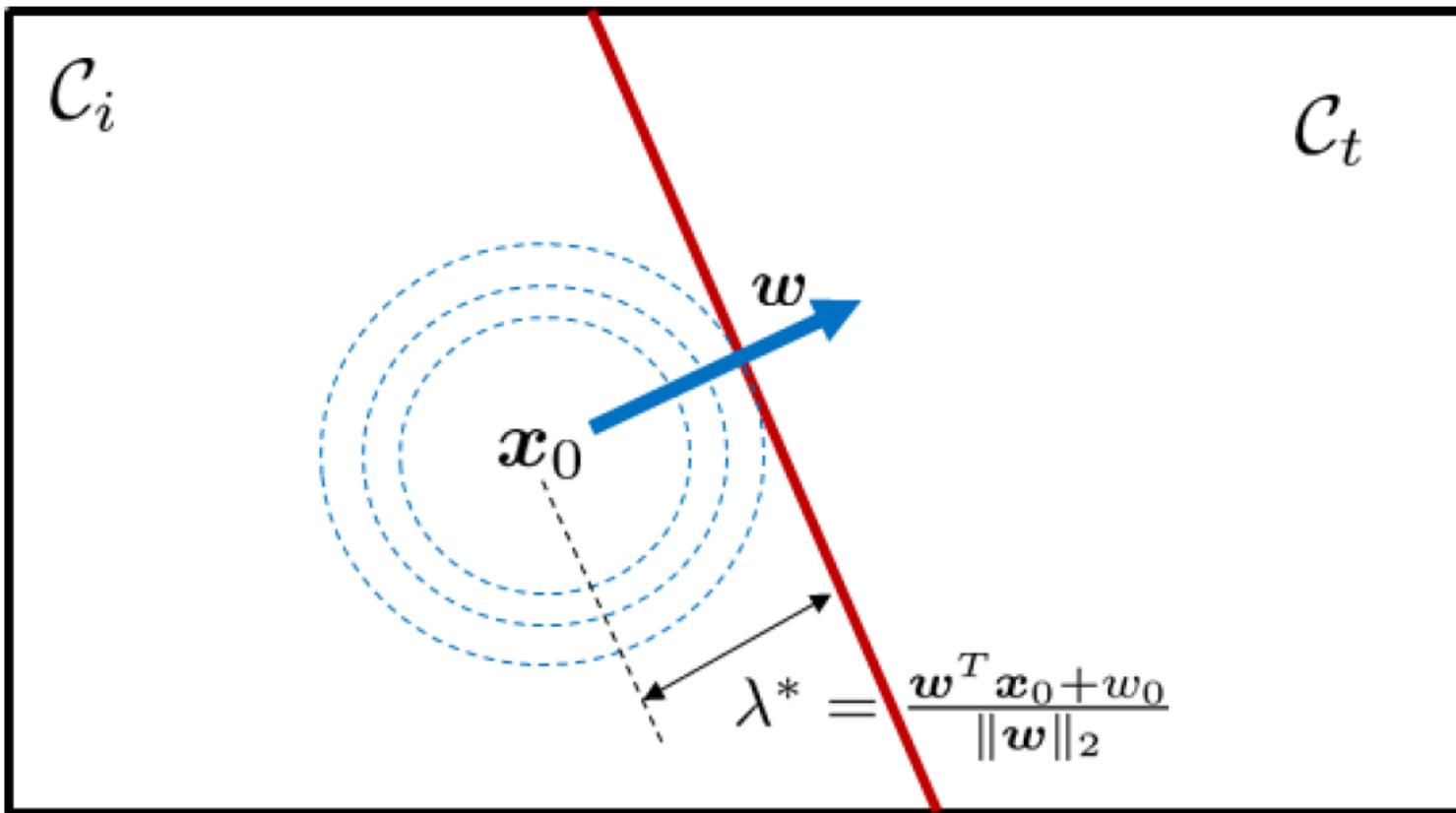
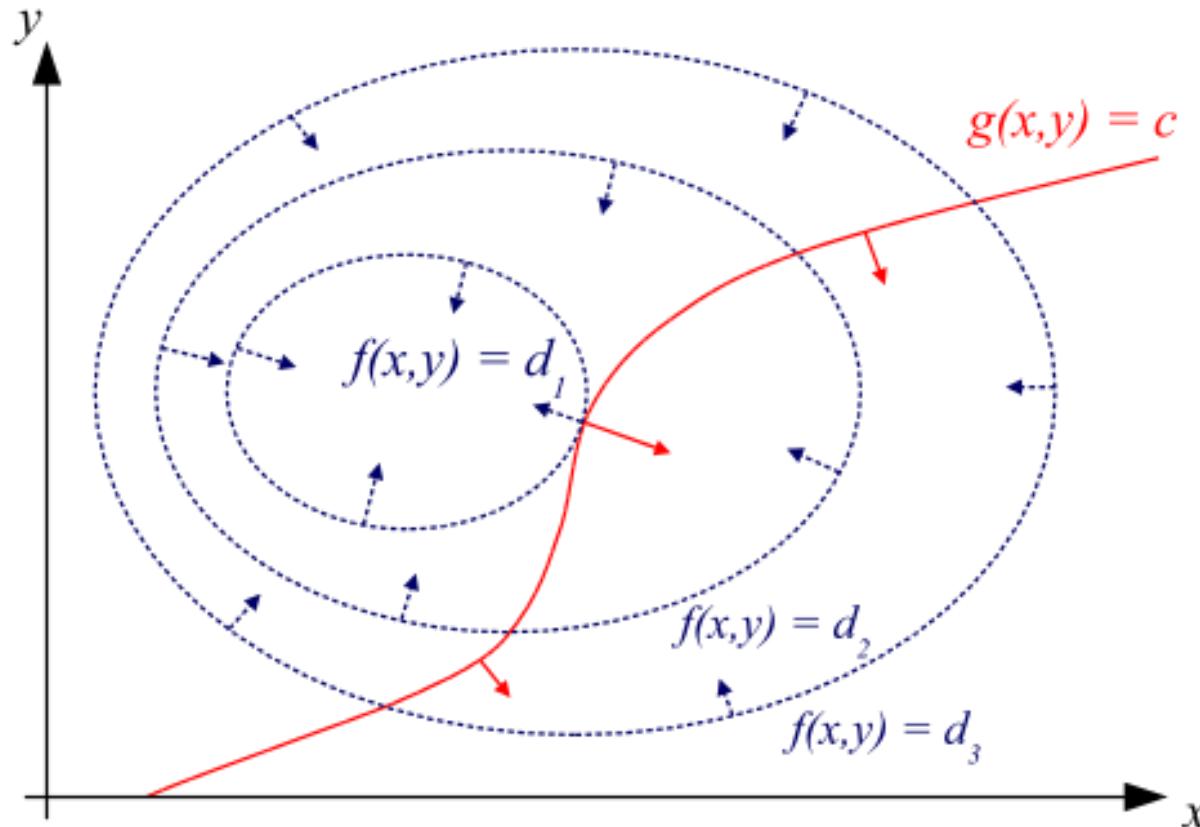


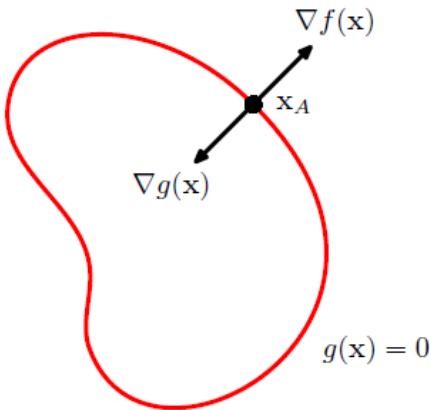
Figure 3.8: Geometry of minimum-norm attack for a two-class linear classifier with objective function $\|\mathbf{x} - \mathbf{x}_0\|^2$. The solution is a projection of the input \mathbf{x}_0 onto the separating hyperplane of the classifier.

Recap: lagrangian multipliers



$$\max f(x, y), \quad s.t. \quad g(x, y) = 0$$

Lagrange Multipliers



Consider the problem:

$$\begin{aligned} & \max_x f(\mathbf{x}) \\ \text{s.t. } & g(\mathbf{x}) = 0 \end{aligned}$$

This is because they are level curves

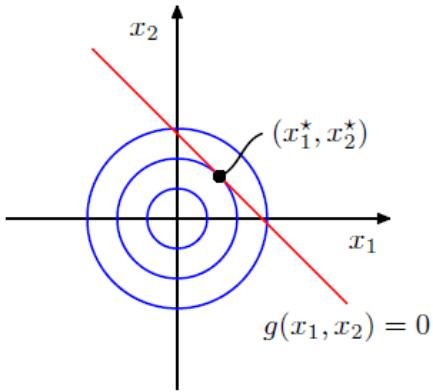
- Points on $g(\mathbf{x}) = 0$ must have $\nabla g(\mathbf{x})$ normal to surface
- A **stationary point** must have no change in f in the direction of the constraint surface, so $\nabla f(\mathbf{x})$ must also be normal to the surface.
 - So there must be some $\lambda \neq 0$ such that $\nabla f(\mathbf{x}) + \lambda \nabla g(\mathbf{x}) = 0$
- Define **Lagrangian**:

$$L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda g(\mathbf{x})$$

This is because they are the same vector up to a multiplicative constant

- Stationary points of $L(\mathbf{x}, \lambda)$ have $\nabla_{\mathbf{x}} L(\mathbf{x}, \lambda) = \nabla f(\mathbf{x}) + \lambda \nabla g(\mathbf{x}) = 0$ and $\nabla_{\lambda} L(\mathbf{x}, \lambda) = g(\mathbf{x}) = 0$

Lagrange Multipliers Example



- Consider the problem

$$\begin{aligned} \max_{\mathbf{x}} f(\mathbf{x}) &= 1 - x_1^2 - x_2^2 \\ \text{s.t.} \quad g(\mathbf{x}) &= x_1 + x_2 - 1 = 0 \end{aligned}$$

- Lagrangian:

$$L(\mathbf{x}, \lambda) = 1 - x_1^2 - x_2^2 + \lambda(x_1 + x_2 - 1)$$

- Stationary points require:

$$\frac{\partial L}{\partial x_1} = -2x_1 + \lambda = 0$$

$$\frac{\partial L}{\partial x_2} = -2x_2 + \lambda = 0$$

$$\frac{\partial L}{\partial \lambda} = x_1 + x_2 - 1 = 0$$

- So stationary point is $(x_1^*, x_2^*) = (\frac{1}{2}, \frac{1}{2})$, $\lambda = 1$

How to calculate the adversarial perturbation?

Proof. The Lagrange multiplier of the constrained optimization is given by

$$\mathcal{L}(\boldsymbol{x}, \lambda) = \frac{1}{2} \|\boldsymbol{x} - \boldsymbol{x}_0\|^2 + \lambda(\boldsymbol{w}^T \boldsymbol{x} + w_0).$$

The solution of the optimization is the saddle point $(\boldsymbol{x}^*, \lambda^*)$ such that $\nabla_{\boldsymbol{x}} \mathcal{L} = 0$ and $\nabla_{\lambda} \mathcal{L} = 0$.

Taking derivative with respect to \boldsymbol{x} and λ yields

$$\begin{aligned}\nabla_{\boldsymbol{x}} \mathcal{L} &= \boldsymbol{x} - \boldsymbol{x}_0 + \lambda \boldsymbol{w} = 0, \\ \nabla_{\lambda} \mathcal{L} &= \boldsymbol{w}^T \boldsymbol{x} + w_0 = 0.\end{aligned}$$

How to calculate the adversarial perturbation?

Multiplying the first equation by \mathbf{w}^T yields

$$0 = \mathbf{w}^T \mathbf{x} - \mathbf{w}^T \mathbf{x}_0 + \lambda \|\mathbf{w}\|^2 = -w_0 - \mathbf{w}^T \mathbf{x}_0 + \lambda \|\mathbf{w}\|^2. \quad (3.27)$$

Thus, the optimal λ is

$$\lambda^* = (\mathbf{w}^T \mathbf{x}_0 + w_0) / \|\mathbf{w}\|^2. \quad (3.28)$$

Correspondingly, the optimal \mathbf{x} is

$$\mathbf{x}^* = \mathbf{x}_0 - \lambda^* \mathbf{w} = \mathbf{x}_0 - \left(\frac{\mathbf{w}^T \mathbf{x}_0 + w_0}{\|\mathbf{w}\|} \right) \frac{\mathbf{w}}{\|\mathbf{w}\|_2}.$$

□

The result of this theorem provides many useful interpretation. First, the search direction is $-\mathbf{w}$ (or $-\mathbf{w}/\|\mathbf{w}\|_2$ for unit norm.) The negative sign is there because $\mathbf{w} = \mathbf{w}_i - \mathbf{w}_t$ is pointing from class \mathcal{C}_t to \mathcal{C}_i , which is opposite to the desired search direction. The search step $\frac{\mathbf{w}^T \mathbf{x}_0 + w_0}{\|\mathbf{w}\|_2}$ measures the distance from the input \mathbf{x}_0 to the target class \mathcal{C}_t . This is coherent to the projection perspective we presented earlier.

Deep fool: non-linear decision boundary

Definition 6 (DeepFool Attack by Moosavi-Dezfooli et al. 2016). *The DeepFool attack for a two-class classification generates the attack by solving the optimization*

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\|^2 \quad \text{subject to} \quad g(\mathbf{x}) = 0, \quad (3.29)$$

where $g(\mathbf{x}) = 0$ is the nonlinear decision boundary separating the two classes.

Corollary 1 (DeepFool Algorithm for Two-Class Problem). *An iterative procedure to obtain the DeepFool attack solution is*

$$\begin{aligned} \mathbf{x}^{(k+1)} &= \underset{\mathbf{x}}{\text{argmin}} \quad \|\mathbf{x} - \mathbf{x}^{(k)}\|^2 \quad \text{subject to} \quad g(\mathbf{x}^{(k)}) + \nabla_{\mathbf{x}}g(\mathbf{x}^{(k)})^T(\mathbf{x} - \mathbf{x}^{(k)}) = 0 \\ &= \mathbf{x}^{(k)} - \left(\frac{g(\mathbf{x}^{(k)})}{\|\nabla_{\mathbf{x}}g(\mathbf{x}^{(k)})\|^2} \right) \nabla_{\mathbf{x}}g(\mathbf{x}^{(k)}). \end{aligned} \quad (3.31)$$

Due to the nonlinearity of $g(\mathbf{x})$, it is generally very difficult to derive a closed-form solution. The numerical procedure to compute the solution can be derived by iteratively updating \mathbf{x} , where each iteration minimizes \mathbf{x} over the first order approximation of $g(\mathbf{x})$:

$$g(\mathbf{x}) \approx g(\mathbf{x}^{(k)}) + \nabla_{\mathbf{x}}g(\mathbf{x}^{(k)})^T(\mathbf{x} - \mathbf{x}^{(k)}),$$

where $\mathbf{x}^{(k)}$ is the k -th iterate of the solution. In other words, we are defining a procedure

$$\mathbf{x}^{(k+1)} = \operatorname{argmin}_{\mathbf{x}} \| \mathbf{x} - \mathbf{x}^{(k)} \|^2 \text{ subject to } g(\mathbf{x}^{(k)}) + \nabla_{\mathbf{x}}g(\mathbf{x}^{(k)})^T(\mathbf{x} - \mathbf{x}^{(k)}) = 0. \quad (3.30)$$

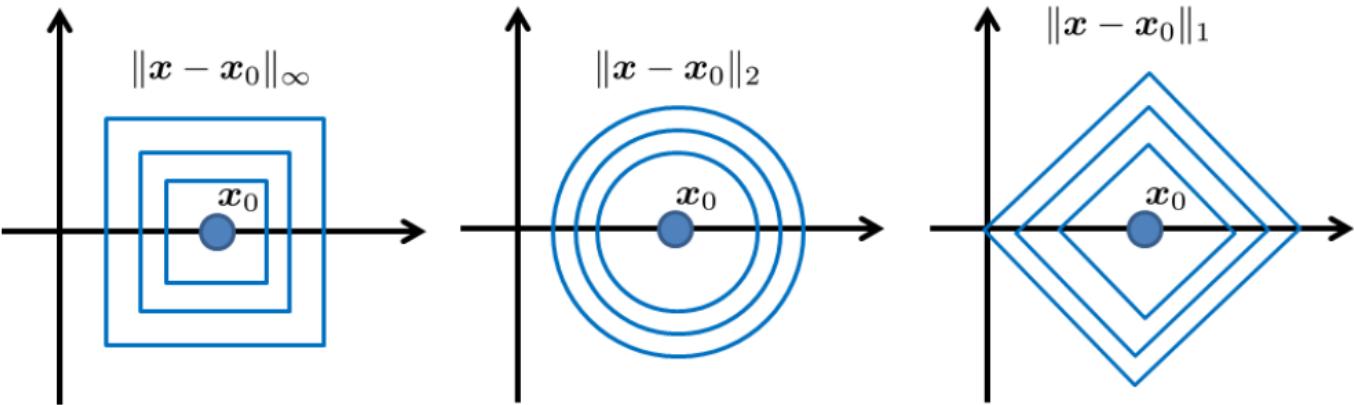
By identifying $\mathbf{w}^{(k)} = \nabla_{\mathbf{x}}g(\mathbf{x}^{(k)})$ and $w_0^{(k)} = g(\mathbf{x}^{(k)}) - \nabla_{\mathbf{x}}g(\mathbf{x}^{(k)})^T\mathbf{x}^{(k)}$, the linearized problem Equation (3.30) becomes

$$\mathbf{x}^{(k+1)} = \operatorname{argmin}_{\mathbf{x}} \| \mathbf{x} - \mathbf{x}^{(k)} \|^2 \text{ subject to } (\mathbf{w}^{(k)})^T \mathbf{x} + w_0^{(k)} = 0$$

The solution can thus be found by using Equation (3.26), yielding

$$\begin{aligned} \mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)} - \left(\frac{(\mathbf{w}^{(k)})^T \mathbf{x}^{(k)} + w_0^{(k)}}{\|\mathbf{w}^{(k)}\|^2} \right) \mathbf{w}^{(k)} \\ &= \mathbf{x}^{(k)} - \left(\frac{g(\mathbf{x}^{(k)})}{\|\nabla_{\mathbf{x}}g(\mathbf{x}^{(k)})\|^2} \right) \nabla_{\mathbf{x}}g(\mathbf{x}^{(k)}). \end{aligned}$$

Exercise!

