

Deep Learning

Fundamentals of Artificial Intelligence

Supervised Learning in Mathematical Terms

map input to output: $y = f(\mathbf{x})$ (estimated: $\hat{f}(\mathbf{x})$)

random variables Y and $\mathbf{X} = (X_1, X_2, \dots, X_p)$ \leftarrow usually high-dimensional

training (curve fitting / parameter estimation):

fit data set of (y_i, \mathbf{x}_i) pairs \rightarrow minimize deviations between y and $\hat{f}(\mathbf{x})$

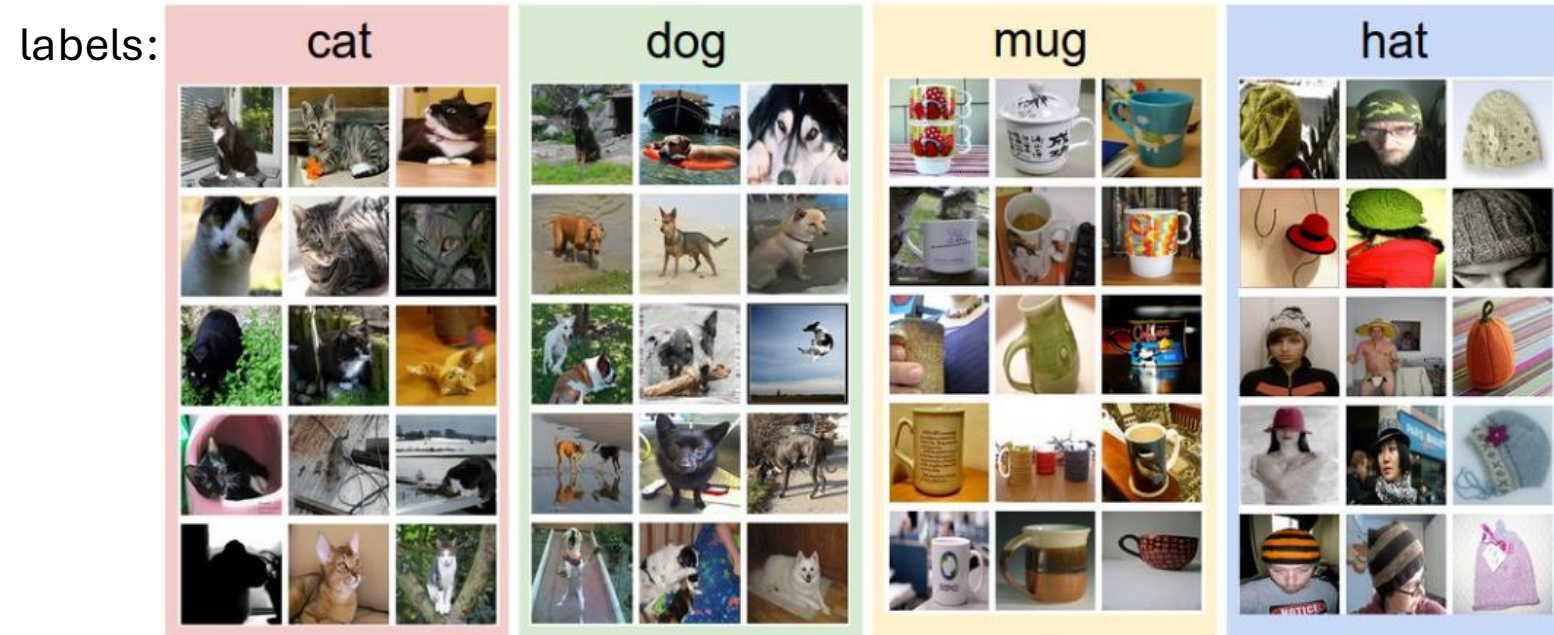
discriminative models: predict conditional density function $p(y|\mathbf{x})$

as opposed to generative models: predict $p(y, \mathbf{x})$ (or just $p(\mathbf{x})$) $\rightarrow \mathbf{x}$ not given, more difficult

Example: Image Classification

training data set:

test data with learned classifier:



$$f\left(\text{img of a cat}\right) = \text{"Cat"}$$

$$f\left(\text{img of a dog}\right) = \text{"Dog"}$$

General Recipe of Statistical Learning

ML algorithm by combining:

- **model** (e.g., linear function & Gaussian distribution)
- **objective function** (e.g., squared residuals)
- **optimization algorithm** (e.g., gradient descent)
- **regularization** (e.g., L2, dropout)

Linear Regression Revisited

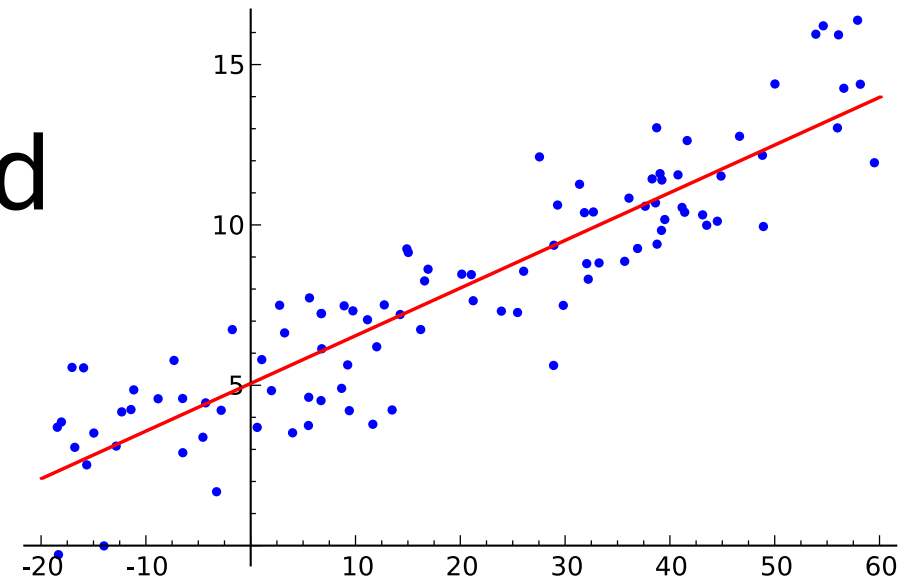
fit: $\hat{f}(\mathbf{x}_i)$

$$y_i = \hat{b} + \sum_{j=1}^p \hat{w}_j x_{ij} + \varepsilon_i$$

predict:

$$\hat{y}_i = E[Y|\mathbf{X} = \mathbf{x}_i] = \hat{f}(\mathbf{x}_i) = \hat{b} + \sum_{j=1}^p \hat{w}_j x_{ij}$$
$$p(y|\mathbf{x}_i) = \mathcal{N}(y; \hat{y}_i, \hat{\sigma}^2)$$

Gaussian mean variance



parameter estimation:
e.g., least squares

parameters to be estimated: $\hat{b}, \hat{\mathbf{w}}$
 $\rightarrow \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(\mathbf{x}_i))^2$
(approximating assumed true b, \mathbf{w}, σ)

Brief Notation

data: vector with p different
feature values for this example i

$$f(\mathbf{x}_i) = \mathbf{w} \cdot \mathbf{x}_i + b$$

vector with slope parameters (one for
each of the p features, but identical
for all examples), aka weights

single intercept
parameter, aka bias

for simplicity of notation, dropping the hats of estimated parameters here

Loss Function

loss function L : expressing deviation between prediction and target

$$L(y_i, \hat{f}(\mathbf{x}_i); \hat{\boldsymbol{\theta}})$$

with $\hat{\boldsymbol{\theta}}$ corresponding to parameters of model $\hat{f}(\mathbf{x})$

e.g., $\hat{b}, \hat{\mathbf{w}}$ in linear regression

prime example: squared residuals for regression problems

$$L(y_i, \hat{f}(\mathbf{x}_i); \hat{\boldsymbol{\theta}}) = \left(y_i - \hat{f}(\mathbf{x}_i; \hat{\boldsymbol{\theta}}) \right)^2$$

