

supervised  
learning

ADULT  
SUPERVISION  
REQUIRED

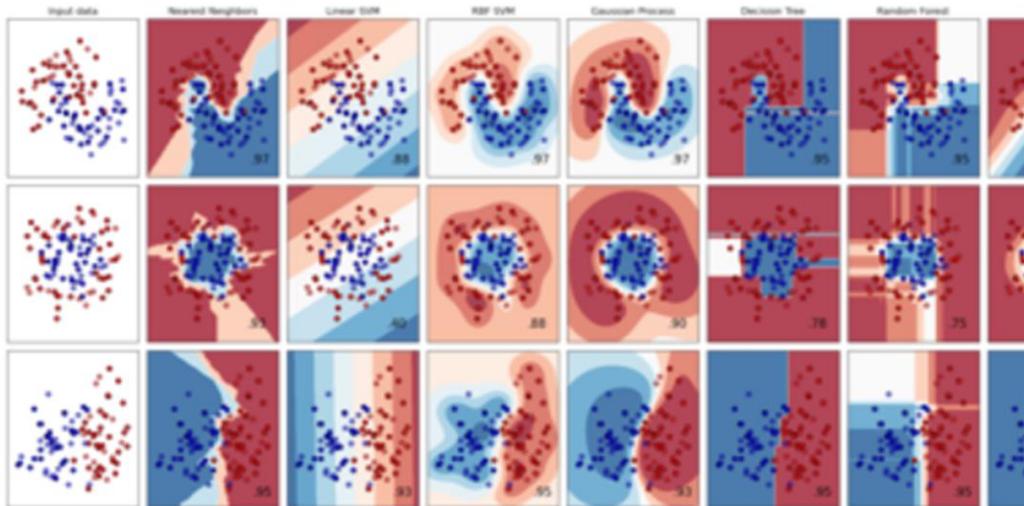
## Classification

any case where  
you know some truth

Identifying which category an object belongs to.

**Applications:** Spam detection, image recognition.

**Algorithms:** Gradient boosting, nearest neighbors, random forest, logistic regression, and more...

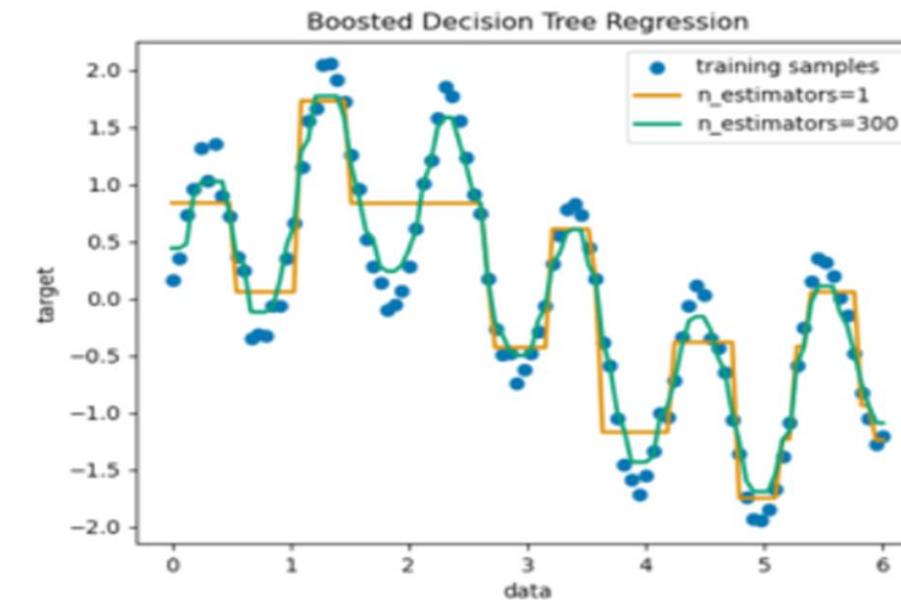


## Regression

Predicting a continuous-valued attribute associated with an object.

**Applications:** Drug response, Stock prices.

**Algorithms:** Gradient boosting, nearest neighbors, random forest, ridge, and more...



carefully curated  
your data

images, spectra,  
measurements,  
whole simulations,  
anything that you  
measured  
and trust dearly

labels

galaxy images  
(or measurements of)  
mass, type, redshift

sims with cosmologies  $X_i$  and  $X_{i+1}$

sim with cosmology  $X_{i+\frac{1}{2}}$

✓ lensed objects  
reconstruction

pairs of similar images  
similarity score

carefully curated  
*your data*

images, spectra,  
measurements,  
whole simulations,  
anything that you  
measured  
and trust dearly

labels

galaxy images  
(or measurements of)  
mass, type, redshift

sims with cosmologies  $X_i$  and  $X_{i+1}$

sim with cosmology  $X_{i+\frac{1}{2}}$

✓ lensed objects  
reconstruction

pairs of similar images  
similarity score

carefully curated  
your data

images, spectra,  
measurements,  
whole simulations,  
anything that you  
measured  
and trust dearly



labels

well-understood  
new data

very similar data  
to your curated  
dataset, except  
you don't have the  
labels!

*carefully curated*  
**your data**

images, spectra,  
measurements,  
whole simulations,  
anything that you  
measured  
and trust dearly

**labels**

*well-understood*  
**new data**

very similar data  
to your curated  
dataset, except  
you don't have the  
labels!

**predicted  
labels**

**machine learning** →

# *supervised learning*

## **regression**

continuous variable

more difficult  
(more room for error)

traces spectrum/evolution

like traditional science

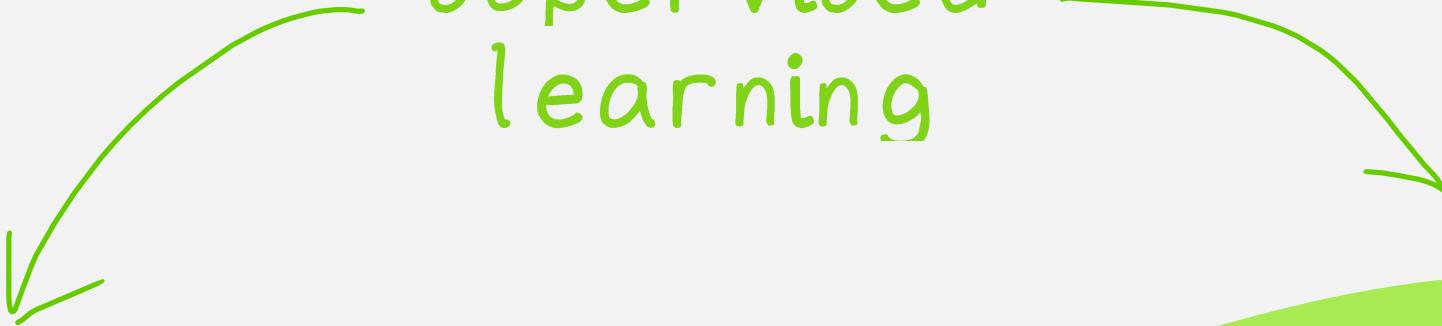
## **classification**

discrete variable

easier to optimize

finds distinct categories/cases

similar to how humans operate?



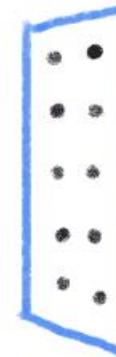
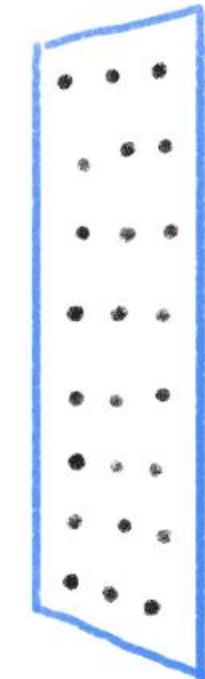
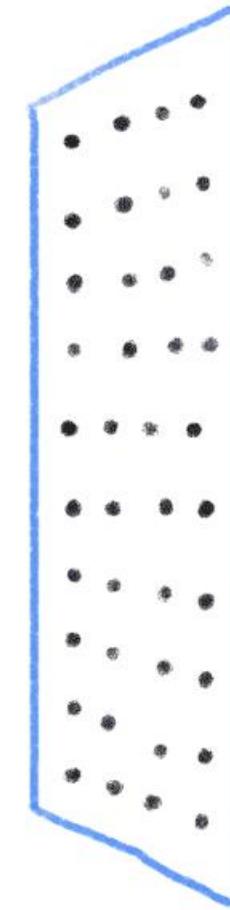
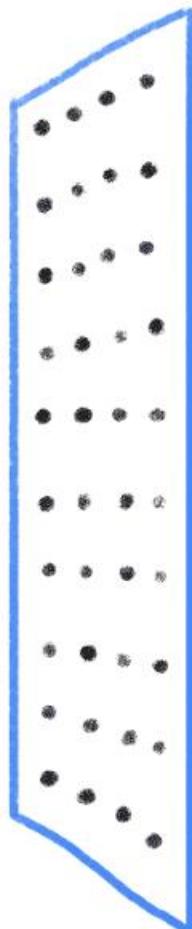
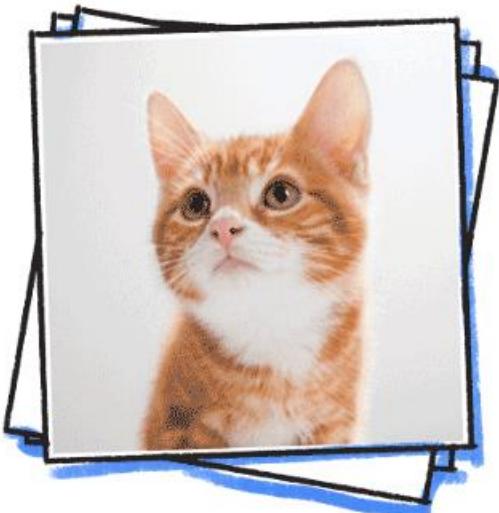
# ***Euclid preparation. Measuring detailed galaxy morphologies for Euclid with Machine Learning***

Euclid Collaboration: B. Aussel<sup>\*1</sup>, S. Kruk<sup>2</sup>, M. Walmsley<sup>3</sup>, M. Huertas-Company<sup>4,5,6,7</sup>, M. Castellano<sup>8</sup>, C. J. Conselice<sup>3</sup>, M. Delli Veneri<sup>9</sup>, H. Domínguez Sánchez<sup>10</sup>, P.-A. Duc<sup>11</sup>, U. Kuchner<sup>12</sup>, A. La Marca<sup>13,14</sup>, B. Margalef-Bentabol<sup>13</sup>, F. R. Marleau<sup>15</sup>, G. Stevens<sup>16</sup>, Y. Toba<sup>17</sup>, C. Tortora<sup>18</sup>, L. Wang<sup>13,14</sup>, N. Aghanim<sup>19</sup>, B. Altieri<sup>2</sup>, A. Amara<sup>20</sup>, S. Andreon<sup>21</sup>, N. Auricchio<sup>22</sup>, M. Baldi<sup>23,22,24</sup>, S. Bardelli<sup>22</sup>, R. Bender<sup>25,26</sup>, C. Bodendorf<sup>25</sup>, D. Bonino<sup>27</sup>, E. Branchini<sup>28,29,21</sup>, M. Brescia<sup>30,18,9</sup>, J. Brinchmann<sup>31</sup>, S. Camera<sup>32,33,27</sup>, V. Capobianco<sup>27</sup>, C. Carbone<sup>34</sup>, J. Carretero<sup>35,36</sup>, S. Casas<sup>37</sup>, S. Cavuoti<sup>18,9</sup>, A. Cimatti<sup>38</sup>, G. Congedo<sup>39</sup>, L. Conversi<sup>40,2</sup>, Y. Copin<sup>41</sup>, F. Courbin<sup>42</sup>, H. M. Courtois<sup>43</sup>, M. Cropper<sup>44</sup>, A. Da Silva<sup>45,46</sup>, H. Degaudenzi<sup>47</sup>, A. M. Di Giorgio<sup>48</sup>, J. Dinis<sup>46,45</sup>, F. Dubath<sup>47</sup>, X. Dupac<sup>2</sup>, S. Dusini<sup>49</sup>, M. Farina<sup>48</sup>, S. Farrens<sup>50</sup>, S. Ferriol<sup>41</sup>, S. Fotopoulou<sup>51</sup>, M. Frailis<sup>52</sup>, E. Franceschi<sup>22</sup>, P. Franzetti<sup>34</sup>, M. Fumana<sup>34</sup>, S. Galeotta<sup>52</sup>, B. Garilli<sup>34</sup>, B. Gillis<sup>39</sup>, C. Giocoli<sup>22,53</sup>, A. Grazian<sup>54</sup>, F. Grupp<sup>25,55</sup>, S. V. H. Haugan<sup>56</sup>, W. Holmes<sup>57</sup>, I. Hook<sup>58</sup>, F. Hormuth<sup>59</sup>, A. Hornstrup<sup>60,61</sup>, P. Hudelot<sup>62</sup>, K. Jahnke<sup>63</sup>,

CAT

(LABLED  
PHOTOS)

DOG

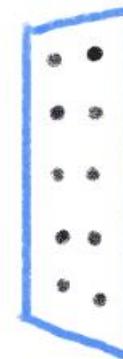
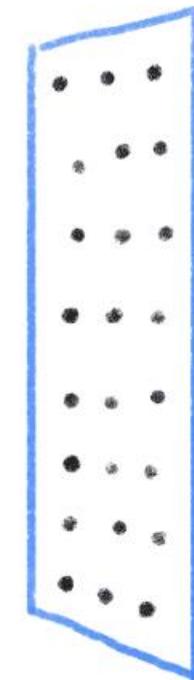
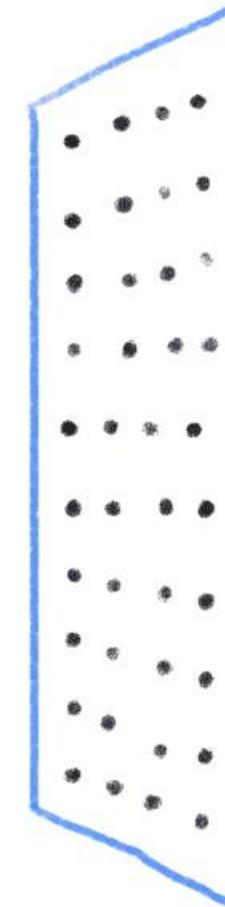
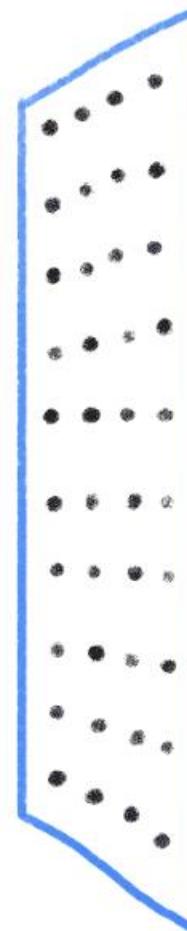
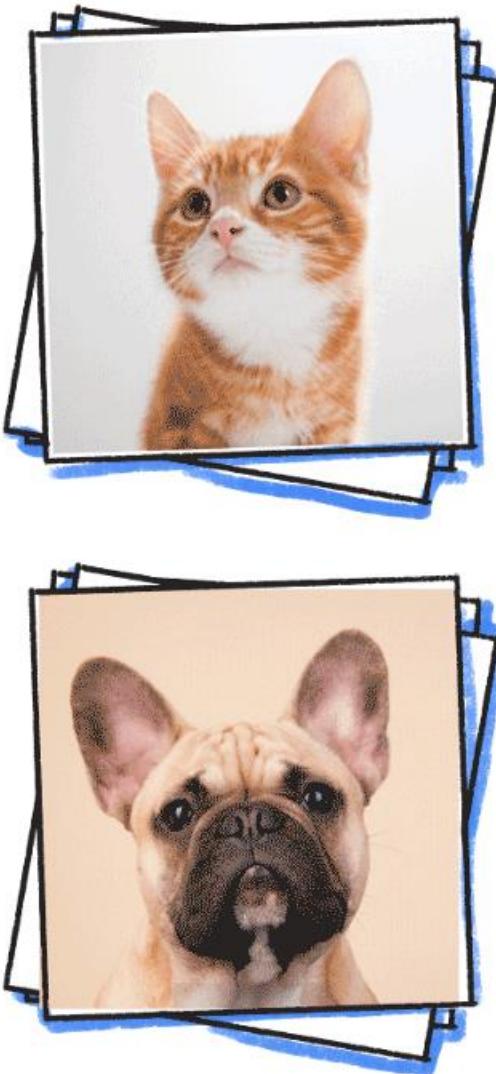


OUTPUT

CAT

(LABLED  
PHOTOS)

DOG



output  $\in (0,1)$

OUTPUT

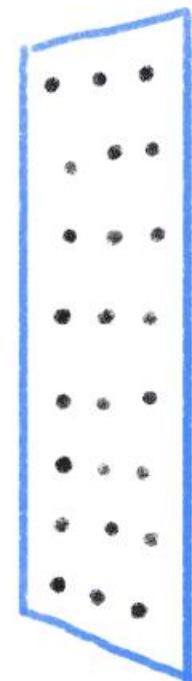
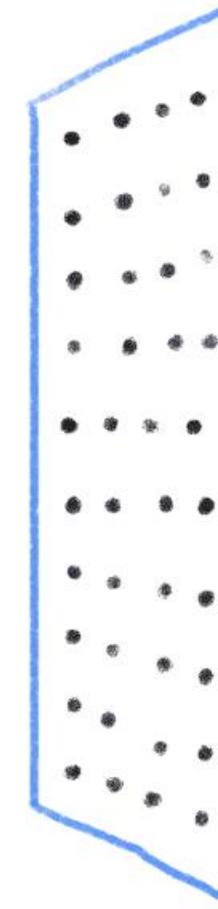
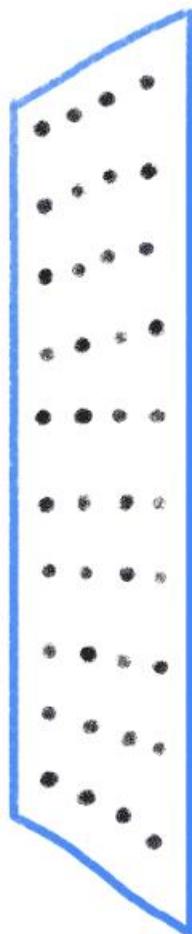
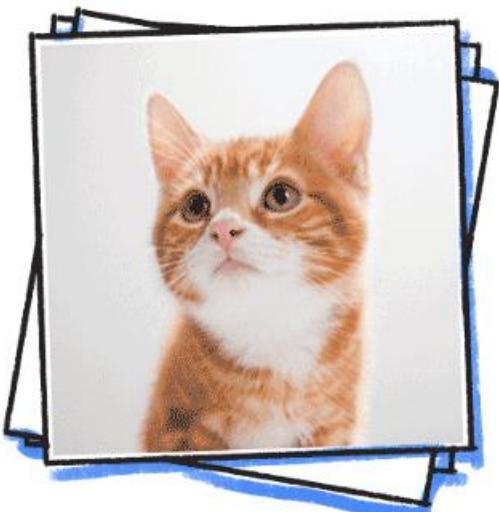
classification

Target label: 0 or 1

CAT

(Labeled  
Photos)

DOG

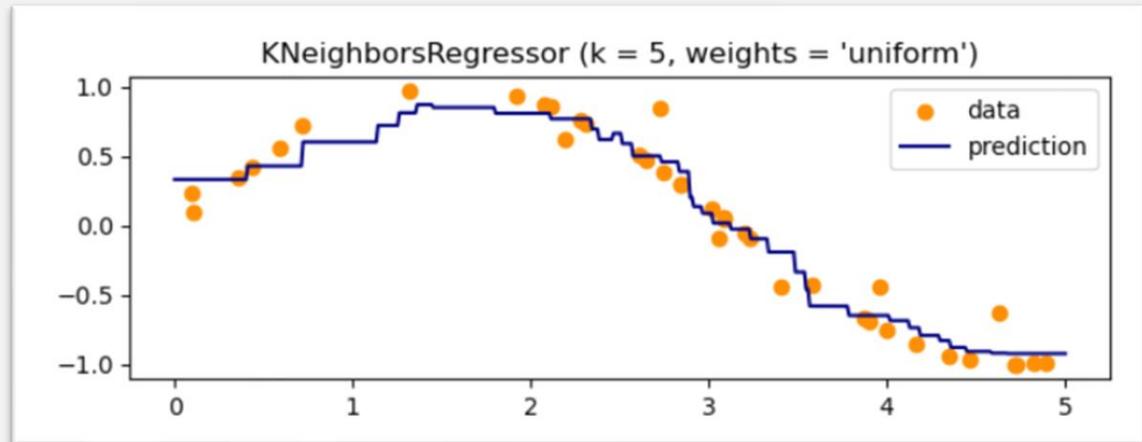
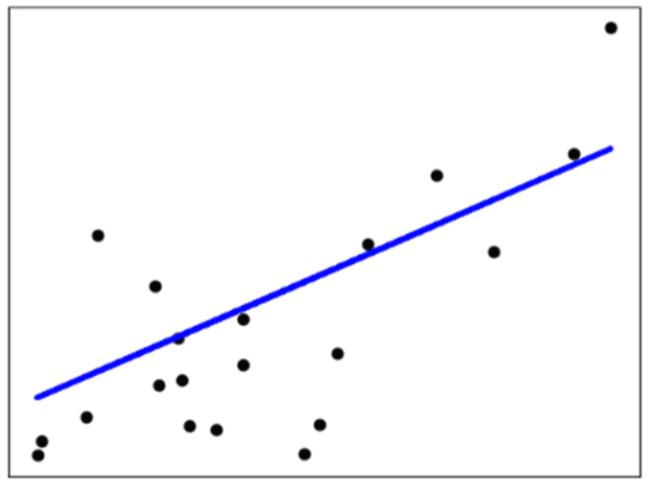