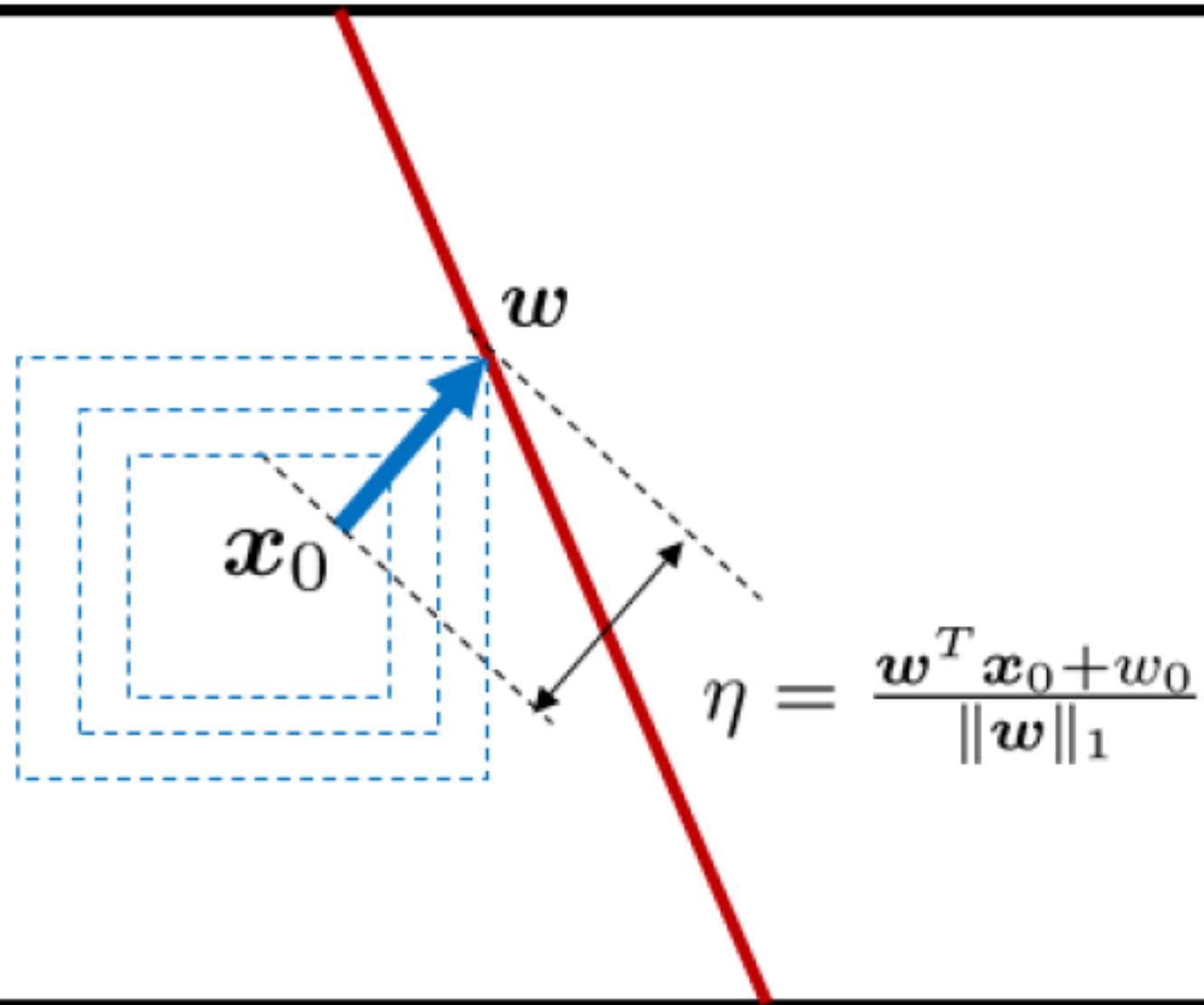


Theorem 4 (Minimum ℓ_∞ Norm Attack for Two-Class Linear Classifier). *The minimum ℓ_∞ norm attack for a two-class linear classifier, i.e.,*

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_0\|_\infty \quad \text{subject to} \quad \mathbf{w}^T \mathbf{x} + w_0 = 0 \quad (3.33)$$

is given by

$$\mathbf{x} = \mathbf{x}_0 - \left(\frac{\mathbf{w}^T \mathbf{x}_0 + w_0}{\|\mathbf{w}\|_1} \right) \cdot \text{sign}(\mathbf{w}). \quad (3.34)$$

\mathcal{C}_i \mathcal{C}_t 

How is minimum ℓ_∞ norm attack related to the maximum allowable ℓ_∞ norm attack? Recall that the maximum allowable attack is given by

$$\underset{\mathbf{x}}{\text{minimize}} \quad \mathbf{w}^T \mathbf{x} + w_0 \quad \text{subject to} \quad \|\mathbf{x} - \mathbf{x}_0\|_\infty \leq \eta. \quad (3.35)$$

By defining $\mathbf{x} = \mathbf{x}_0 + \mathbf{r}$, we have $\mathbf{w}^T \mathbf{x} + w_0 = \mathbf{w}^T \mathbf{x}_0 + \mathbf{w}^T \mathbf{r} + w_0$. By defining $b_0 = -(\mathbf{w}^T \mathbf{x}_0 + w_0)$, the optimization can be rewritten as

$$\underset{\mathbf{r}}{\text{minimize}} \quad \mathbf{w}^T \mathbf{r} - b_0 \quad \text{subject to} \quad \|\mathbf{r}\|_\infty \leq \eta. \quad (3.36)$$

Using the Holder's inequality (the negative side), we can show that

$$\mathbf{w}^T \mathbf{r} \geq -\|\mathbf{r}\|_\infty \|\mathbf{w}\|_1 \geq -\eta \|\mathbf{w}\|_1.$$

The lower bound of $\mathbf{w}^T \mathbf{r}$ is attained when $\mathbf{r} = -\eta \cdot \text{sign}(\mathbf{w})$, because

$$\mathbf{w}^T \mathbf{r} = -\eta \mathbf{w}^T \text{sign}(\mathbf{w}) = -\eta \sum_{i=1}^d w_i \text{sign}(w_i) = -\eta \sum_{i=1}^d |w_i| = -\eta \|\mathbf{w}\|_1.$$

Therefore, the attacked data point \mathbf{x} is $\mathbf{x} = \mathbf{x}_0 - \eta \cdot \text{sign}(\mathbf{w})$. This leads to the following theorem.

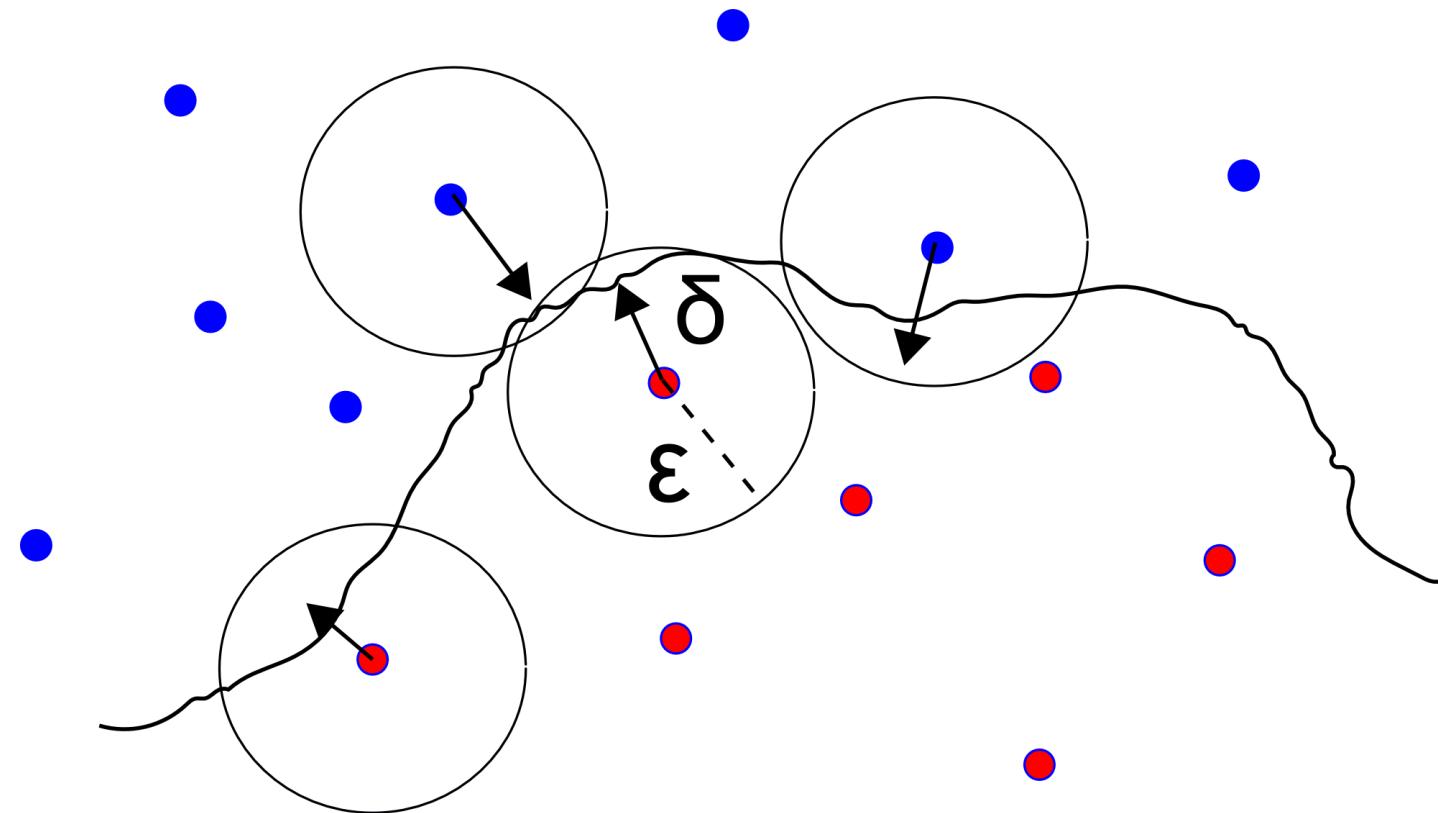
Theorem 5 (Maximum Allowable ℓ_∞ Norm Attack of Two-Class Linear Classifier). *The maximum allowable ℓ_∞ norm attack for a two-class linear classifier, i.e.,*

$$\underset{\mathbf{x}}{\text{minimize}} \quad \mathbf{w}^T \mathbf{x} + w_0 \quad \text{subject to} \quad \|\mathbf{x} - \mathbf{x}_0\|_\infty \leq \eta. \quad (3.37)$$

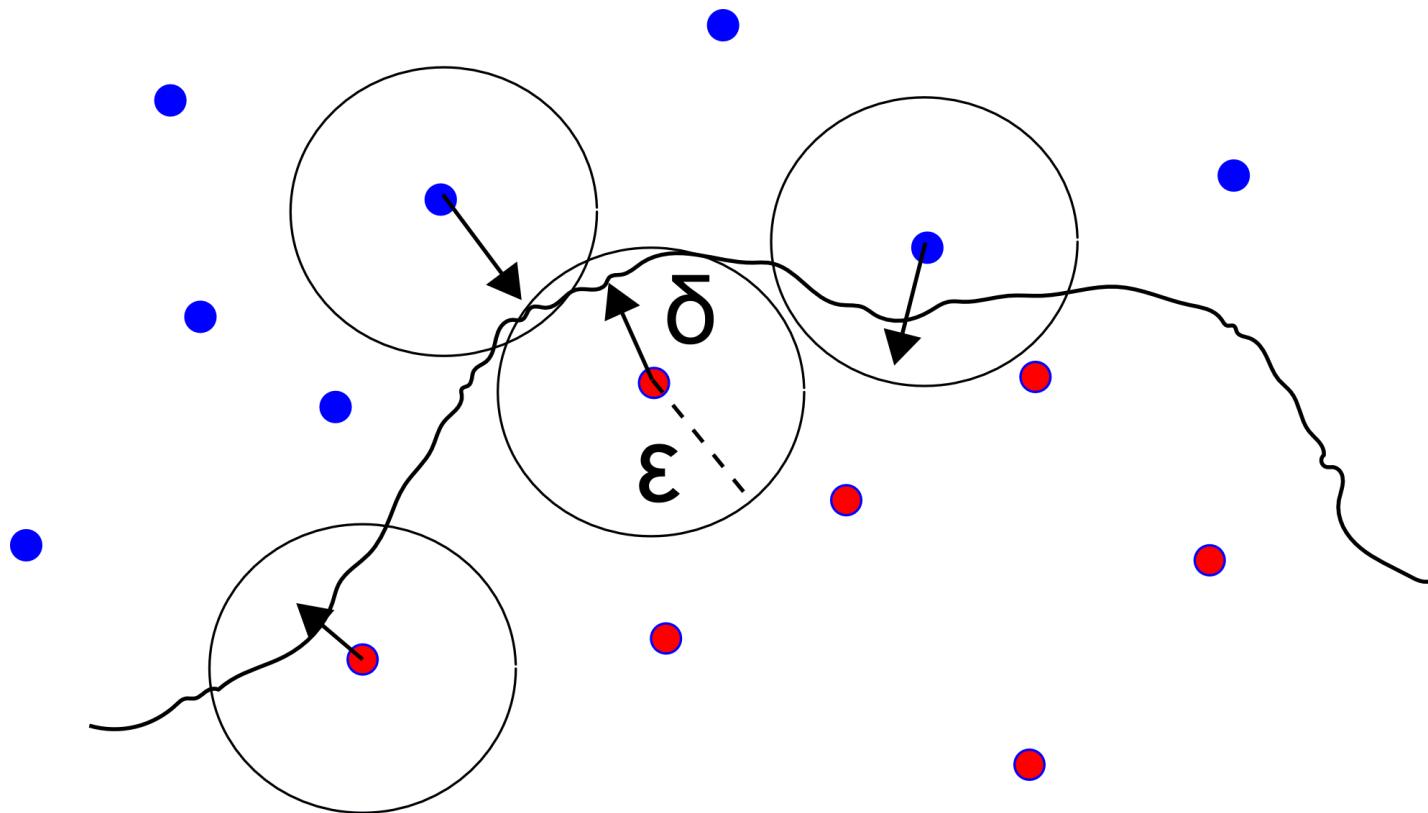
is given by

$$\mathbf{x} = \mathbf{x}_0 - \eta \cdot \text{sign}(\mathbf{w}). \quad (3.38)$$

How can we "fix" the adversarial weakness?
Adversarial training

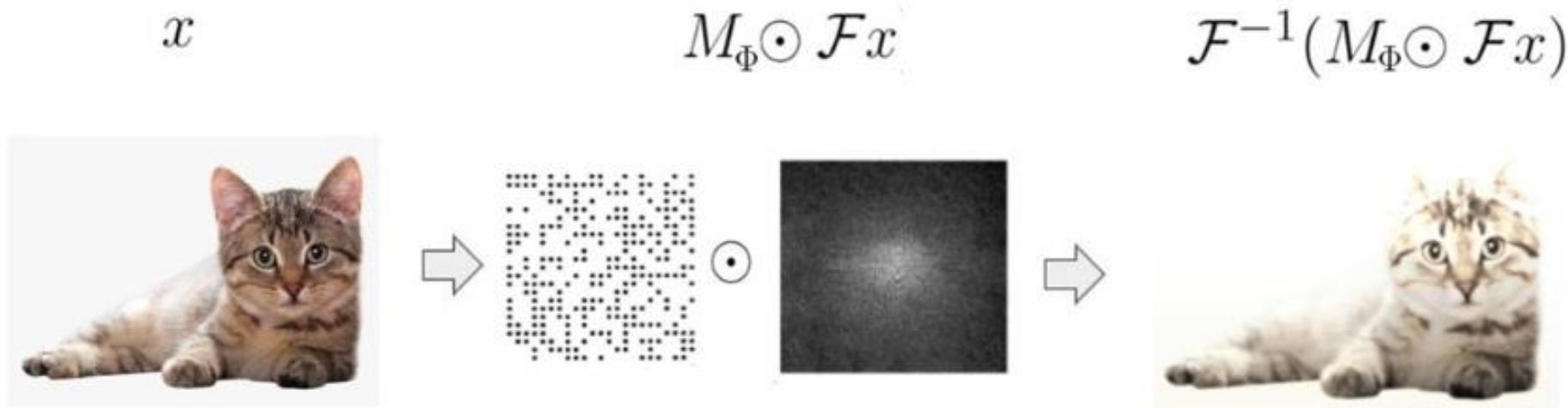


We need to modify the boundaries such that the adv examples are correctly classified

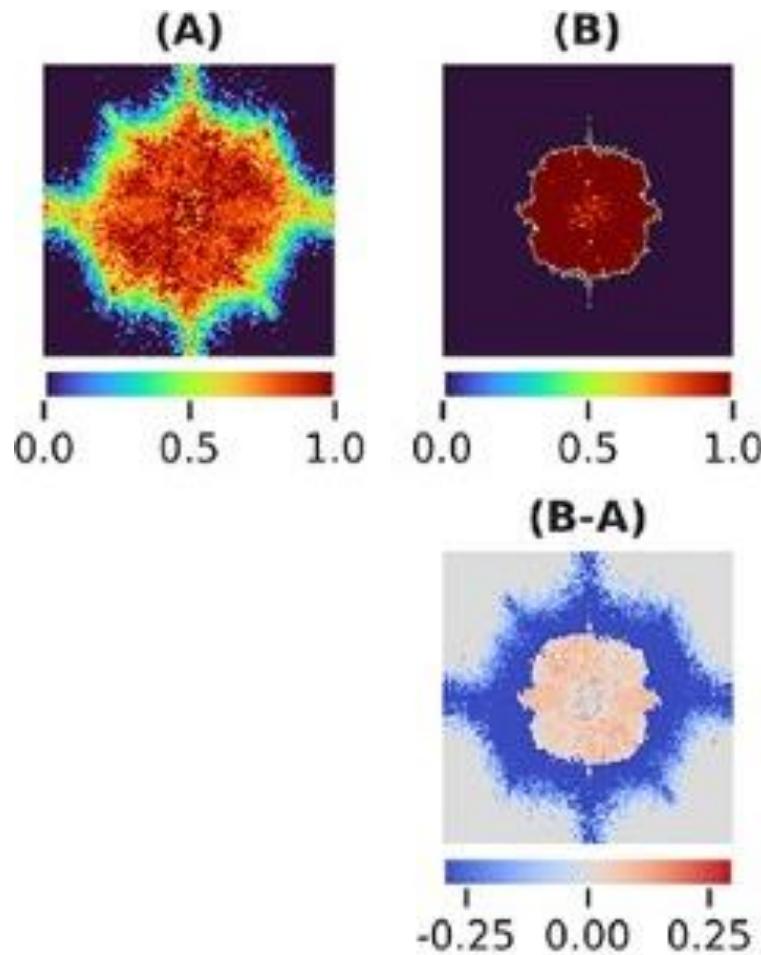


Origin of adversarial weakness and
adversarial robustness in Fourier space

Adversarial attacks in Fourier space



$$M_{\Phi}(\lambda, p) = \operatorname{argmin}_{M_{\Phi}} \sum_{x \in \mathcal{X}_V} e^{[\mathcal{L}(\Phi(\bar{x}), y) - \mathcal{L}(\Phi(x), y)]^2} + \lambda \|M_{\Phi}\|_p, \quad \lambda \in \mathbb{R}_+,$$



Before and after adversarial training: what is your conclusion in terms of frequency bias of the Adversarial attacks?

Class 6

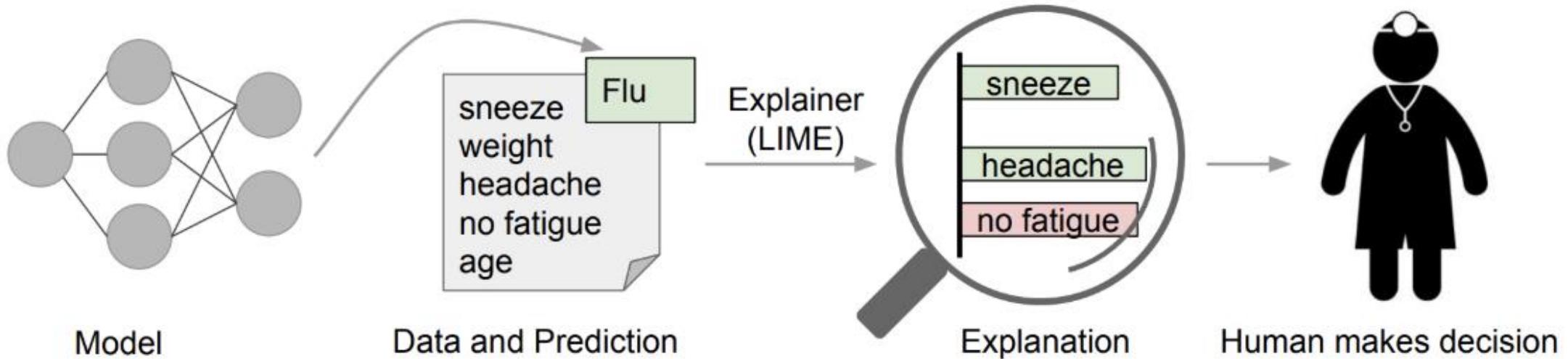
Interpretability

- Problems of annotation artifacts, biases and adversarial attacks
- What does it mean for a model to be interpretable?**
- No mathematical definition of interpretability.
 - *"Interpretability is the degree to which a human can understand the cause of a decision."* (*)
 - *"Interpretability is the degree to which a human can consistently predict the model's result"* (**)
 - Interpret on the level of **predictions**
 - Why did the model make certain predictions?
 - Interpret on the level of **components**
 - How does the model make predictions?