Cost Function

averaging losses over (empirical) training data set:

$$J(\hat{\theta}) = \frac{1}{n} \sum_{i=1}^n L(y_i, \hat{f}(\mathbf{x}_i); \hat{\theta})$$

cost function to be minimized according to model parameters $\hat{\theta}$

→ objective function

Cost Minimization

minimize training costs $J(\hat{\boldsymbol{\theta}})$ according to model parameters $\hat{\boldsymbol{\theta}}$:

$$\nabla_{\hat{\boldsymbol{\theta}}} J(\hat{\boldsymbol{\theta}}) = 0$$

for mean squared error (aka least squares method):

$$\nabla_{\hat{\boldsymbol{\theta}}} \frac{1}{n} \sum_{i=1}^n \left(y_i - \hat{f}(\mathbf{x}_i; \hat{\boldsymbol{\theta}}) \right)^2 = 0$$

Gradient Descent

typically no closed-form solution to minimization of cost function: $\nabla_{\hat{\theta}} J(\hat{\theta}) = 0$

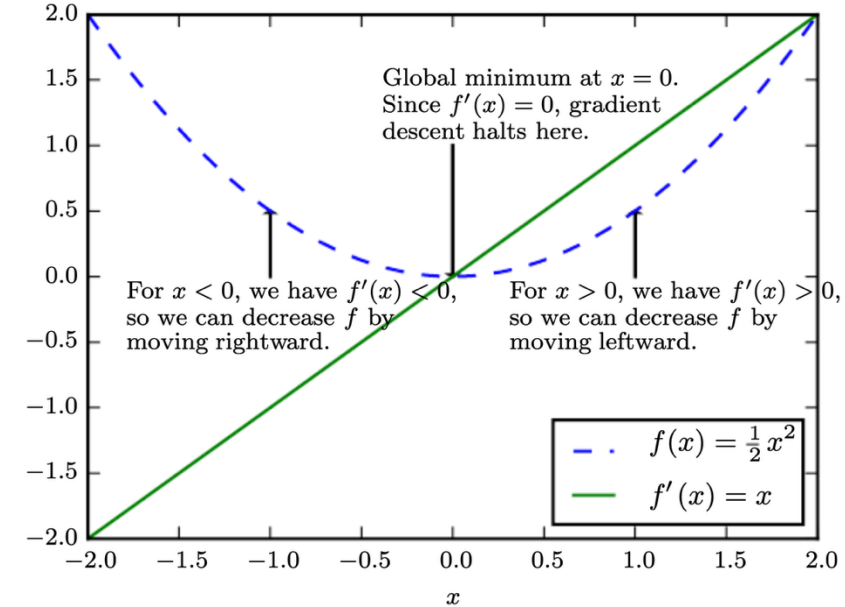
→ need for numerical methods

popular choice: gradient descent

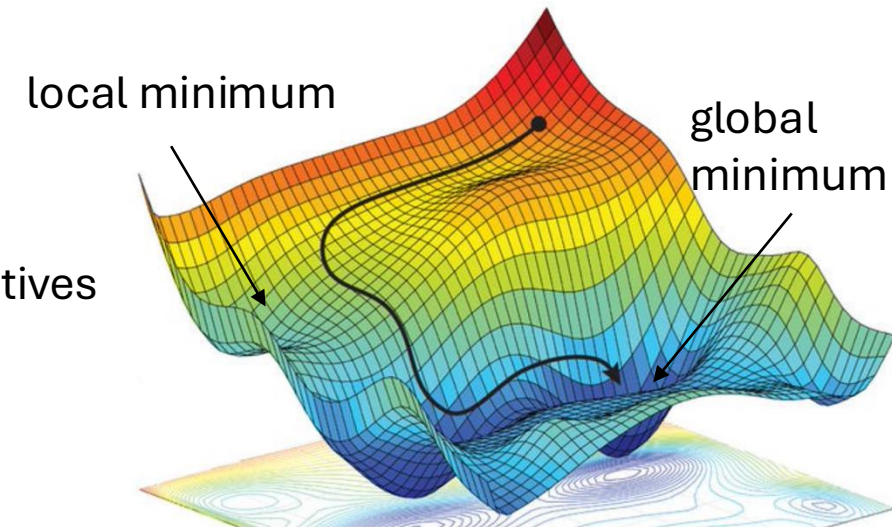
iteratively moving in direction of steepest descent
(negative gradient): $\hat{\theta} \leftarrow \hat{\theta} - \eta \nabla_{\hat{\theta}} J(\hat{\theta})$

step size
(learning rate)

vector containing all partial derivatives



[source](#)



L^2 Regularization (Ridge Regression)

add **penalty term** on parameter L^2 norm to cost function

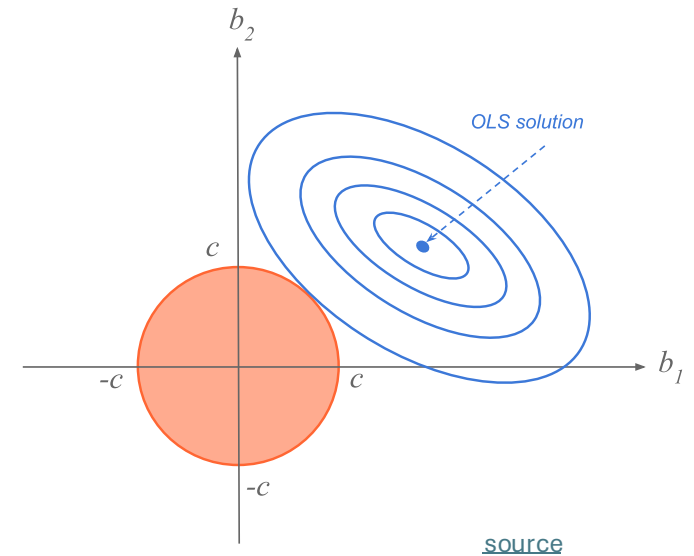
$$\tilde{J}(\hat{\theta}) = J(\hat{\theta}) + \lambda \cdot \hat{\theta}^T \hat{\theta}$$

hyperparameter controlling
strength of regularization

aka **weight decay** in neural networks

corresponds to imposing constraint on parameters to lie in L^2 region (size controlled by λ)