OUTPUT

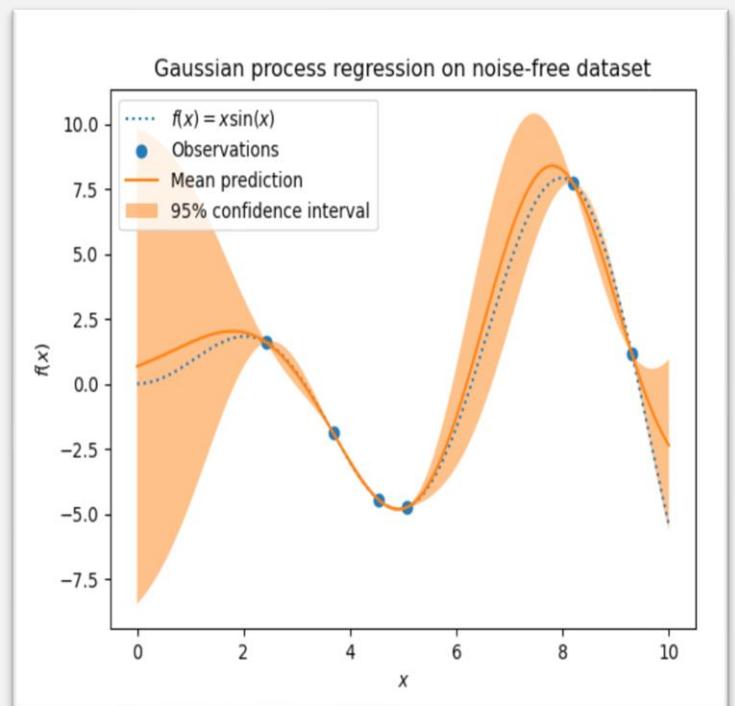
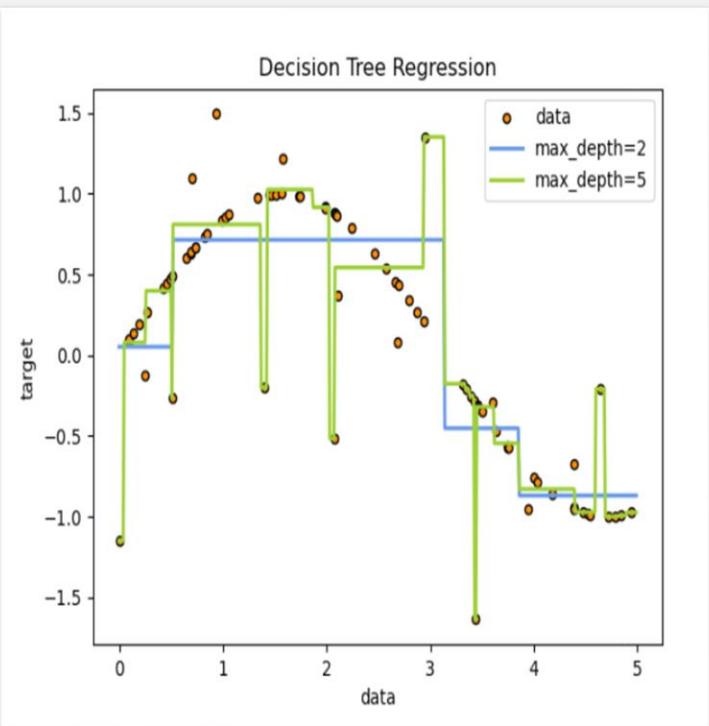
Target label:  $p(\text{cat}) \in (0,1)$

regression

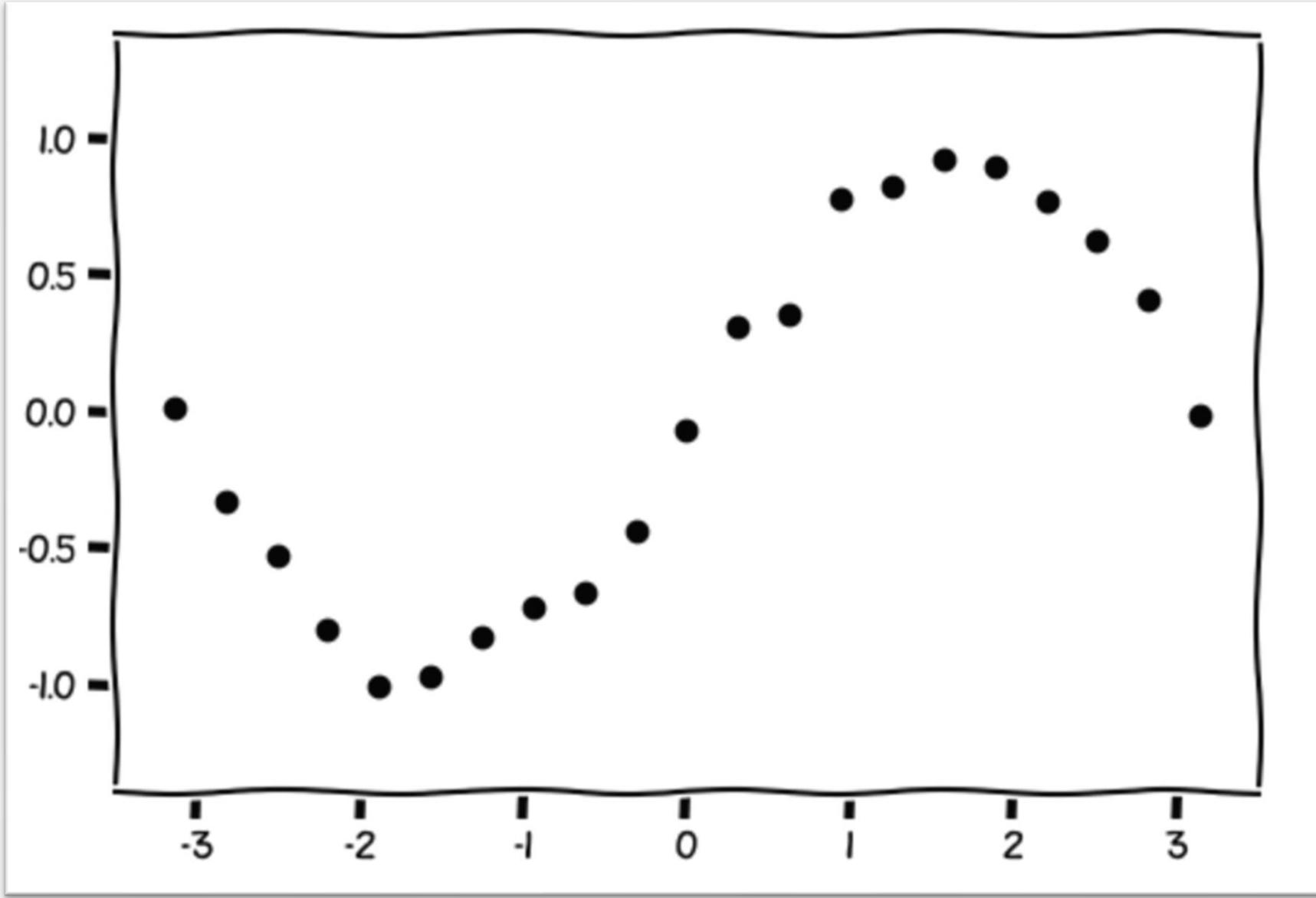
regression

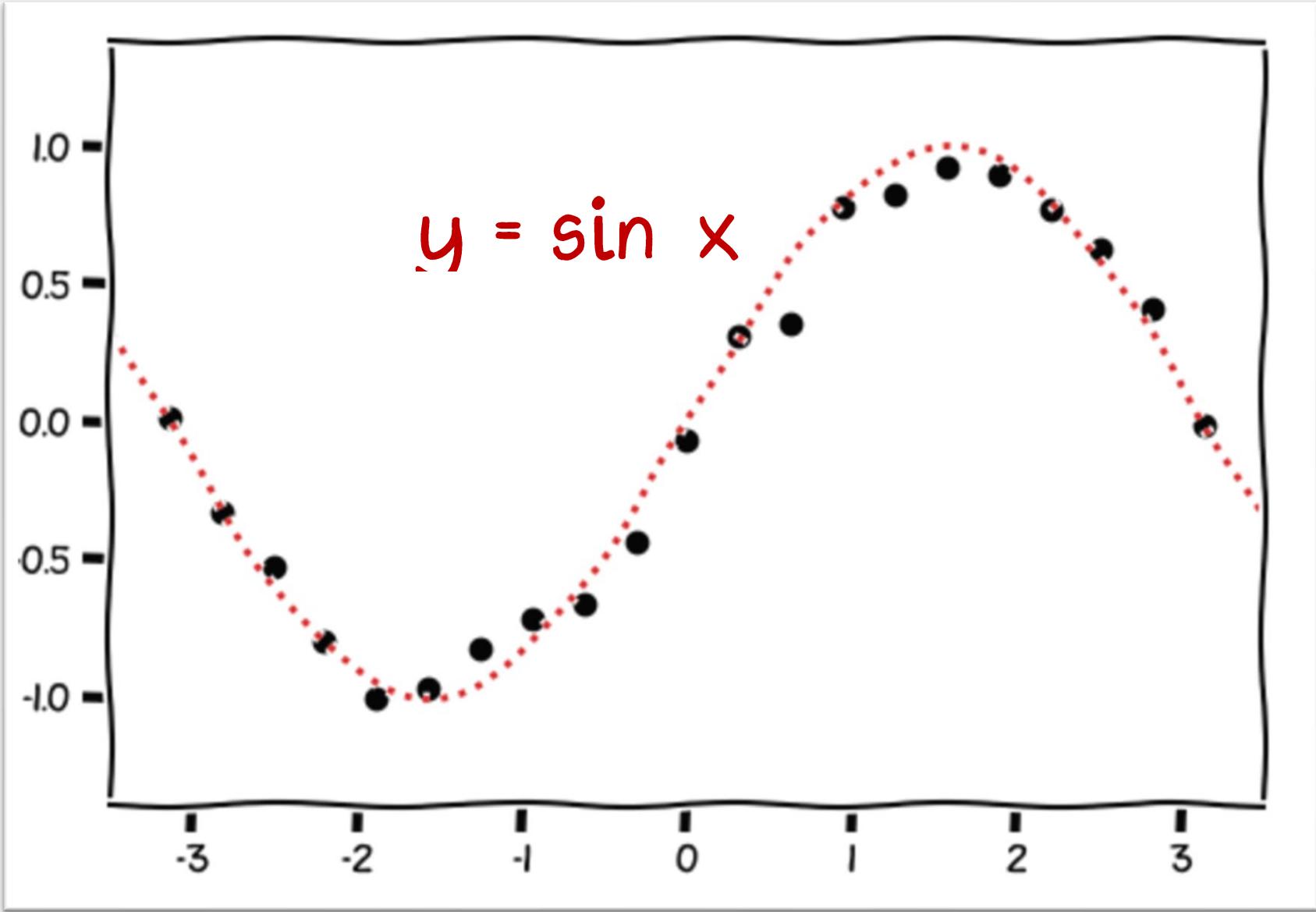


regression  
given  $\vec{x}$  predict  $\vec{y} = f(\vec{x})$



linear  
regression





# linear regression

$$\vec{y} = \alpha_1 f_1(\vec{x}) + \alpha_2 f_2(\vec{x}) + \cdots + \alpha_n f_n(\vec{x})$$

# linear regression

$$\vec{y} = \alpha_1 f_1(\vec{x}) + \alpha_2 f_2(\vec{x}) + \cdots + \alpha_n f_n(\vec{x})$$

you can do whatever you want to your dependent variables,  
but the coefficients must be linear

$$y = a x_1^b$$

$$y = e^{a+bx}$$

$$y = e^a x$$

$$y = ax + b$$

$$ax^2 + by^2 = 1$$

$$y = (ax_1 + bx_2)^2$$

$$y = a \sin x_1^2 + \frac{b}{x_2}$$

$$y = ax + a^2 x^2$$

$$y = \sin ax$$

- **Linearity.** This means that the mean of the response variable is a [linear combination](#) of the parameters (regression coefficients) and the predictor variables. Note that this assumption is much less restrictive than it may at first seem. Because the predictor variables are treated as fixed values (see above), linearity is really only a restriction on the parameters. The predictor variables themselves can be arbitrarily transformed, and in fact multiple copies of the same underlying predictor variable can be added, each one transformed differently. This technique is used, for example, in [polynomial regression](#), which uses linear regression to fit the response variable as an arbitrary [polynomial](#) function (up to a given degree) of a predictor variable. With this much flexibility, models such as polynomial regression often have "too much power", in that they tend to [overfit](#) the data. As a result, some kind of