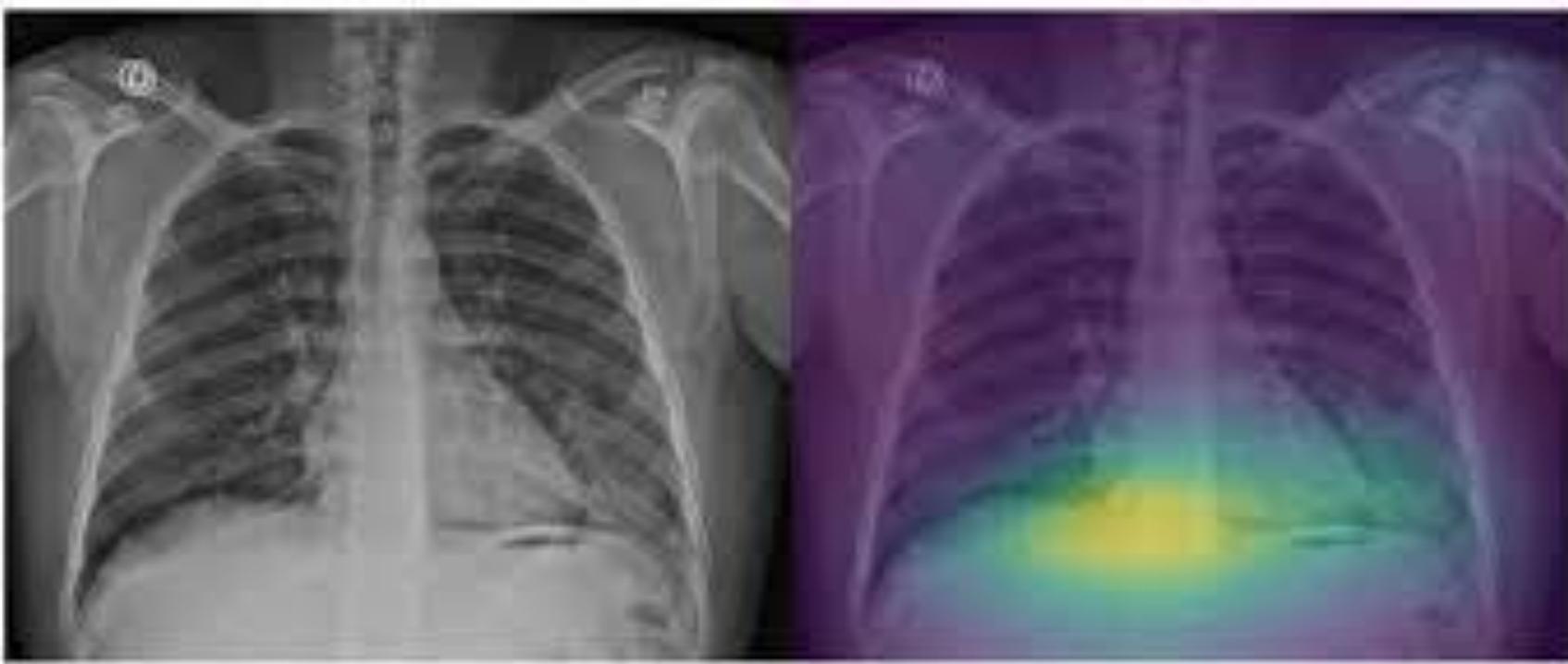
^{*}Miller, Tim. "Explanation in artificial intelligence: Insights from the social sciences." arXiv Preprint arXiv:1706.07269. (2017)

^{**} Kim, Been, Rajiv Khanna, and Oluwasanmi O. Koyejo. "Examples are not enough, learn to criticize! Criticism for interpretability." Advances in Neural Information Processing Systems (2016)

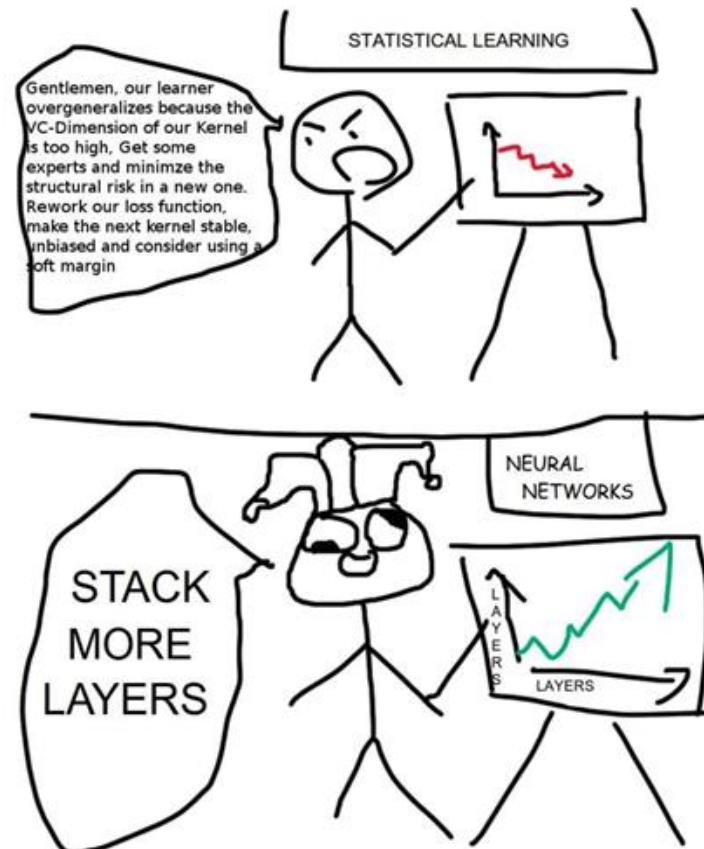
Why interpretability is important?



Why interpretability is important?



Why interpretability is important?

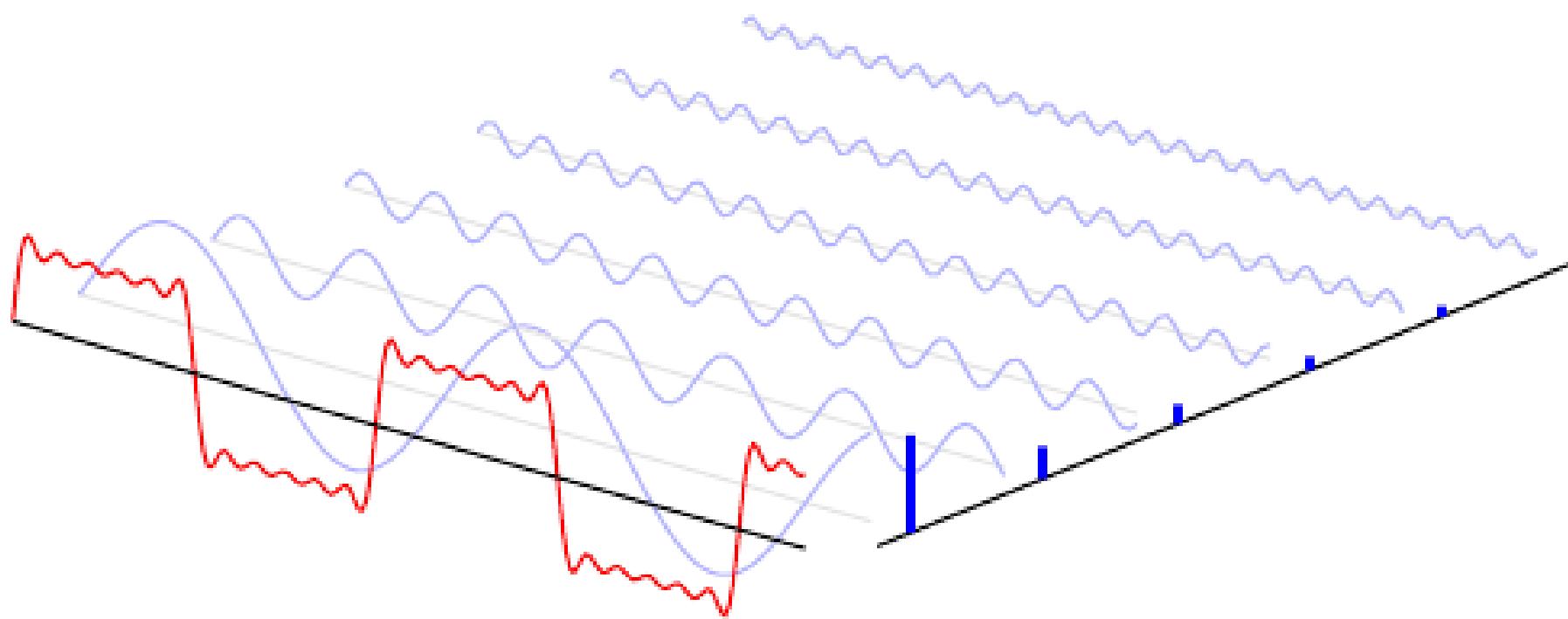


[Reddit; source unknown]

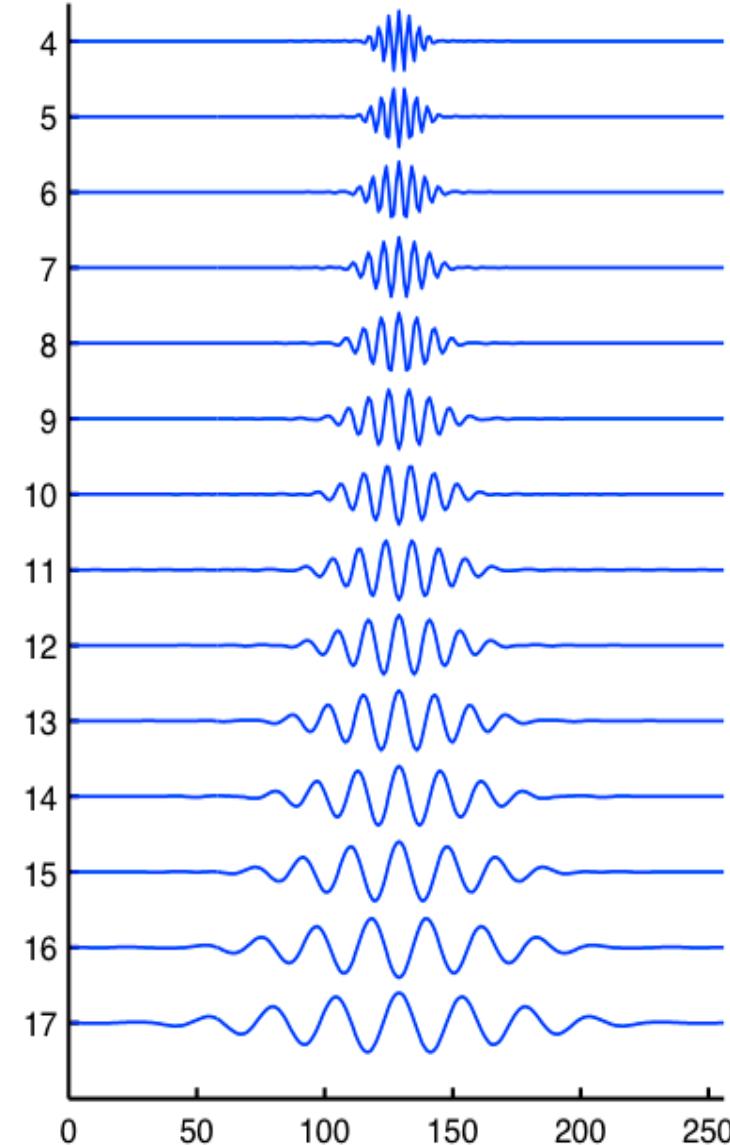


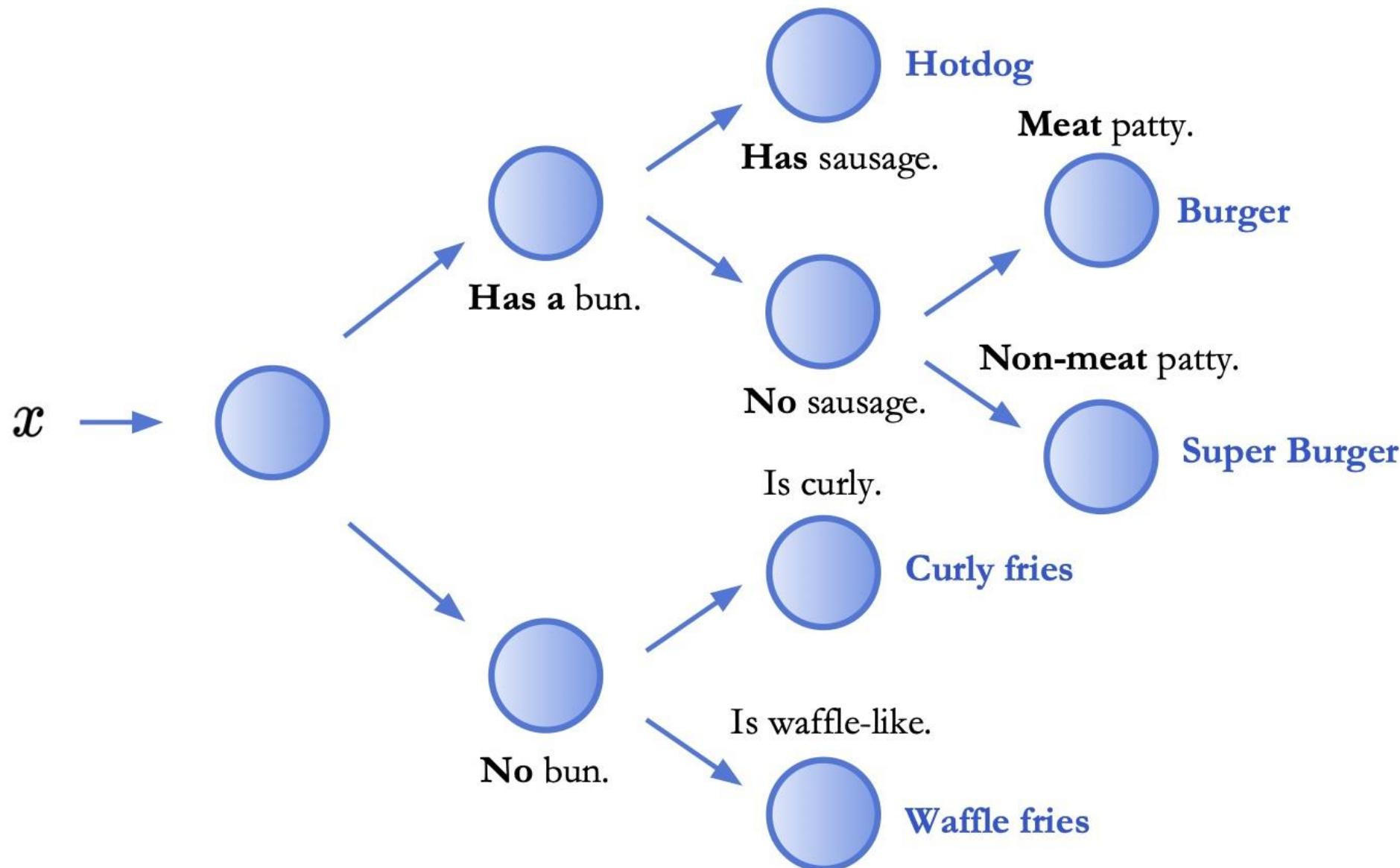
[xkcd.com]

Designed features: Fourier transform



Designed
features:
Wavelet
transform





Can we extract a similar interpretability looking at learned convolutional filters?

Models and Hypothesis space

- Linear separable

$$f(x) = w^T x$$

We learn the separator on data

- Non-linear separable

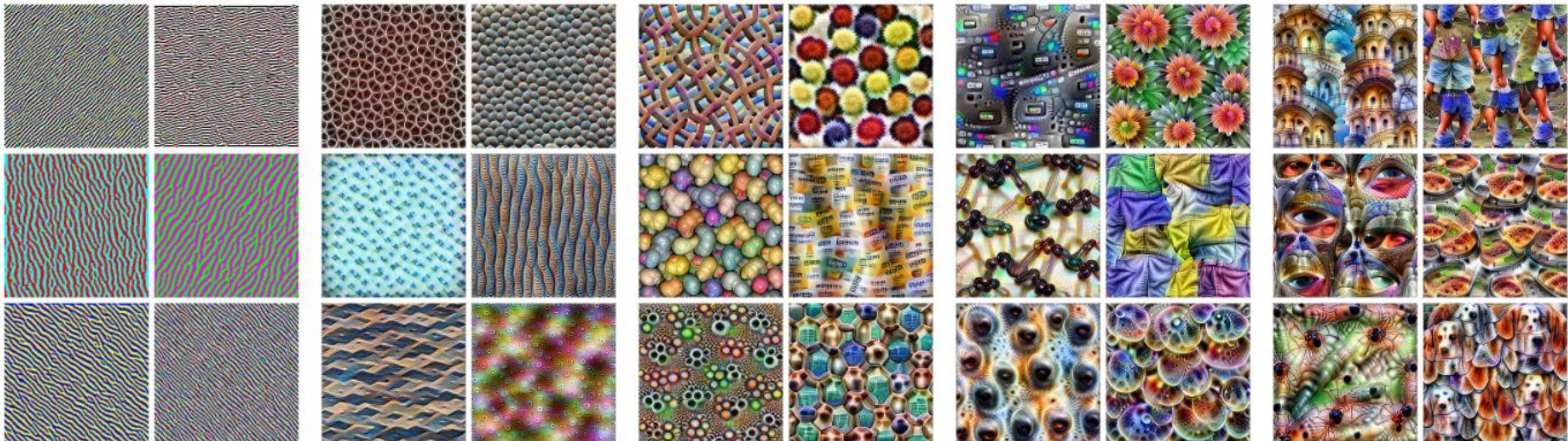
$$f(x) = w^T \phi(x)$$

We learn the separator on a
priori hand-crafted features

What about Anns?

What images/patterns maximally activate a neuron?

How neural networks build up their understanding of images



Edges (layer conv2d0)

Textures (layer mixed3a)

Patterns (layer mixed4a)

Parts (layers mixed4b & mixed4c)

Objects (layers mixed4d & mixed4e)

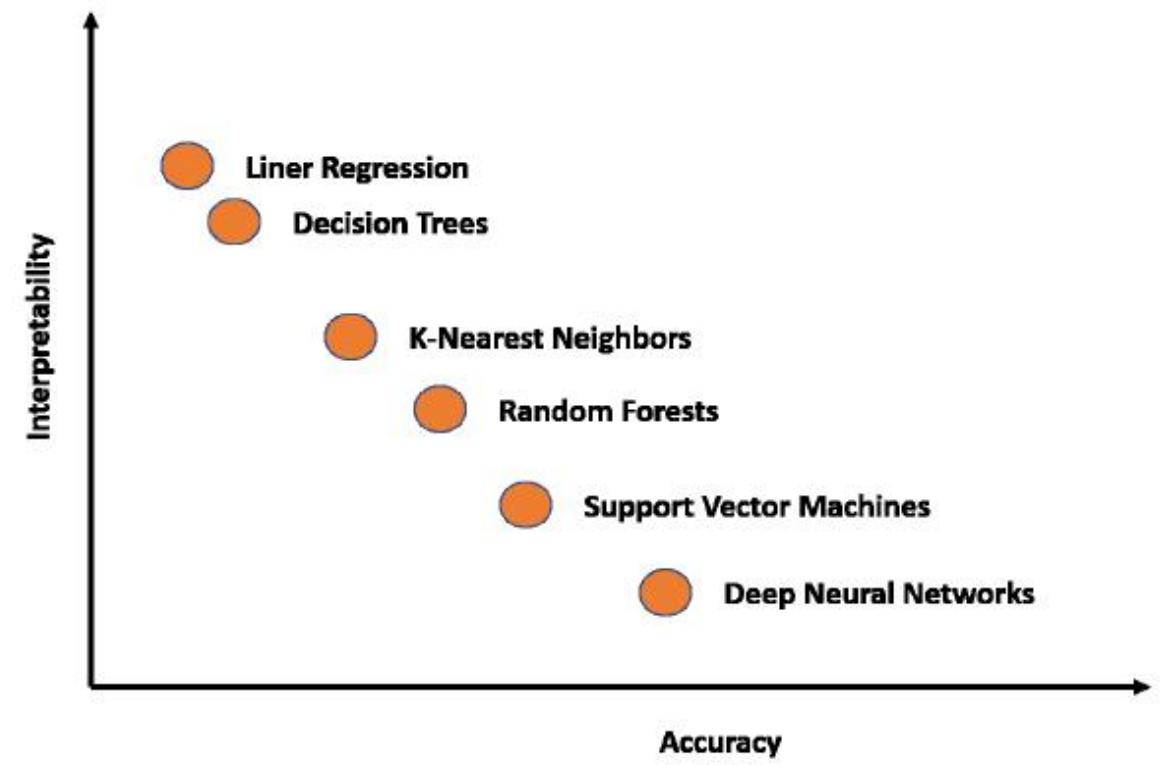
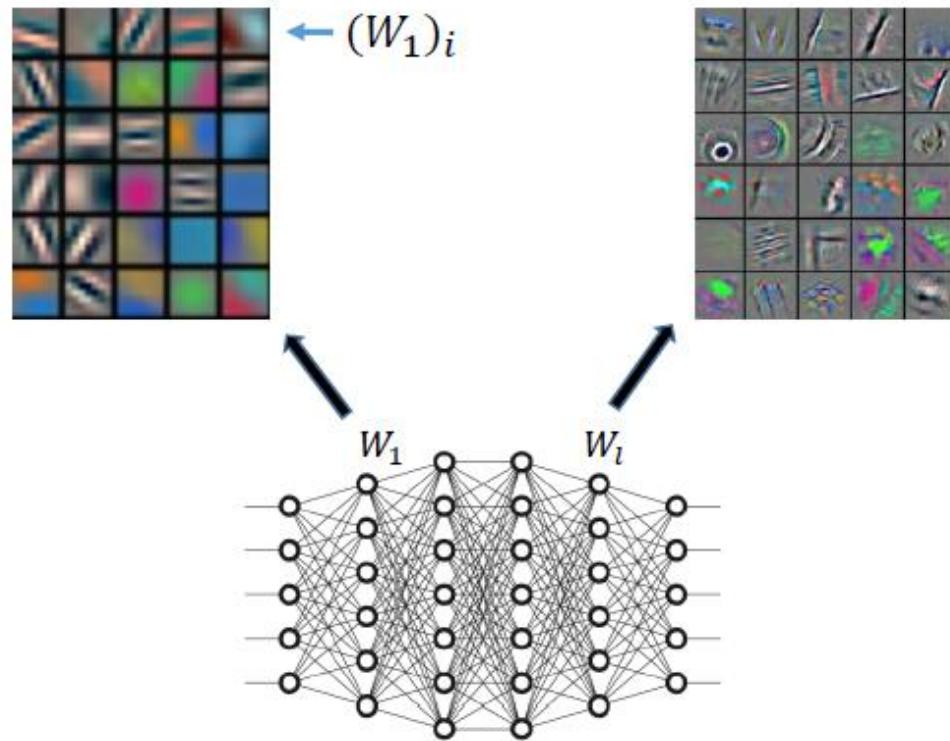
How are those images calculated?