goal: reduce overfitting

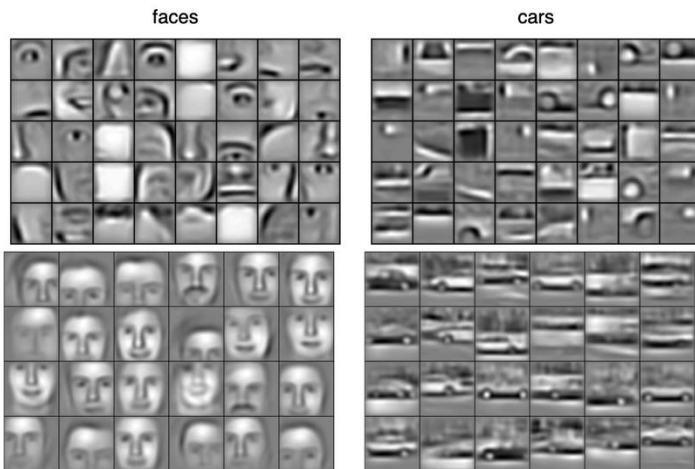


Ladder of Generalization

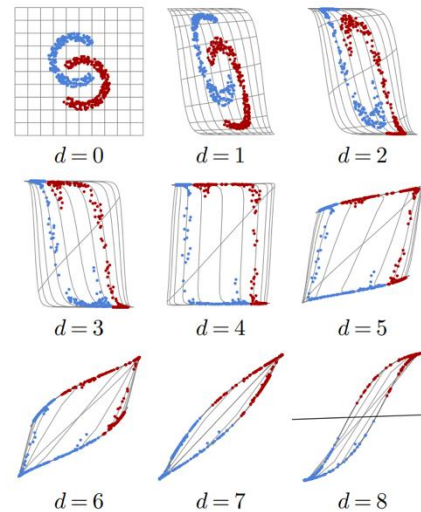
classic ML: feature engineering

deep learning: feature learning

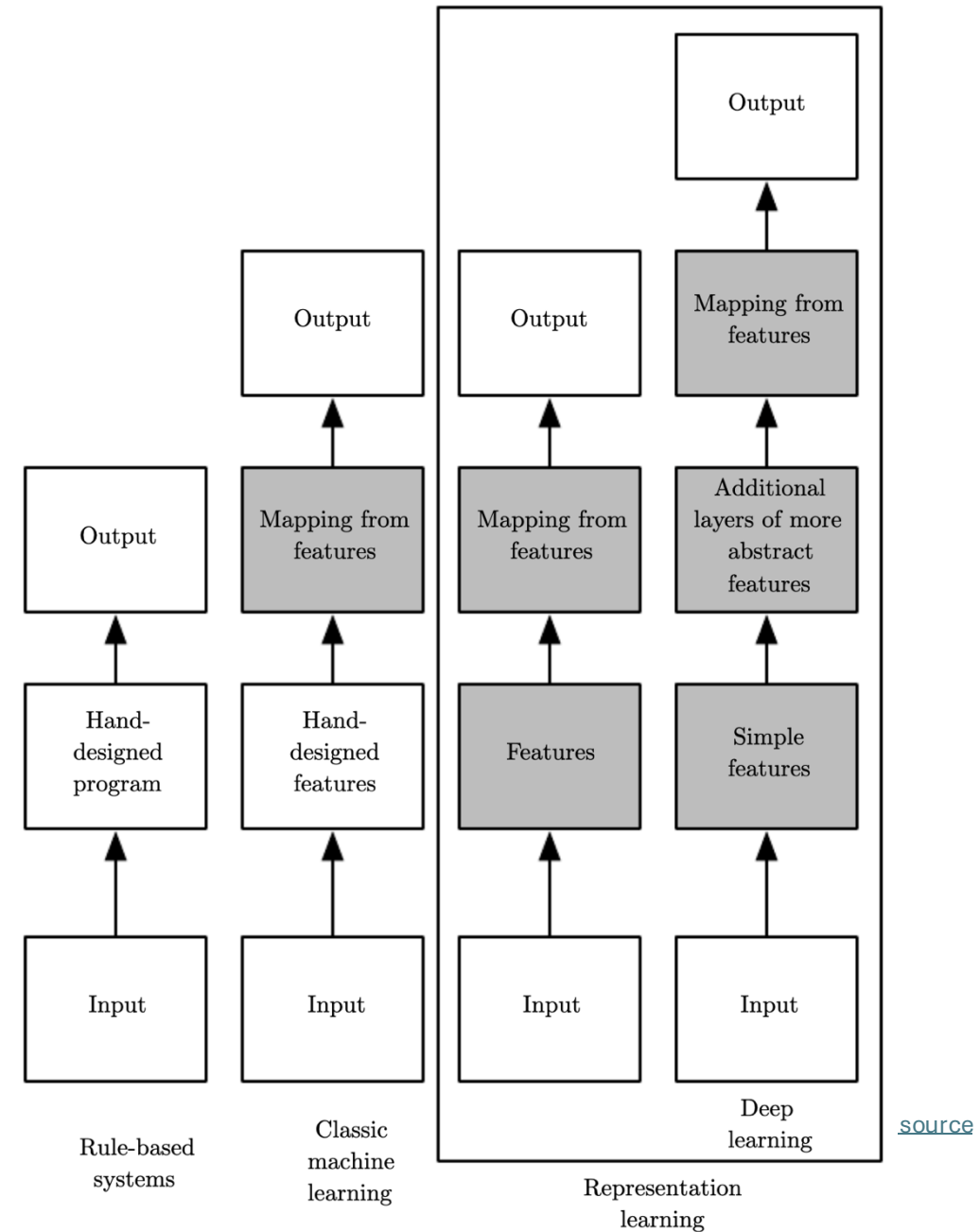
(hierarchy of concepts learned from raw data in deep graph with many layers)



[source](#)

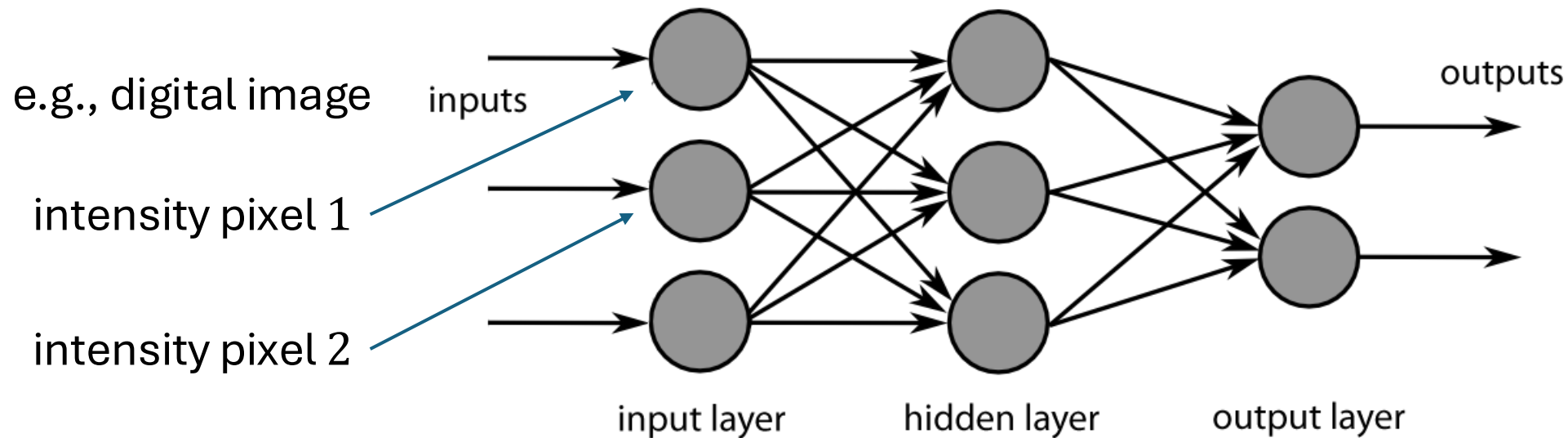


[source](#)



Neural Networks

idea: powerful algorithm by combining many linear building blocks



each node exchanges information only with directly connected nodes
→ parallel computation

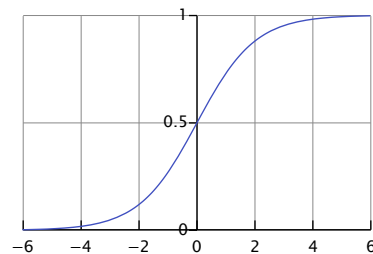
Artificial Neuron

linear model with parameters called weights \mathbf{w} (including bias, not shown for simplicity)

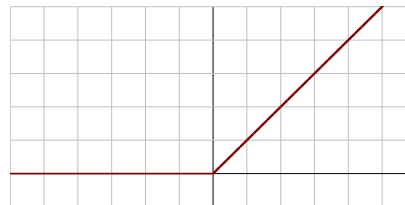
non-linear via (differentiable) activation function on top of linear model

examples:

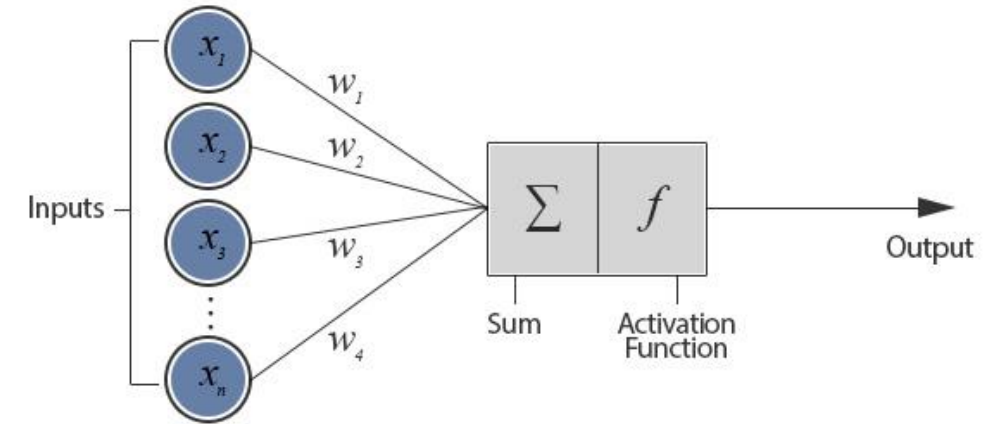
- sigmoid



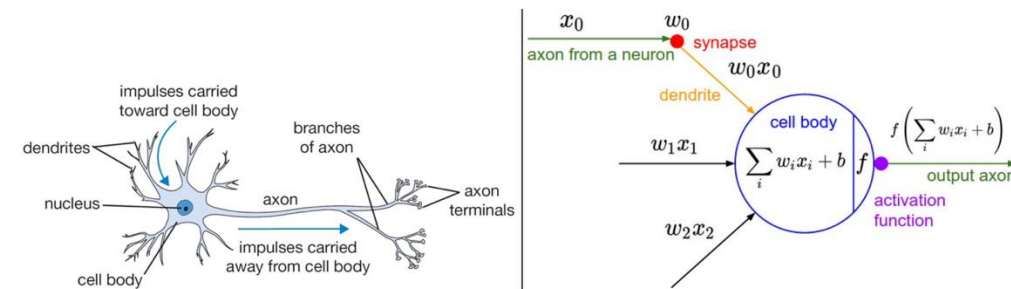
- ReLU (Rectified Linear Unit)



artificial neuron (node in neural network):