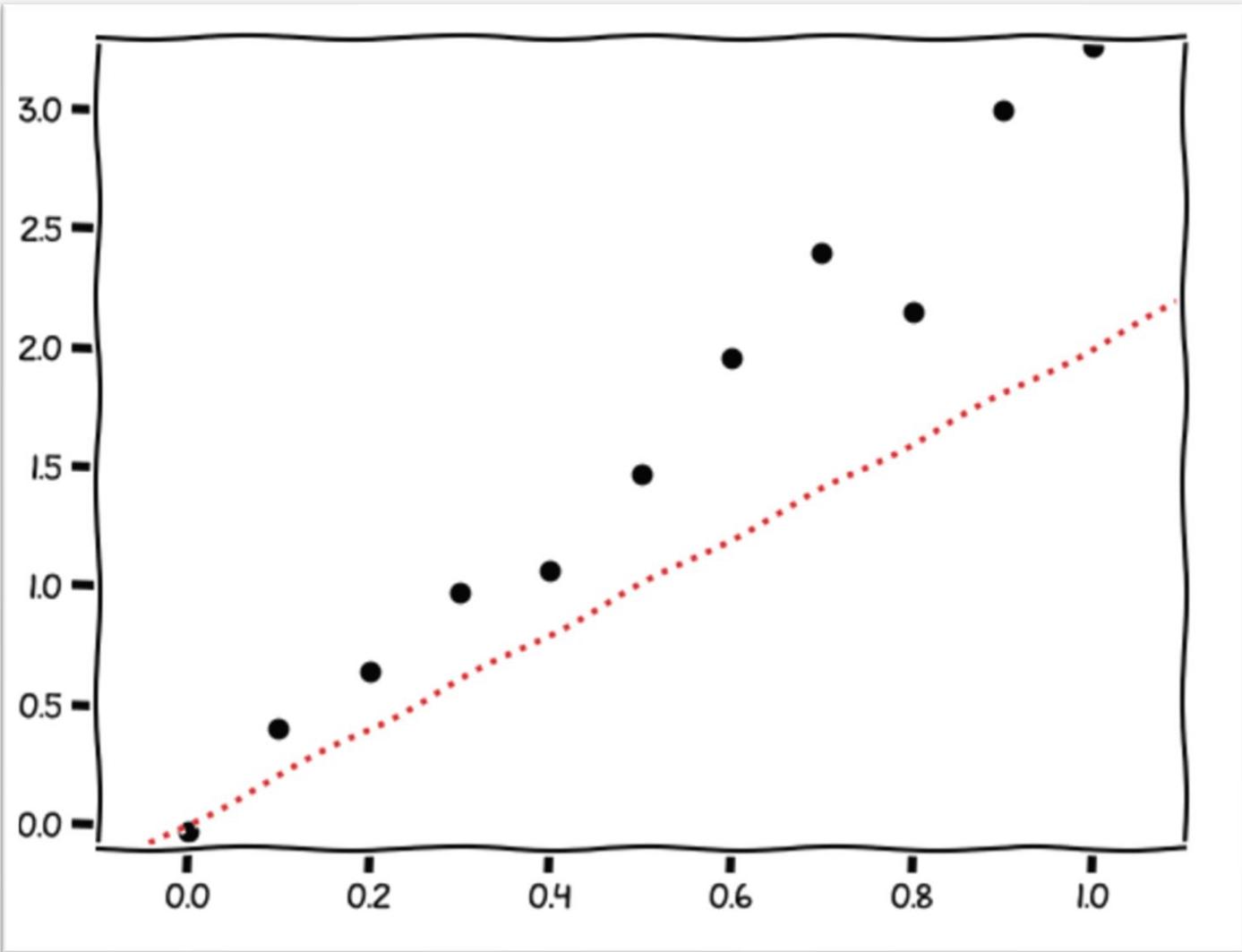
# linear regression

fitting a line  $y = ax$

1. Try some guess  $a_0$
2. Calculate loss

RMSE = 0.5

```
from sklearn.metrics import mean_squared_error
```



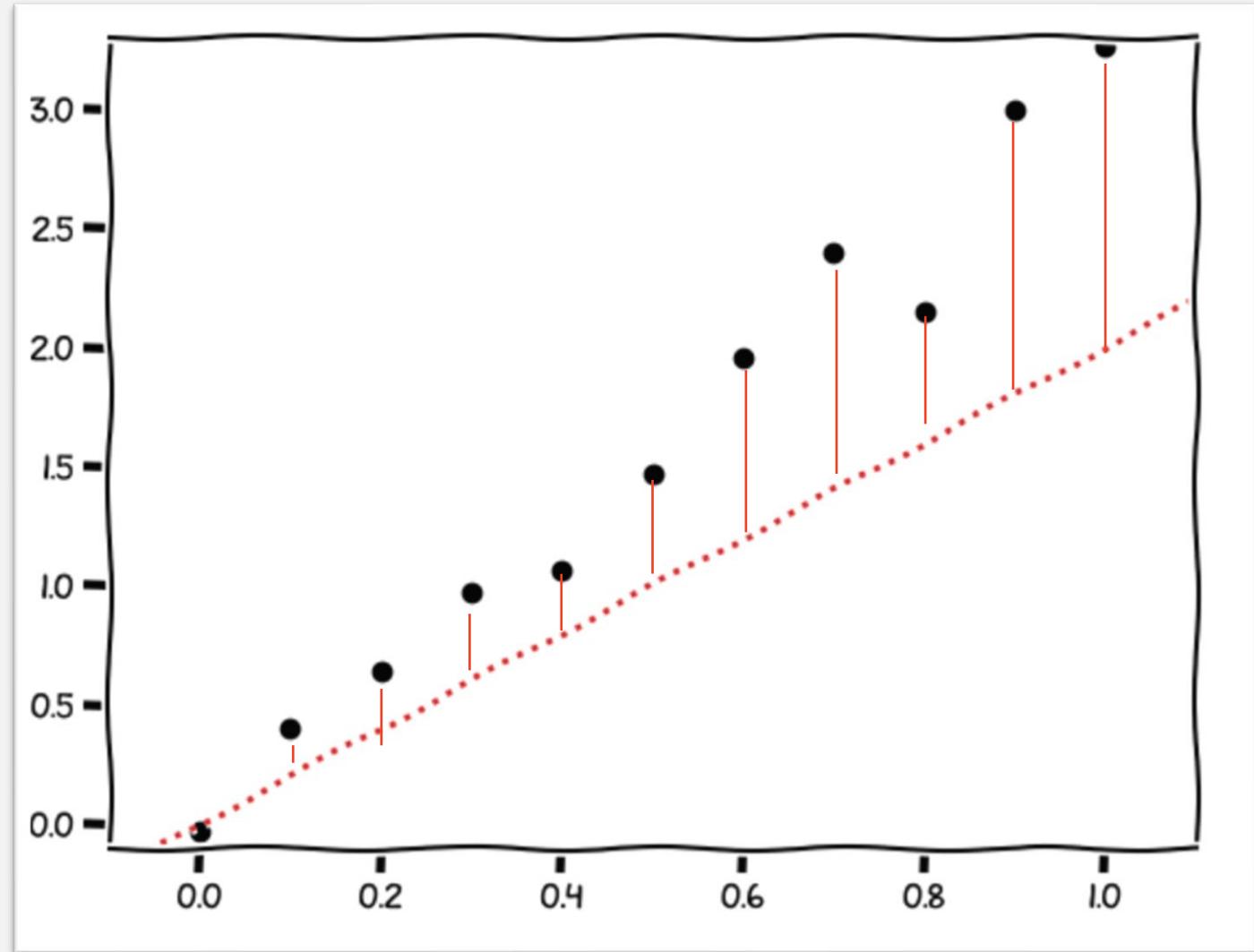
# linear regression

fitting a line  $y = ax$

1. Try some guess  $a_0$
2. Calculate loss

RMSE = 0.5

```
from sklearn.metrics import mean_squared_error
```



# linear regression

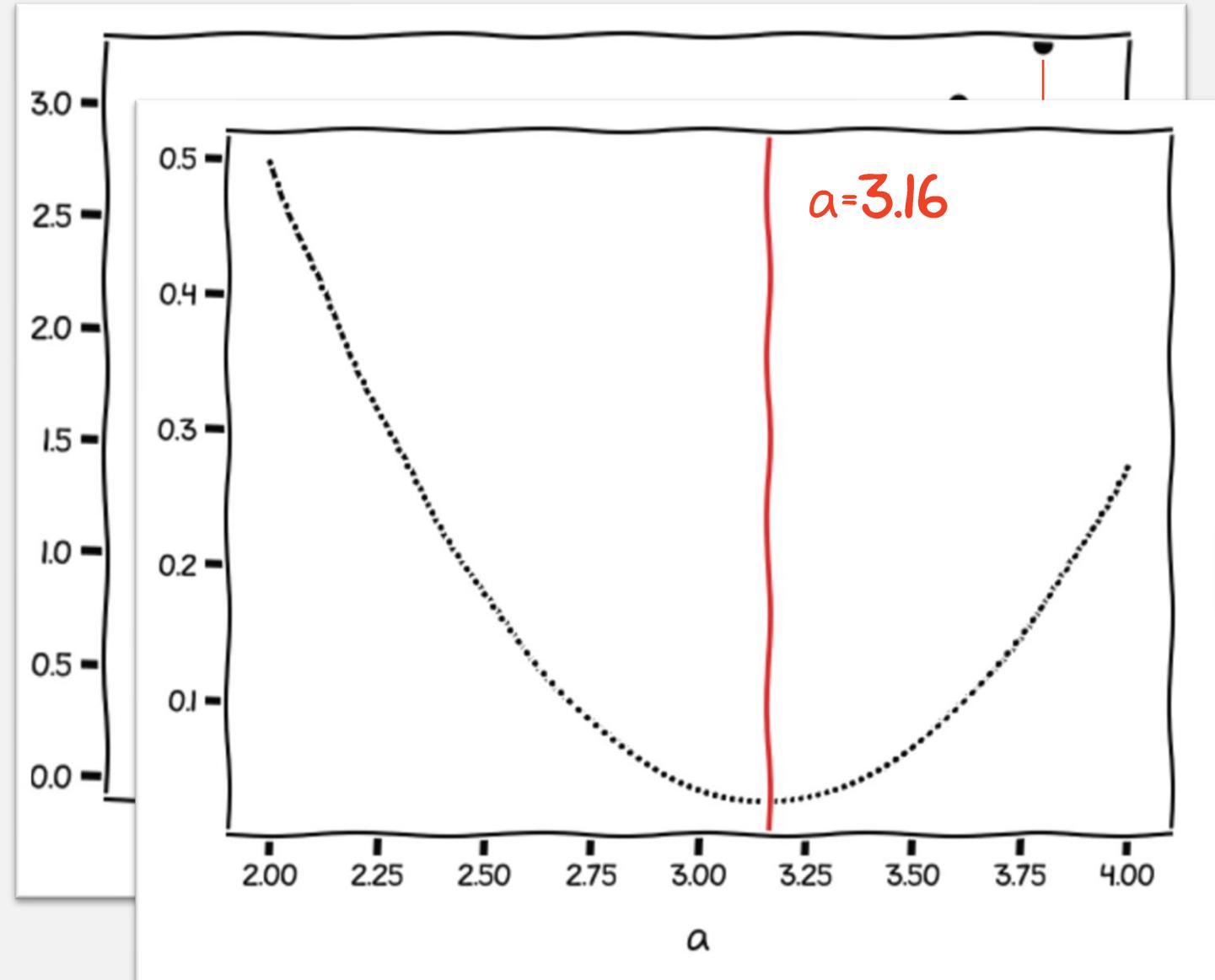
fitting a line  $y = ax$

1. Try some guess  $a_0$
2. Calculate loss

RMSE = 0.5

```
from sklearn.metrics import mean_squared_error
```

3. Find the minimum of the loss function  
E.g., through gradient descent



# linear regression

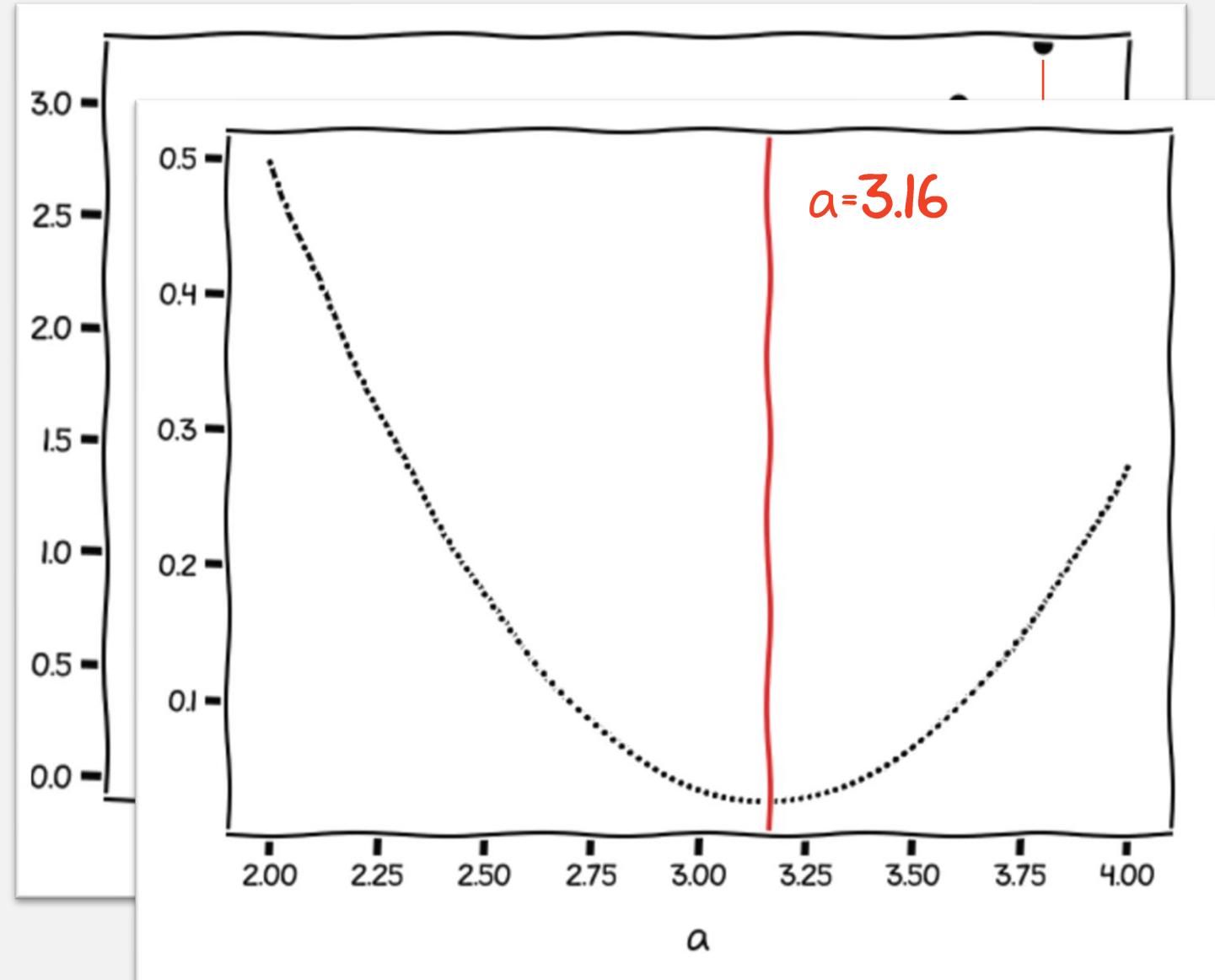
fitting a line  $y = ax$

1. Try some guess  $a_0$
2. Calculate loss

RMSE = 0.5

```
from sklearn.metrics import mean_squared_error
```

3. Find the minimum of the loss function  
E.g., through gradient descent



# linear regression

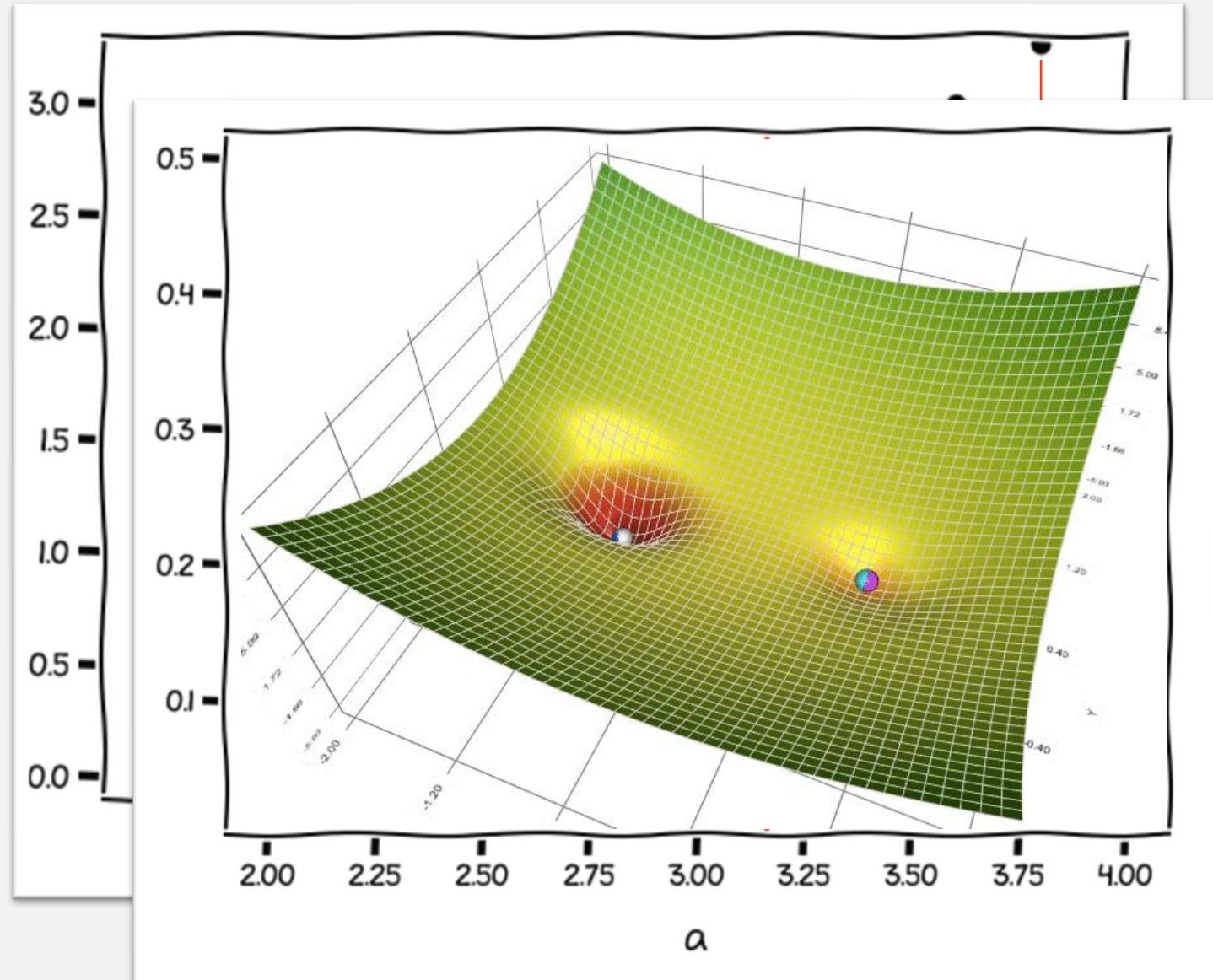
fitting a line  $y = ax$

1. Try some guess  $a_0$
2. Calculate loss

RMSE = 0.5

```
from sklearn.metrics import mean_squared_error
```

3. Find the minimum of the loss function  
E.g., through gradient descent



# linear regression

fitting a line  $y = ax$

1. Try some guess  $a_0$
2. Calculate loss

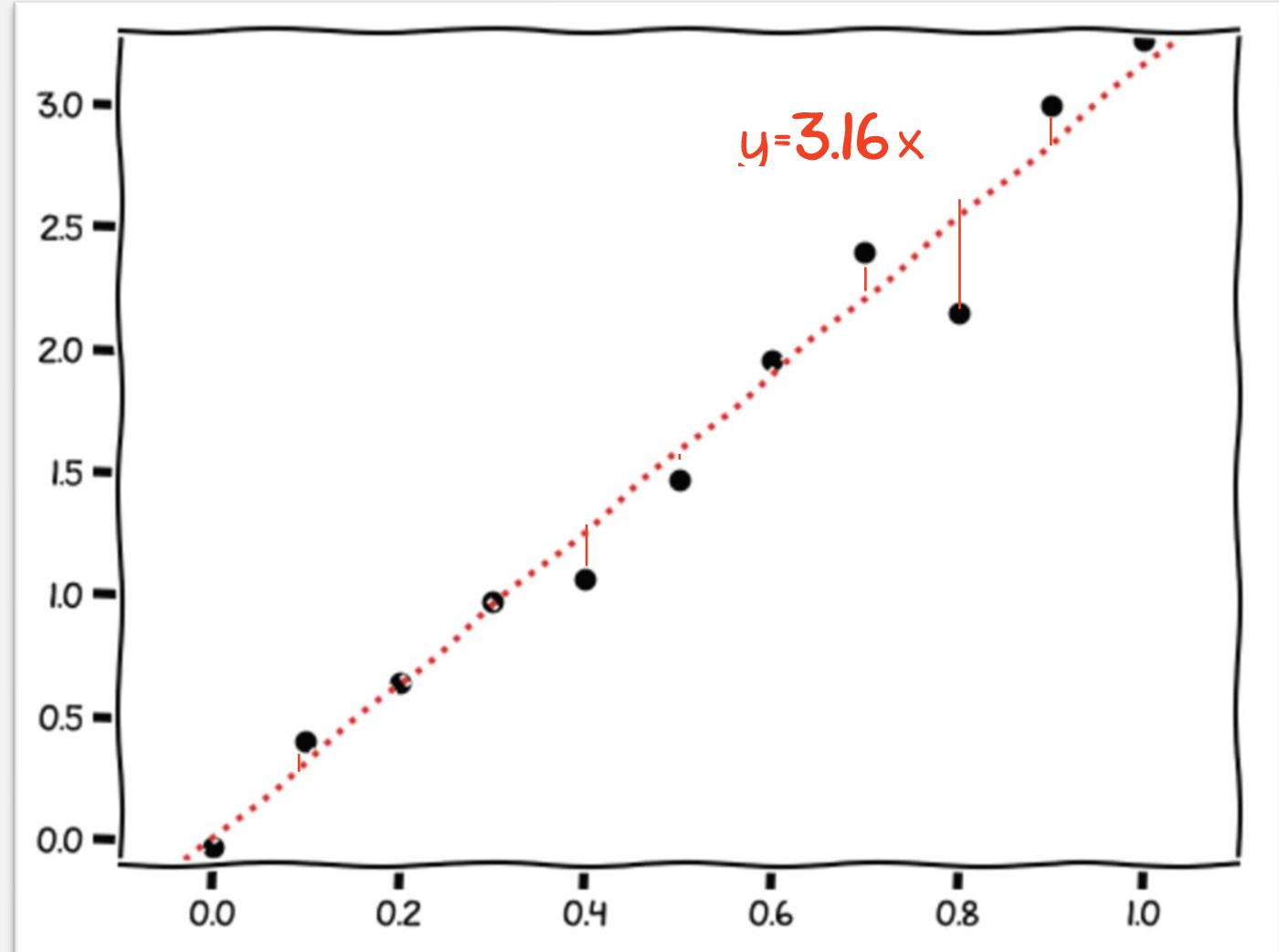
RMSE = 0.5

```
from sklearn.metrics import mean_squared_error
```

3. Find the minimum of the loss function

E.g., through gradient descent

4. Celebrate!



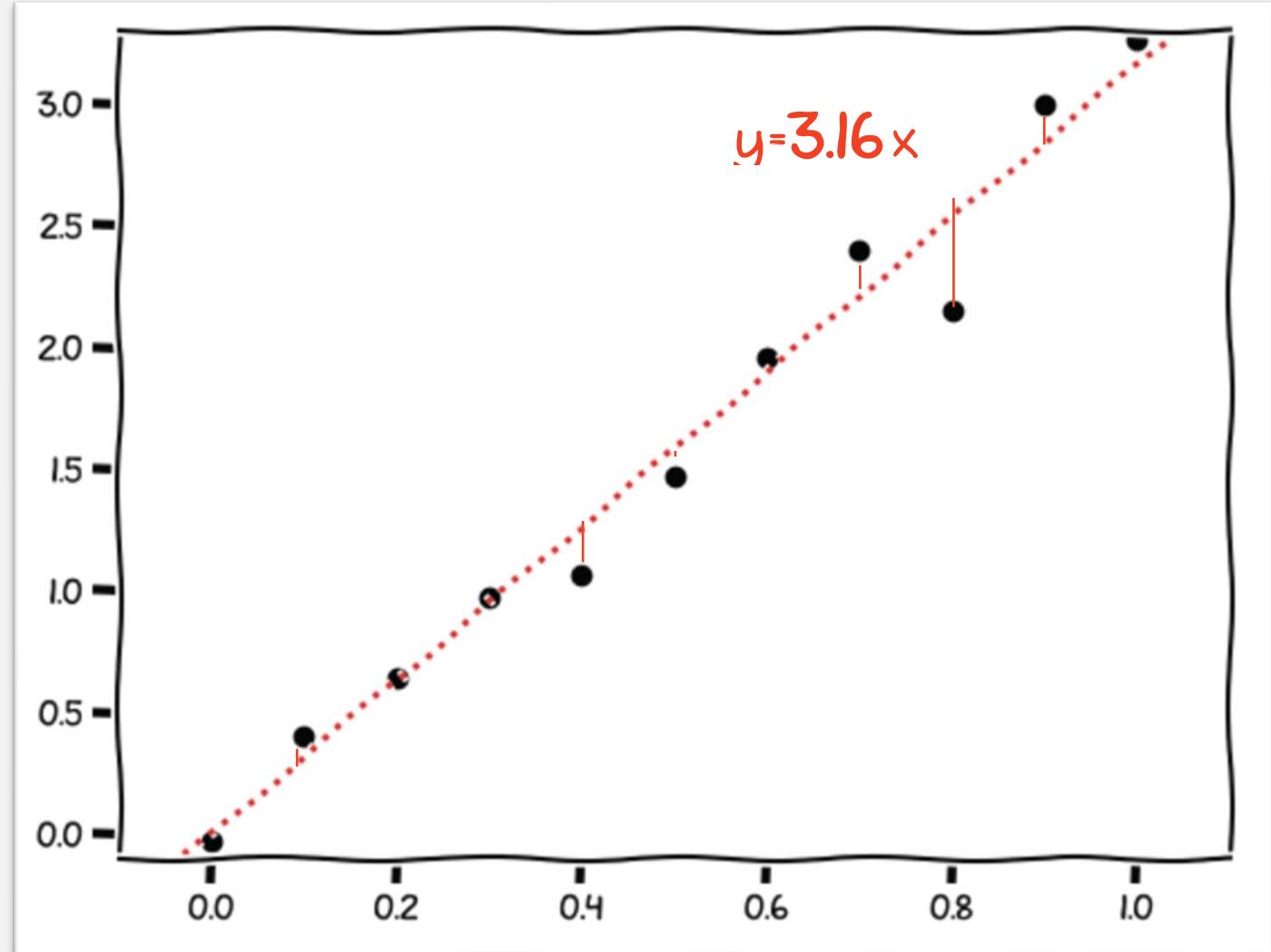
`sklearn.linear_model.LinearRegression`

# linear regression

fitting a line  $y = ax$

This assumes y depends on x!

If you want a *correlation* you will get a different answer



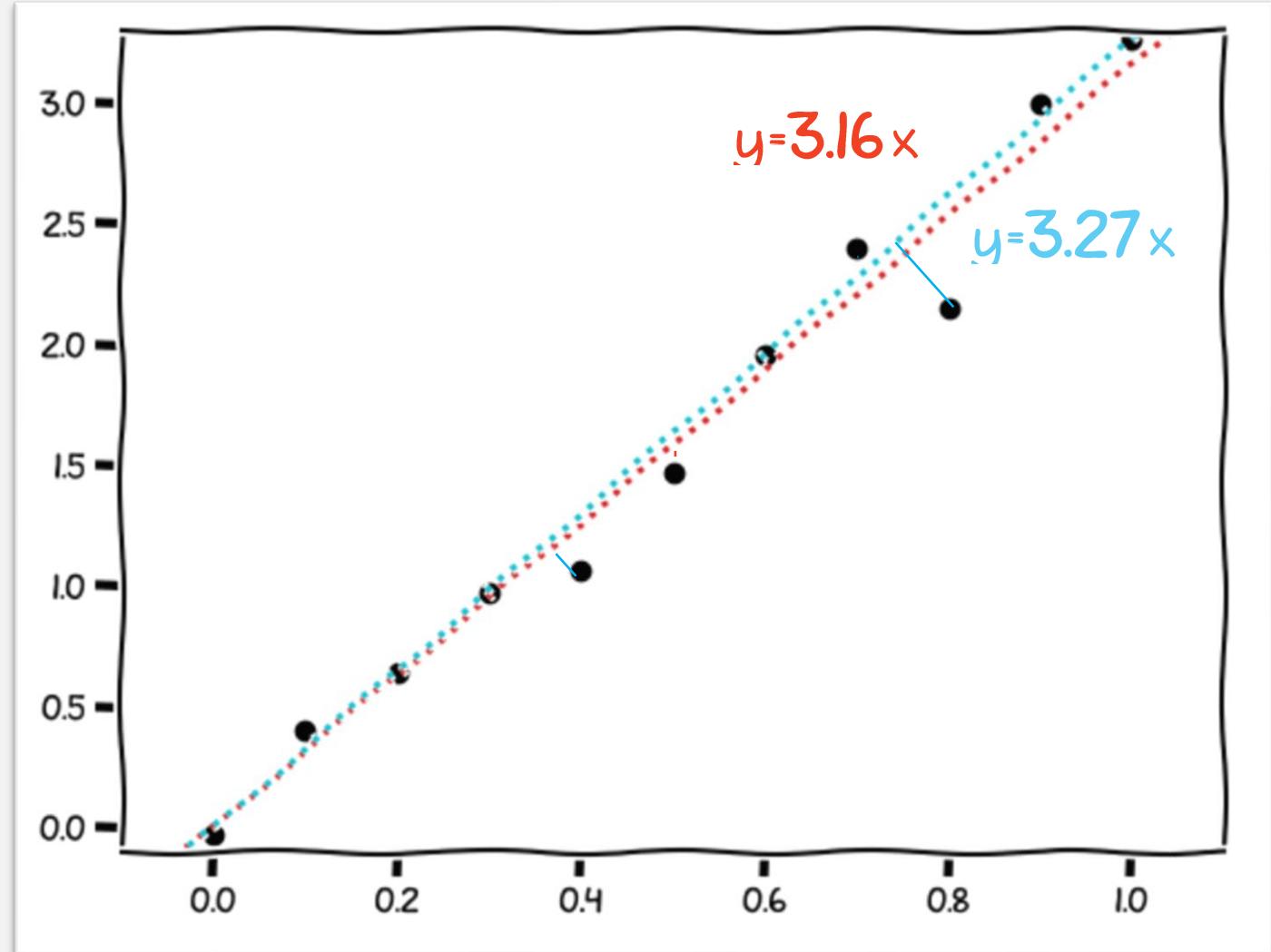
`sklearn.linear_model.LinearRegression`

# linear regression

fitting a line  $y = ax$

This assumes y depends on x!

If you want a *correlation* you will get a different answer

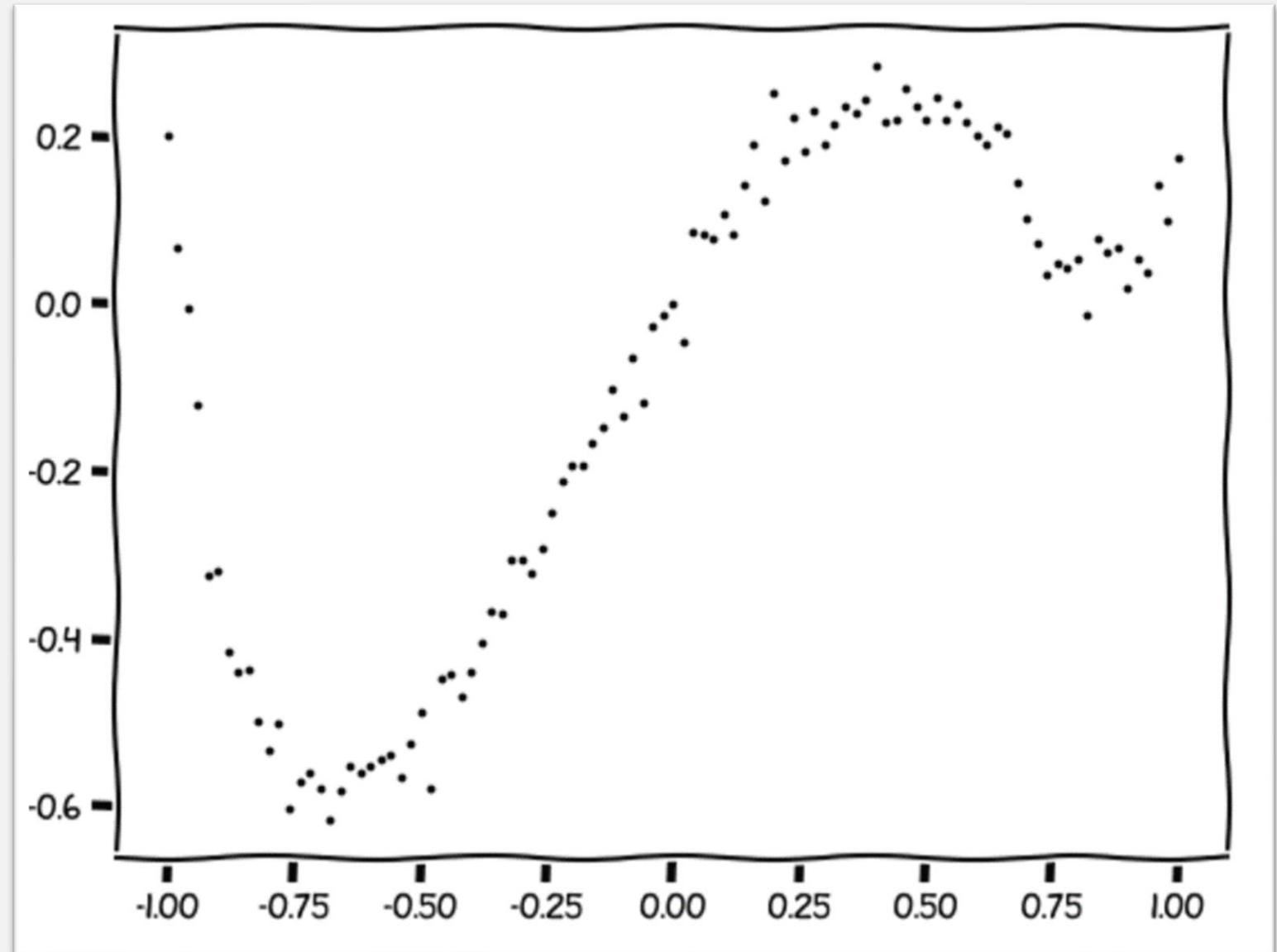


`sklearn.linear_model.LinearRegression`  
`sklearn.decomposition.PCA`



`sklearn.linear_model.Lasso`

regression

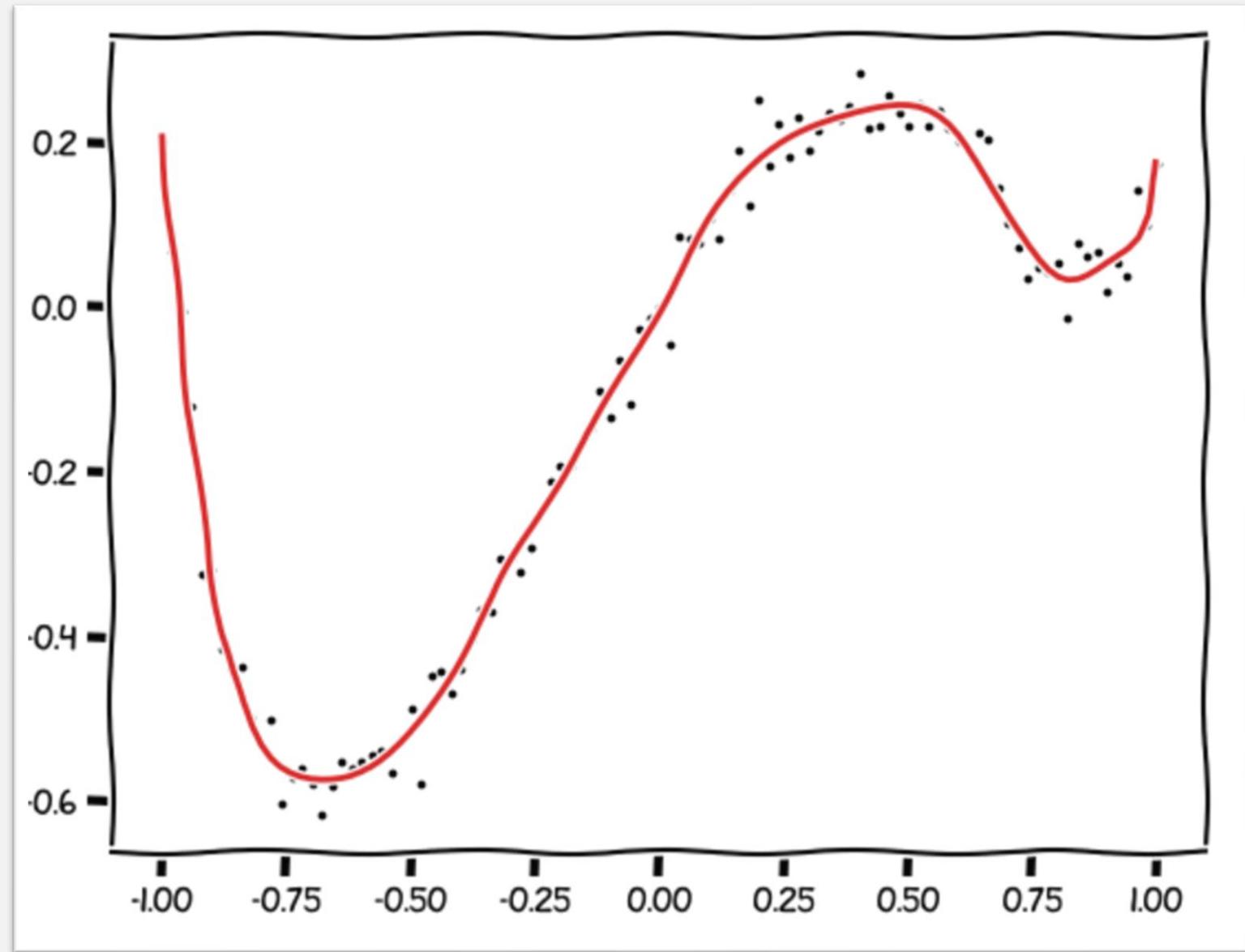


regression

$$y = \Sigma($$

-0.002	$x$
1	$x^2$
0.05	$x^3$
-4	$x^4$
-10	$x^5$
28	$x^6$
63	$x^7$
-100	$x^8$
-178	$x^9$
173	$x^{10}$
256	$x^{11}$
-140	$x^{12}$
-181	$x^{13}$
43	$x^{14}$
50	$x^{15}$