$$\arg \max_x f_i(x)$$

Single neuron response

$$x_{t+1} = x_t + \eta \nabla_x f_i(x)$$

Interpretable features: images

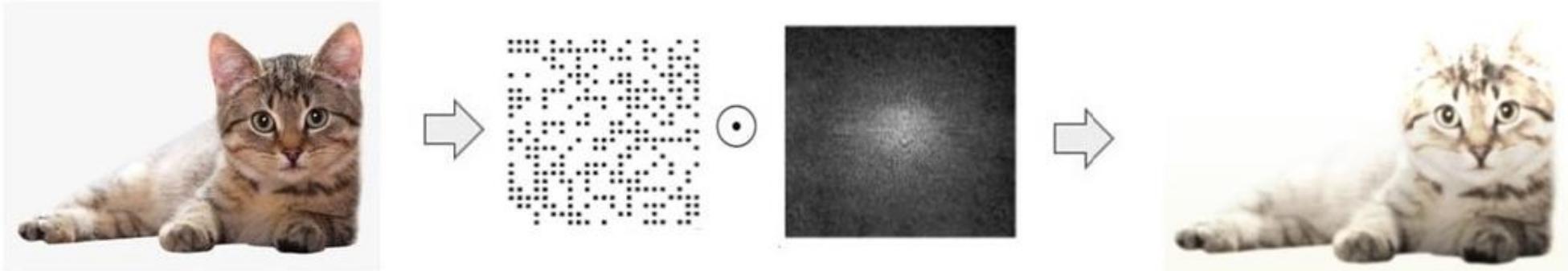


How does a network make a decision? Fourier
important features.

x

$M_\Phi \odot \mathcal{F}x$

$\mathcal{F}^{-1}(M_\Phi \odot \mathcal{F}x)$

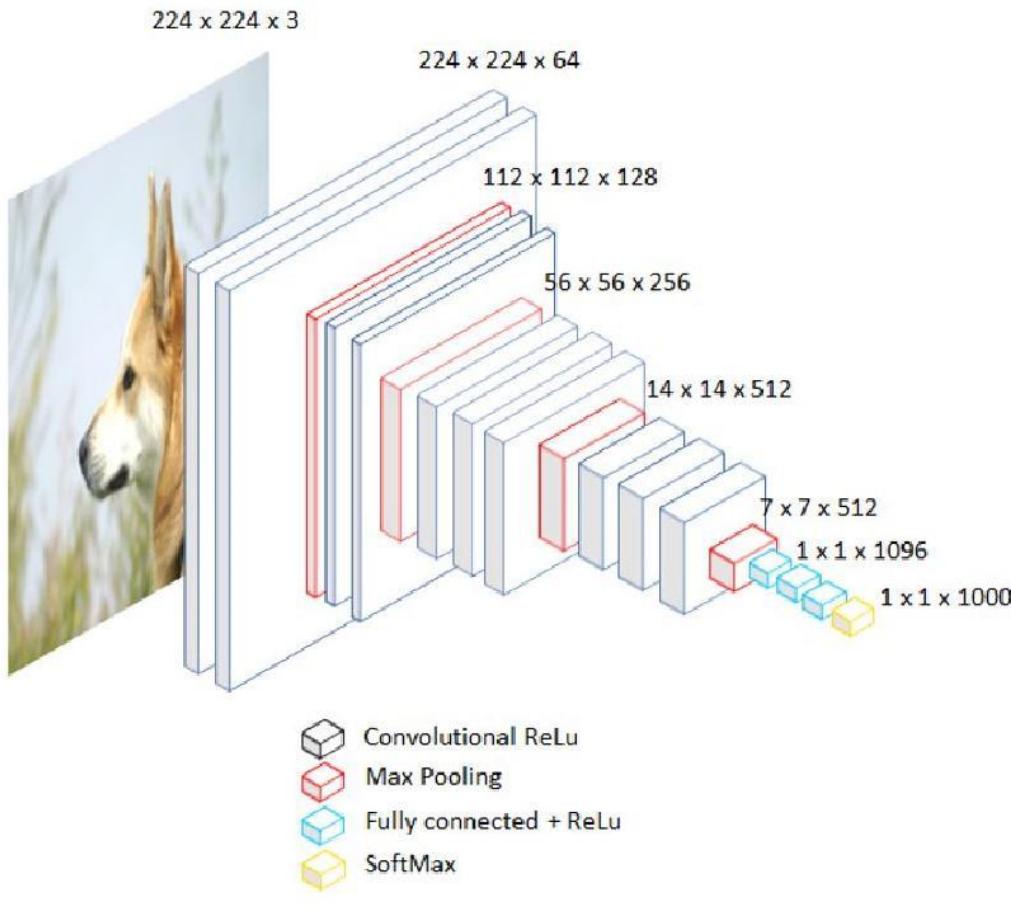


$$M_\Phi(\lambda, p) = \operatorname{argmin}_{M_\Phi} \sum_{x \in \mathcal{X}_V} e^{[\mathcal{L}(\Phi(\bar{x}), y) - \mathcal{L}(\Phi(x), y)]^2} + \lambda \|M_\Phi\|_p, \quad \lambda \in \mathbb{R}_+$$

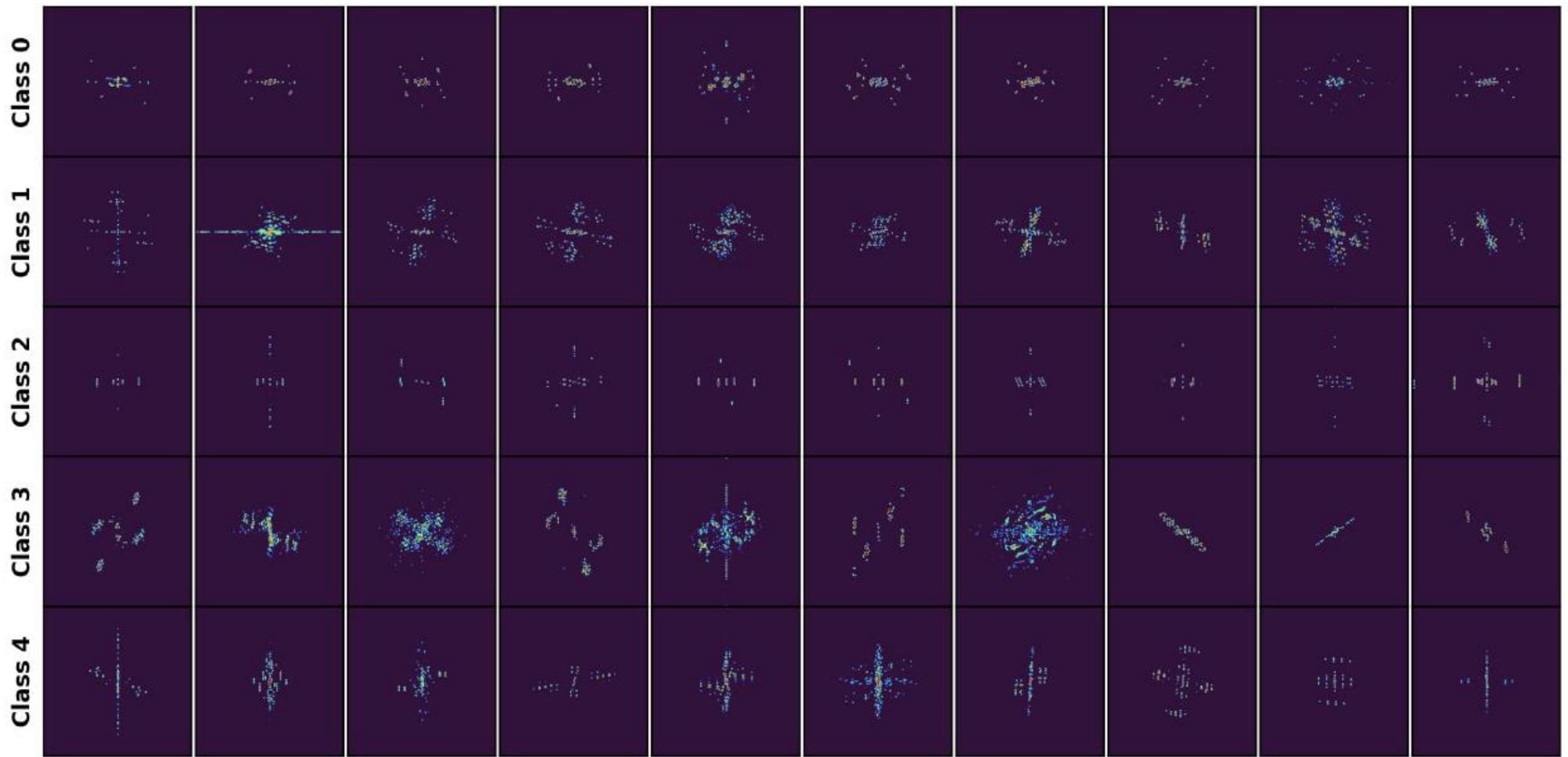
Data and Model: Imagenette and VGG11



```
lbl_dict = dict(  
    n01440764='trench',  
    n02102040='English springer',  
    n02979186='cassette player',  
    n03000684='chain saw',  
    n0028079='church',  
    n03394916='French horn',  
    n03417042='garbage truck',  
    n03425413='gas pump',  
    n03445777='golf ball',  
    n03888257='parachute'  
)
```



6, 644 image/label pairs: training set. 1, 934 pairs: validation set. For simplicity, we used grayscale versions.



(A)

(B)

(C)

What is your conclusion?

Or texture?

IMAGENET-TRAINED CNNS ARE BIASED TOWARDS
TEXTURE; INCREASING SHAPE BIAS IMPROVES
ACCURACY AND ROBUSTNESS

Robert Geirhos
University of Tübingen & IMPRS-IS
robert.geirhos@bethgelab.org

Patricia Rubisch
University of Tübingen & U. of Edinburgh
p.rubisch@sms.ed.ac.uk

Claudio Michaelis
University of Tübingen & IMPRS-IS
claudio.michaelis@bethgelab.org

Matthias Bethge*
University of Tübingen
matthias.bethge@bethgelab.org

Felix A. Wichmann*
University of Tübingen
felix.wichmann@uni-tuebingen.de

Wieland Brendel*
University of Tübingen
wieland.brendel@bethgelab.org



(a) Texture image
81.4% **Indian elephant**
10.3% indri
8.2% black swan

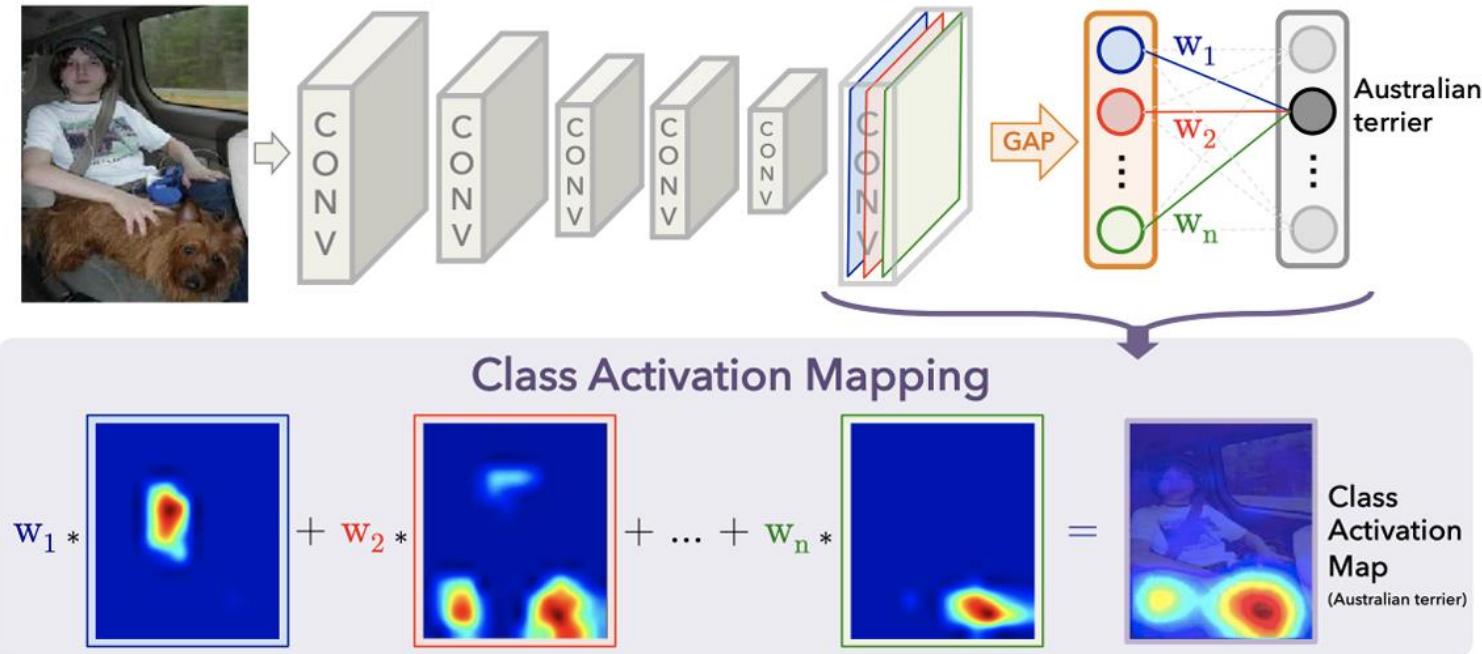


(b) Content image
71.1% **tabby cat**
17.3% grey fox
3.3% Siamese cat



(c) Texture-shape cue conflict
63.9% **Indian elephant**
26.4% indri
9.6% black swan

Where is a ANN looking at? Class activation mapping



Class activation mapping

- $f_k(x, y)$: Activation of unit k in spatial location (x, y)
- $F^k = \sum_{x,y} f_k(x, y)$: Result of global average pooling
- $S_c = \sum_k w_k^c F_k$: input to Softmax layer for class c
- $M_c(x, y) = \sum_k w_k^c f_k(x, y)$: CAM for class c

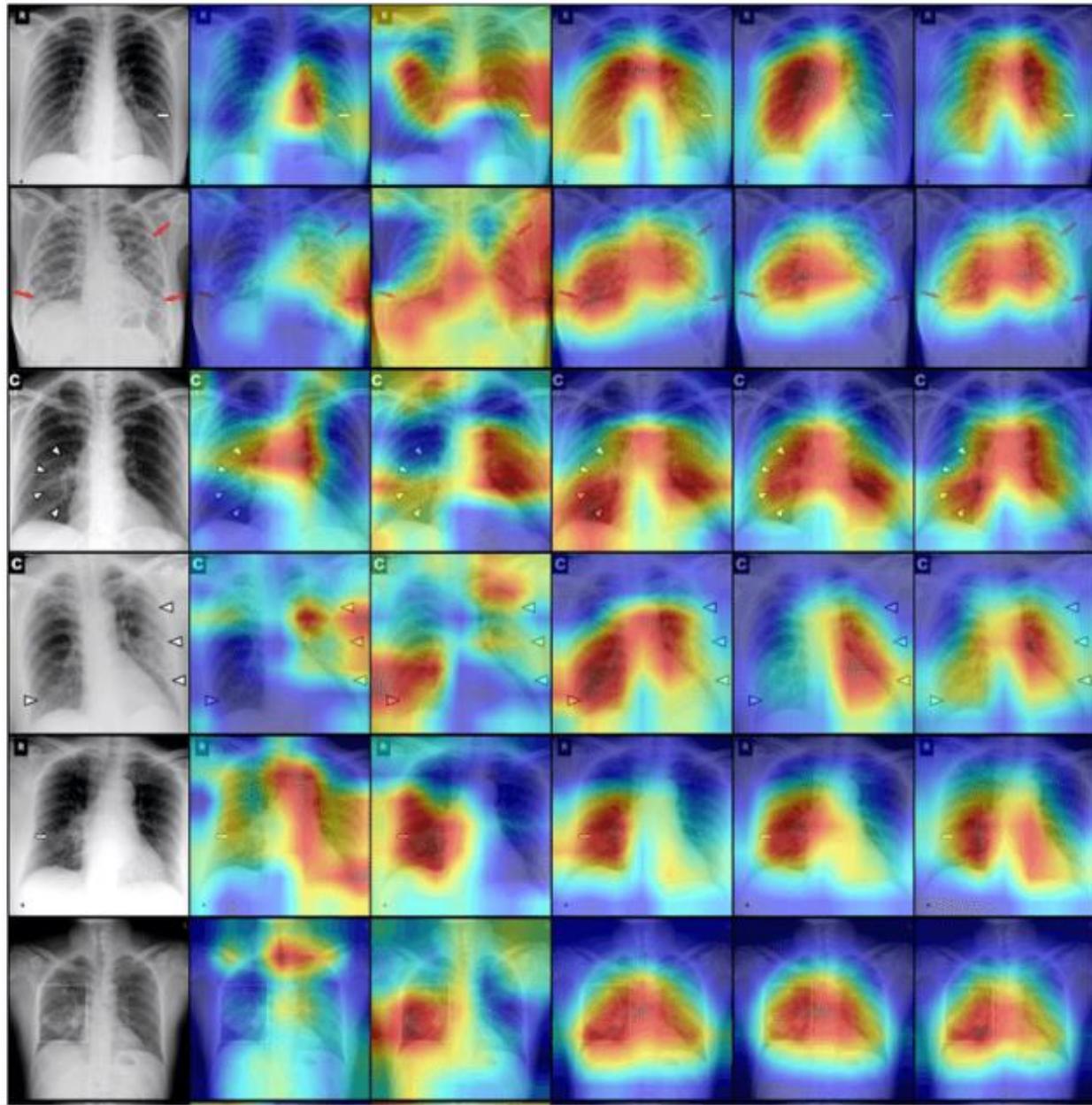
ResNet50

DenseNet121

DANN

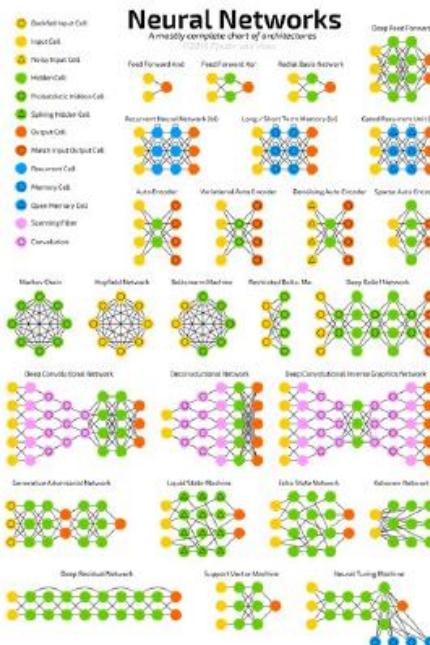
PADA

SODA



What is your conclusion?

Interpretable operations



Navier-Stokes Equations

Continuity Equation

$$\nabla \cdot \vec{V} = 0$$

Momentum Equations

$$\rho \frac{D\vec{V}}{Dt} = -\nabla p + \rho \vec{g} + \mu \nabla^2 \vec{V}$$

Total derivative
Pressure gradient
Body force term
Diffusion term
 $\rho \left[\frac{\partial V}{\partial t} + (\vec{V} \cdot \nabla) V \right]$
Change of velocity with time
Connective term
Fluid flows in the direction of largest change in pressure.
External forces, that act on the fluid (gravitational force or electromagnetic).
For a Newtonian fluid, viscosity operates as a diffusion of momentum.