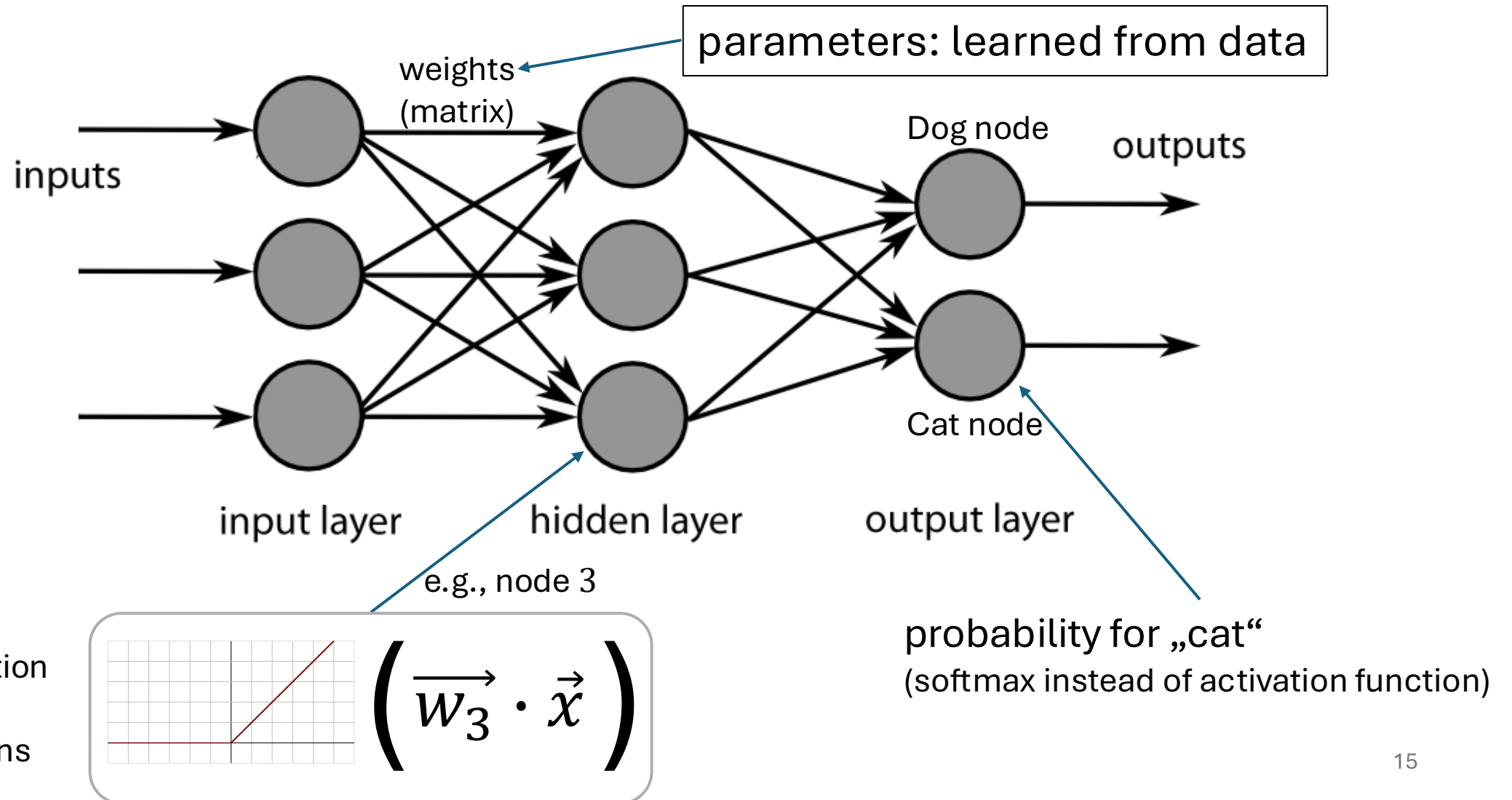


inspired from biological neurons:



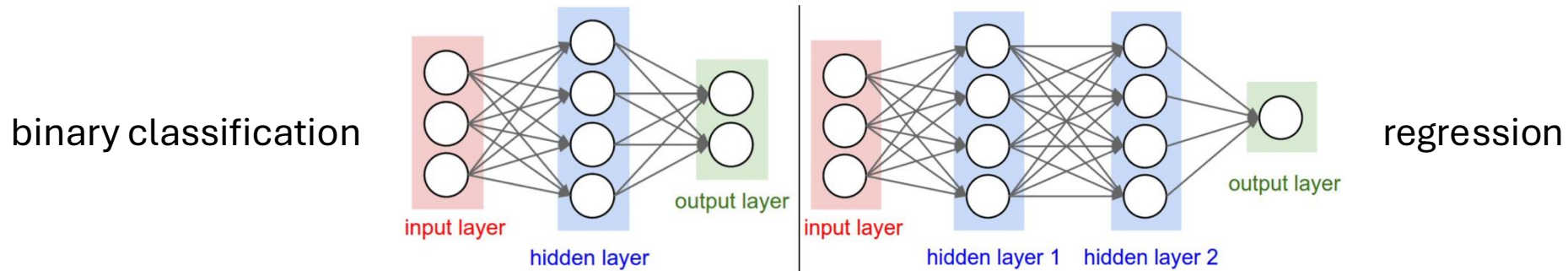
[source](#)

Neural Network as Nonlinear Function



Multi-Layer Perceptron (MLP)

fully-connected feed-forward network with at least one hidden layer
(universal function approximator)



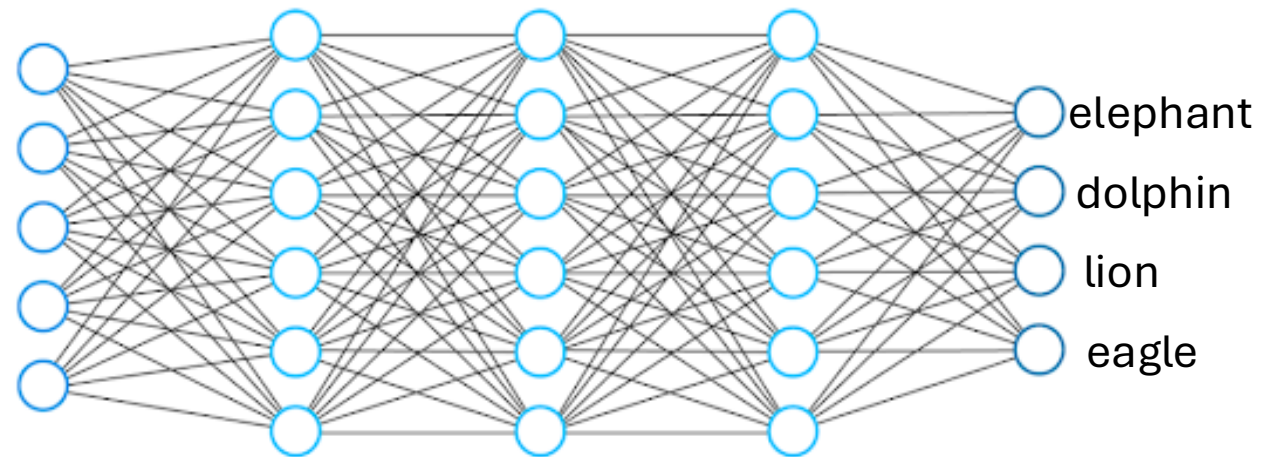
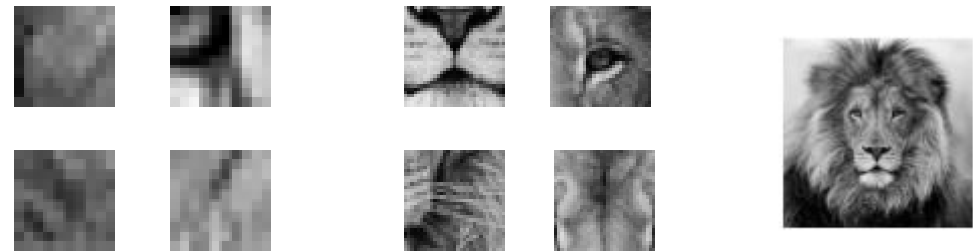
toward deep learning: add hidden layers

more layers (depth) more efficient than just more nodes (width):
less parameters needed for same function complexity