)

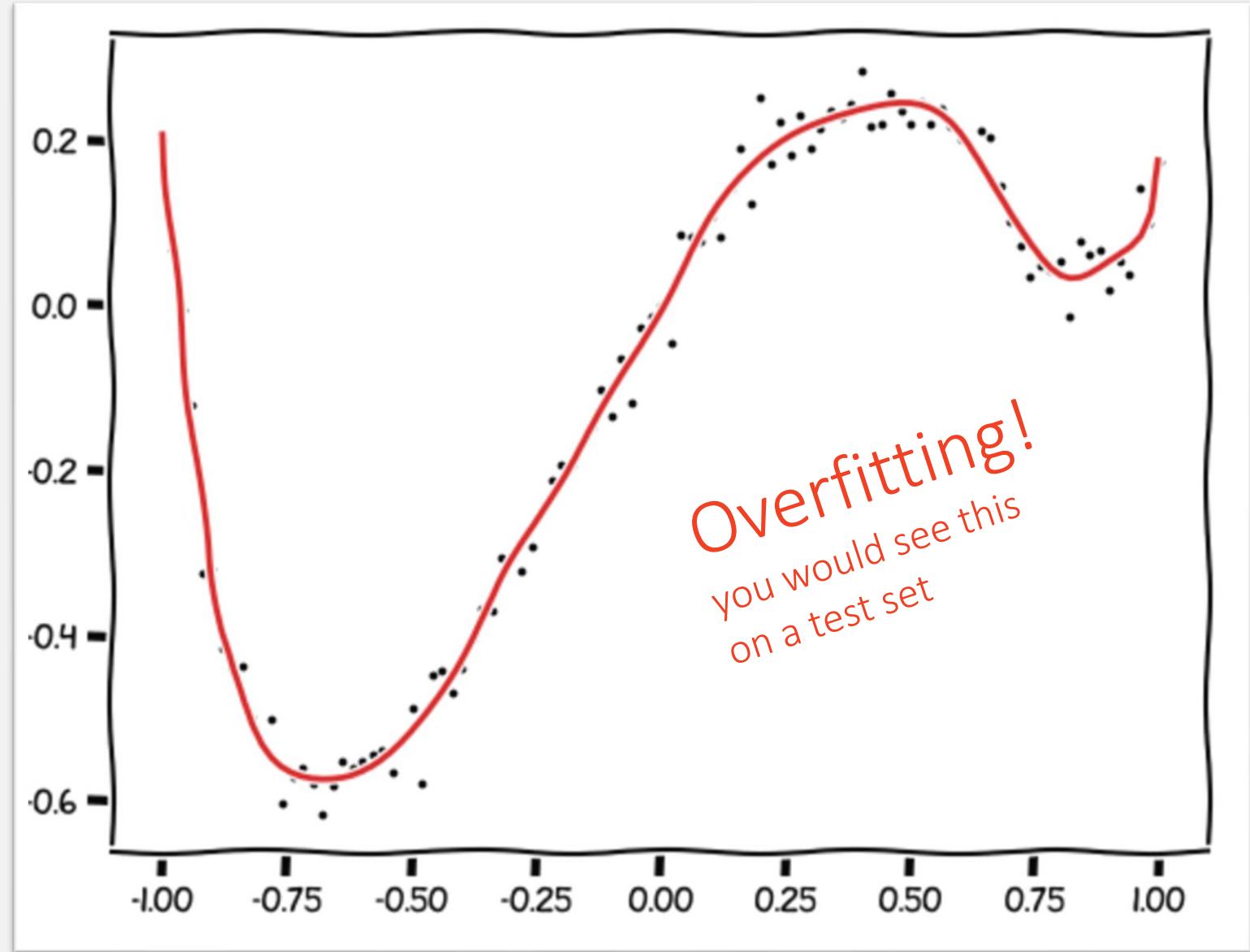


regression

$$y = \Sigma($$

-0.002	$x$
1	$x^2$
0.05	$x^3$
-4	$x^4$
-10	$x^5$
28	$x^6$
63	$x^7$
-100	$x^8$
-178	$x^9$
173	$x^{10}$
256	$x^{11}$
-140	$x^{12}$
-181	$x^{13}$
43	$x^{14}$
50	$x^{15}$

)



# lasso regression

“least absolute shrinkage and selection operator”

is a regularization technique

$$\mathcal{L} = \underbrace{\frac{|y - \vec{\alpha} \cdot \vec{x}|^2}{N}}_{\text{least squares loss}} + \lambda \sum |\vec{\alpha}|$$

linear coefficients

lasso parameter  
(lagrange multiplier)

Our first fancy loss function!

# lasso regression

“least absolute shrinkage and selection operator”

is a regularization technique

$$\mathcal{L} = \underbrace{\frac{|y - \vec{\alpha} \cdot \vec{x}|^2}{N}}_{\text{least squares loss}} + \lambda \sum |\vec{\alpha}|$$

linear coefficients  
lasso parameter  
(lagrange multiplier)

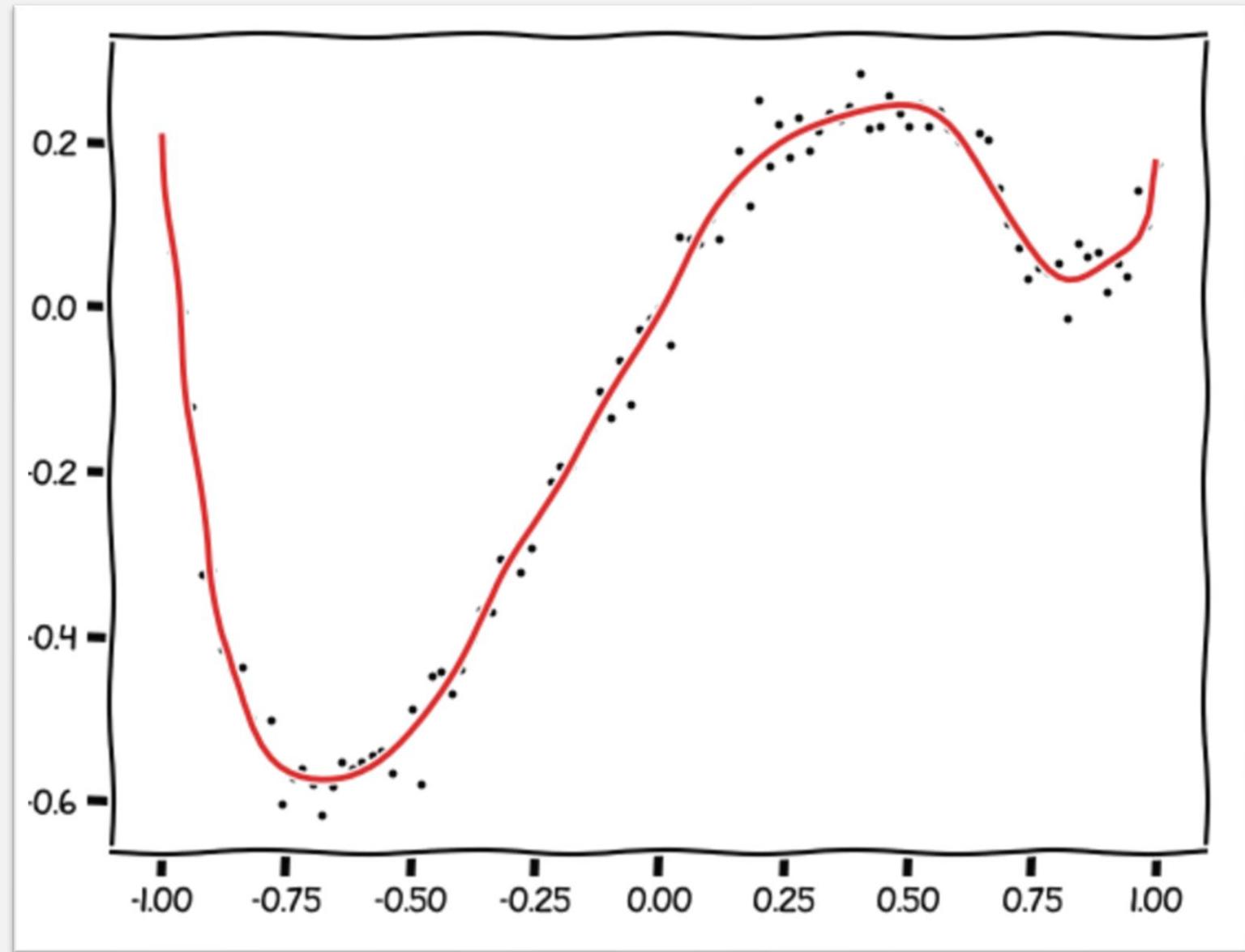


This thing is called  
alpha in scikit-learn  
to make your life  
harder!

Our first fancy loss function!

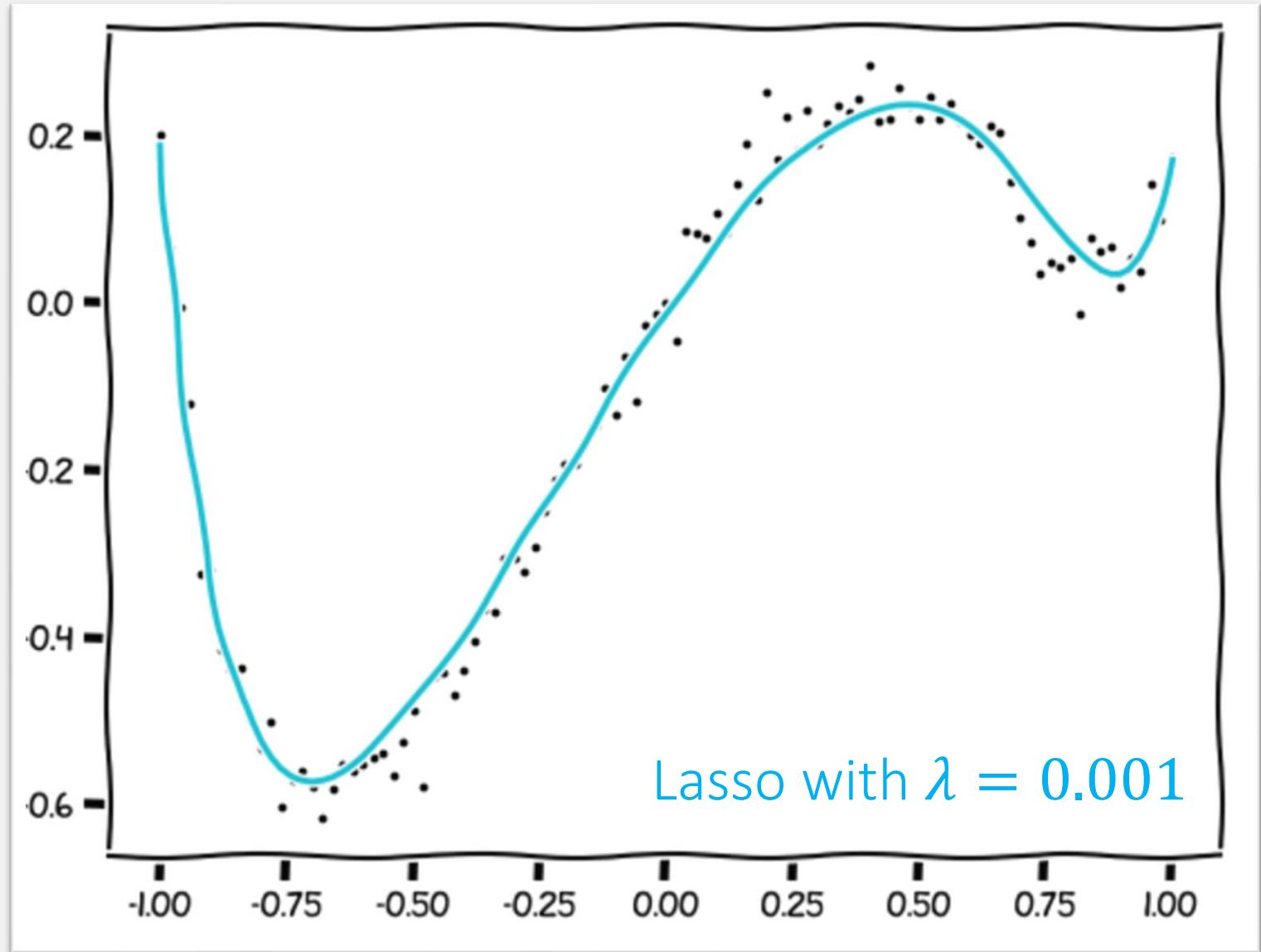
$$y = \Sigma($$

-0.002	$x$
1	$x^2$
0.05	$x^3$
-4	$x^4$
-10	$x^5$
28	$x^6$
63	$x^7$
-100	$x^8$
-178	$x^9$
173	$x^{10}$
256	$x^{11}$
-140	$x^{12}$
-181	$x^{13}$
43	$x^{14}$
50	$x^{15}$

 $)$ 

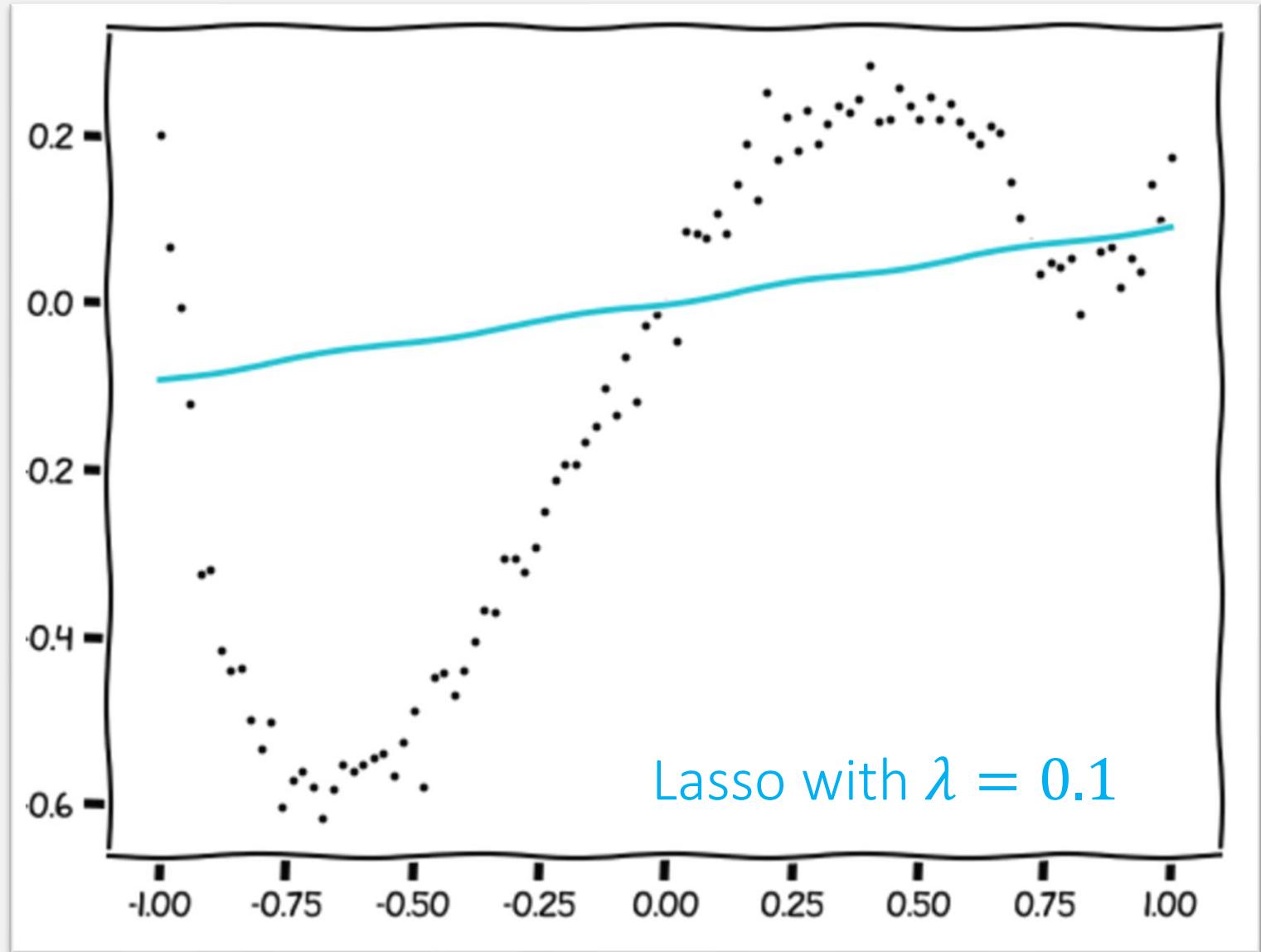
$$y = \Sigma($$

-0.01	$x$
0.9	$x^2$
-0.4	$x^3$
-0.6	$x^4$
0	$x^5$
-0.3	$x^6$
0	$x^7$
0	$x^8$
0	$x^9$
0	$x^{10}$
0.6	$x^{11}$
0	$x^{12}$
0.03	$x^{13}$
0	$x^{14}$
0	$x^{15}$

 $)$ 

$$y = \Sigma($$

-0.01	$x$
0.09	$x^2$
-0.4	$x^3$
-0.6	$x^4$
0	$x^5$
-0.3	$x^6$
0	$x^7$
0	$x^8$
0	$x^9$
0	$x^{10}$
0.6	$x^{11}$
0	$x^{12}$
0.03	$x^{13}$
0	$x^{14}$
0	$x^{15}$

 $)$ 

# lasso regression

“least absolute shrinkage and selection operator”

is a regularization technique

## pros

- very easy (my usual pro)
- good at reducing # of dimensions and selecting important features
- easy to implement to other techniques: just change your loss
- Helps to avoid overfitting – amazing!

## cons

- Need to tune  $\lambda$  – see `sklearn.linear_model.LassoCV`
- Not stable: different resamples will output different coefficients
- Has no knowledge of complexity by default
- All of this is worse if coefficients are correlated

# ridge regression

Aka Tikhonov regularization

## another regularization technique

$$\mathcal{L} = \underbrace{\frac{|y - \vec{\alpha} \cdot \vec{x}|^2}{N}}_{\text{least squares loss}} + \lambda \sum |\vec{\alpha}|^2$$

linear coefficients  
ridge parameter  
(lagrange multiplier)

# ridge regression

Aka Tikhonov regularization

## another regularization technique

$$\mathcal{L} = \frac{|y - \vec{\alpha} \cdot \vec{x}|^2}{N} + \lambda \sum |\vec{\alpha}|^2$$

linear coefficients

**pros**

more stable than Lasso

**cons**

does not set parameters to 0 – not good for feature selection  
If you ask me why I can probably explain but I don't want to

classification



decision trees!