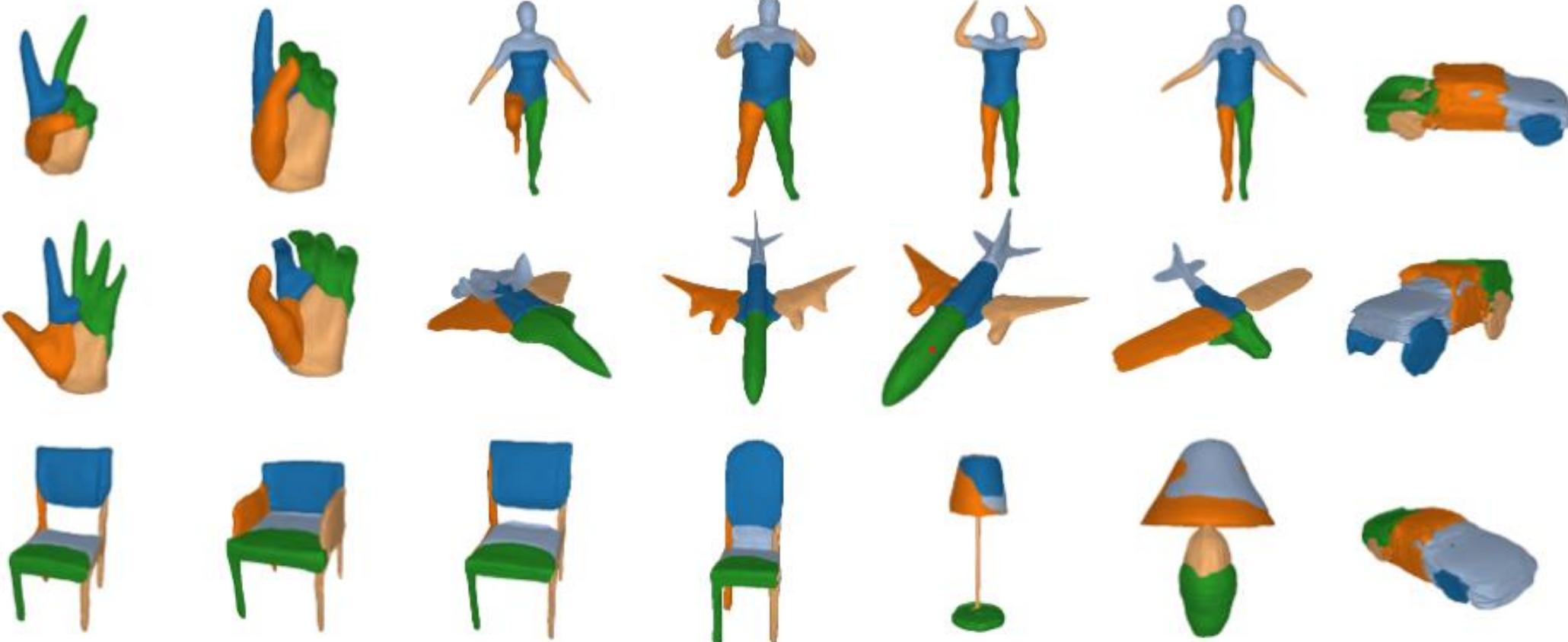
Suppose we learnt a φ_W from input/output to a physical system: $y_i = \varphi_W(x_i)$

Approximate φ_W with a dictionary of known operators:

$$\arg \min_{\alpha} \sum_i \|\varphi_W(\mathbf{x}_i) - (\{\partial_t, \partial_x, \dots\} \alpha)(\mathbf{x}_i)\|_2 + \|\alpha\|_0$$

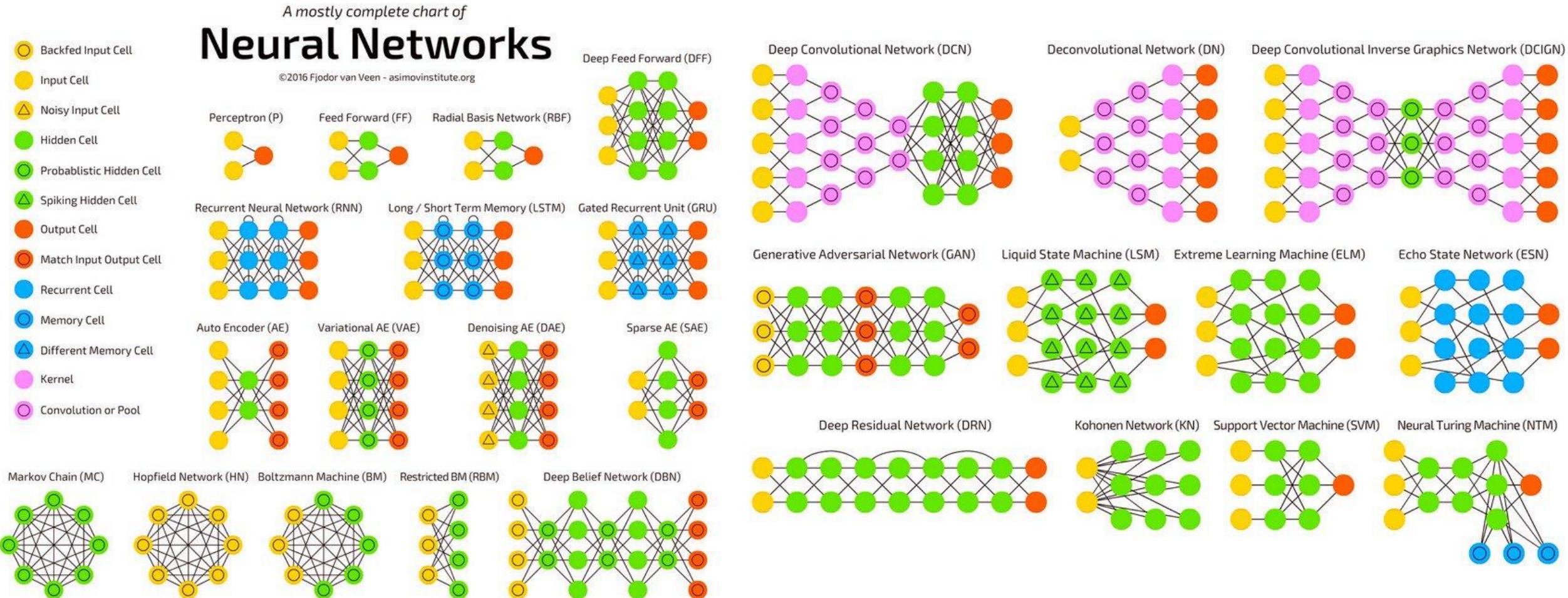
How do we understand a visual scene?

Interpretable parts



Other architectures and beyond images

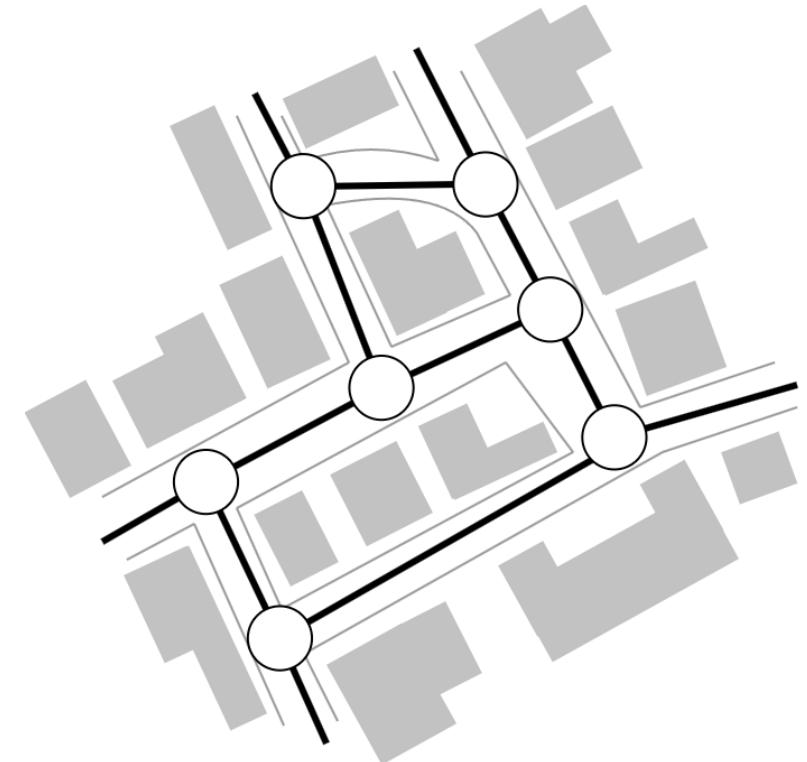
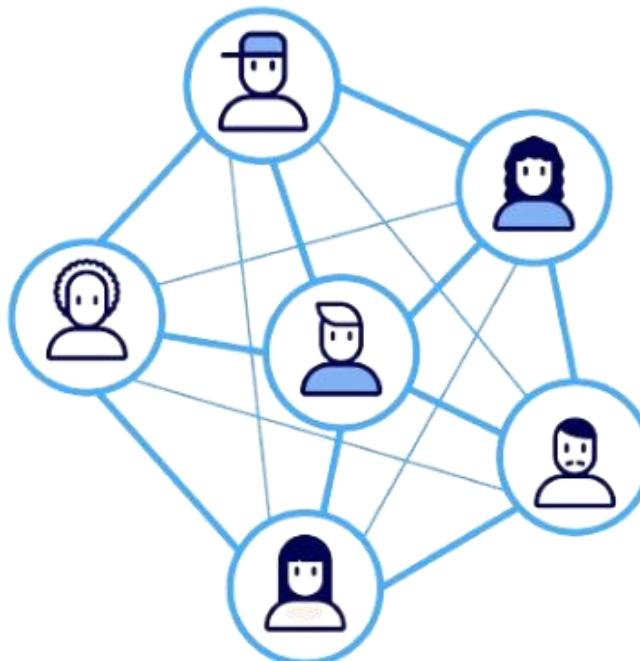
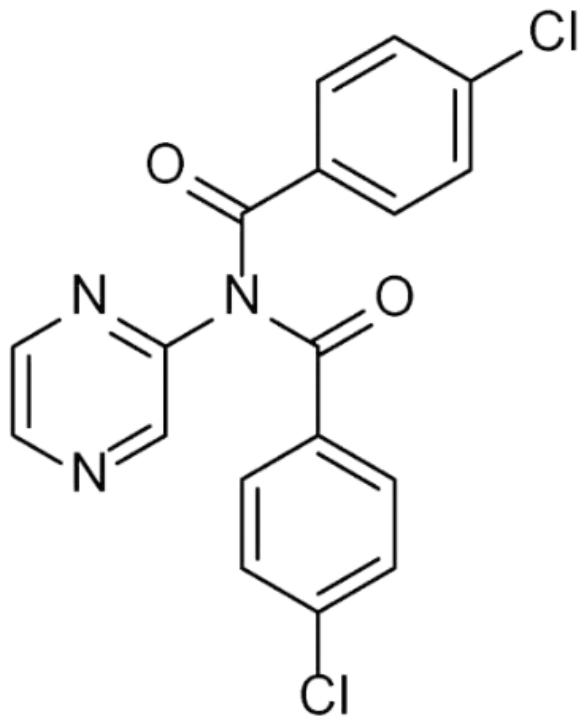
A zoo of networks



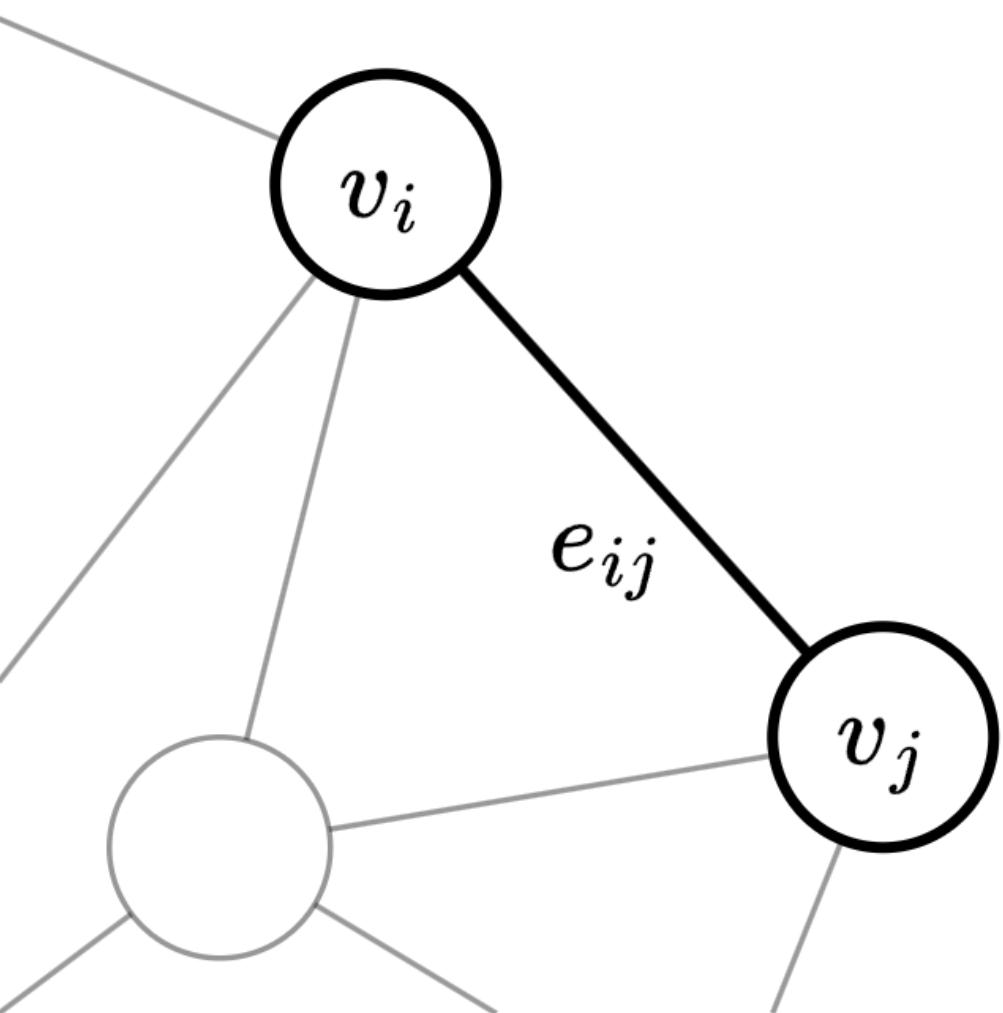
<https://modelzoo.co/>

Short intro to graph neural networks

Graph structures in the real world

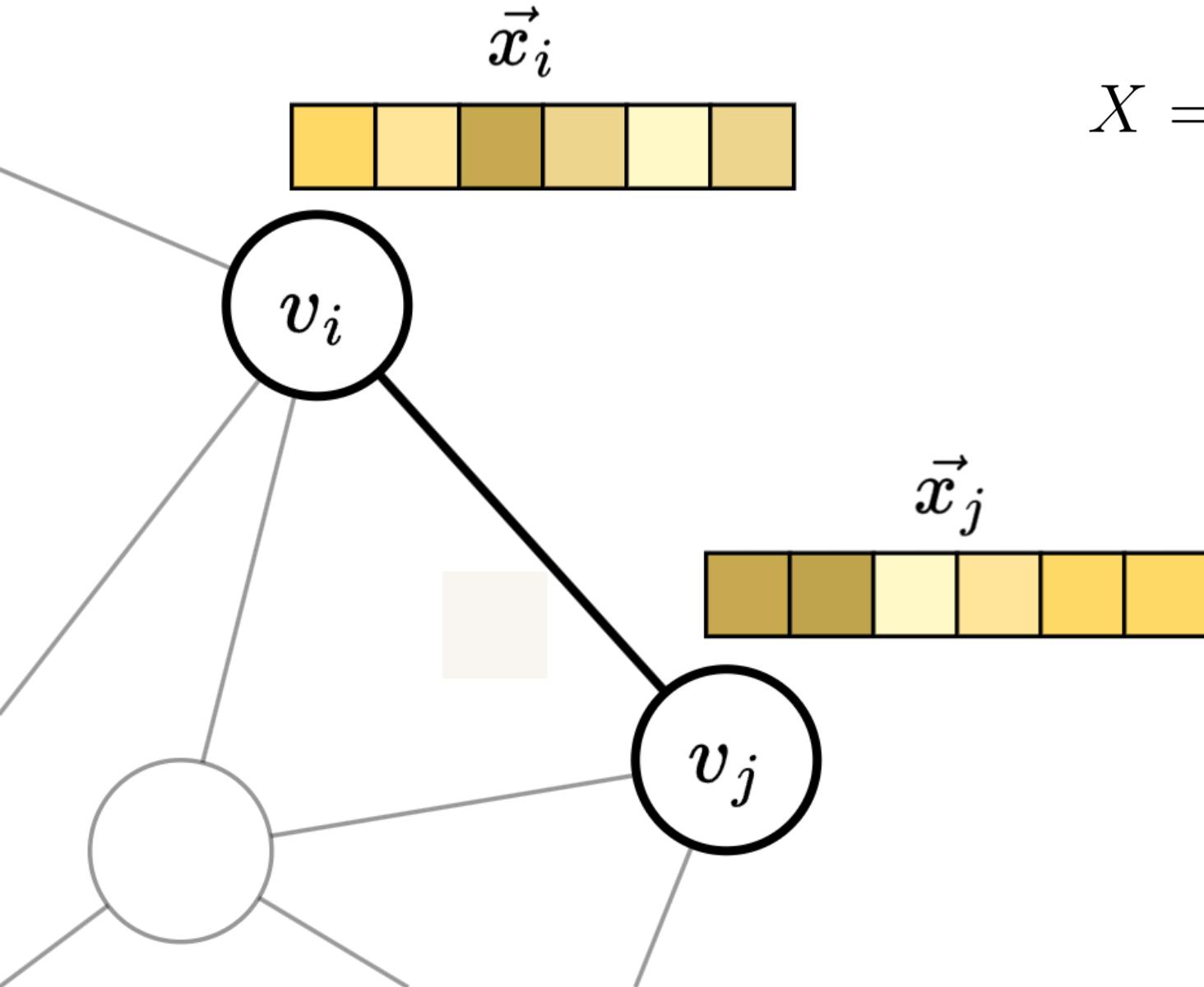


$$\mathcal{G} = (V, E)$$



Adjacency matrix of \mathcal{G}

$$A_{ij} = \begin{cases} 0 & \text{if } (i, j) \notin E \\ 1 & \text{if } (i, j) \in E \end{cases}$$



$$X = \begin{pmatrix} \vec{x}_1 \\ \vdots \\ \vec{x}_i \\ \vdots \\ \vec{x}_n \end{pmatrix}$$

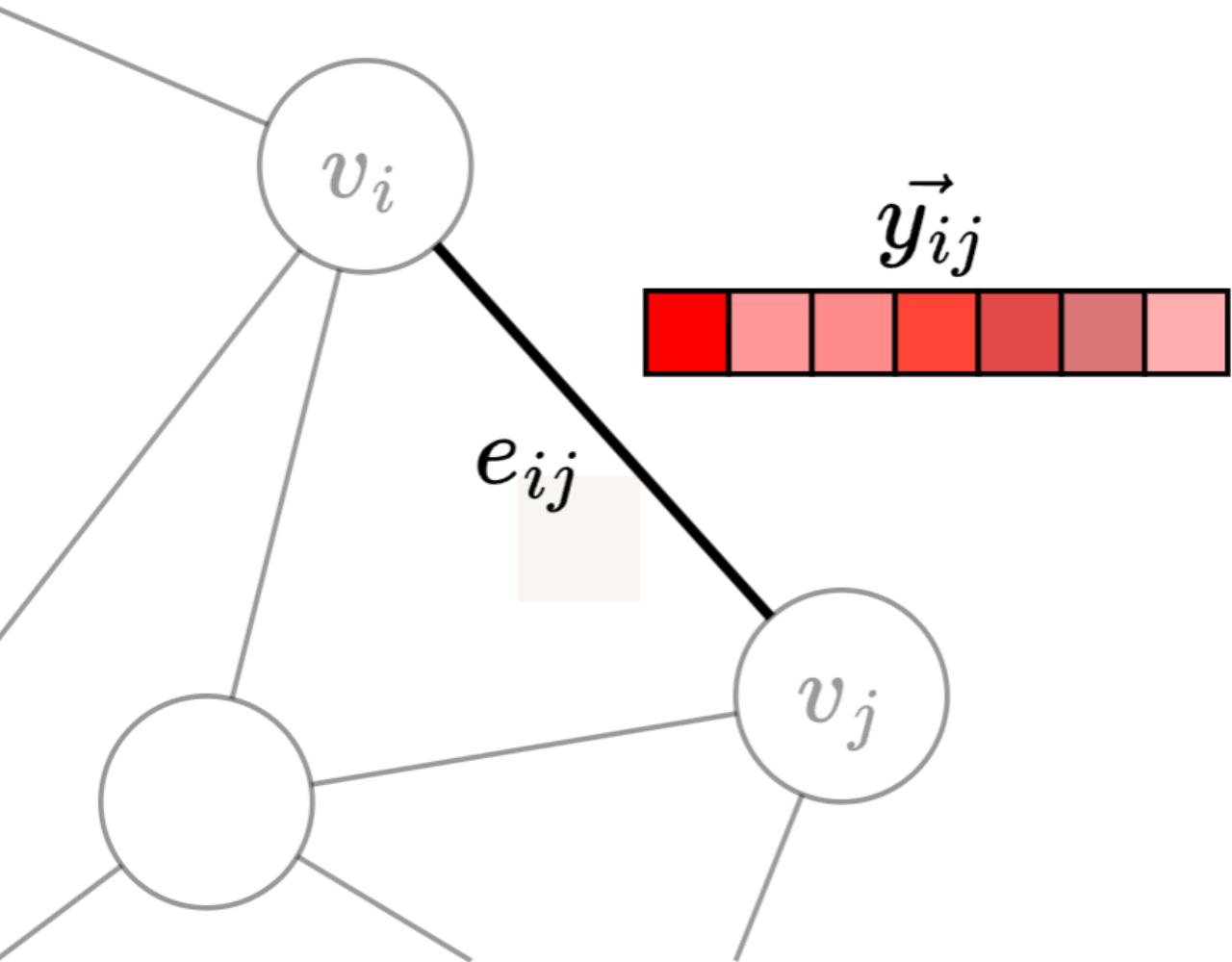
node features of \mathcal{G}

for example:

- properties of the atom
- data of the account
- number/type of crossroads

$$Y_{ij} = \vec{y}_{ij}$$

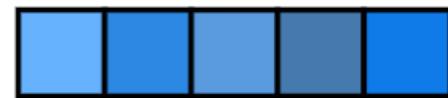
edge features of \mathcal{G}



for example:

- type of bond
- date of friend request, relationship
- average traffic, number of car parks

$$\vec{g}$$



global features of \mathcal{G}

for example:

- known properties
- year
- day of the week

Tasks on graphs

Global tasks

tasks on the **whole graph**.