Hierarchical Representation

deep learning:
several hidden layers

increasing abstraction
→



Classification: Logistic Regression

binary classification: $y = 1$ or $y = 0$

→ predict probabilities p_i for $y = 1$

- logit (log-odds) as link function
- Y following Bernoulli distribution

$$\text{logit}(E[Y|\mathbf{X} = \mathbf{x}_i]) = \ln\left(\frac{p_i}{1 - p_i}\right) = \mathbf{w} \cdot \mathbf{x}_i + b$$

Multi-Classification: Softmax

for classification in K categories:

- use “score” function with K outputs

$$f_k(\mathbf{x}_i) = \mathbf{w}_k \cdot \mathbf{x}_i + b_k$$

- and convert into probabilities by means of softmax function

$$\frac{\exp(f_k(\mathbf{x}_i))}{\sum_{l=1}^K \exp(f_l(\mathbf{x}_i))}$$

Classification Networks

cross-entropy loss (one output node for each class k):

$$L(y, \hat{y}) = - \sum_{k=1}^K y_k \log \hat{y}_k$$

using softmax on output nodes:

$$\hat{y}_k = \frac{e^{o_k}}{\sum_{l=1}^K e^{o_l}}$$

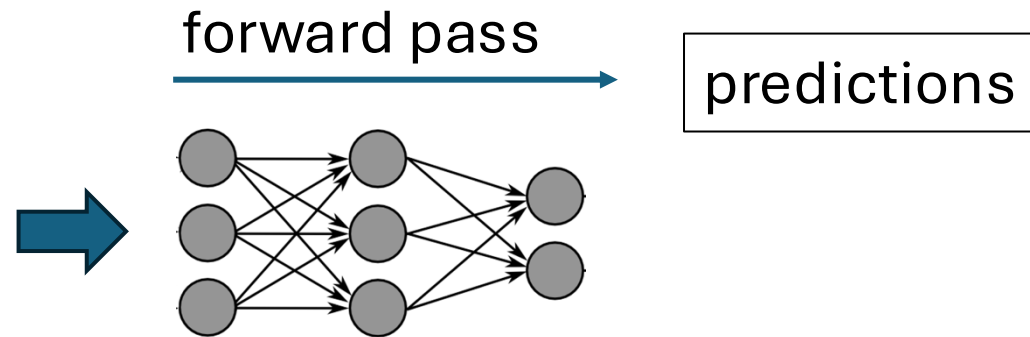
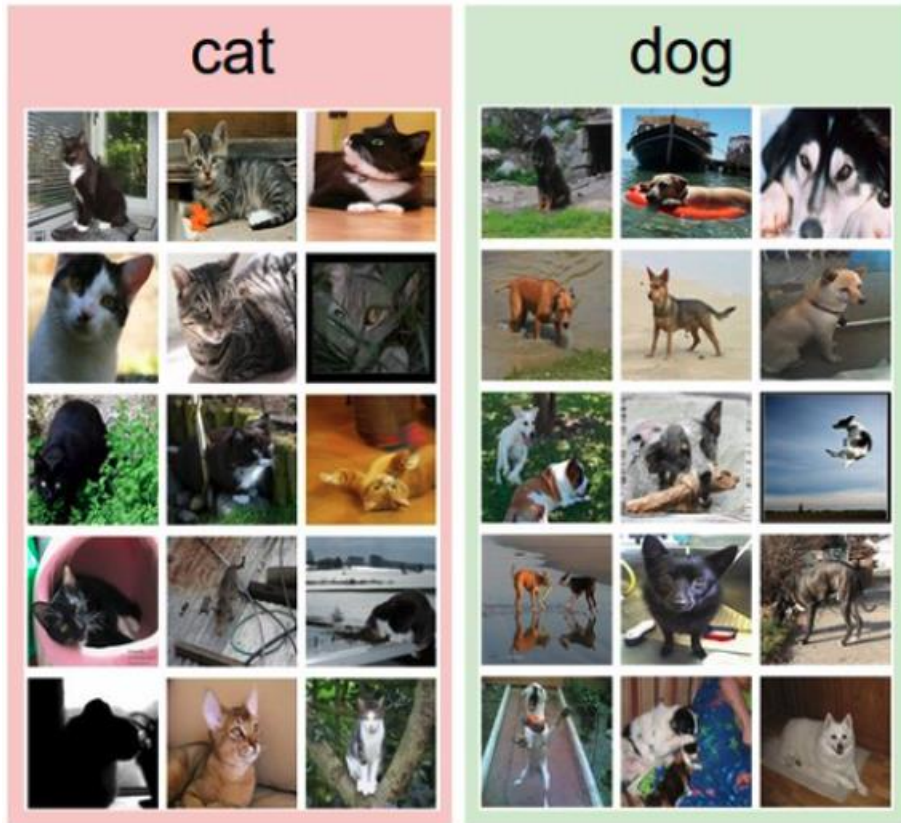
→ prediction of probabilities

regression networks:

- squared error loss:
 $L(y, \hat{y}) = (y - \hat{y})^2$
- just one output node
- no function on output
→ prediction of real numbers

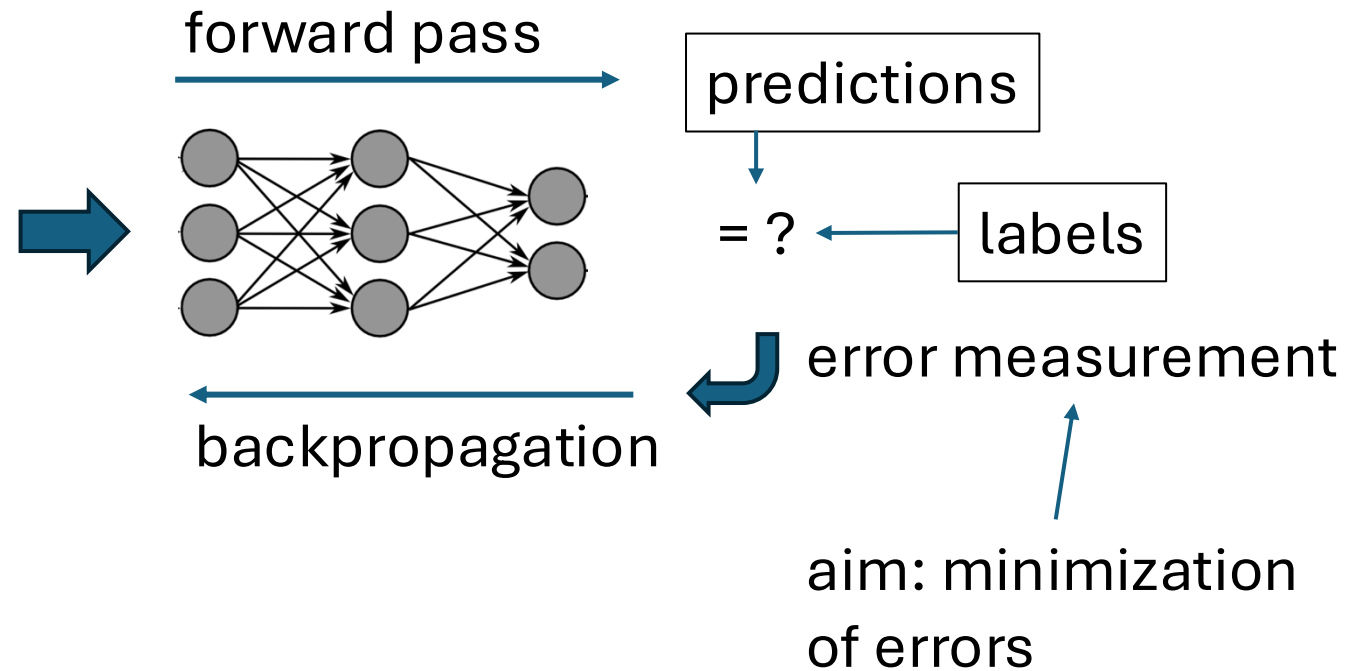
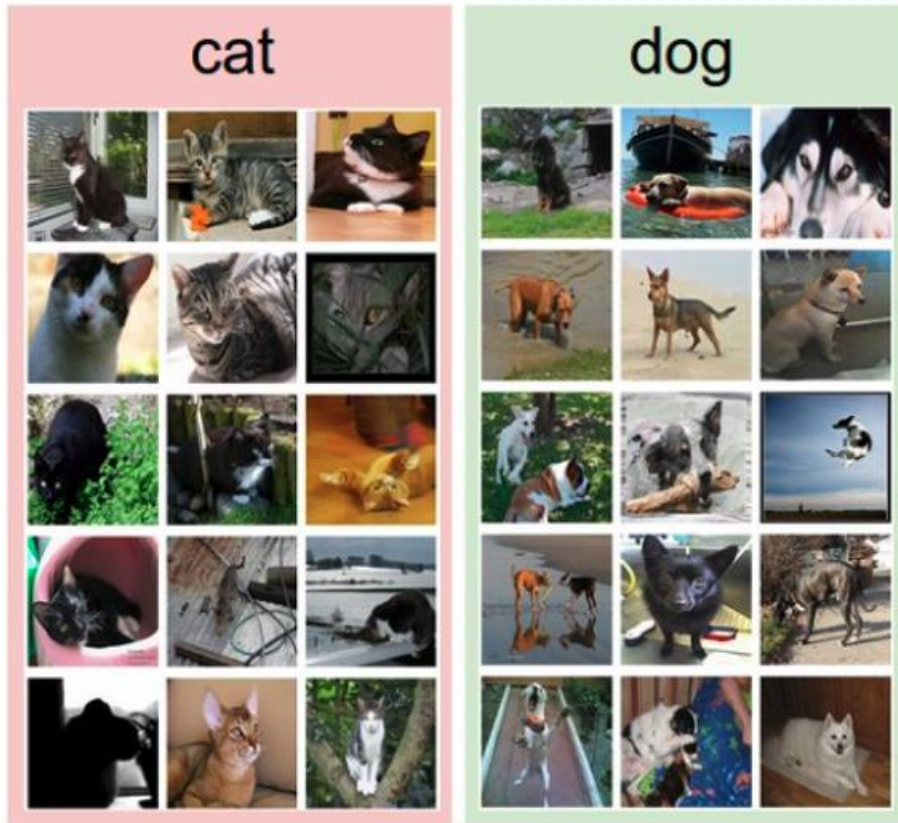
Training with Labeled Images