`sklearn.tree.DecisionTreeClassifier`

lots of input  
parameters  
can be  
heterogeneous,  
unnormalized,  
continuous,  
categorical,  
whatever really



single  
categorical  
label



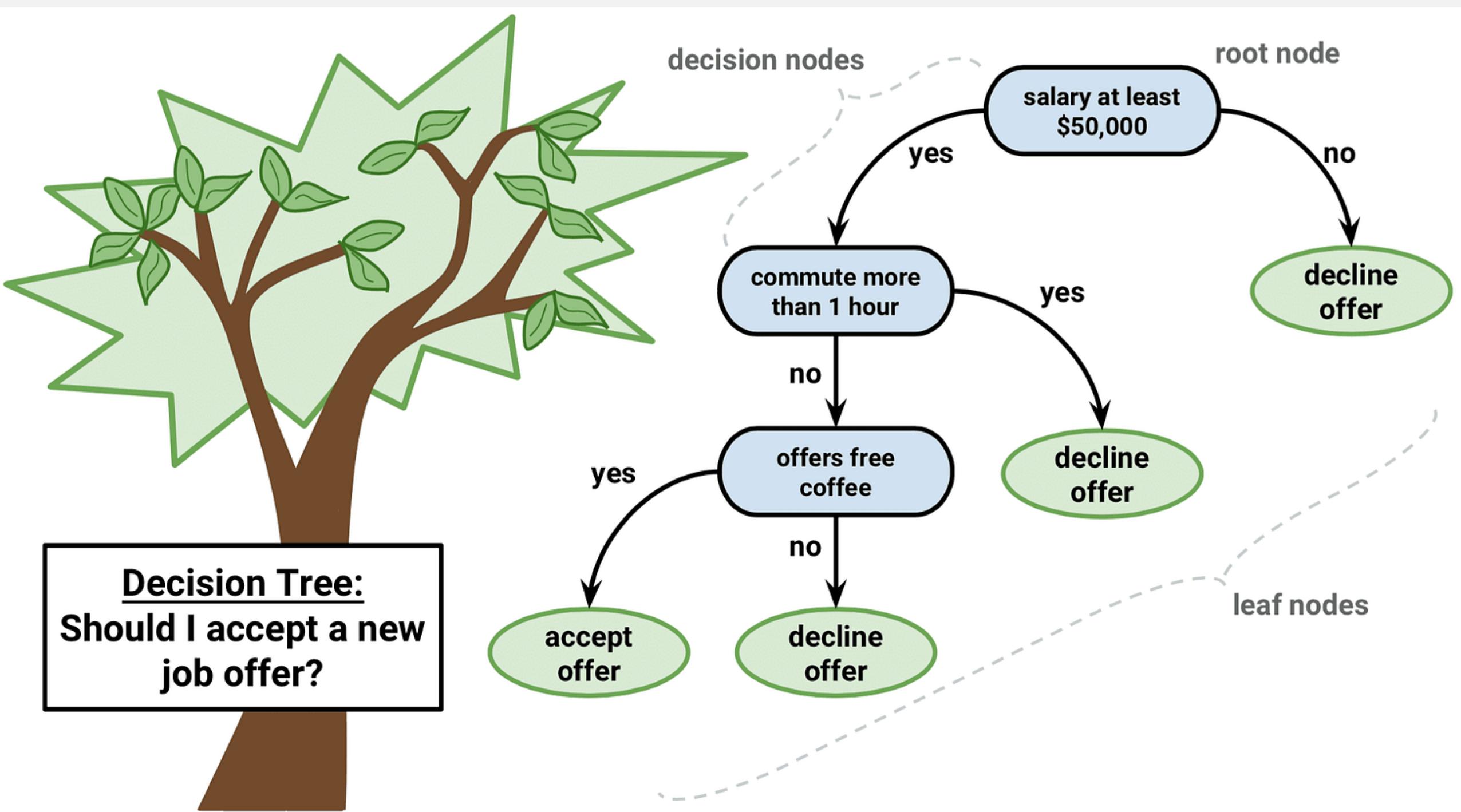
lots of input  
parameters

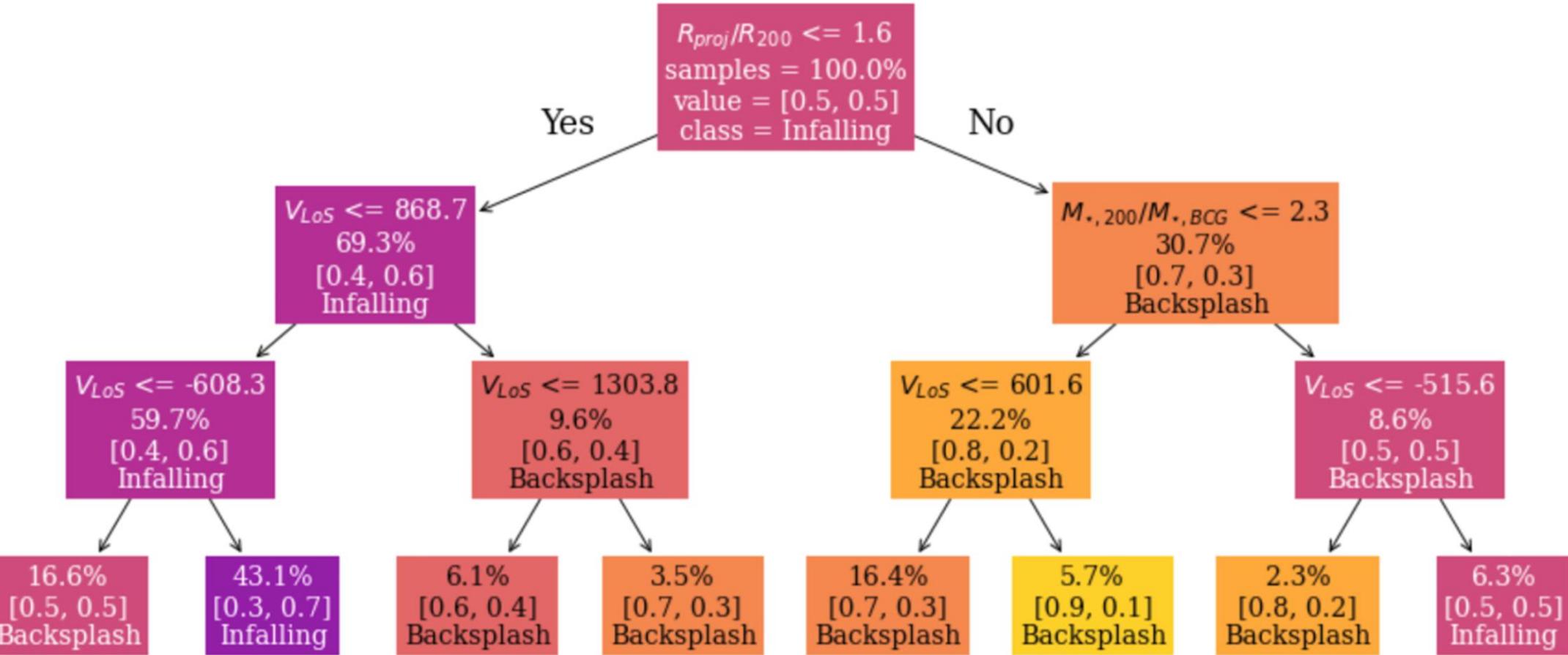
can be  
heterogeneous,  
unnormalized,  
continuous,  
categorical,  
whatever really

the reason they  
are great

single  
categorical  
label



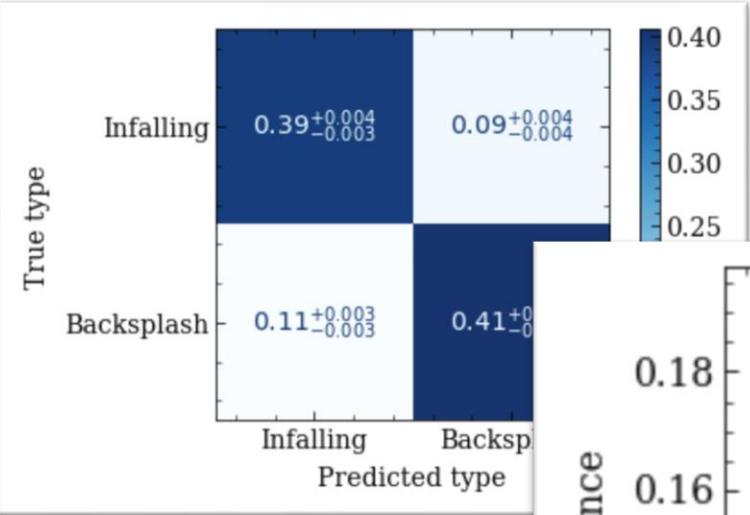




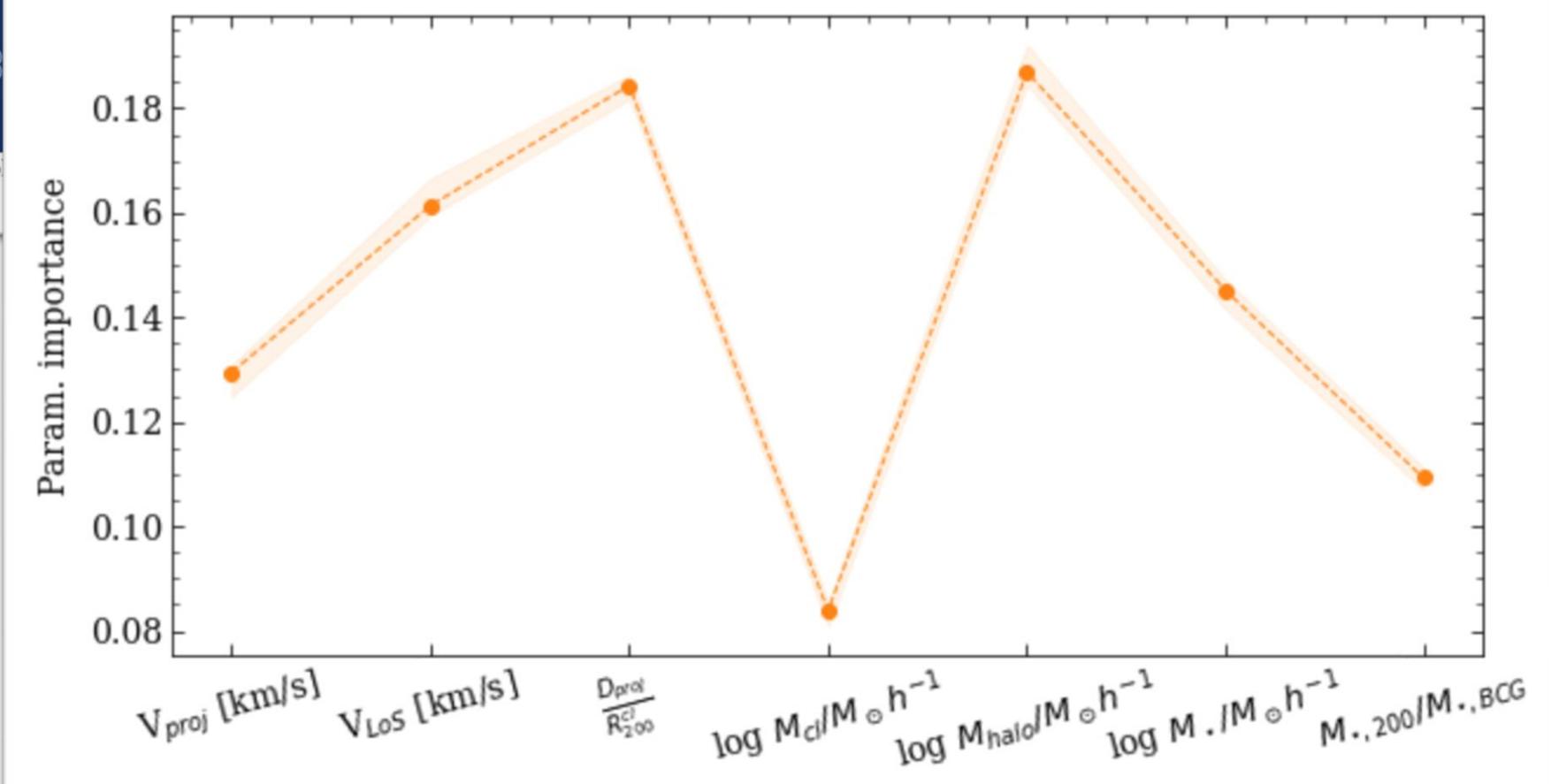
Make nice interpretable figures!

If your trees are simple

`sklearn.tree.plot_tree` ↗



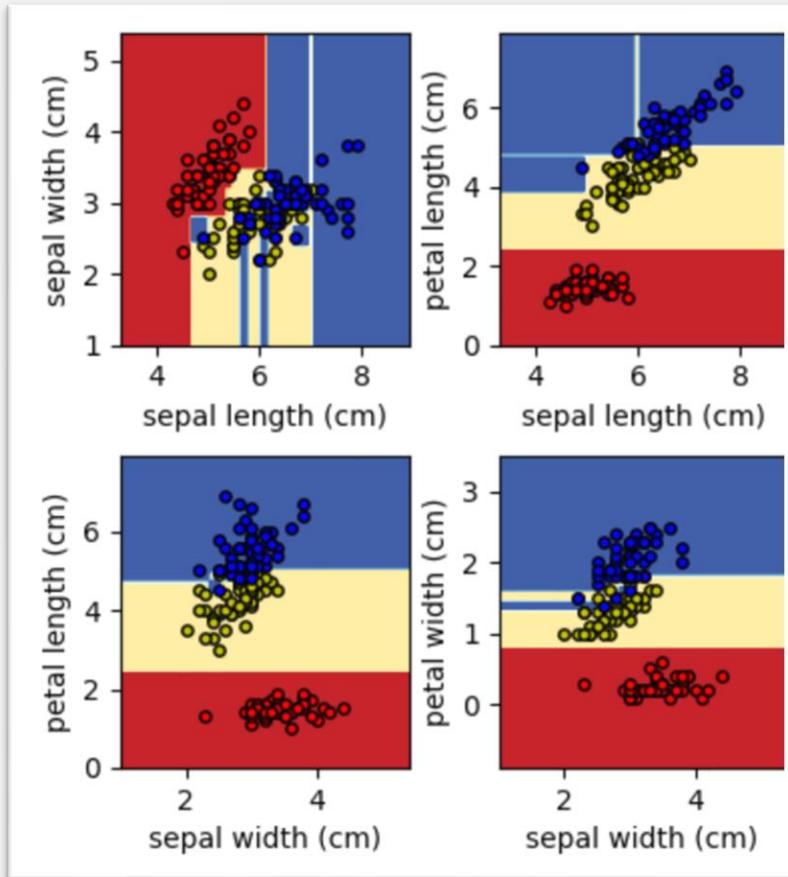
`sklearn.metrics.confusion_matrix`  
`sklearn.metrics.ConfusionMatrixDisplay`



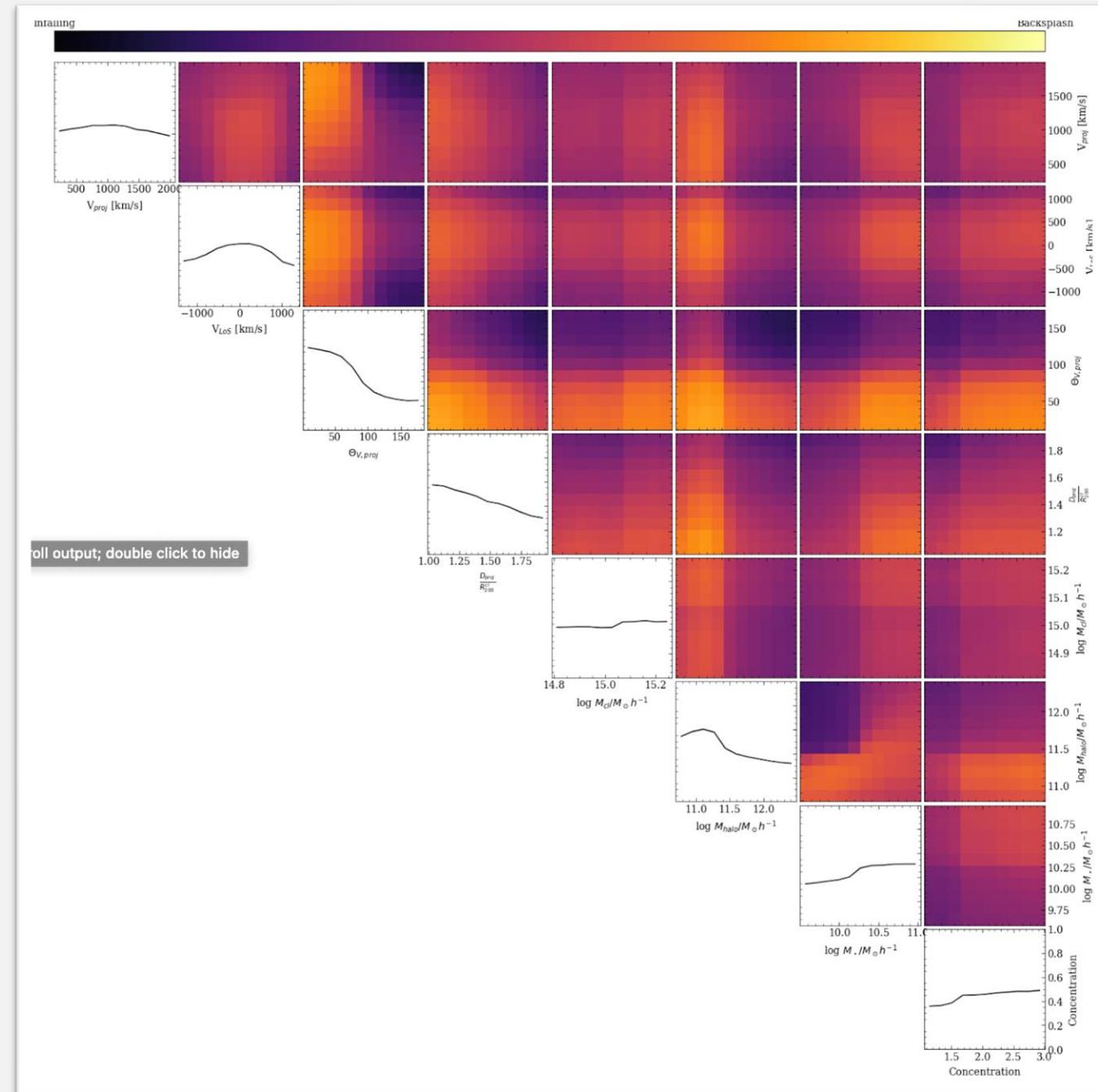
***property feature\_importances\_***

Get great  
results & see  
which parameters  
matter most!

# Show non-linear decision boundaries!



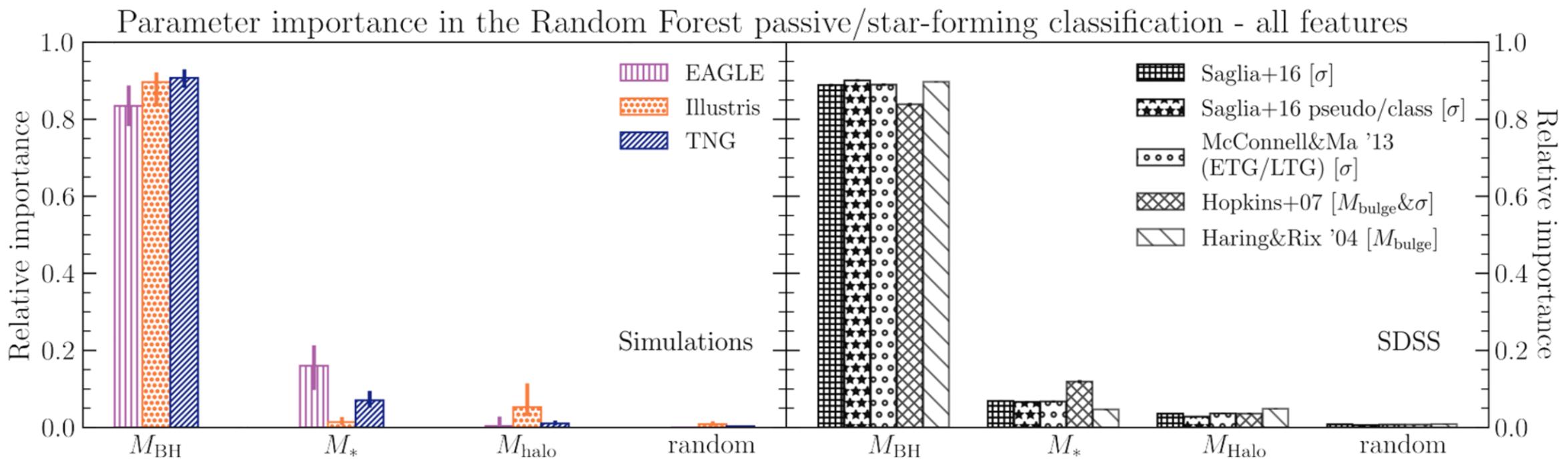
`sklearn.inspection.DecisionBoundaryDisplay`



# Can even prove causation over correlation?

Piotrowska et al., 2022

14 *J. M. Piotrowska et al.*



<https://arxiv.org/pdf/2112.07672.pdf>

lesson 2



decision trees!

`sklearn.tree.DecisionTreeClassifier`

## pros

Trees!  
Extremely interpretable  
Can capture non-linear relationships  
Feature importances – causation vs correlation?  
Works with basically any input data w/out worries

## cons

Very unstable  
Require *a lot* of tuning  
Not good on imbalanced datasets (like most things)  
Cannot do anything except “yes” or “no” classification  
Neural networks beat them on accuracy



## pros

Trees!  
Extremely interpretable  
Can capture non-linear relationships  
Feature importances – causation vs correlation?  
Works with basically any input data w/out worries

## cons

~~Very unstable~~  
~~Require a lot of tuning~~  
Not good on imbalanced datasets (like most things)  
~~Cannot do anything except “yes” or “no” classification~~  
~~Neural networks beat them on accuracy???~~maybe



## pros

Trees!  
Extremely interpretable  
Can capture non-linear relationships  
Feature importances – causation vs correlation  
Works with basically any input data w/o modification

## cons

~~Very unstable~~  
~~Require a lot of tuning~~  
Not good on imbalanced datasets (~~like cancer~~)  
~~Cannot do anything except “yes” or “no”~~  
~~Neural networks beat them on accuracy~~



random forests



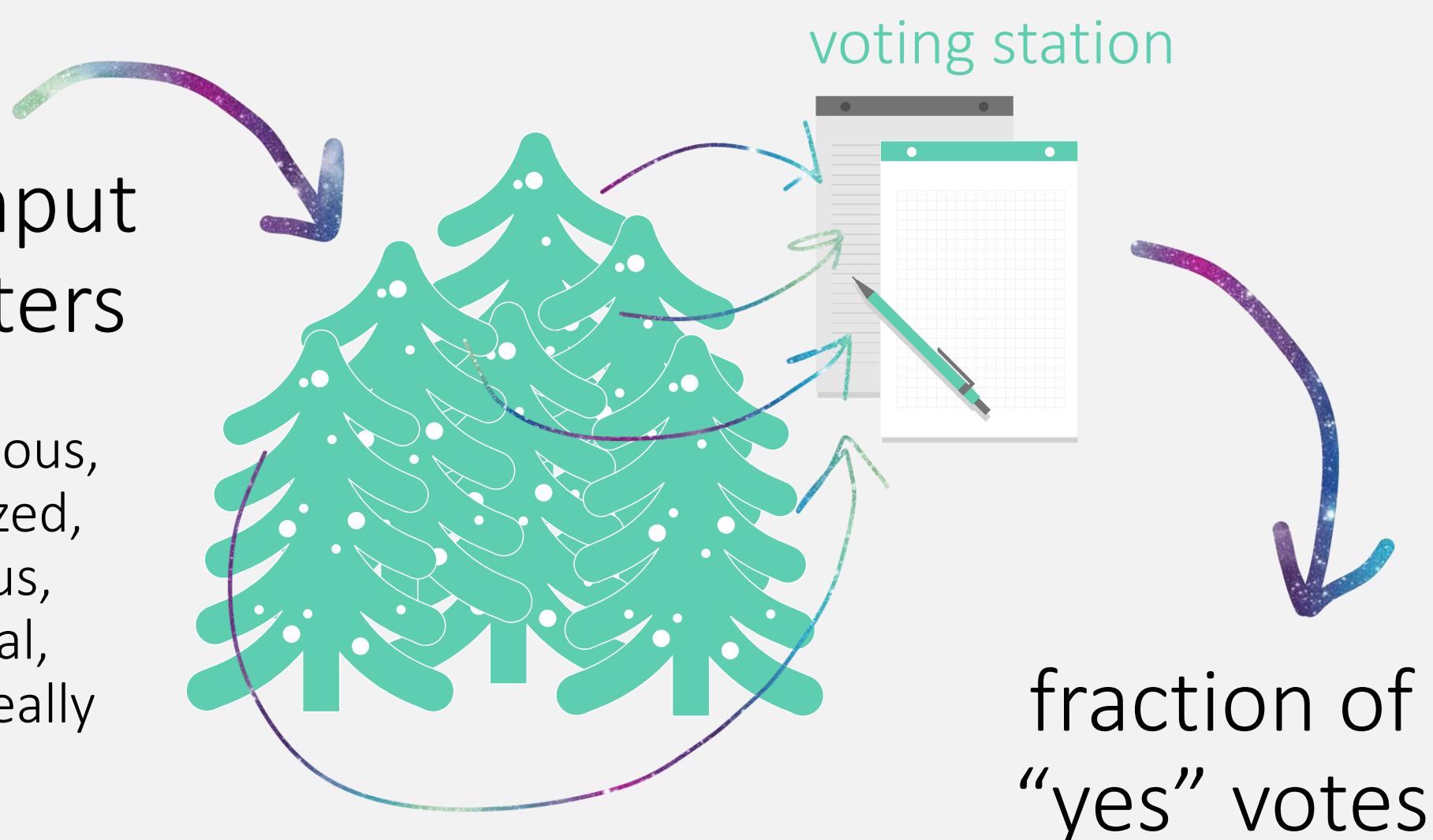
lots of input  
parameters  
can be  
heterogeneous,  
unnormalized,  
continuous,  
categorical,  
whatever really