Both classification and regression.

Examples:

- predicting properties of a molecule
- predicting the hobby shared by all the people
- predicting the likelihood of a car accident

Tasks on graphs

Node tasks

tasks at the level of the **single node**.
Both classification and regression.

Examples:

- predicting if an atom is part of an aromatic ring
- predicting political stance of each users on a topic
- classifying intersections based on predicted traffic

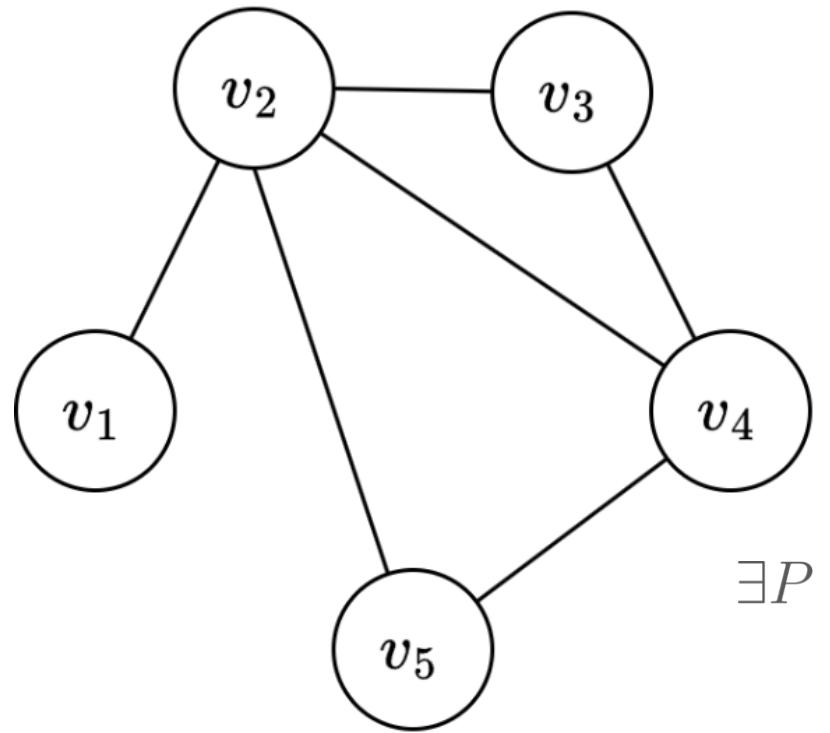
Tasks on graphs

tasks at the level of the **single edge**.
Both classification and regression.

Examples:

- predicting which bond is easier to break
- predicting how likely two users will become friends
- predicting how busy a road will be in an hour

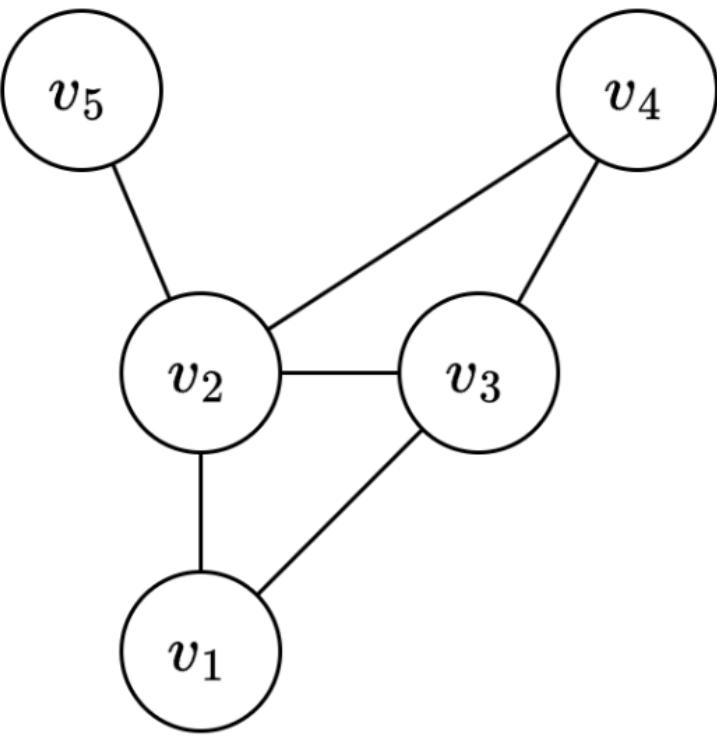
Edge tasks



$\Leftarrow \Rightarrow$

$\exists P$ permutation matrix s.t.

$$PA_1P^{-1} = A_2$$



$$A_1 = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 1 & 1 \\ 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 & 1 \\ 0 & 1 & 0 & 1 & 0 \end{bmatrix}$$

$$A_2 = \begin{bmatrix} 0 & 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 1 \\ 1 & 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \end{bmatrix}$$

$\forall P$ permutation matrix

invariance in the context of global tasks

$$F_G(A, X, Y, \vec{g}) = F_G(PAP^{-1}, PX, PYP^{-1}, \vec{g}) \in \mathbb{R} \text{ or } \in \mathbb{R}^d$$

equivariance in the context of node or edge tasks

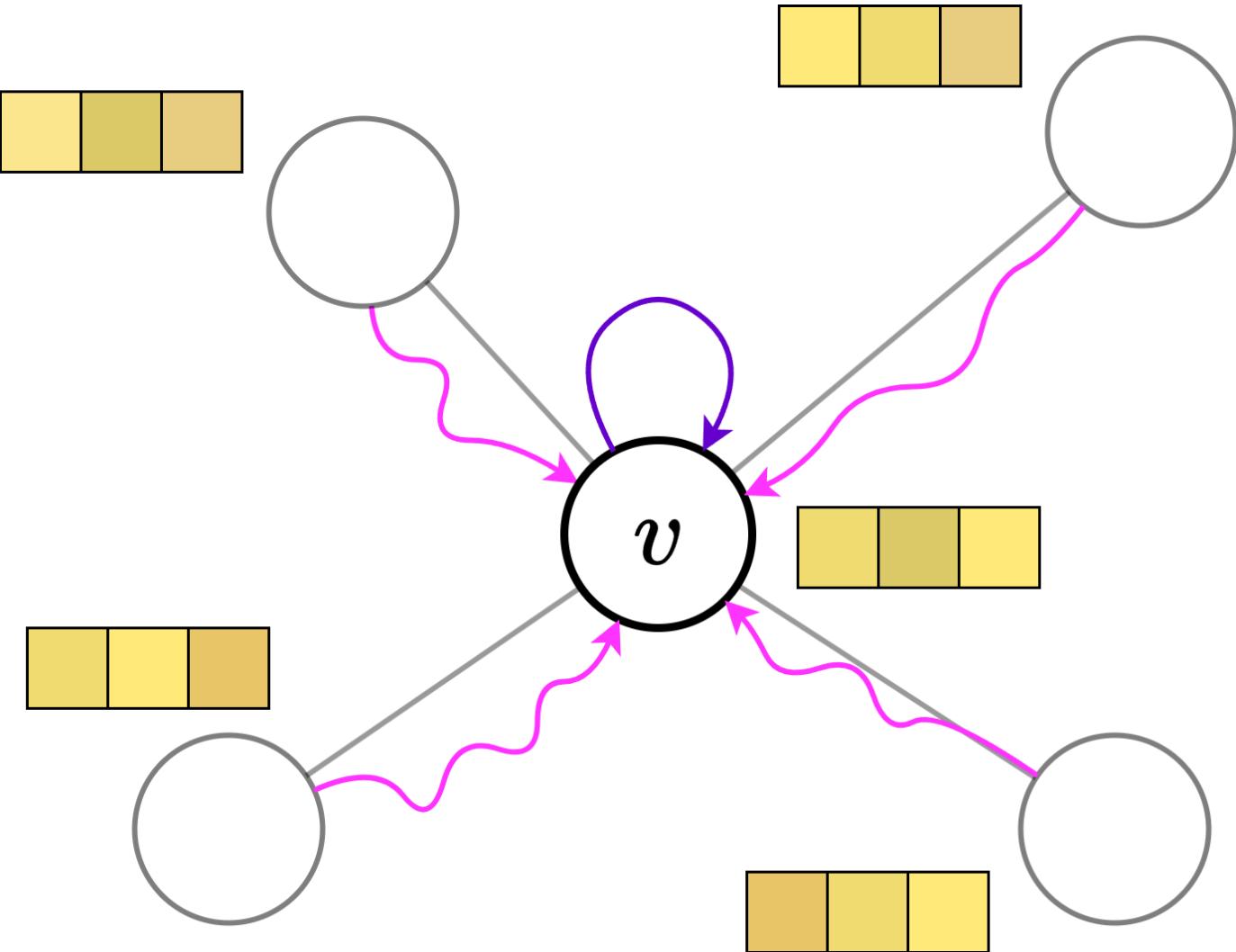
$$F_V(A, X, Y, \vec{g}) = P^{-1}F_V(PAP^{-1}, PX, PYP^{-1}, \vec{g})$$
$$\in \mathbb{R}^n \text{ or } \in \mathbb{R}^{n \cdot d}$$

$$F_E(A, X, Y, \vec{g}) = P^{-1}F_E(PAP^{-1}, PX, PYP^{-1}, \vec{g})P$$
$$\in \mathcal{M}_{\mathbb{R}}(n, n) \text{ or } \in \mathcal{M}_{\mathbb{R}^d}(n, n)$$

in the context of node/global tasks

input: A, X

- **Aggregating** over \mathcal{N}_v
- **Updating** \vec{x}_v



Convolutional approach

$$\vec{x}_v = \phi \left(\vec{x}_v, \bigoplus_{u \in \mathcal{N}_v} c_{vu} \psi(\vec{x}_u) \right)$$

Update

$$\phi(\vec{x}, \vec{z}) = \sigma(W\vec{x} + U\vec{z} + \vec{b})$$

learnable

$$\vec{x}_v = \phi \left(\vec{x}_v, \bigoplus_{u \in \mathcal{N}_v} c_{vu} \psi(\vec{x}_u) \right)$$

coefficient
in the case of summation
can be

$$c_{uv} = \frac{1}{\sqrt{\deg(u) \cdot \deg(v)}}$$

must be a **permutation invariant operator**

- mean
- max
- sum

transformation of
the features

$$\psi(\vec{x}) = W\vec{x} + \vec{b}$$

learnable

Attentional approach

$$\vec{x}_v = \phi \left(\vec{x}_v, \bigoplus_{u \in \mathcal{N}_v} a(\vec{x}_v, \vec{x}_u) \psi(\vec{x}_u) \right)$$



coefficients are **learned**
via NN and softmax
normalized

$$a(\vec{x}_v, \vec{x}_u) \in \mathbb{R}$$

Message-passing approach

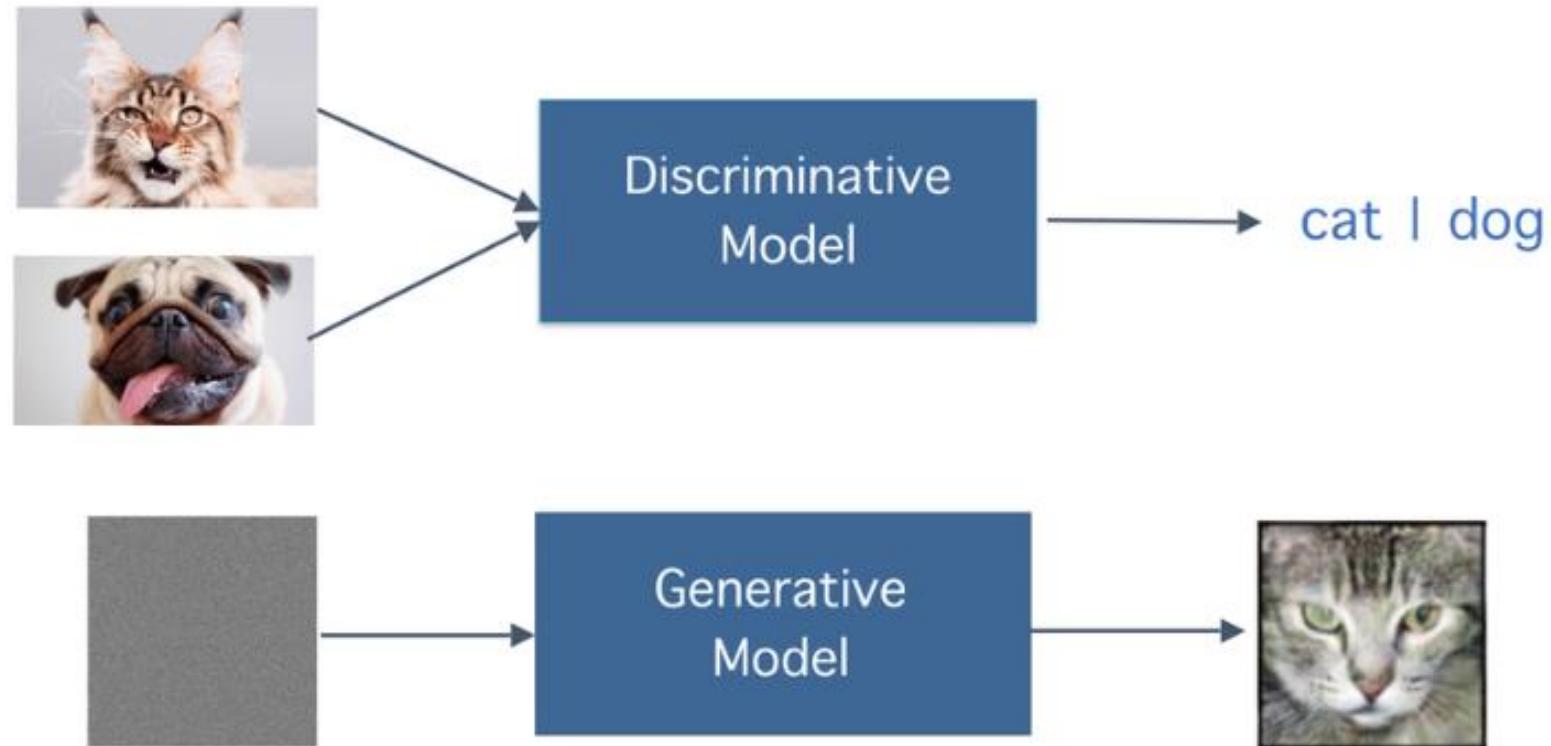
$$\vec{x}_v = \phi \left(\vec{x}_v, \bigoplus_{u \in \mathcal{N}_v} \psi(\vec{x}_v, \vec{x}_u) \right)$$

ψhere is a **learnable message function**

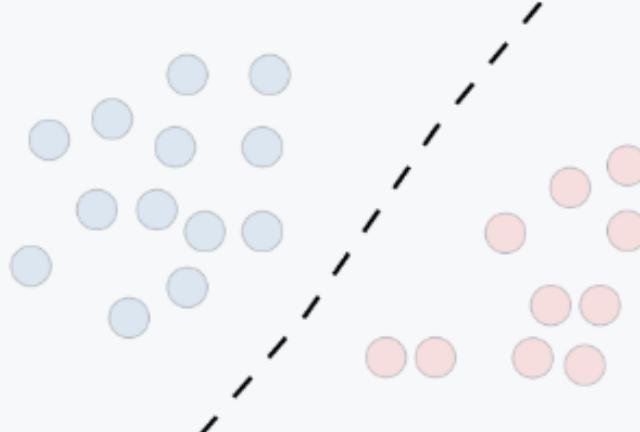
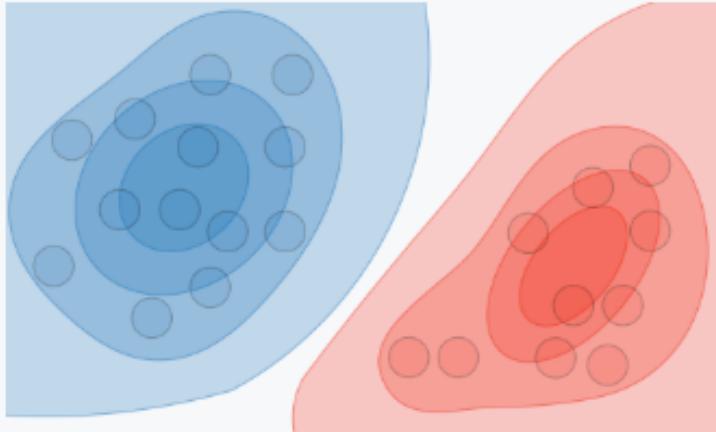
$$\psi(\vec{x}_v, \vec{x}_u) \in \mathbb{R}^m$$

Short intro to generative models

Generative vs discriminative model



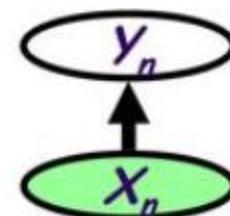
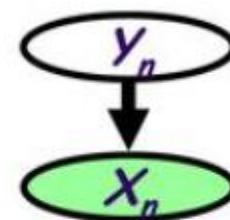
Generative vs discriminative model

	Discriminative model	Generative model
Goal	Directly estimate $P(y x)$	Estimate $P(x y)$ to then deduce $P(y x)$
What's learned	Decision boundary	Probability distributions of the data
Illustration		

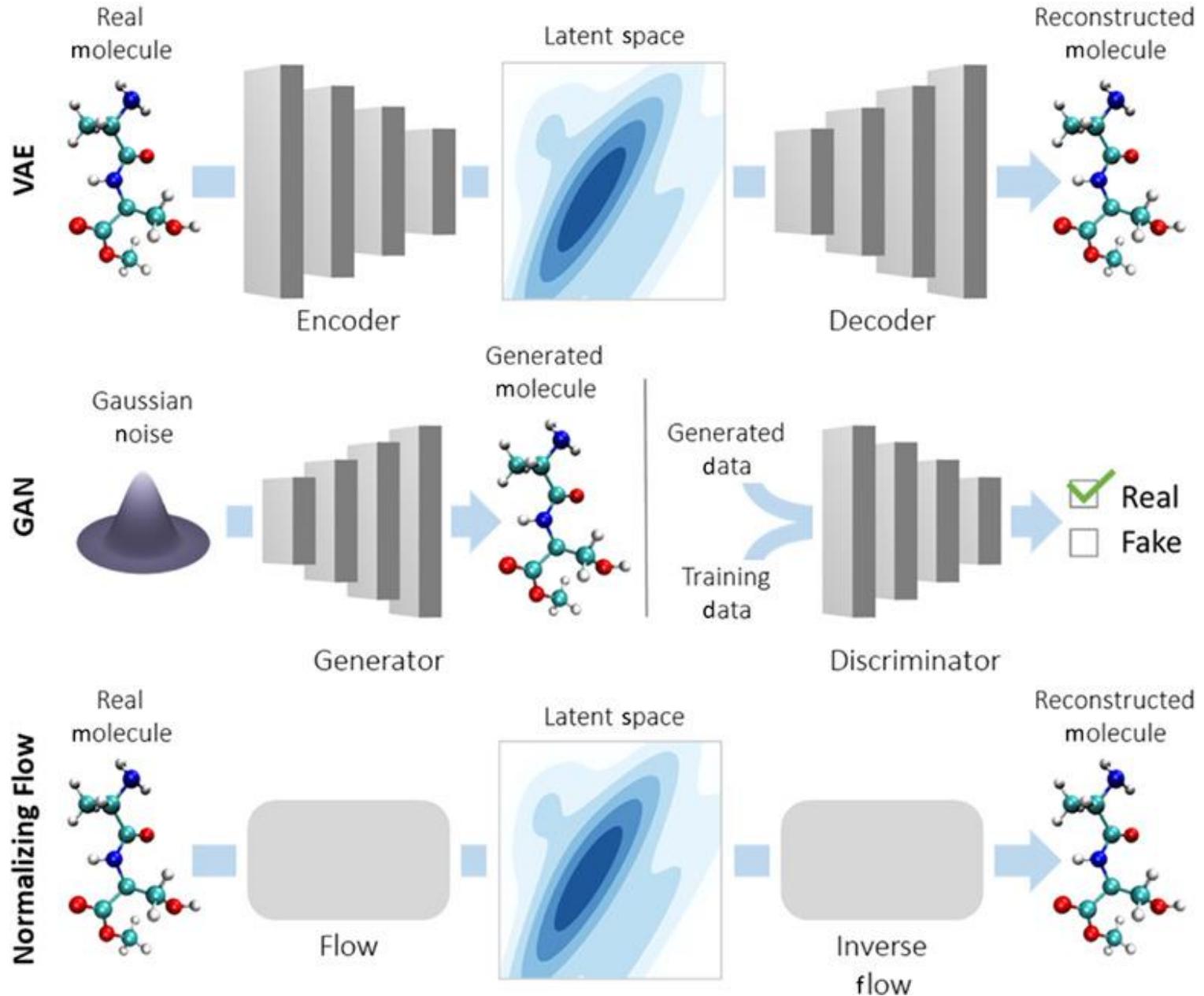
Some examples of discriminative models?