label:



Training with Labeled Images

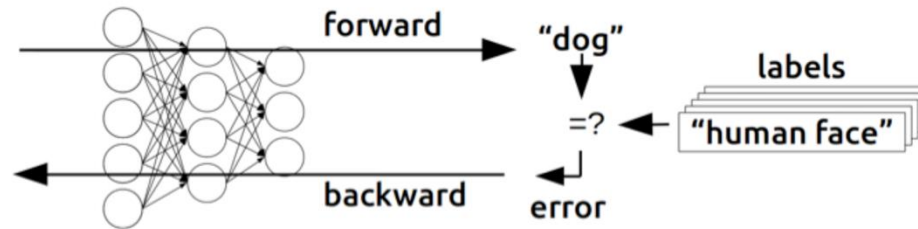
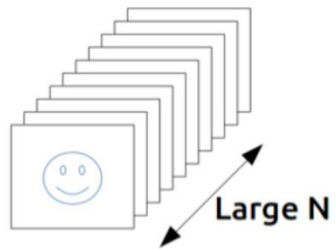
label:



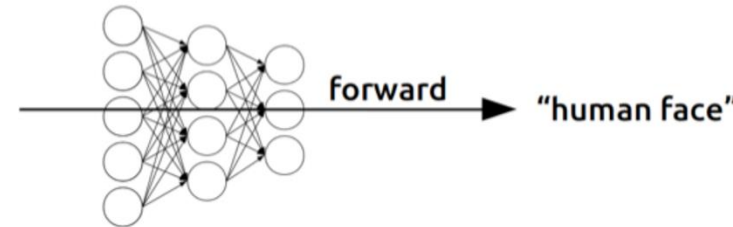
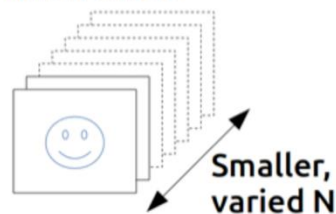
Find Gradients for (Deep) Neural Networks

backpropagation of errors through network layers via chain rule of calculus (enables learning of deep neural networks)

Training



Inference



[source](#)

forward pass:

- fix current weights
- compute predictions

backward pass:

- compute errors
- calculate gradients (backpropagation of errors)
- update weights accordingly (gradient descent)

Training: Iterative Adjustment of Weights

```
for epoch in range(epochs):  
    for X, y in mini_batches:  
        pred = model(X)  
        cost = loss(pred, y)  
        cost.backward()  
        optimizer.step()
```

quantification of error
(here: cross entropy)

chain rule to calculate
gradient (direction) of
cost function with
respect to weights

one step in direction of steepest
descent: (stochastic) gradient descent

Example WOLOG

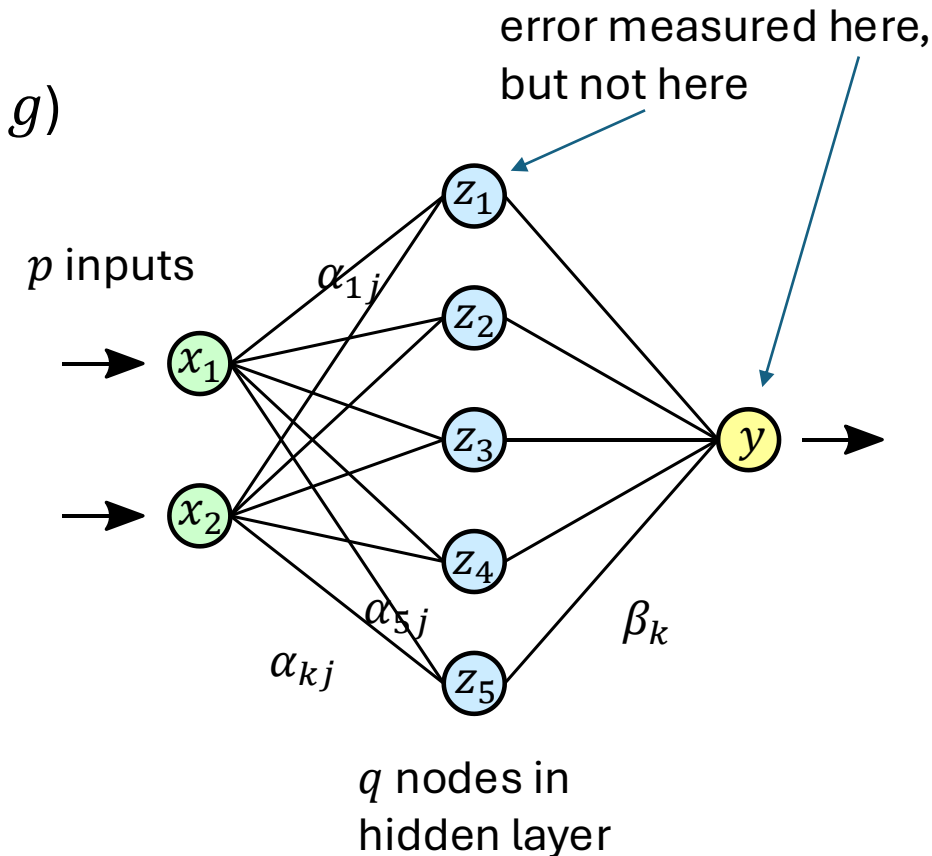
- regression (squared error loss, identity output function g)
- with one hidden layer ($\hat{\mathbf{w}}: \hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\beta}}$)

$$\hat{y}_i = \hat{f}(\mathbf{x}_i; \hat{\mathbf{w}}) = g(\mathbf{z}_i; \hat{\boldsymbol{\beta}}) = \sum_{k=0}^q \hat{\beta}_k z_{ik}$$
$$z_{ik} = h(\mathbf{x}_i; \hat{\boldsymbol{\alpha}}_k) = h\left(\sum_{j=0}^p \hat{\alpha}_{kj} x_{ij}\right)$$

activation function

cost function:

$$J(\hat{\mathbf{w}}) = \sum_{i=1}^n L_i(y_i, \hat{f}(\mathbf{x}_i); \hat{\mathbf{w}}) = \sum_{i=1}^n \left(y_i - \hat{f}(\mathbf{x}_i; \hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\beta}})\right)^2$$