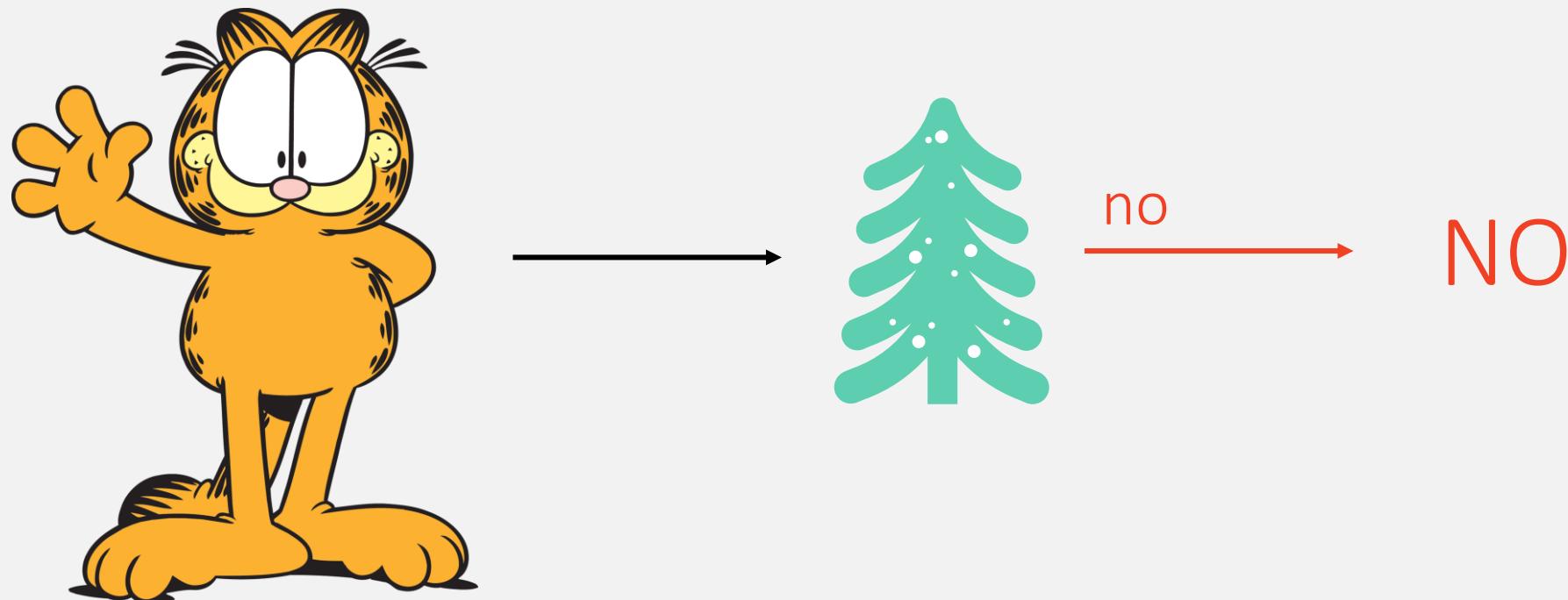
single  
categorical  
label



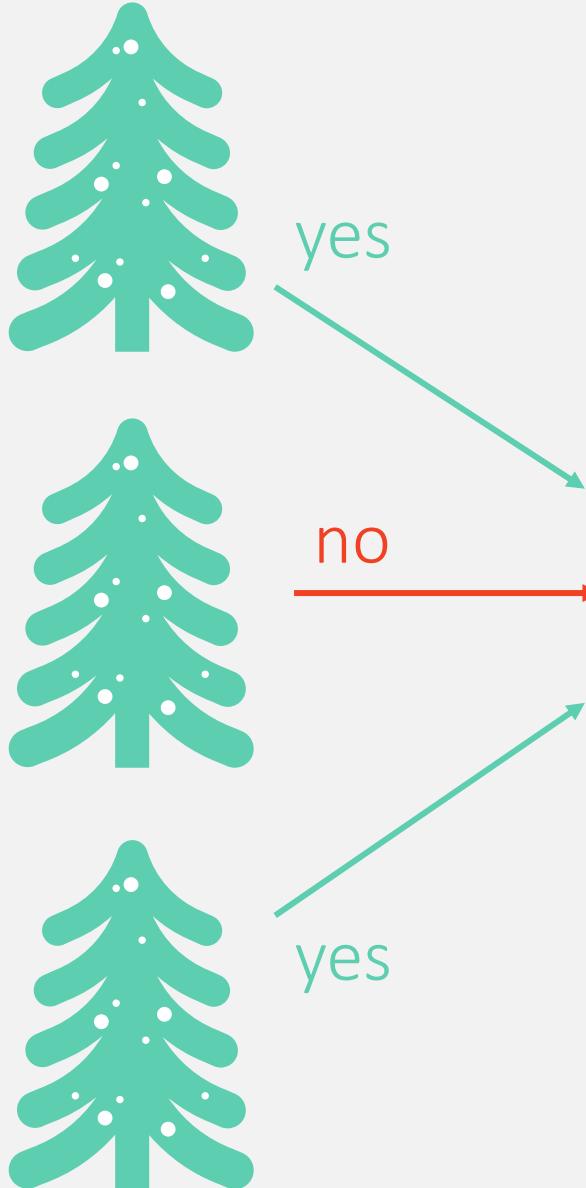
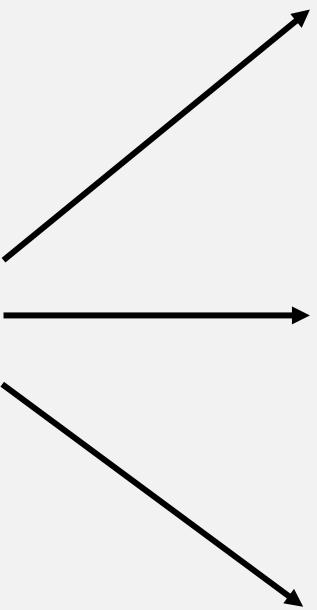
lots of input  
parameters  
can be  
heterogeneous,  
unnormalized,  
continuous,  
categorical,  
whatever really



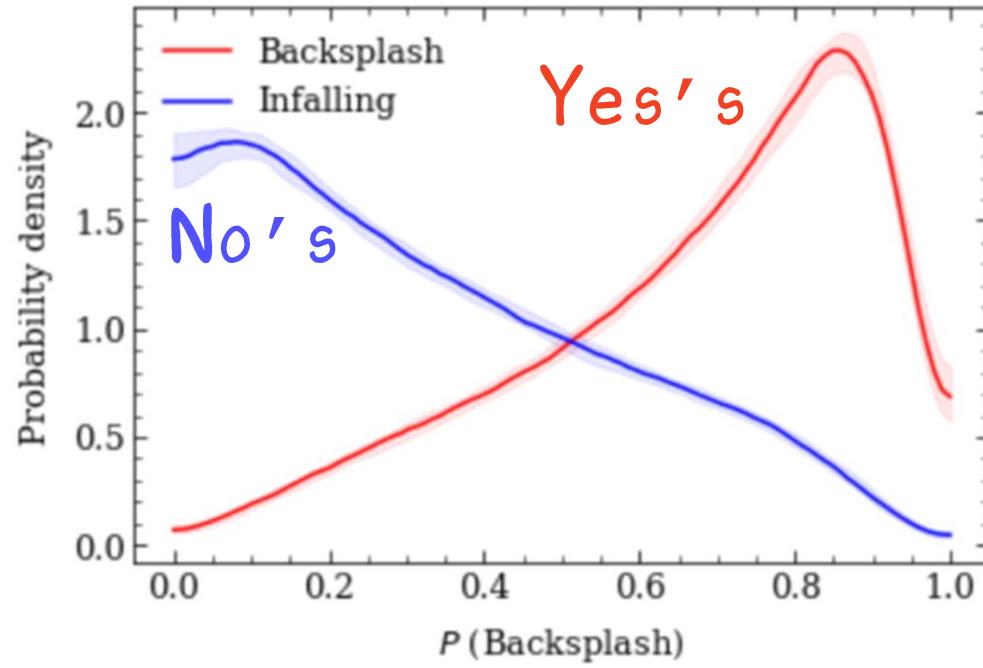
is this a cat?



is this a cat?



$p(\text{cat})=0.67$

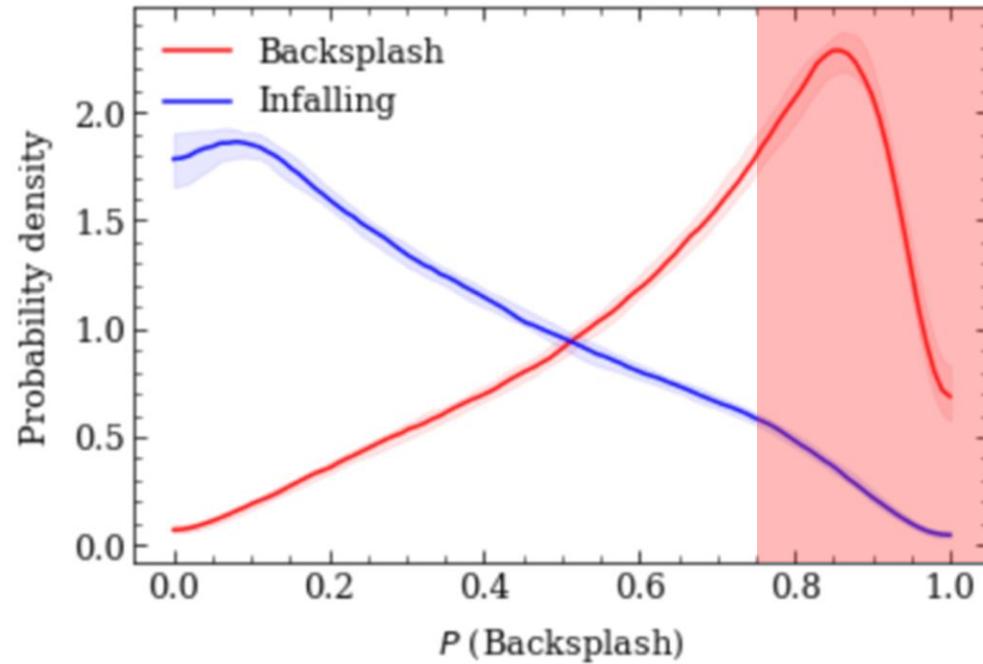


**Probability density distribution** of **infalling** and **backsplash** galaxies, as a function of **predicted backsplash probability**.

Infalling galaxies are more likely to be misclassified as backsplash than vice versa

$p(\text{yes})$

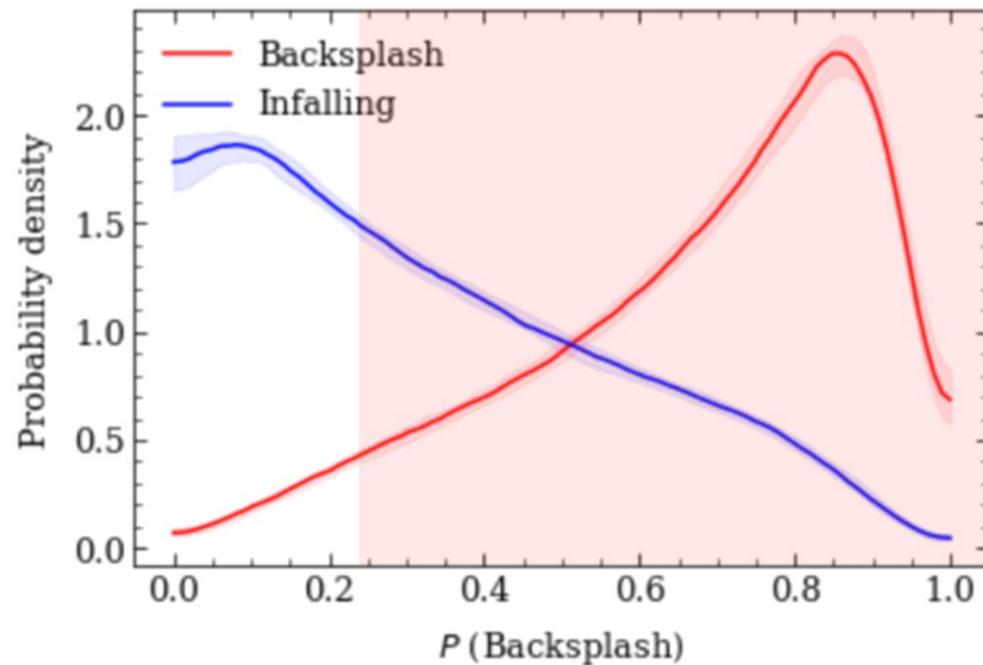
*pure sample*



**Probability density distribution** of infalling and backsplash galaxies, as a function of **predicted backsplash probability**.

Infalling galaxies are more likely to be misclassified as backsplash than vice versa

complete sample



**Probability density distribution** of infalling and backsplash galaxies, as a function of **predicted backsplash probability**.

Infalling galaxies are more likely to be misclassified as backsplash than vice versa

# Galaxy Zoo

20 citizen scientists classify each galaxy,  
giving  $p(\text{smooth})$



ⓘ You should sign in!

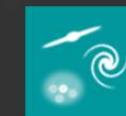
## TASK

## TUTORIAL

Is the galaxy simply smooth and rounded, with no sign of a disk?



Smooth



Features or Disk



Star, Artifact, or Bad Zoom

Is the galaxy simply smooth and rounded, with no sign of a disk?



Smooth



Features or Disk



Star or Artifact

How rounded is it?



Round

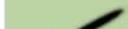


In Between



Cigar Shaped

Could this be a disk viewed edge-on?



Yes - Edge On Disk



No - Something Else

Does the galaxy have a bulge at its centre? If so, what shape?



Rounded



Boxy



No bulge

Is there a bar feature through the centre of the galaxy?



No Bar



Weak Bar



Strong Bar

How many spiral arms are there?



1



2



3



4



More than 4



Can't tell

How tightly wound do the spiral arms appear?



Tight



Medium



Loose

Is there any sign of a spiral arm pattern?



Yes



No

Is the galaxy merging or disturbed?



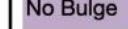
Merging



Major Disturbance



Minor Disturbance



None

Do you see any of these rare features?



Ring



Lens or arc



Irregular



Dust lane



Overlapping



Something Else

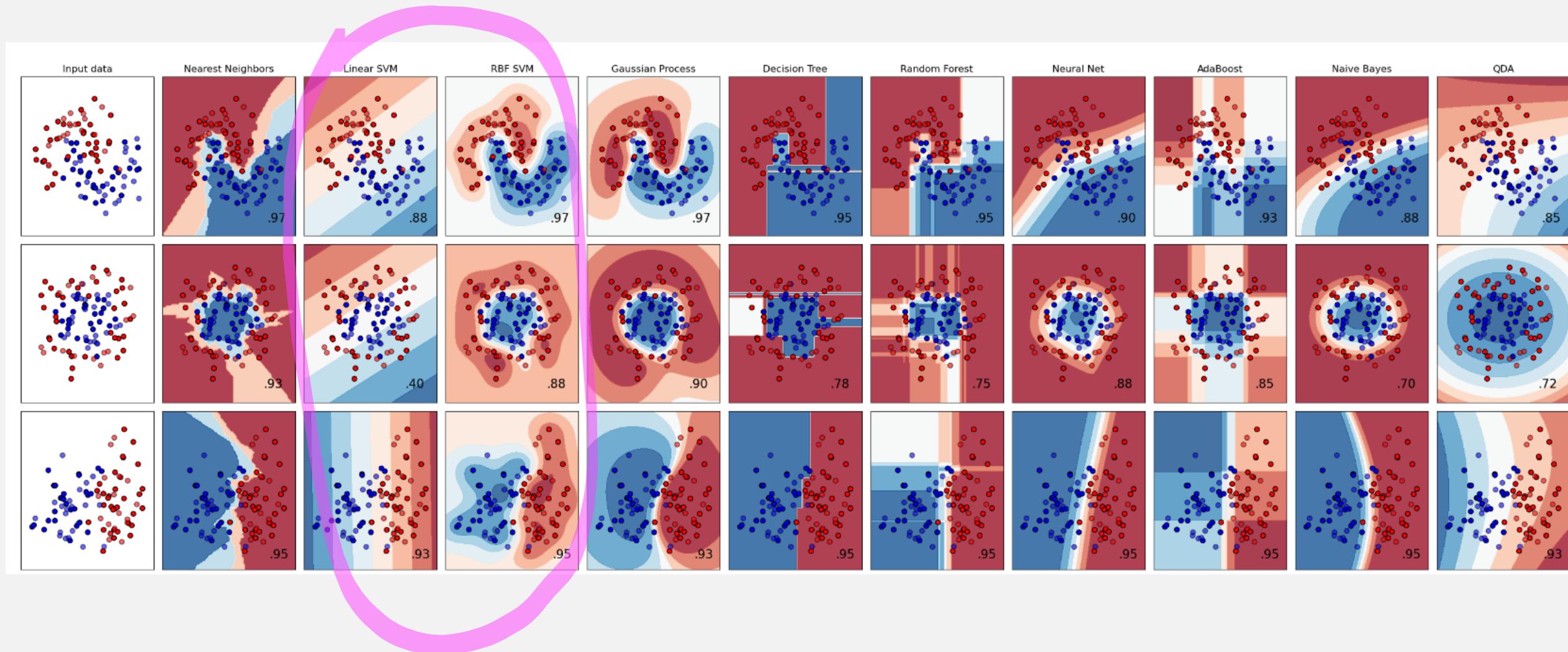


Nothing Unusual

END

A random forest = a natural loss function!

# Support Vector Machines



`sklearn.svm.SVC`

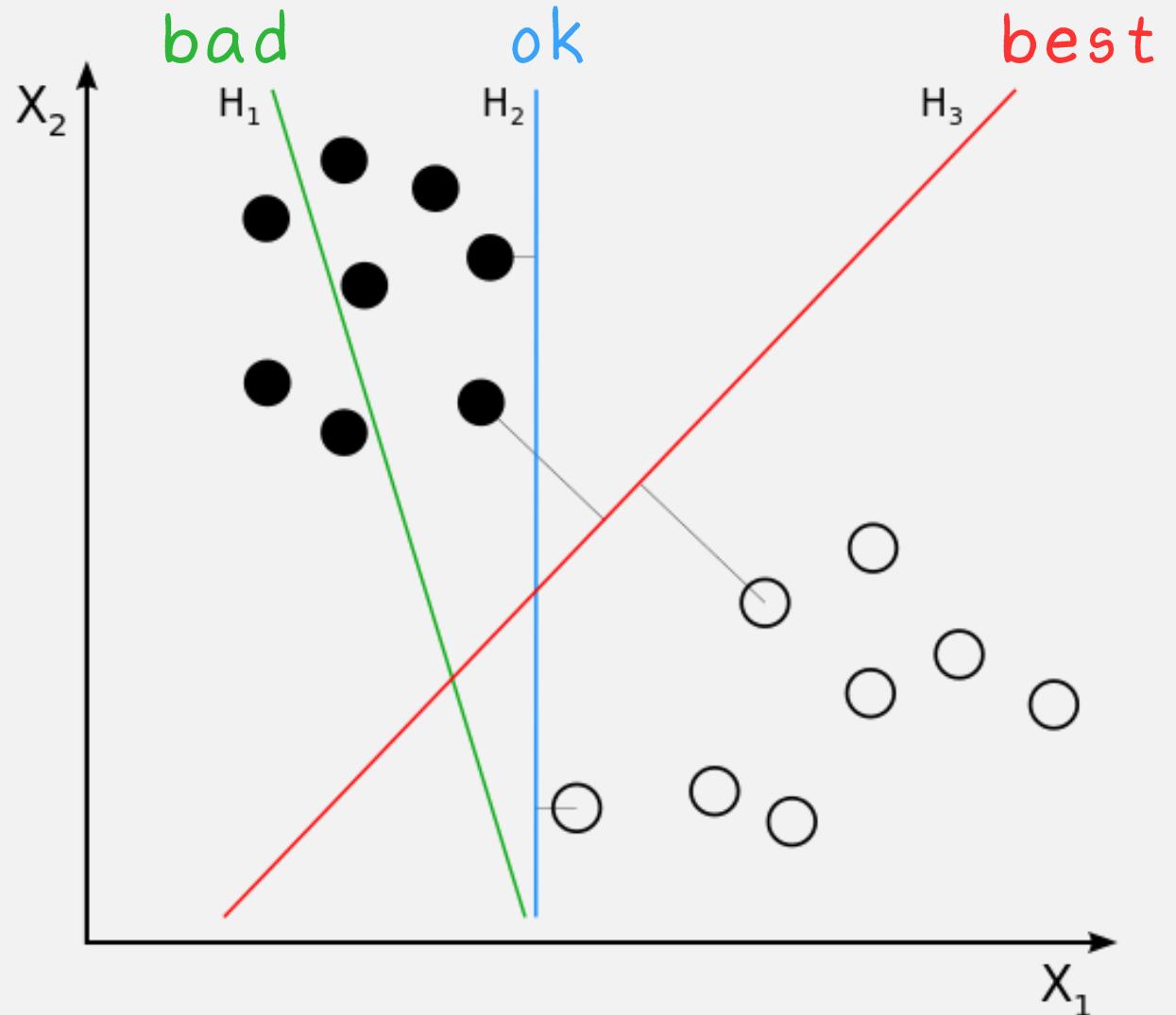
# SVMs (SVCs)

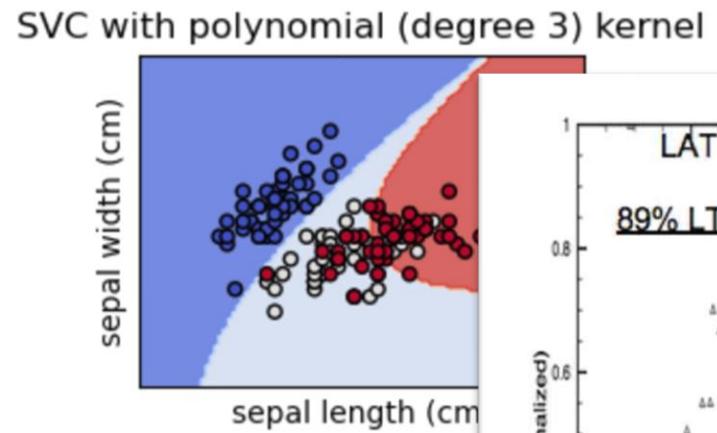
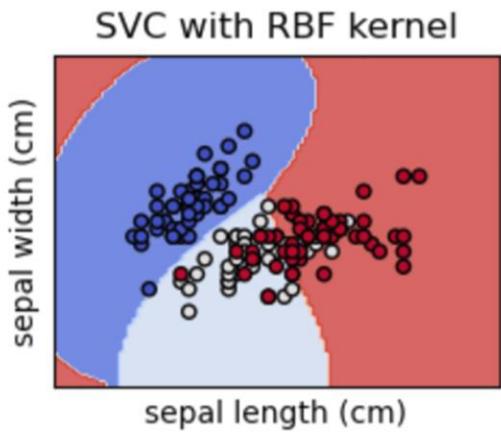
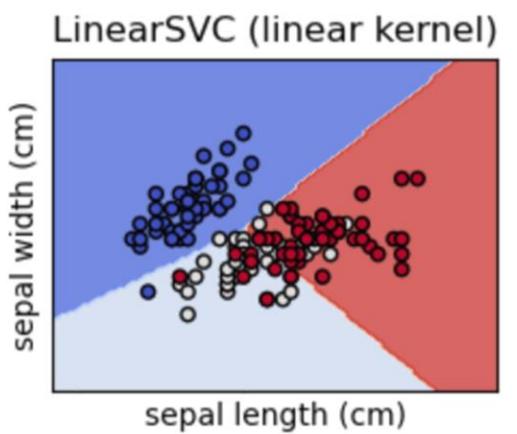
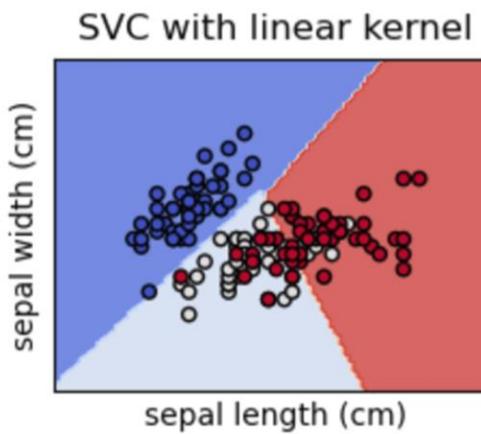
Think of them as linear regression  
but for classification

Don't quote me on this

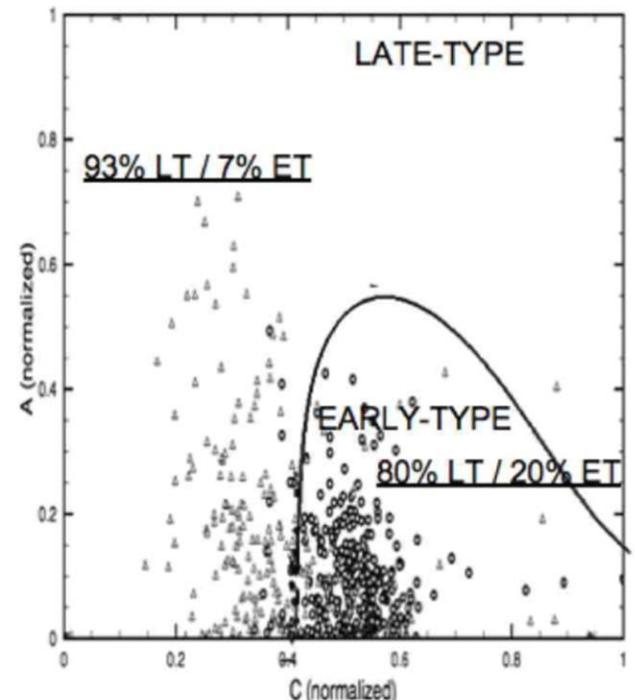
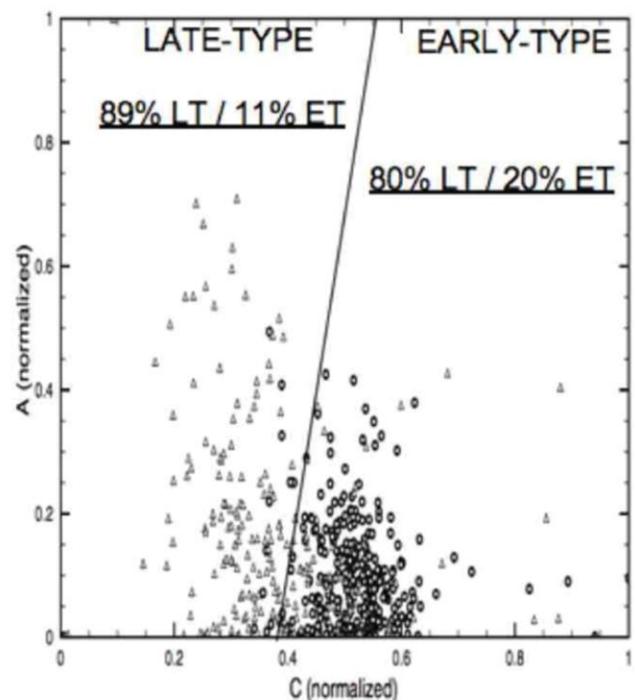
Fits some function (**kernel** or  
**hyperplane**) to get the **maximum**  
separation between classes

Distance to the **nearest** point on  
each side is maximized





Huertas-Company et al. 2008



# SVMs (SVCs)

## pros

very versatile (with RBF/polynomial kernels)  
give you a meaningful formula you can tell other people  
good in many-dimensional spaces  
good for small datasets (even if  $N_{\text{dim}} > N_{\text{samples}}$ )  
multi-class classifier!