Classifiers

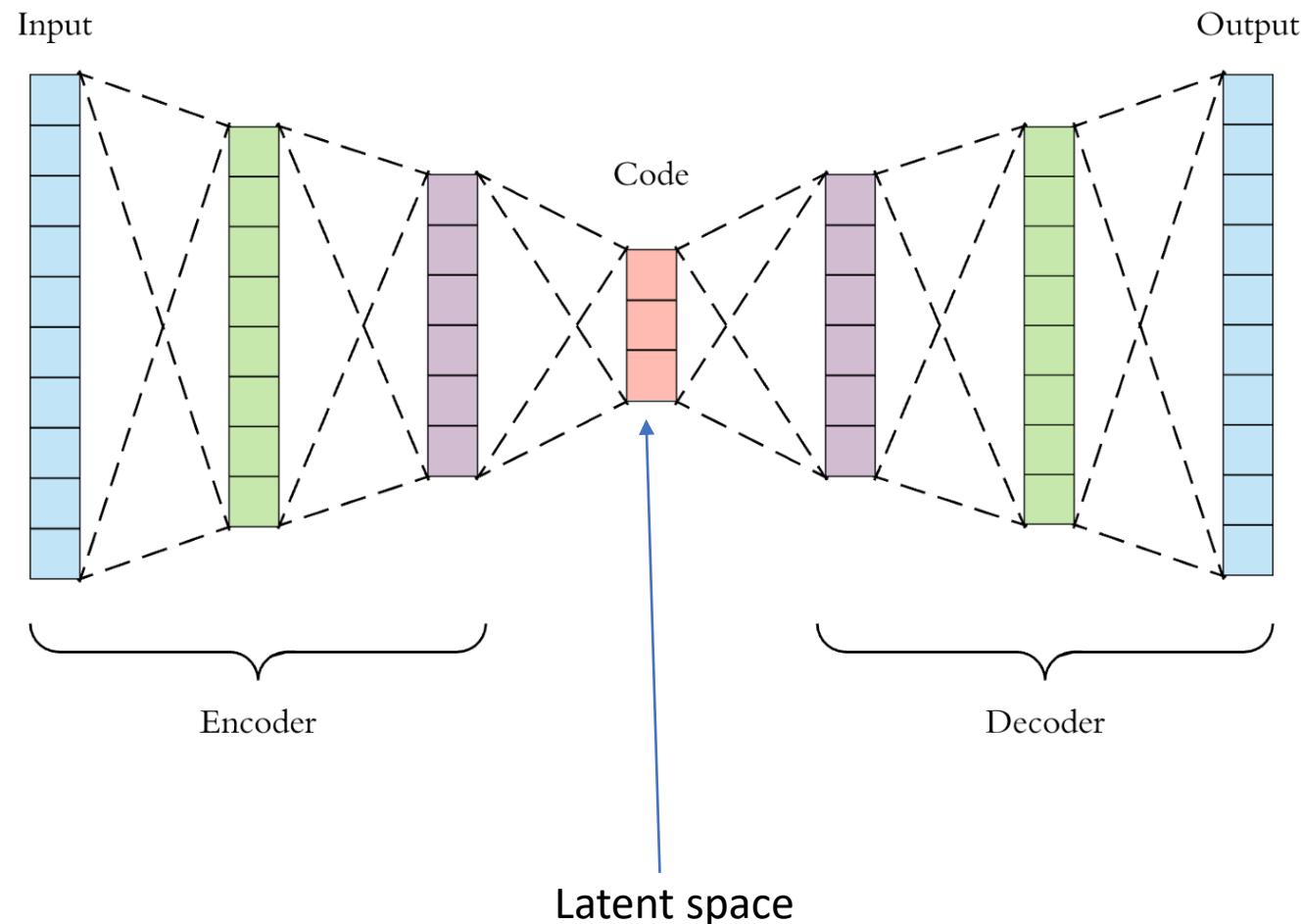
- Goal: Wish to learn $f: X \rightarrow Y$, e.g., $P(Y|X)$
- Generative classifiers (e.g., Naïve Bayes):
 - Assume some functional form for $P(X|Y), P(Y)$
This is a '**generative**' model of the data!
 - Estimate parameters of $P(X|Y), P(Y)$ directly from training data
 - Use Bayes rule to calculate $P(Y|X=x)$
- Discriminative classifiers (e.g., logistic regression)
 - Directly assume some functional form for $P(Y|X)$
This is a '**discriminative**' model of the data!
 - Estimate parameters of $P(Y|X)$ directly from training data



Types of generative models



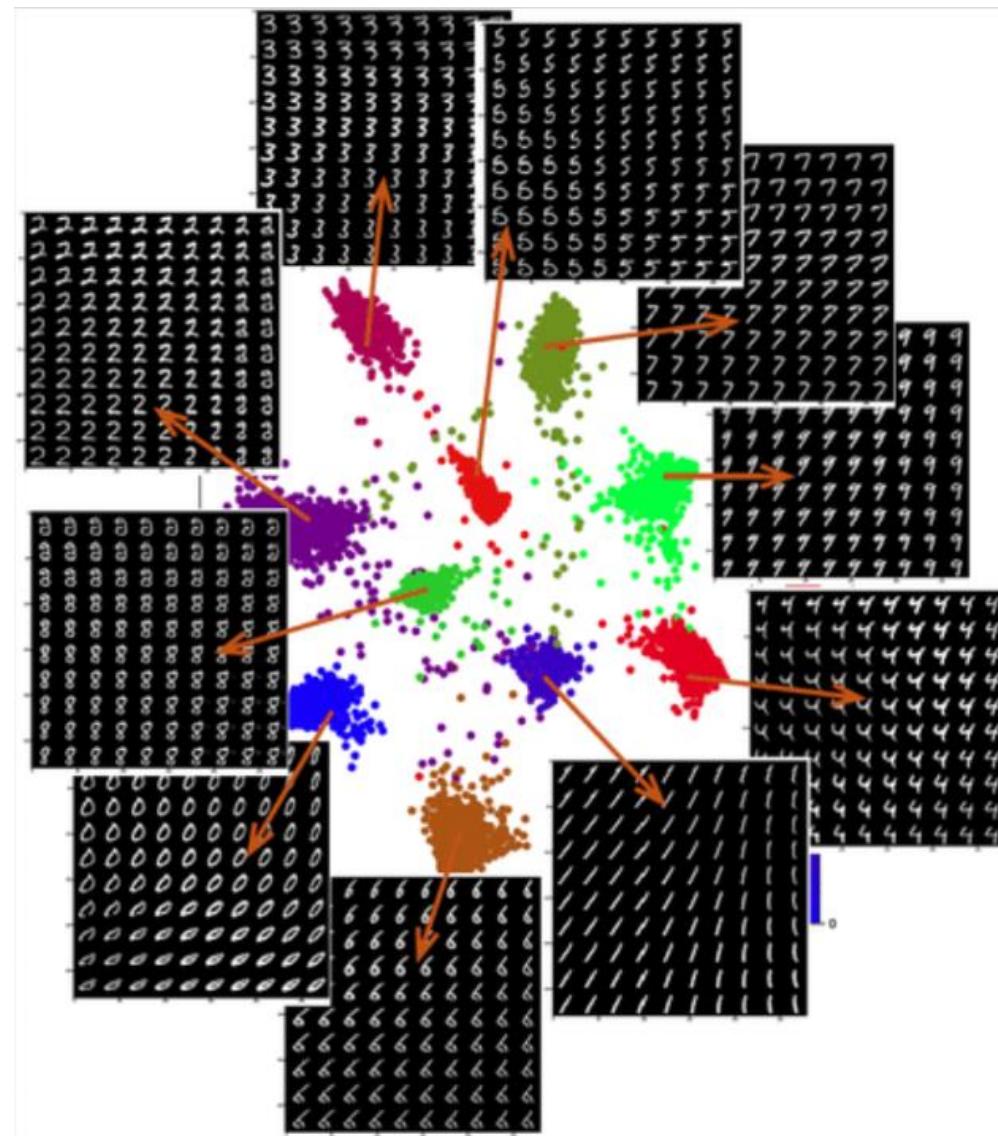
Autoencoders



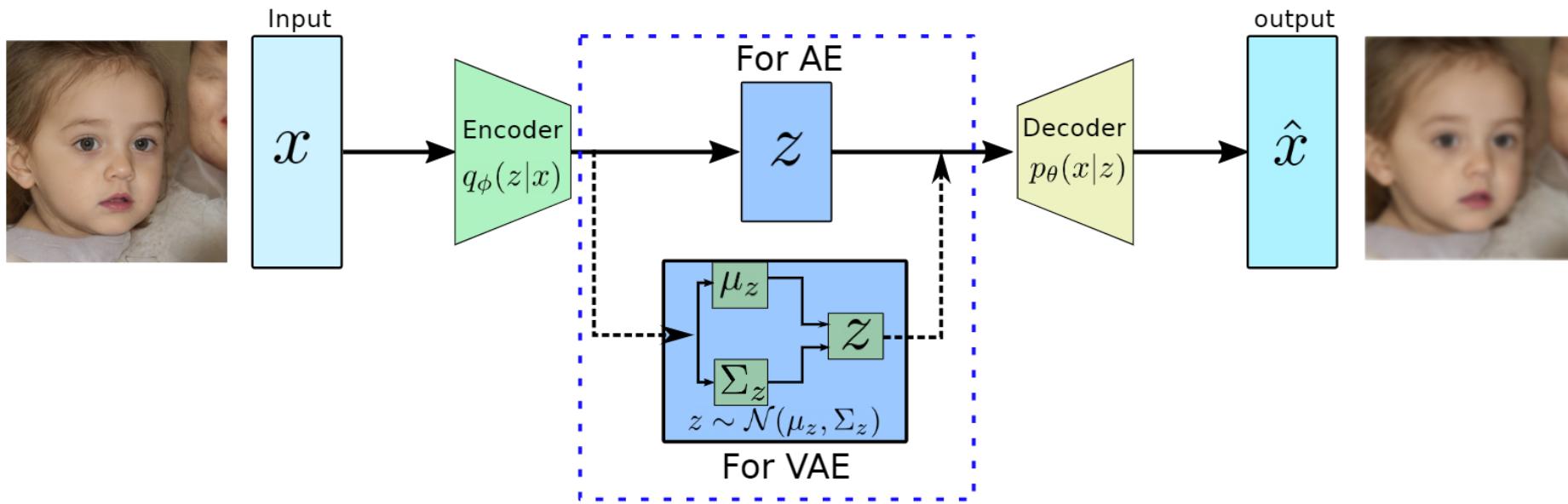
How would you write the loss of an autoencoder?

How would you use an autoencoder as a generative model?

Latent space for MNIST

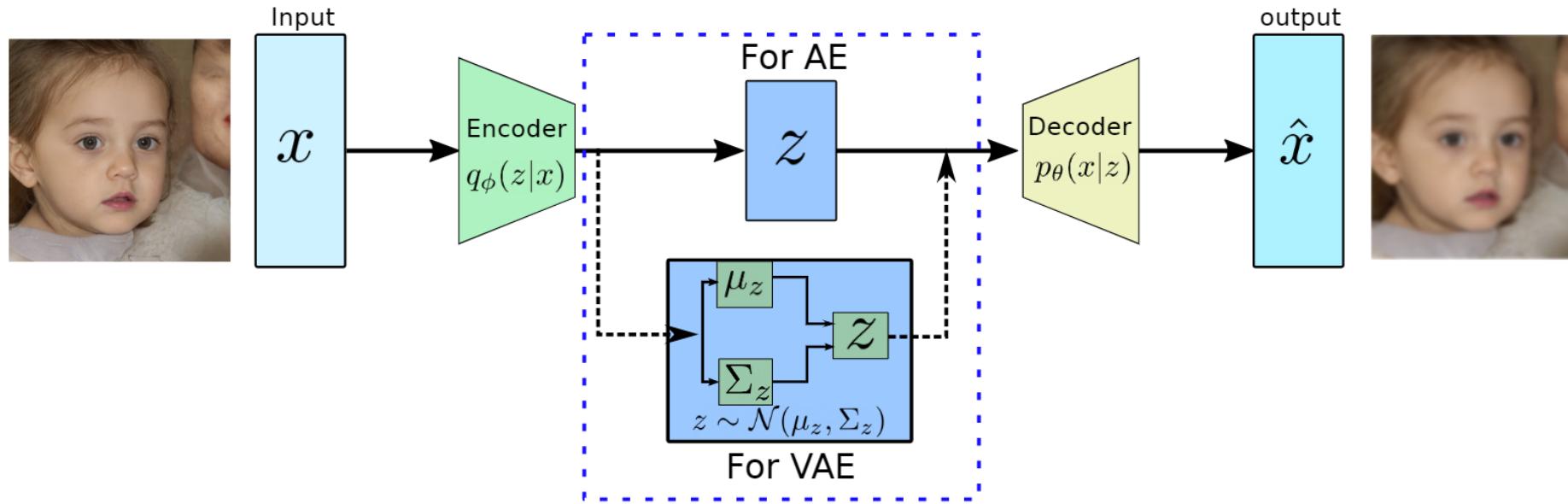


Variational autoencoder



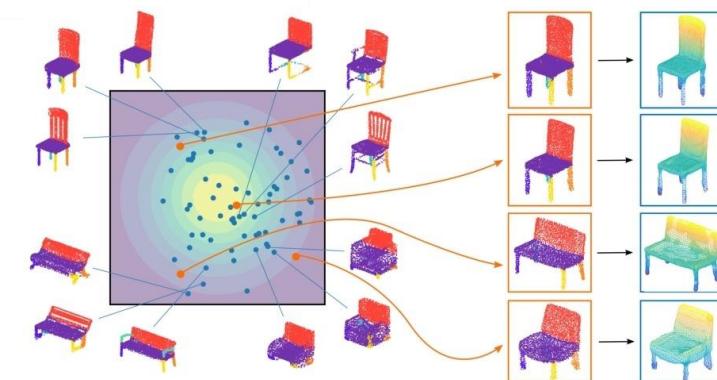
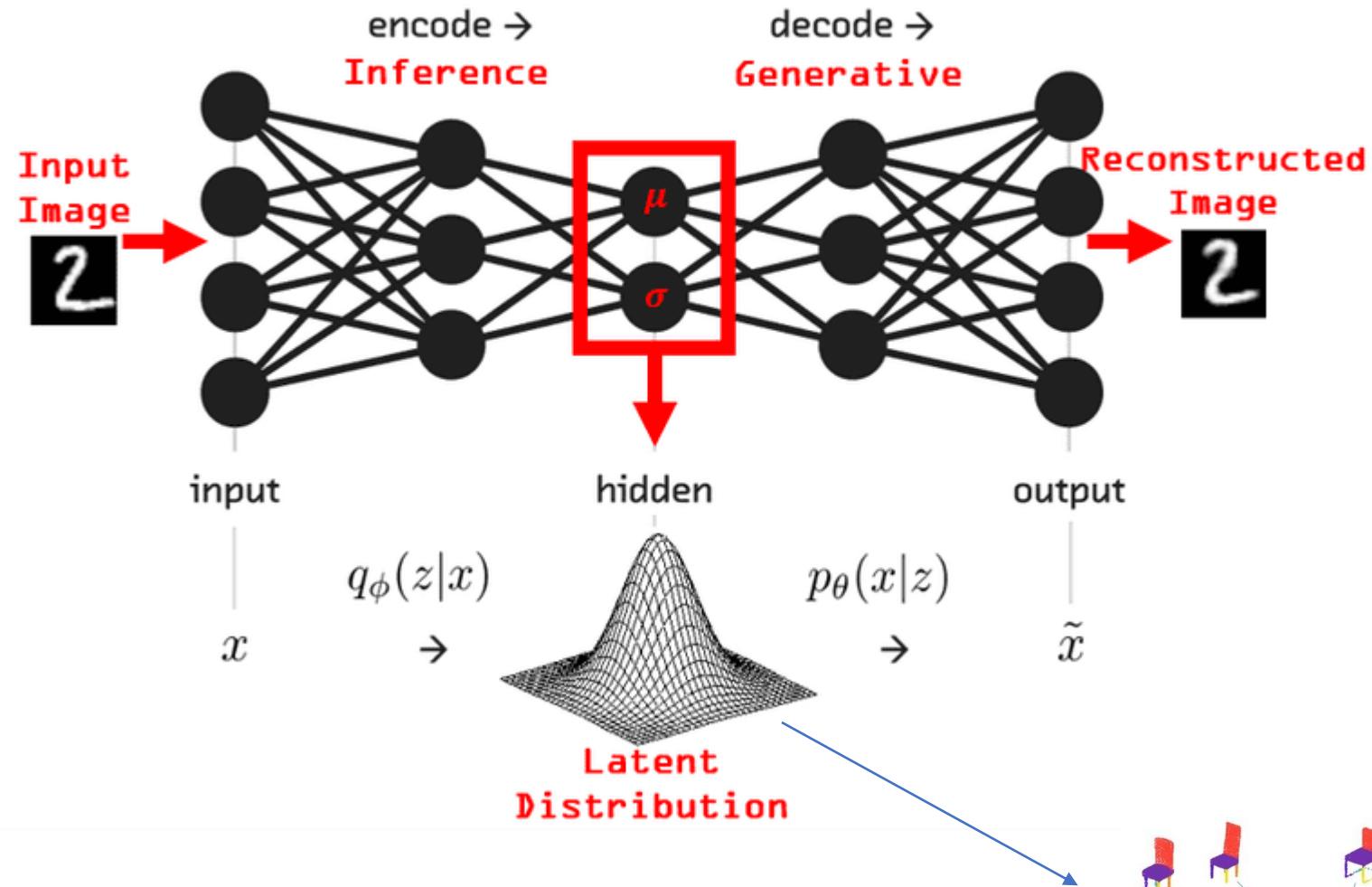
Instead of outputting the vectors in the latent space, the encoder of VAE outputs **parameters of a pre-defined distribution in the latent space for every input**. The VAE then imposes a constraint on this latent distribution forcing it to be a normal distribution.

Variational autoencoder

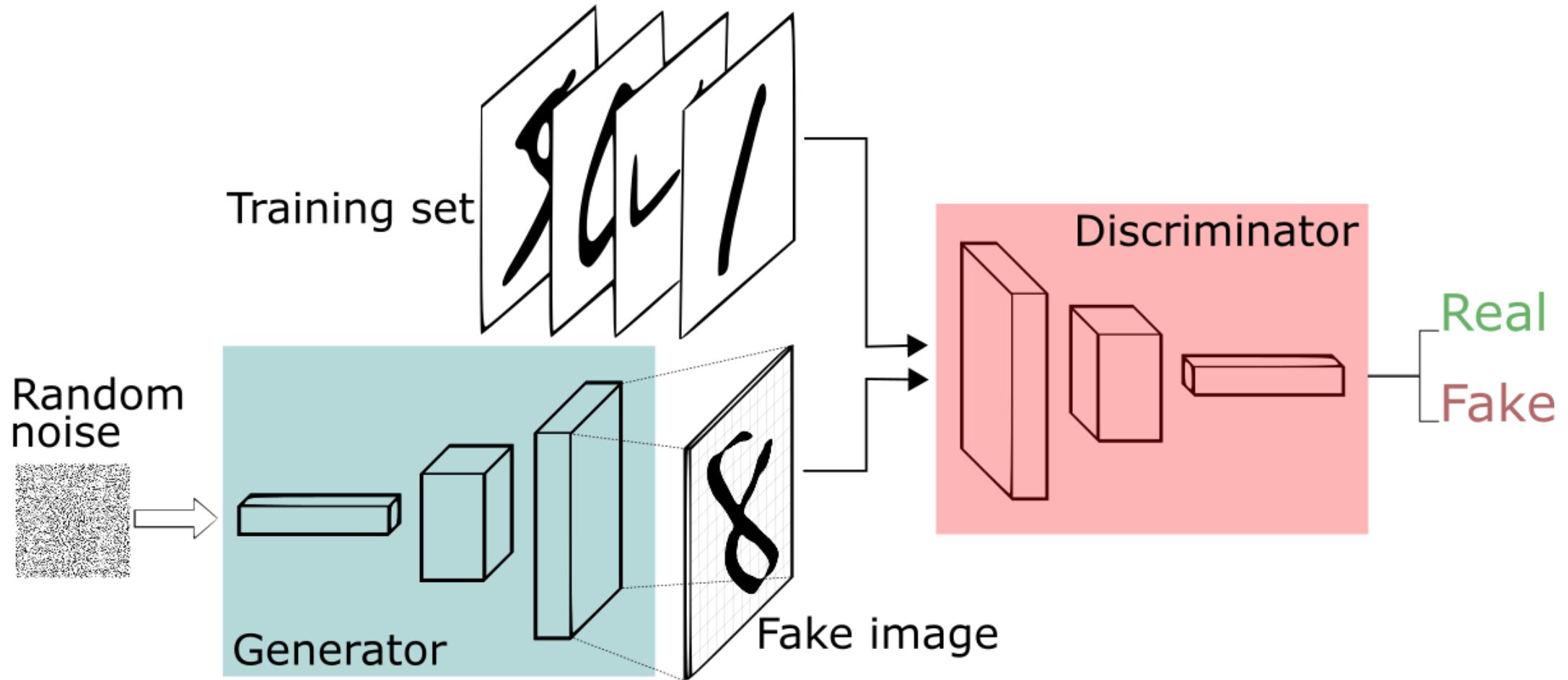


The latent vector is sampled from the encoder-generated distribution before feeding it to the decoder. This random sampling makes it difficult for backpropagation to happen for the encoder since we can't trace back errors due to this random sampling. Hence we use a **reparameterization trick** to model the sampling process which makes it possible for the errors to propagate through the network. The latent vector z is represented as a function of the encoder's output.

$$z = \mu_x + \sigma_x \epsilon, \quad \epsilon \sim \mathcal{N}(0, I)$$



Generative adversarial networks



Gan Loss

gradient ascent

predict well on real images
=> want probability close to 1

predict well on fake images
=> want probability close to 0

$$\nabla_{\mathbf{W}_D} \frac{1}{n} \sum_{i=1}^n \left[\underbrace{\log D(\mathbf{x}^{(i)})}_{\text{predict well on real images}} + \underbrace{\log (1 - D(G(\mathbf{z}^{(i)})))}_{\text{predict well on fake images}} \right]$$

Discriminator objective in the neg. log-likelihood (binary cross entropy) perspective:

Real images, $y = 1$

$$\mathcal{L}(\mathbf{w}) = \boxed{-y^{(i)} \log (\hat{y}^{(i)})} - (1 - y^{(i)}) \log (1 - \hat{y}^{(i)})$$

Want, $\hat{y} = 1$

Fake images, $y = 0$

$$\mathcal{L}(\mathbf{w}) = -y^{(i)} \log (\hat{y}^{(i)}) \boxed{- (1 - y^{(i)}) \log (1 - \hat{y}^{(i)})}$$

Want, $\hat{y} = 0$

<https://poloclub.github.io/ganlab/>

<https://pianalytix.com/6-eye-catching-applications-of-gans/>

Can we use NNs to process text?