Example WOLOG

gradients:

$$\frac{\partial L_i}{\partial \hat{\beta}_k} = -2 \left(y_i - \hat{f}(\mathbf{x}_i; \hat{\mathbf{w}}) \right) z_{ik} = \delta_i z_{ik}$$
$$\frac{\partial L_i}{\partial \hat{\alpha}_{kj}} = -2 \left(y_i - \hat{f}(\mathbf{x}_i; \hat{\mathbf{w}}) \right) \hat{\beta}_k h'_k \left(\sum_{j=0}^p \hat{\alpha}_{kj} x_{ij} \right) x_{ij} = s_{ik} x_{ij}$$

backpropagation equations: $s_{ik} = h'_k \left(\sum_{j=0}^p \hat{\alpha}_{kj} x_{ij} \right) \hat{\beta}_k \delta_i$

use errors of later layers to calculate errors of earlier ones (avoiding redundant calculations of intermediate terms)

Using Gradients for Iterative Learning

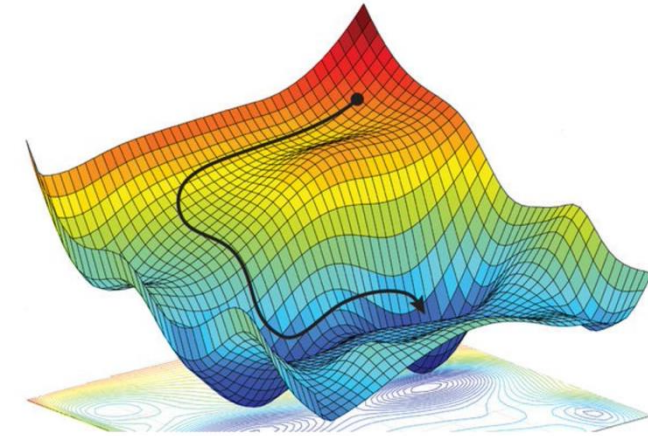
use gradients found via backpropagation to iteratively update weights (gradient descent):

$$\hat{\beta}_k^{(r+1)} = \hat{\beta}_k^{(r)} - \eta_r \sum_{i=1}^n \frac{\partial L_i}{\partial \hat{\beta}_k^{(r)}}$$

$$\hat{\alpha}_{kj}^{(r+1)} = \hat{\alpha}_{kj}^{(r)} - \eta_r \sum_{i=1}^n \frac{\partial L_i}{\partial \hat{\alpha}_{kj}^{(r)}}$$

- adaptive learning rate η_r : learning rate often adjusted per iteration
- weight initialization: choose small random weights as starting values to break symmetry (typically using some heuristics)

Stochastic Gradient Descent (SGD)



weight updates: $\hat{\mathbf{w}} \leftarrow \hat{\mathbf{w}} - \eta \nabla_{\hat{\mathbf{w}}} J(\hat{\mathbf{w}})$

step size η cost function $J(\hat{\mathbf{w}})$

options:

- $J(\hat{\mathbf{w}}) = \frac{1}{n} \sum_{i=1}^n J_i(\hat{\mathbf{w}})$ batch gradient descent (whole training dataset)
- $J(\hat{\mathbf{w}}) = J_i(\hat{\mathbf{w}})$ stochastic gradient descent (single example)
- $J(\hat{\mathbf{w}}) = \frac{1}{m} \sum_{i=1}^m J_i(\hat{\mathbf{w}})$ mini-batch stochastic gradient descent

mini-batch size \rightarrow hyperparameter: tradeoffs between runtime & memory, convergence & generalization

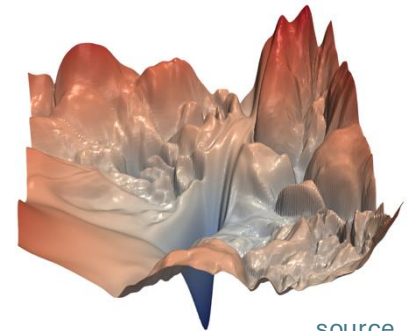
random fluctuations of SGD help to escape saddle points

How to Train Deep Neural Networks Effectively?

optimization and regularization difficult

- local vs global minima, saddle points, easily overfitting
- many hyperparameters to tune

typical loss surface:



[source](#)

but many methods to get it working in practice (despite partly patchy theoretical understanding)

- optimization: activation and loss functions, weight initialization, adaptive learning rate (e.g., Adam), learning rate schedulers
- explicit regularization: weight decay, dropout, data augmentation
- implicit regularization: early stopping, batch/layer normalization, SGD

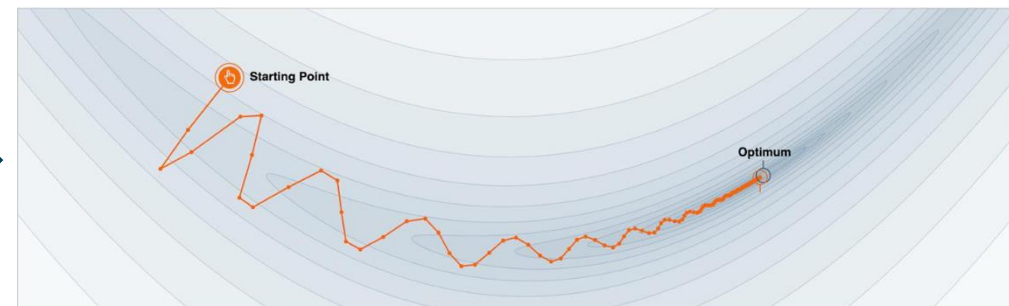
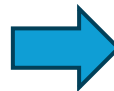
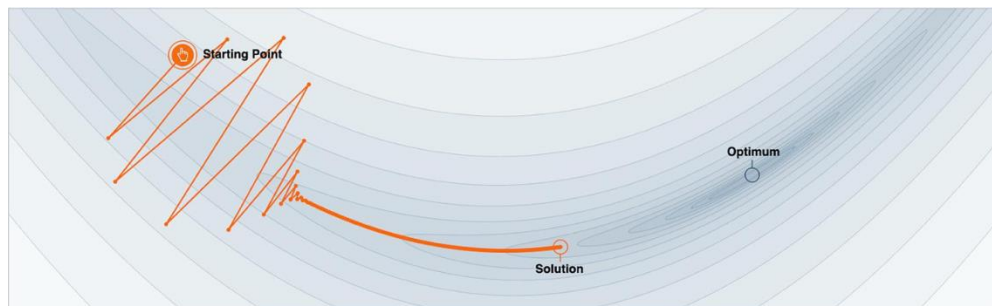
Gradient Descent with Momentum

algorithm:

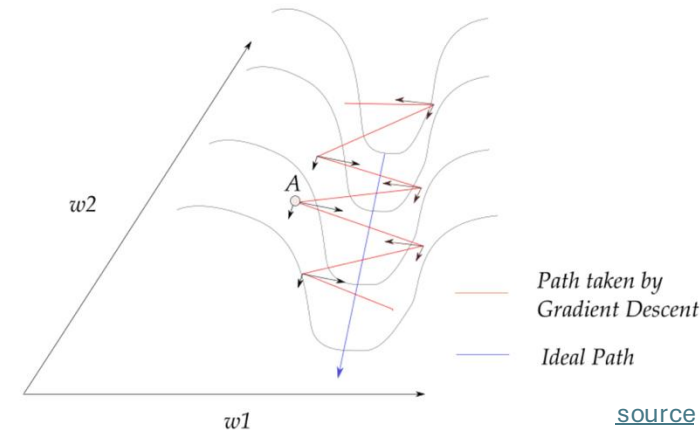
- estimate gradient: $\mathbf{g} = \nabla_{\hat{\mathbf{w}}} J(\hat{\mathbf{w}})$
- compute velocity update: $\mathbf{v} \leftarrow \alpha \mathbf{v} - \eta \mathbf{g}$
- apply parameter update: $\hat{\boldsymbol{\theta}} \leftarrow \hat{\boldsymbol{\theta}} + \mathbf{v}$

hyperparameter ($0 < \alpha < 1$)
specifying exponential decay

→ no bouncing around of parameters, escape from local minima and saddle points



Adaptive Learning Rate



better convergence by adapting learning rate for each weight per iteration:
lower/higher learning rates for weights with large/small updates

popular methods ($g_{\hat{w}}$ denotes component of gradient for individual weight \hat{w}):