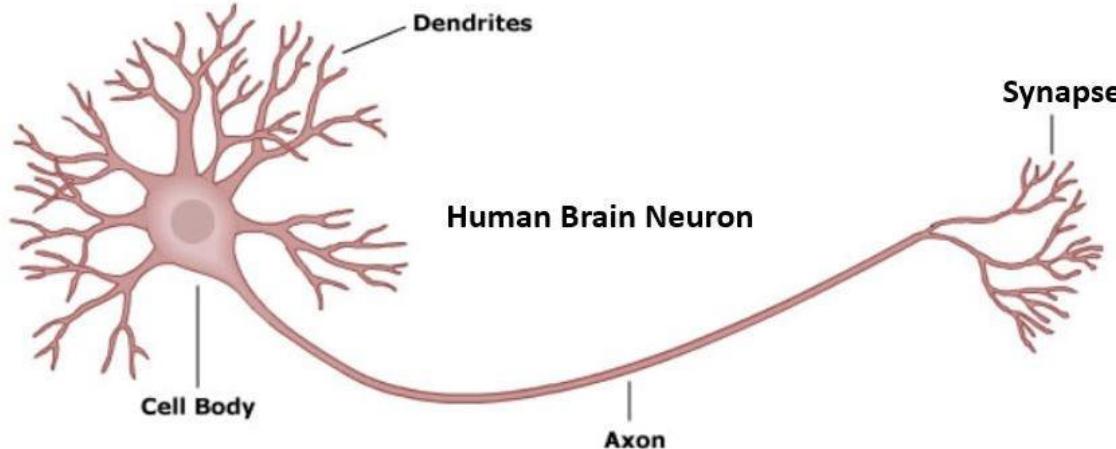
## cons

prone to overfitting (like most things in life)  
can be slow  
not great if classes are very overlapping (e.g., Gaussian mixtures)  
need to be careful if classes are not balanced

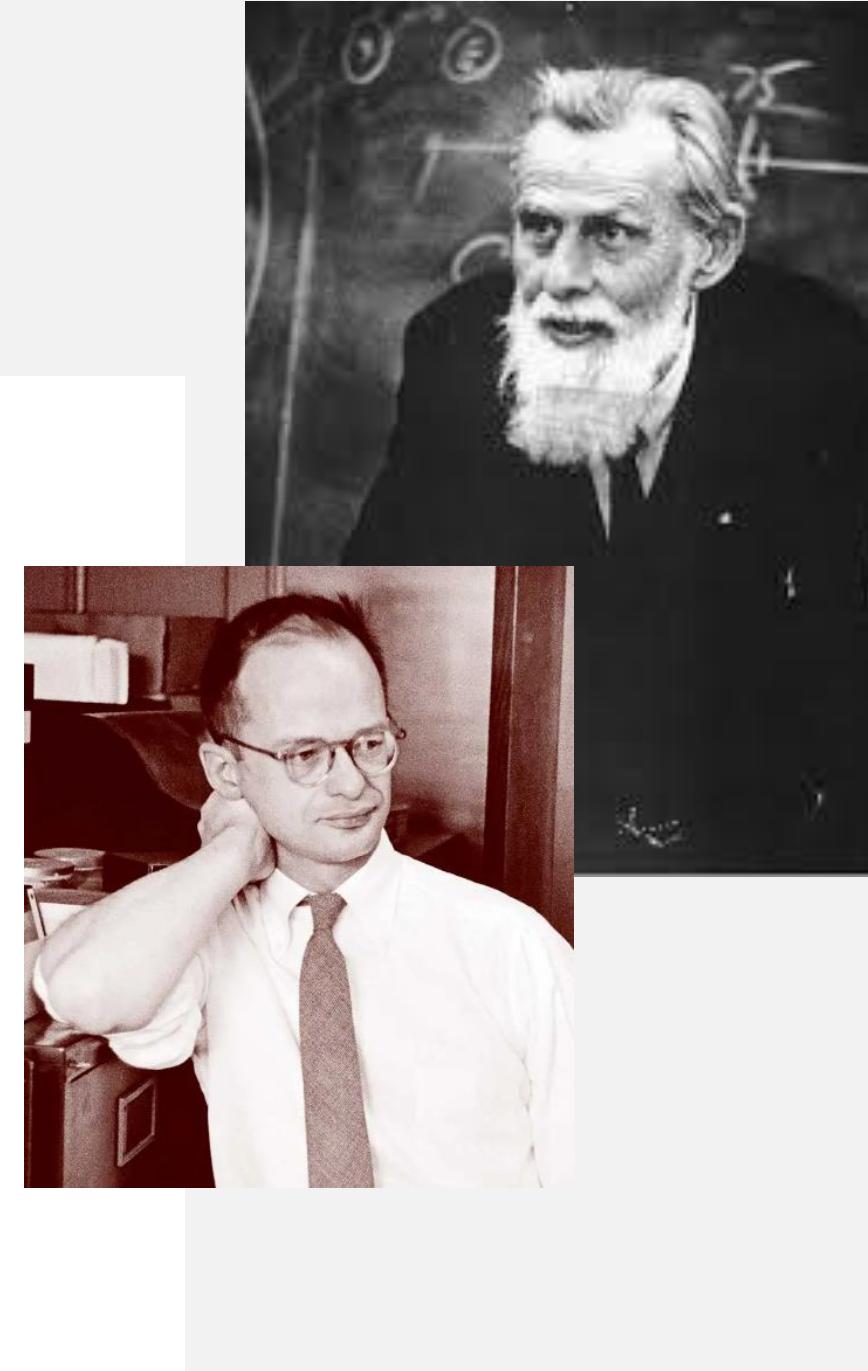
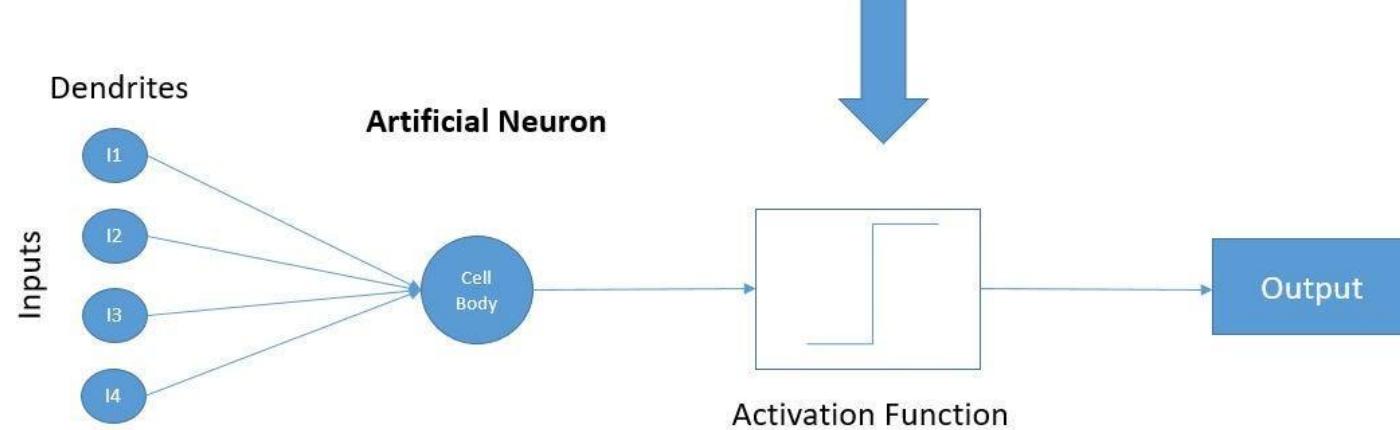
Deep Learning

1943 McCulloch & Pitts

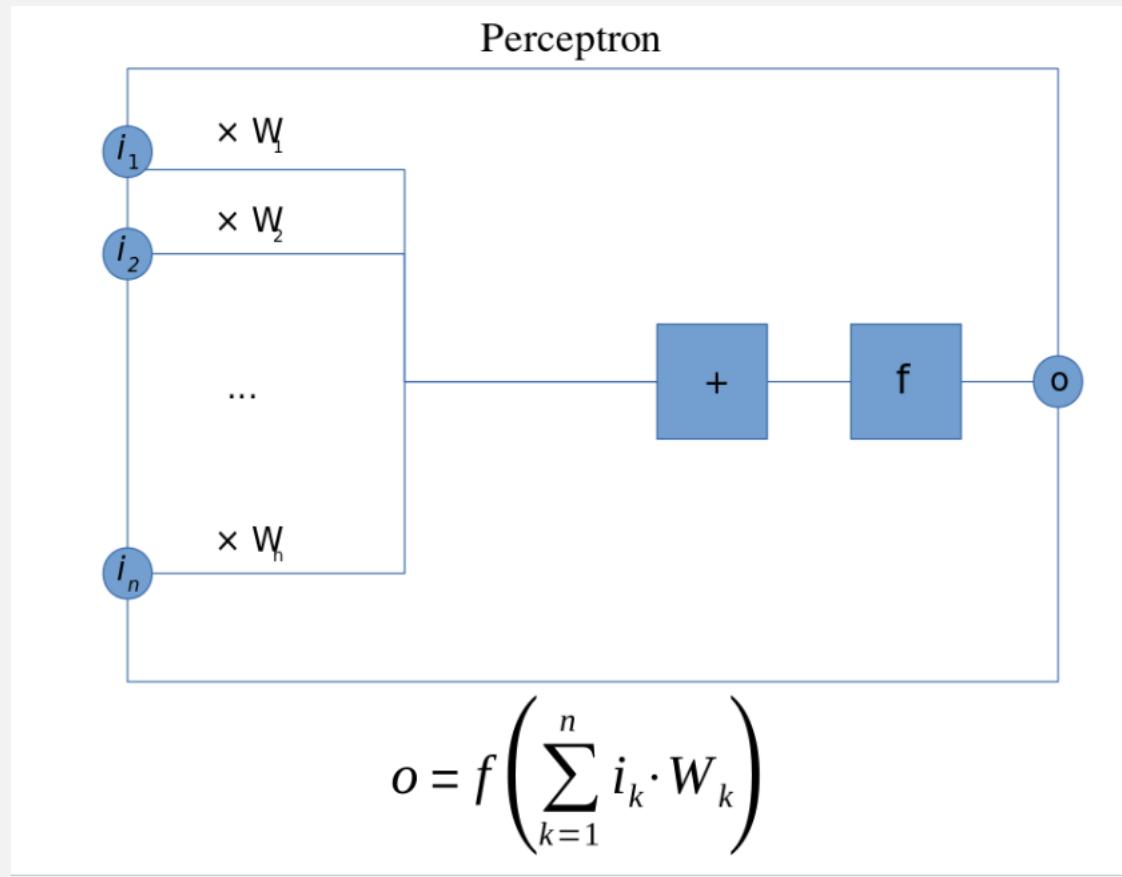
First mathematical model  
of an artificial neuron



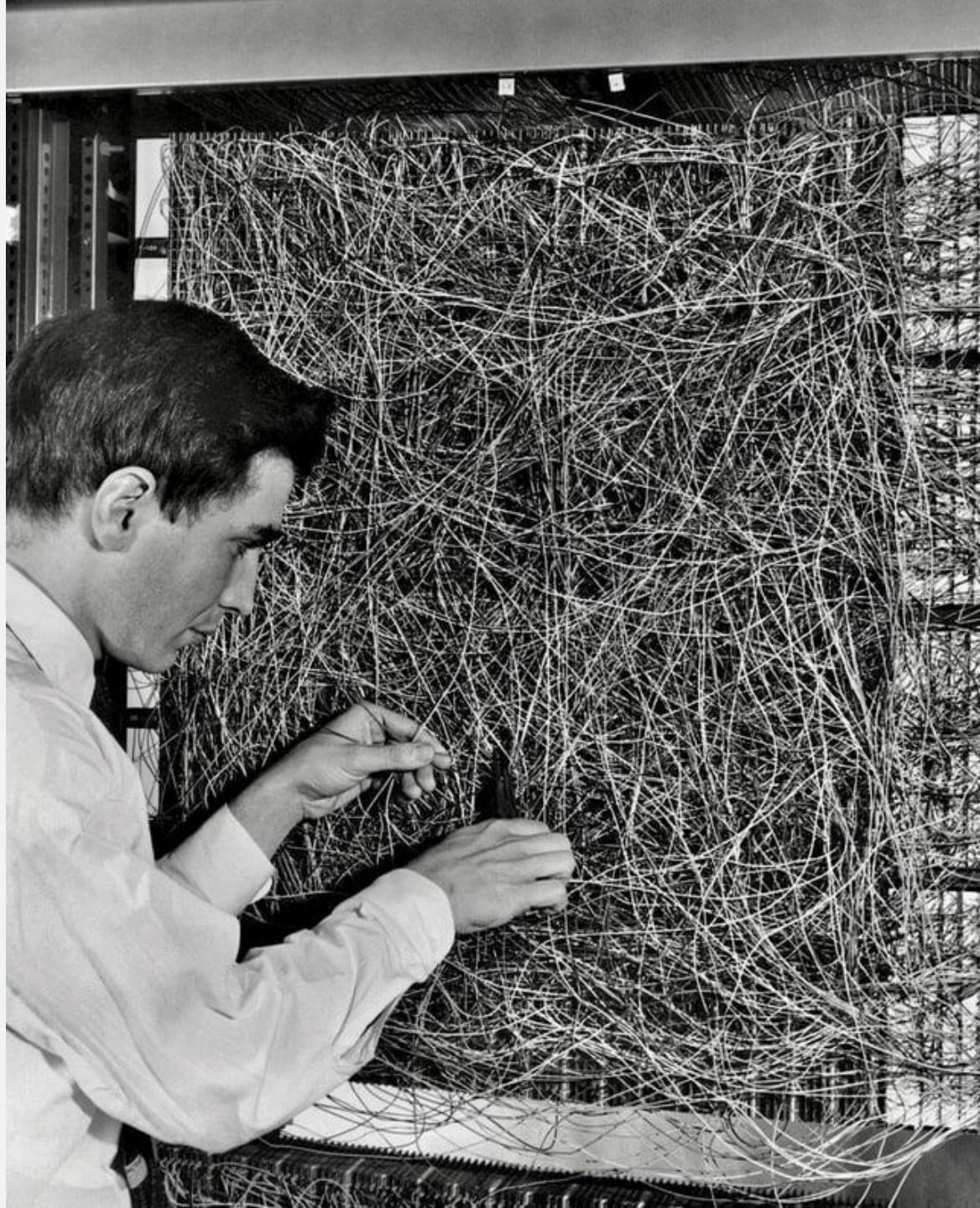
Human Brain Neuron



# 1958 Frank Rosenblatt Perceptron

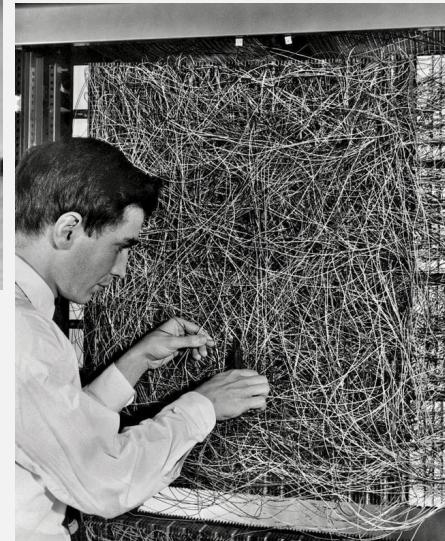


`sklearn.neural_network.MLPClassifier`



## 1958 Minsky & Papert Perceptrons

Minsky and Papert write a book called “Perceptrons”



They show that a single layer of these devices was unable to learn some critical Mathematical functions, such as XOR

In the same book, they also show that using multiple layers of these devices would allow these limitations to be addressed

# ML Breakthroughs

1986

Back-propagation

1989

Universal Approximation  
Theorem

2000+

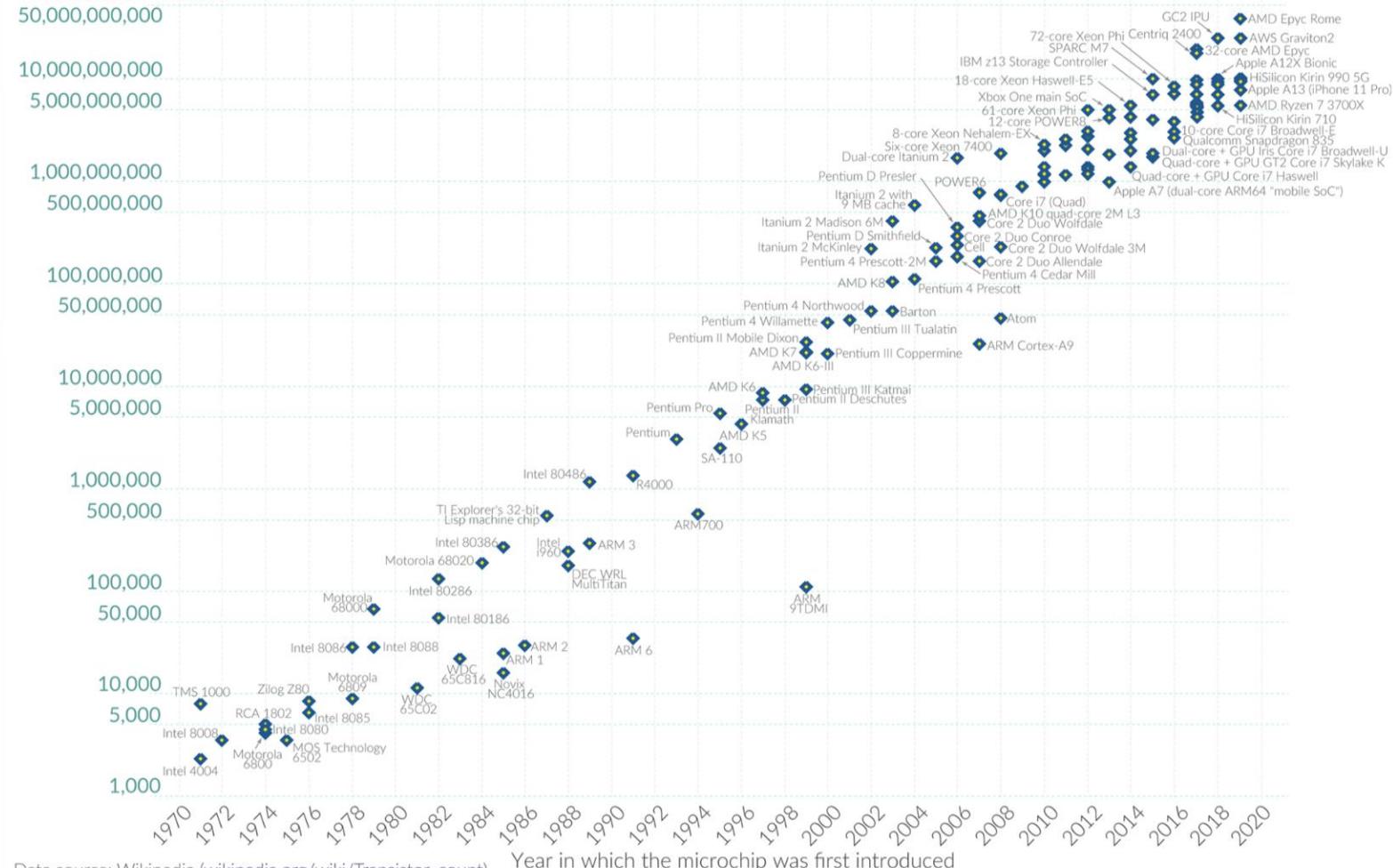
Moore's Law

## Moore's Law: The number of transistors on microchips doubles every two years

Moore's law describes the empirical regularity that the number of transistors on integrated circuits doubles approximately every two years. This advancement is important for other aspects of technological progress in computing – such as processing speed or the price of computers.