Representing words as discrete symbols

In traditional NLP, we regard words as discrete symbols:
hotel, conference, motel – a localist representation

one 1, the rest 0's



Words can be represented by one-hot vectors:

hotel = [0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0]

motel = [0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0]

Vector dimension = number of words in vocabulary (e.g., 500,000)

Challenge: How to compute similarity of two words?

Representing words by their context

Distributional hypothesis: words that occur in similar contexts tend to have similar meanings



J.R.Firth 1957

- “You shall know a word by the company it keeps”
- One of the most successful ideas of modern statistical NLP!

...government debt problems turning into banking crises as happened in 2009...

...saying that Europe needs unified banking regulation to replace the hodgepodge...

...India has just given its banking system a shot in the arm...

These context words will represent *banking*.

Distributional hypothesis

“tejuino”



C1: A bottle of ____ is on the table.

C2: Everybody likes ____.

C3: Don't have ____ before you drive.

C4: We make ____ out of corn.

Distributional hypothesis

C1: A bottle of ____ is on the table.

C2: Everybody likes ____.

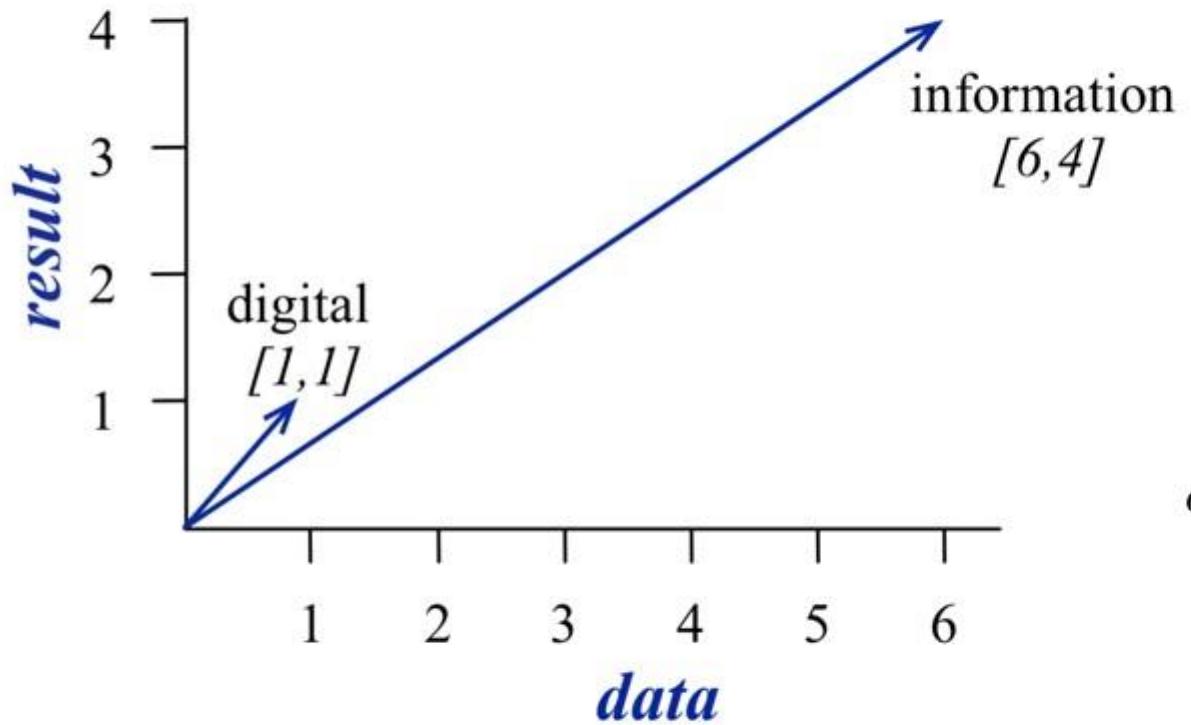
C3: Don't have ____ before you drive.

C4: We make ____ out of corn.

	C1	C2	C3	C4
tejuino	1	1	1	1
loud	0	0	0	0
motor-oil	1	0	0	0
tortillas	0	1	0	1
choices	0	1	0	0
wine	1	1	1	0

“words that occur in similar contexts tend to have similar meanings”

Words as vectors

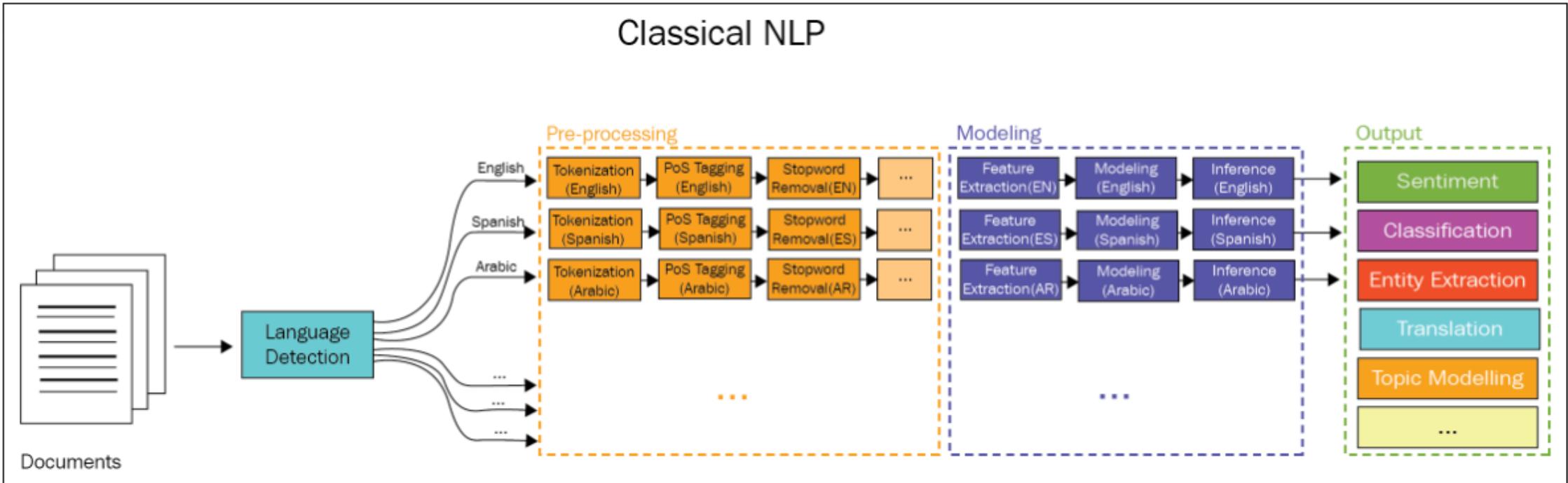


$$\cos(\mathbf{u}, \mathbf{v}) = \frac{\mathbf{u} \cdot \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|}$$

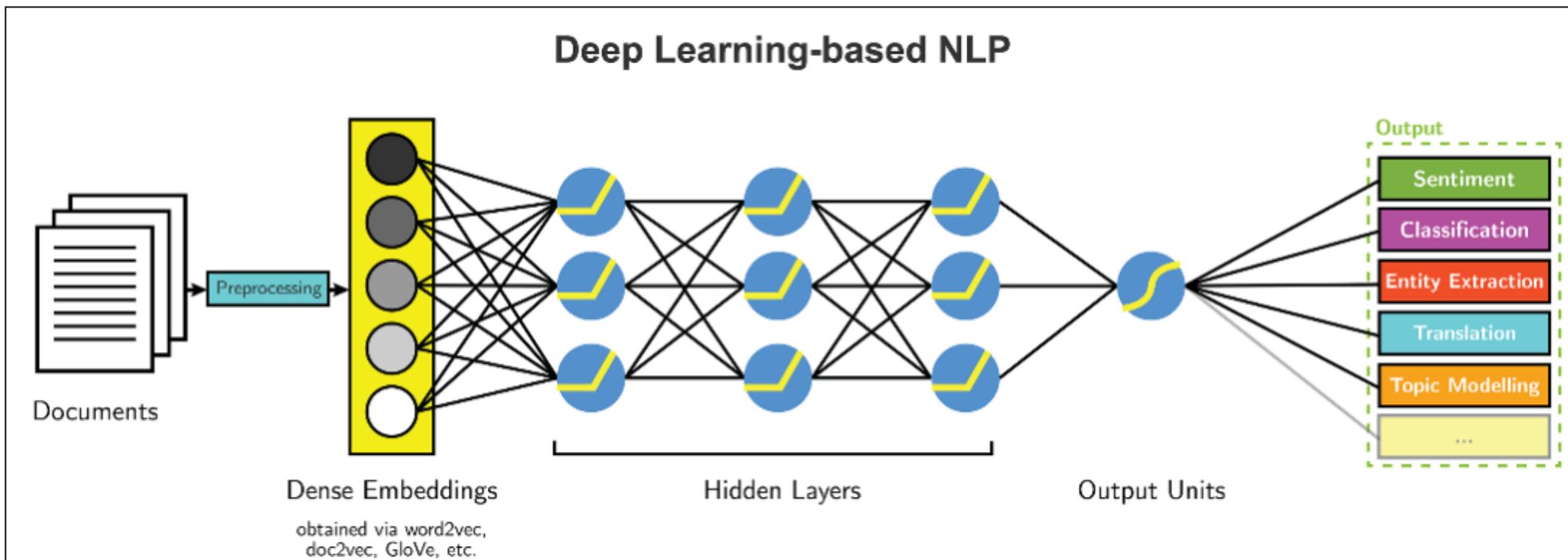
$$\cos(\mathbf{u}, \mathbf{v}) = \frac{\sum_{i=1}^V u_i v_i}{\sqrt{\sum_{i=1}^V u_i^2} \sqrt{\sum_{i=1}^V v_i^2}}$$

What is the range of $\cos(\cdot)$?

Classical NLP



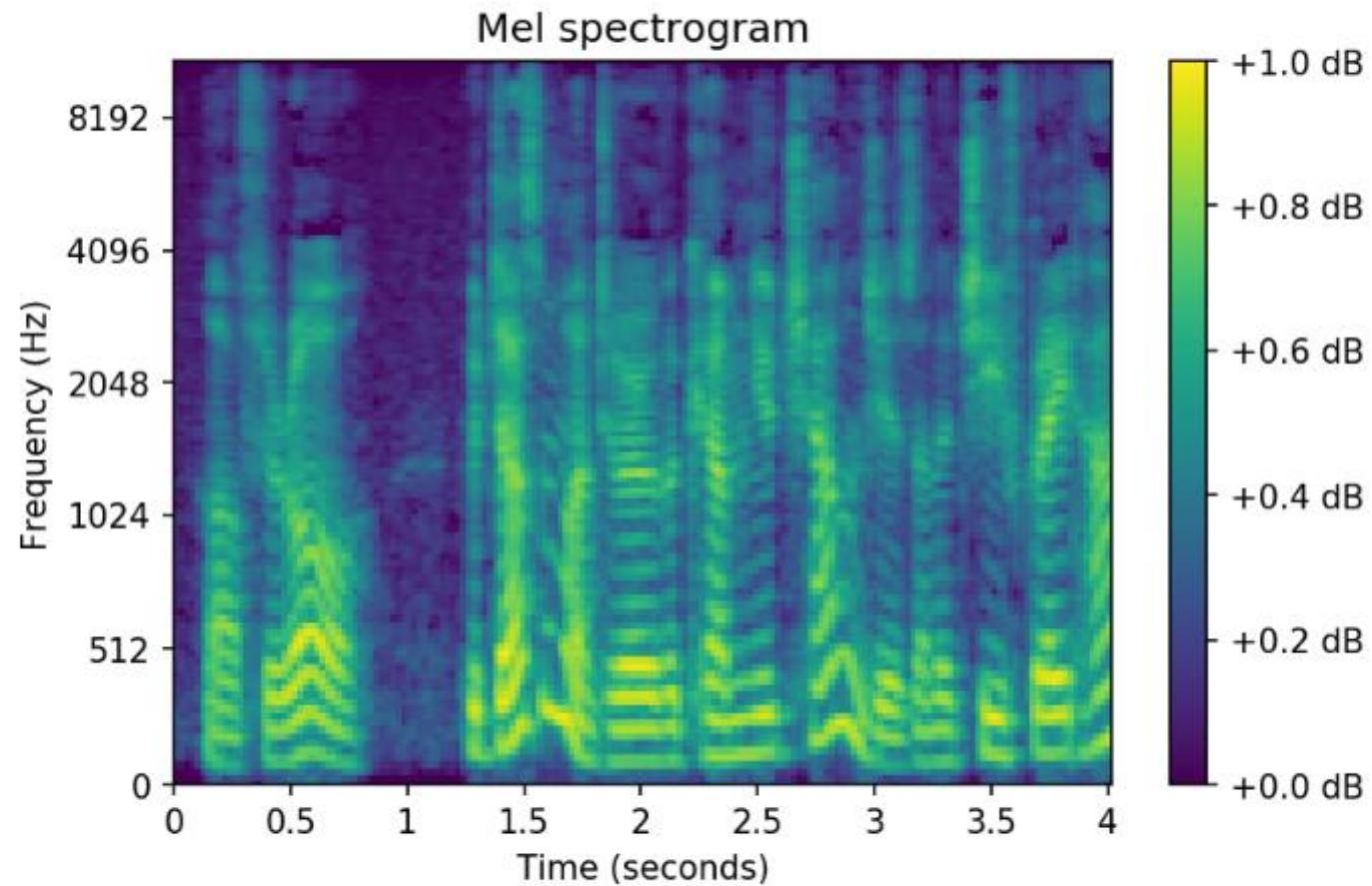
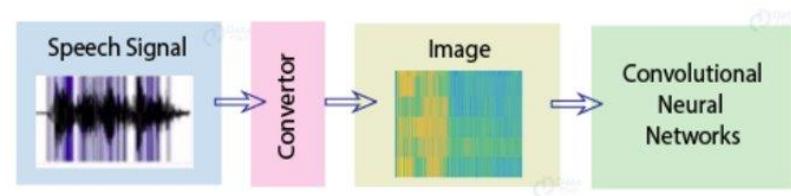
Feed-forward neural LM is a standard feedforward network that takes as input at time t a representation of some number of previous words ($w_{t-1}, w_{t-2} \dots$) and outputs probability distribution over possible next words

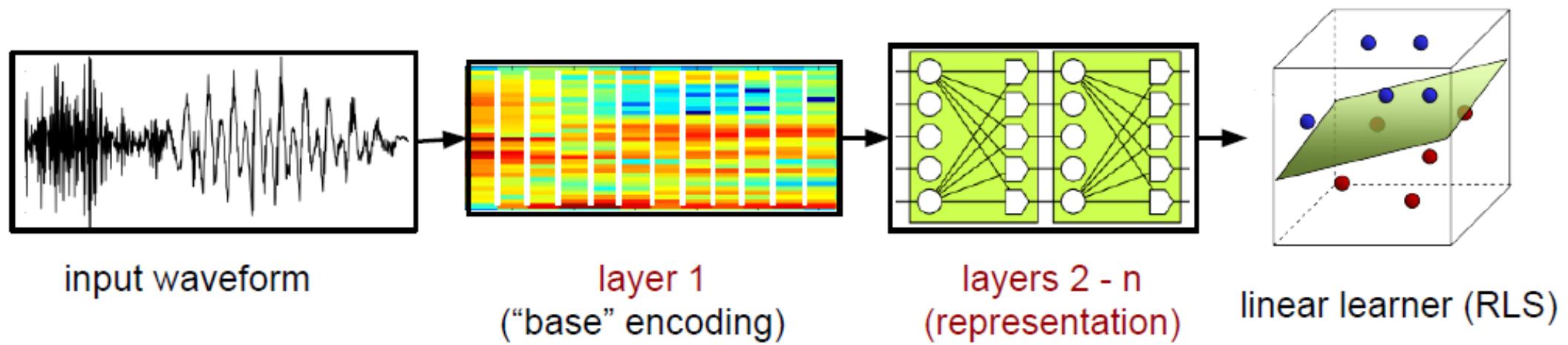
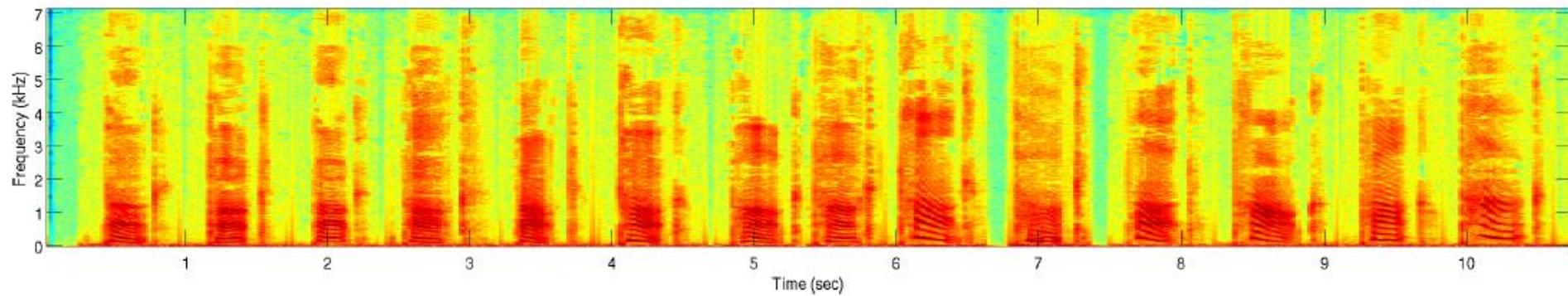
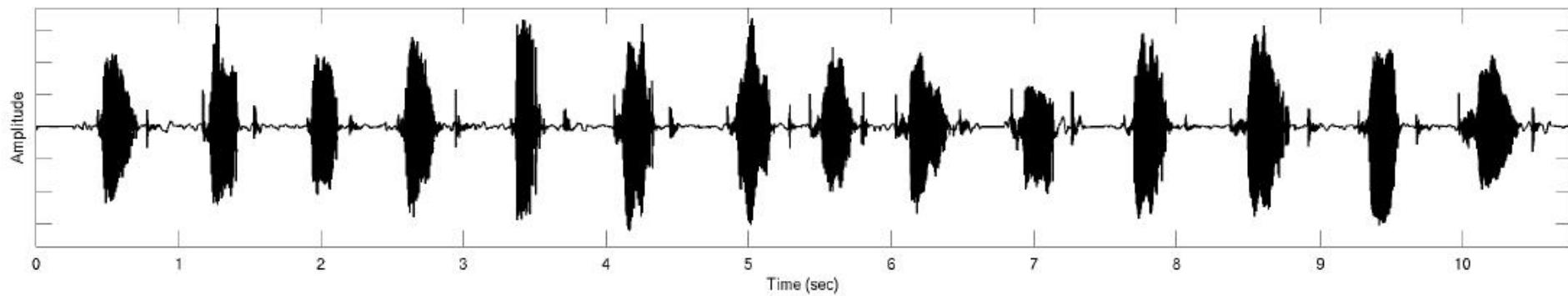


Transformers and large language models

What about sound?

Mel Spectrogram





Class 7

Loose hands, questions