- Adagrad: $\hat{w} \leftarrow \hat{w} - \frac{\eta}{\sqrt{\sum_{\tau=1}^t g_{\hat{w},\tau}^2}} g_{\hat{w}}$ with t, τ denoting current and past iterations
(issue: sum in denominator grows with more iterations \rightarrow can get stuck)
- RMSProp: $\hat{w} \leftarrow \hat{w} - \frac{\eta}{\sqrt{v(\hat{w})}} g_{\hat{w}}$ with $v(\hat{w}) \leftarrow \gamma v(\hat{w}) + (1 - \gamma) g_{\hat{w}}^2$
- Adam (Adaptive Moment Optimization): combines RMSProp with momentum

Learning Rate Schedulers

learning rate η in gradient descent typically set to small value (e.g., 0.001)

use momentum and adaptive learning rates to improve optimization

learning rate scheduling as further alternative (can be used on top):

- decay by a certain factor every given number of epochs
- reduce when a metric has stopped improving
- smooth decay by cosine annealing (optionally cyclic with restarts)
- ...

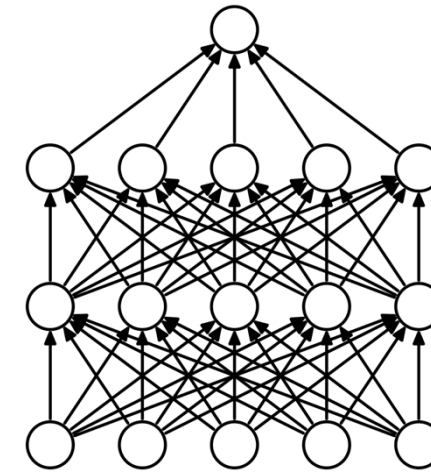
Dropout in Neural Networks

goal: prevent overfitting of large neural networks

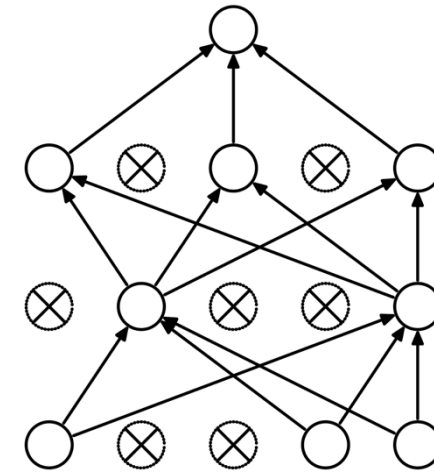
randomly drop non-output nodes (along with their connections) during training (not prediction)

→ adaptability: regularizing each hidden node to perform well regardless of which other hidden nodes are in the model

- for each mini-batch, randomly sample independent binary masks for the nodes
- destroying extracted features rather than input values



(a) Standard Neural Net



(b) After applying dropout.

[source](#)

dropout for 2D structures (such as images) usually drops entire channels instead of individual nodes (because locality nullifies the effect of standard dropout)

Batch Normalization

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

Parameters to be learned: γ, β

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

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

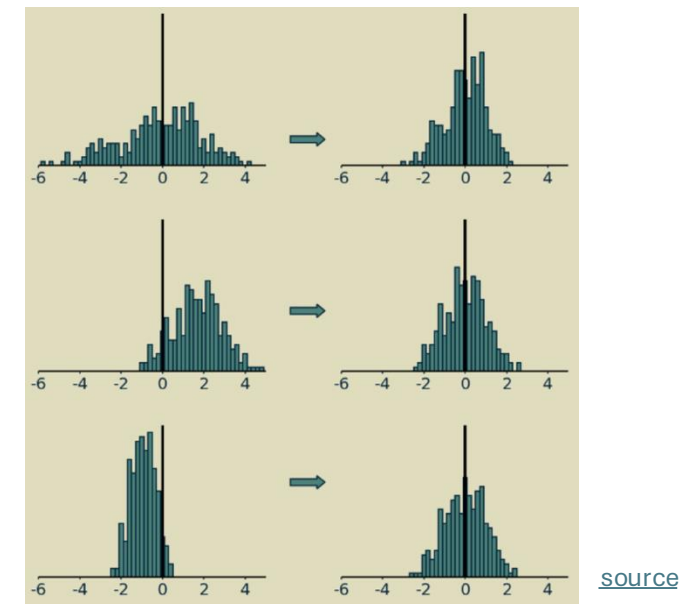
$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

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

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

Algorithm 1: Batch Normalizing Transform, applied to activation x over a mini-batch.

[source](#)



[source](#)

adaptive reparameterization of inputs to network layer (before or after activation)
independently for each input/feature
(not to confuse with weight normalization)

to maintain expressive power (optional):
introduce parameters γ, β (learned together with weights via backpropagation)

Benefits from Batch Normalization

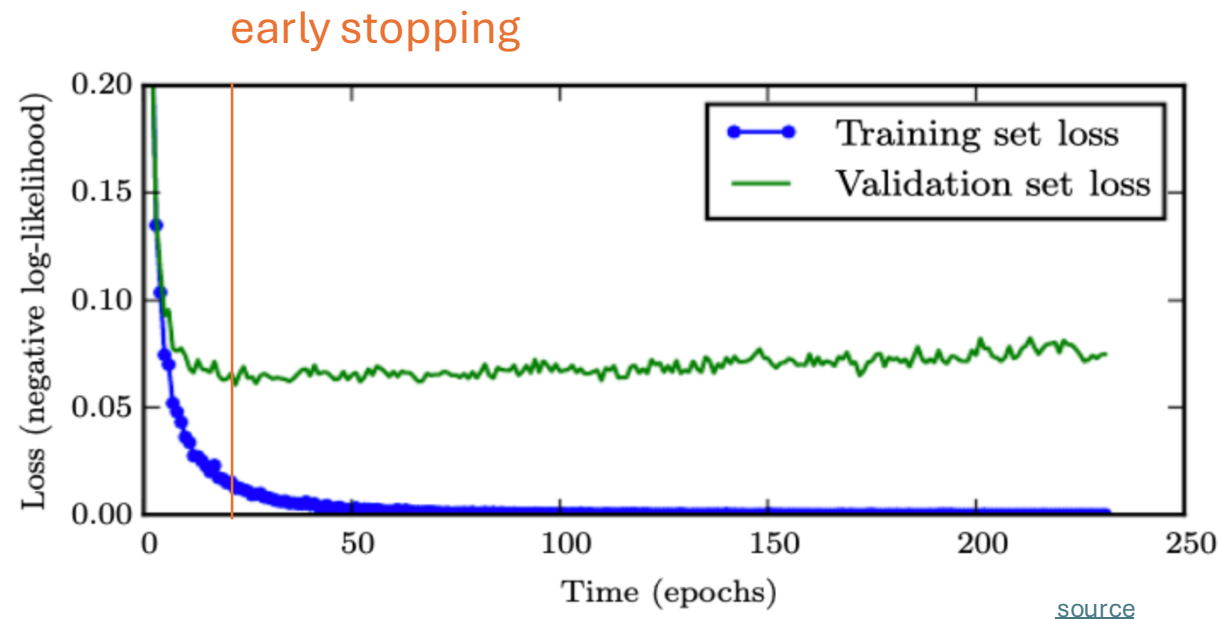
- allows higher learning rates
- reduces importance of weight initialization
- alleviates vanishing/exploding gradients
- (implicit) regularization effect: introducing noise

reason why batch normalization improves optimization still controversial

most plausible explanation: smoothening of loss landscape

Early Stopping

loss independently measured on validation set
halting training when overfitting begins to occur



Tabular vs Unstructured Data

deep learning methods dominate applications on unstructured data (like text or images), but not necessarily on tabular data

typical characteristics of tabular data difficult to handle for deep learning:

- irregular patterns in target function (neural networks require piecewise continuous targets)
- uninformative features
- non-rotationally invariant data (linear combinations of features misrepresent the information)

tree-based models (e.g., gradient boosting) can naturally deal with these situations