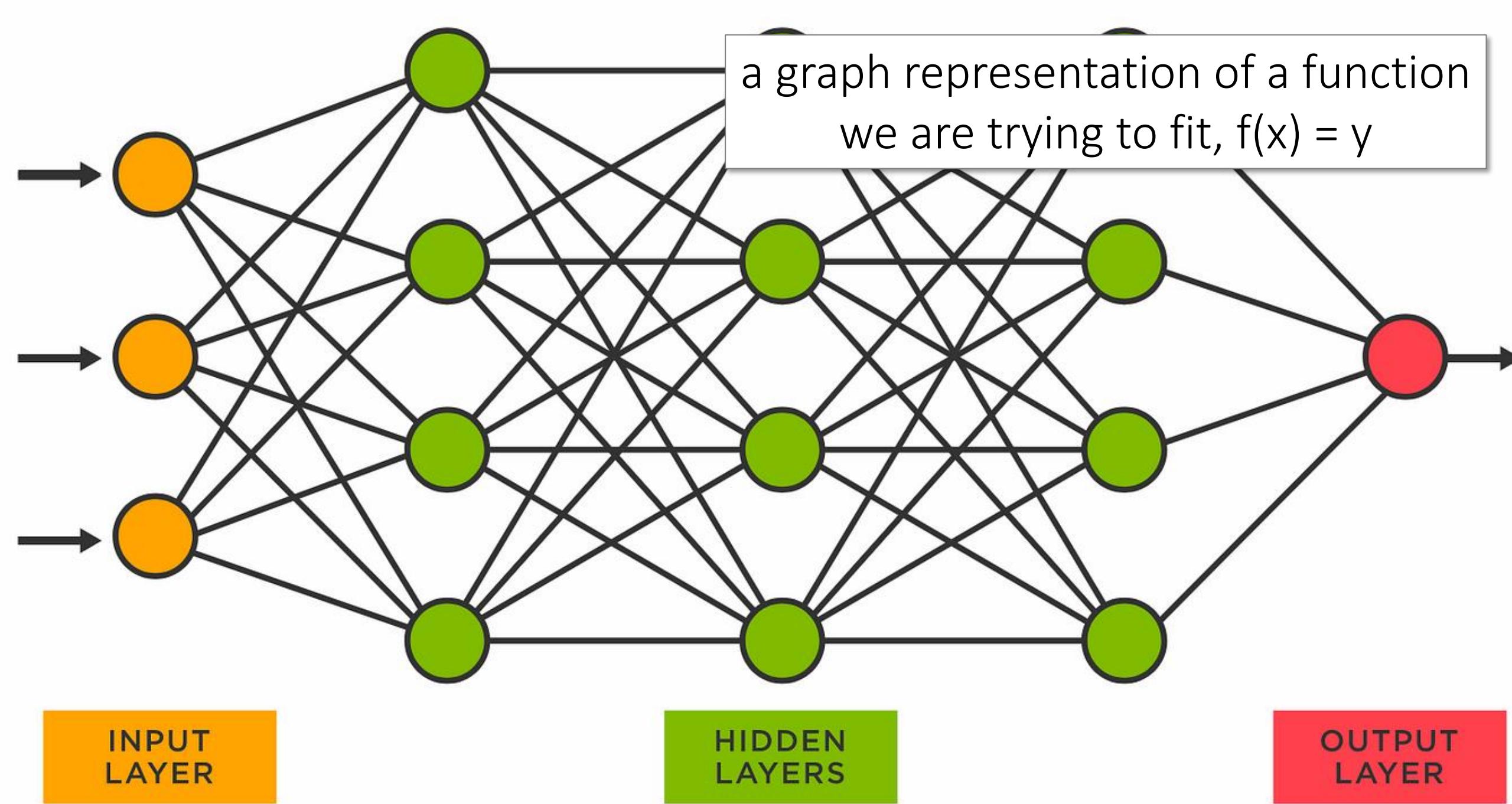
Our World  
in Data

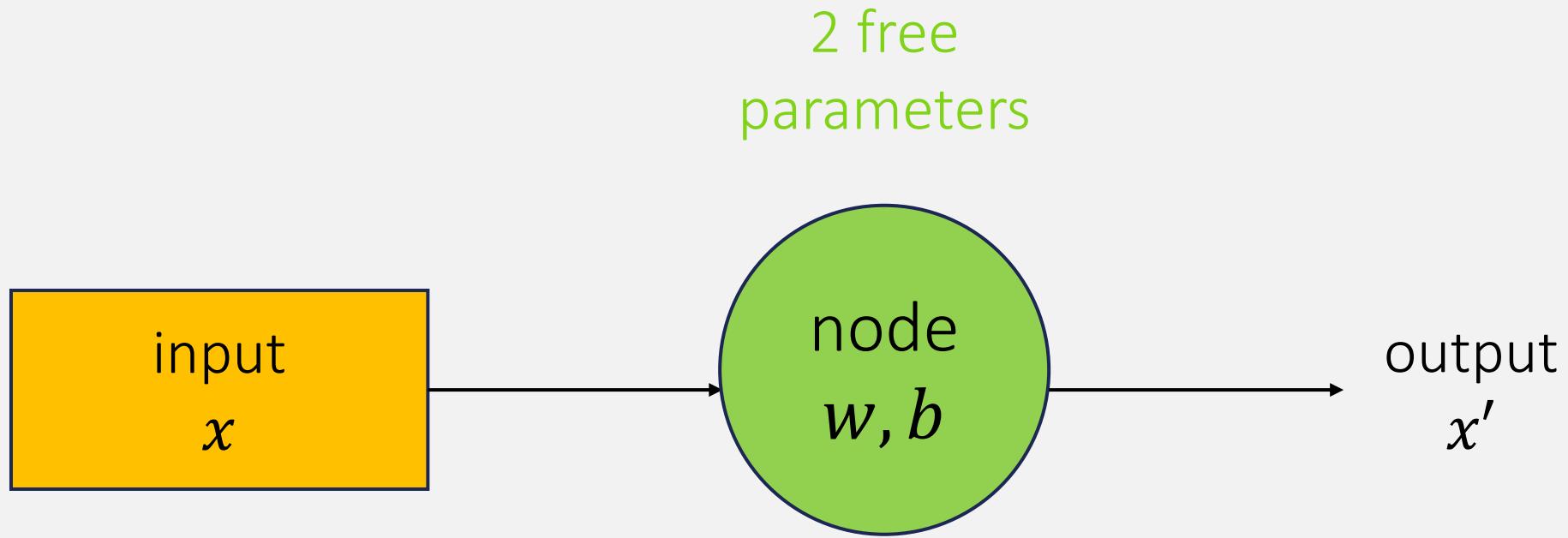
### Transistor count



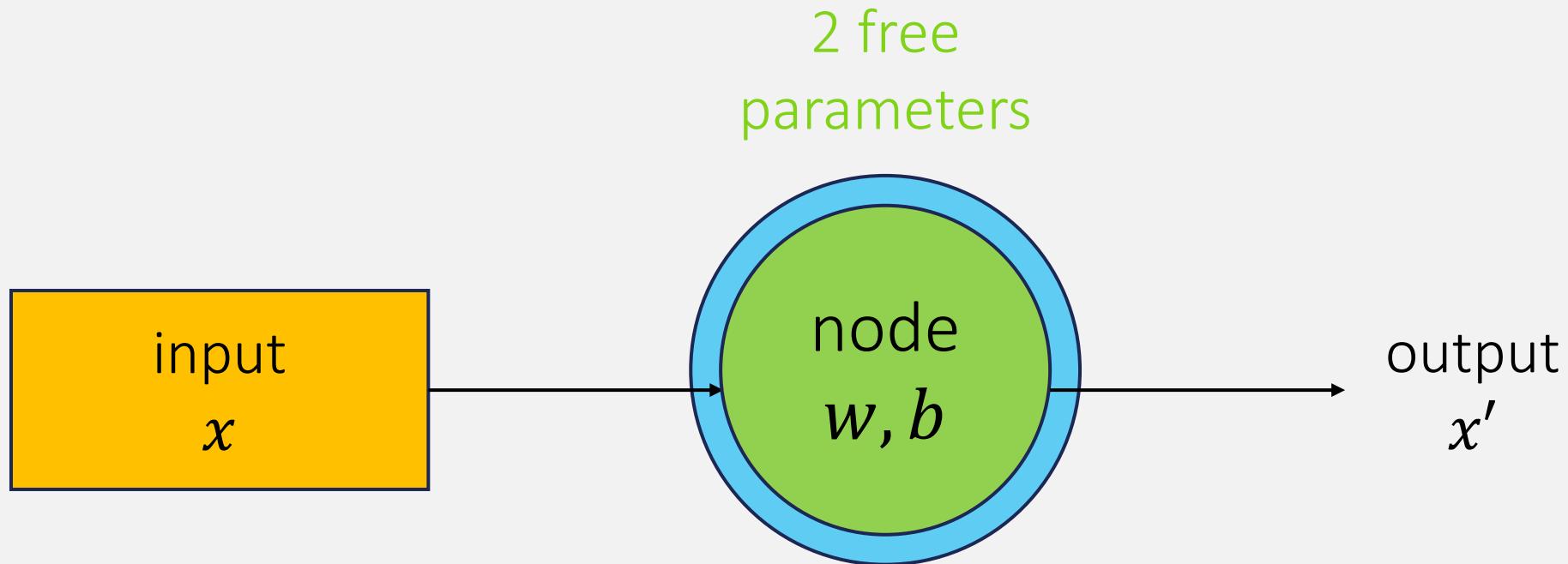
# Machine Learning 101

## deep<sub>ish</sub> neural nets



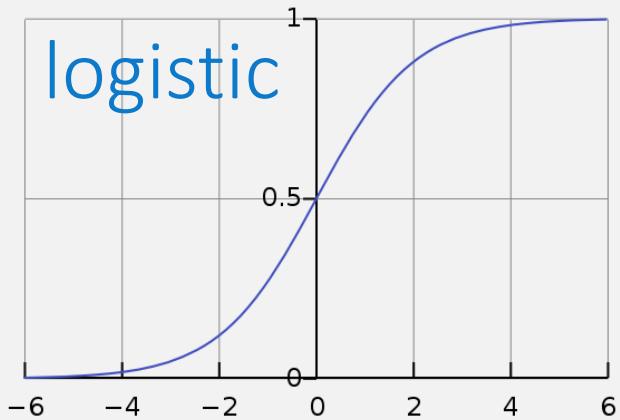


$$x' = \textcolor{green}{wx} + b$$

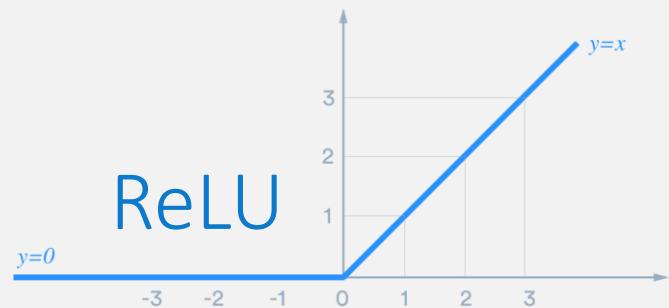
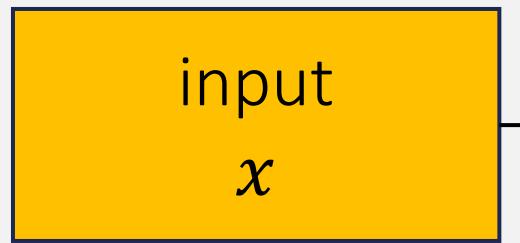
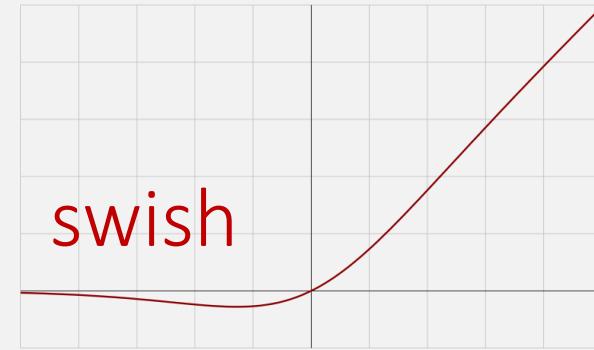


$$x' = f(wx + b)$$

activation function

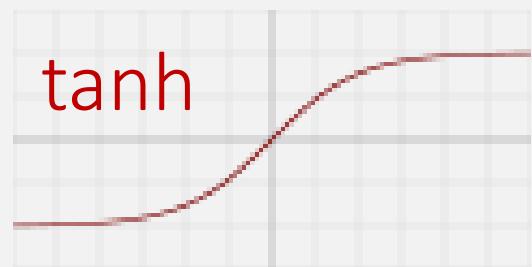


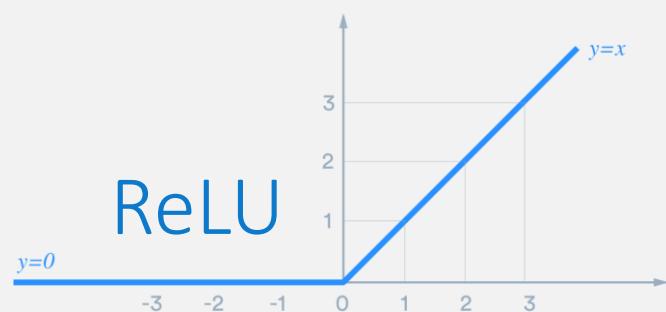
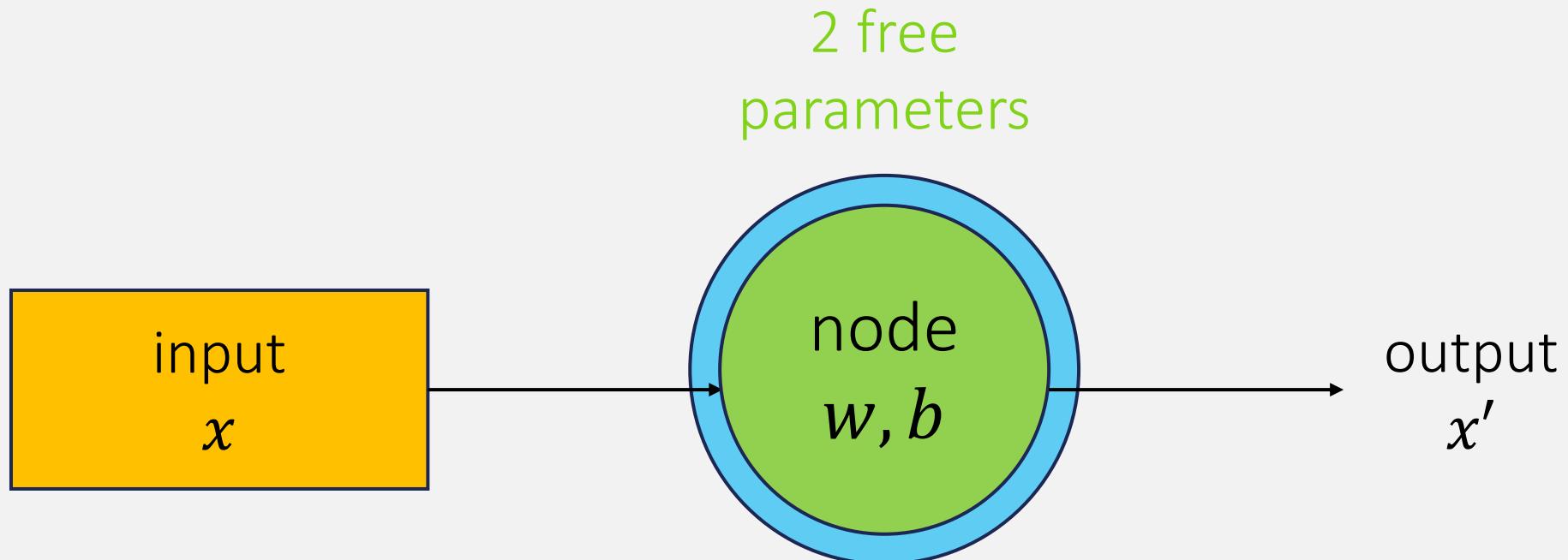
2 free  
parameters



$$x' = f(wx + b)$$

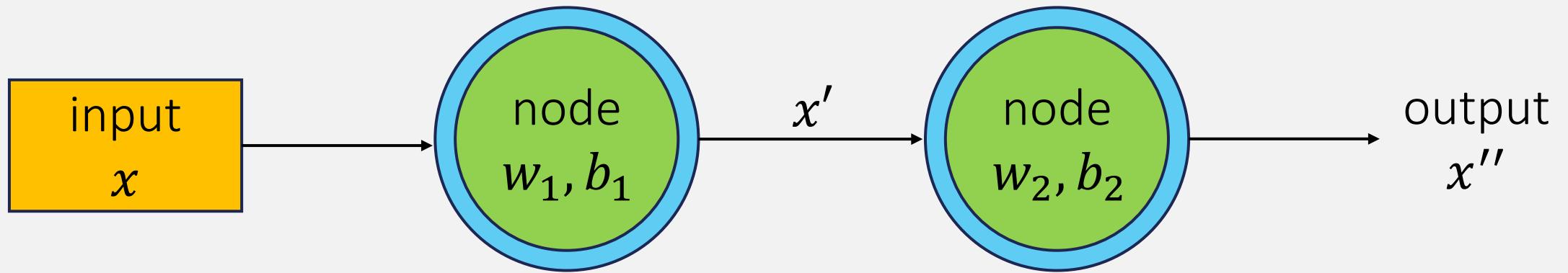
activation function



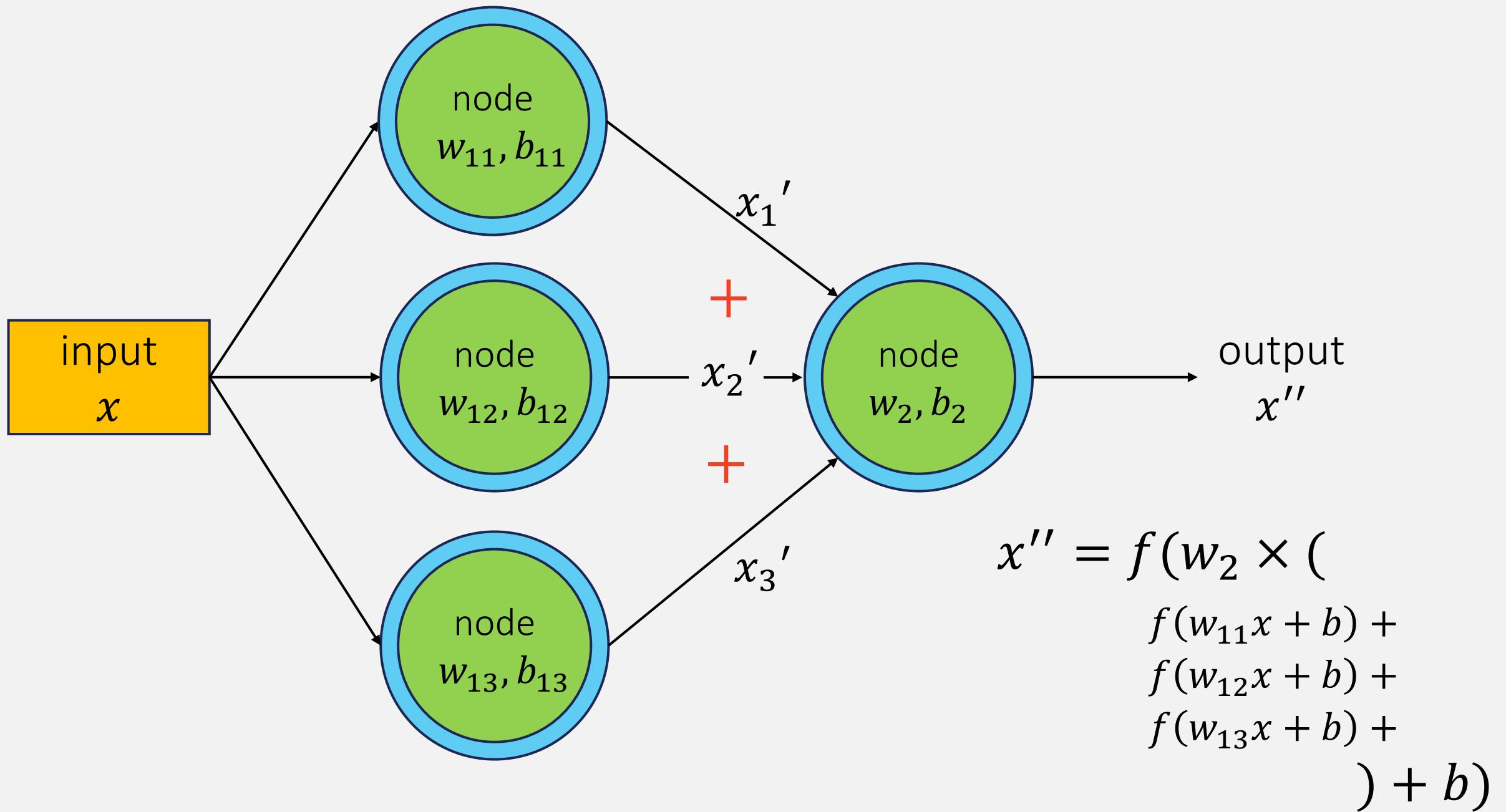


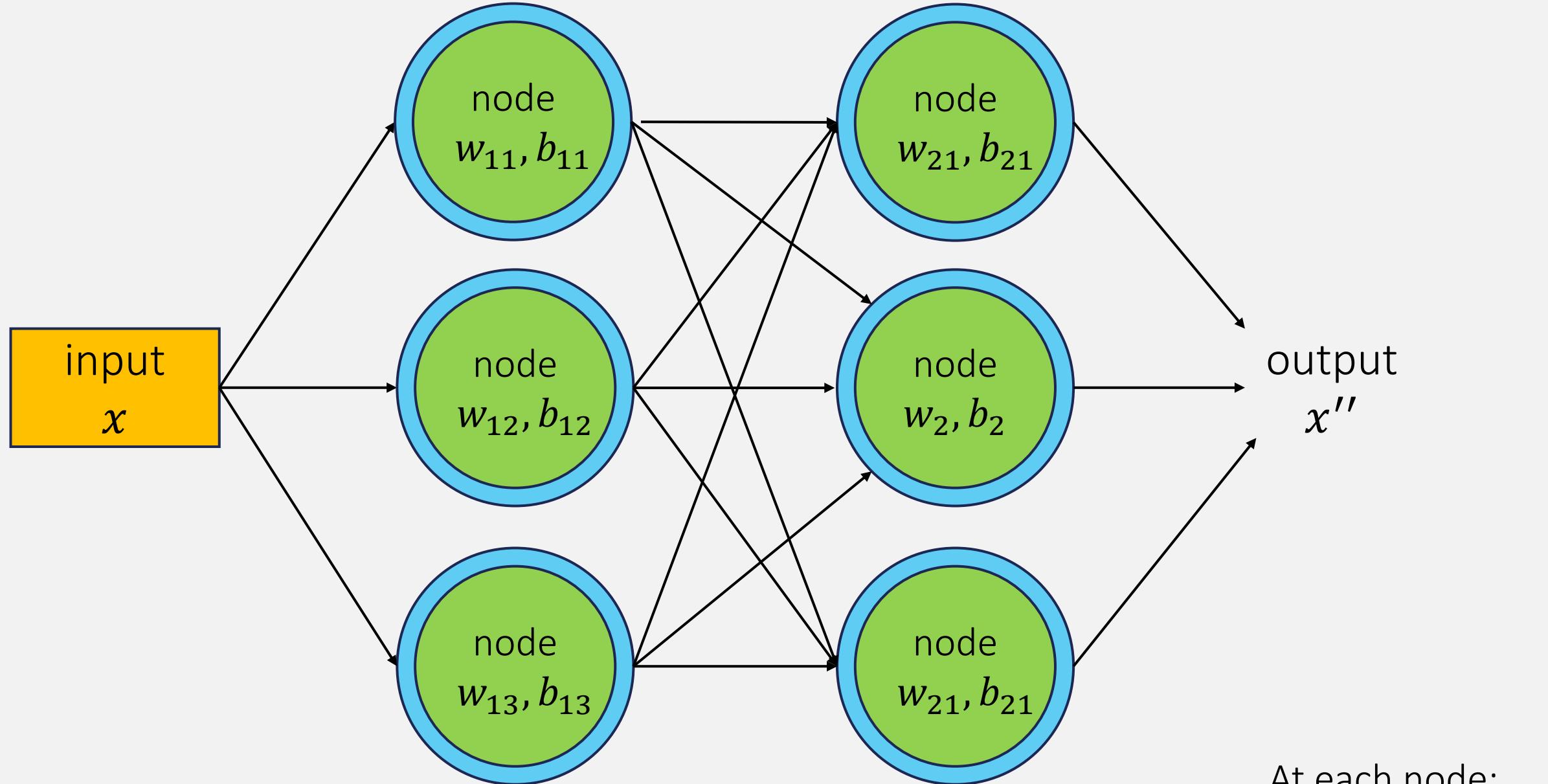
$$x' = f(wx + b)$$

activation function

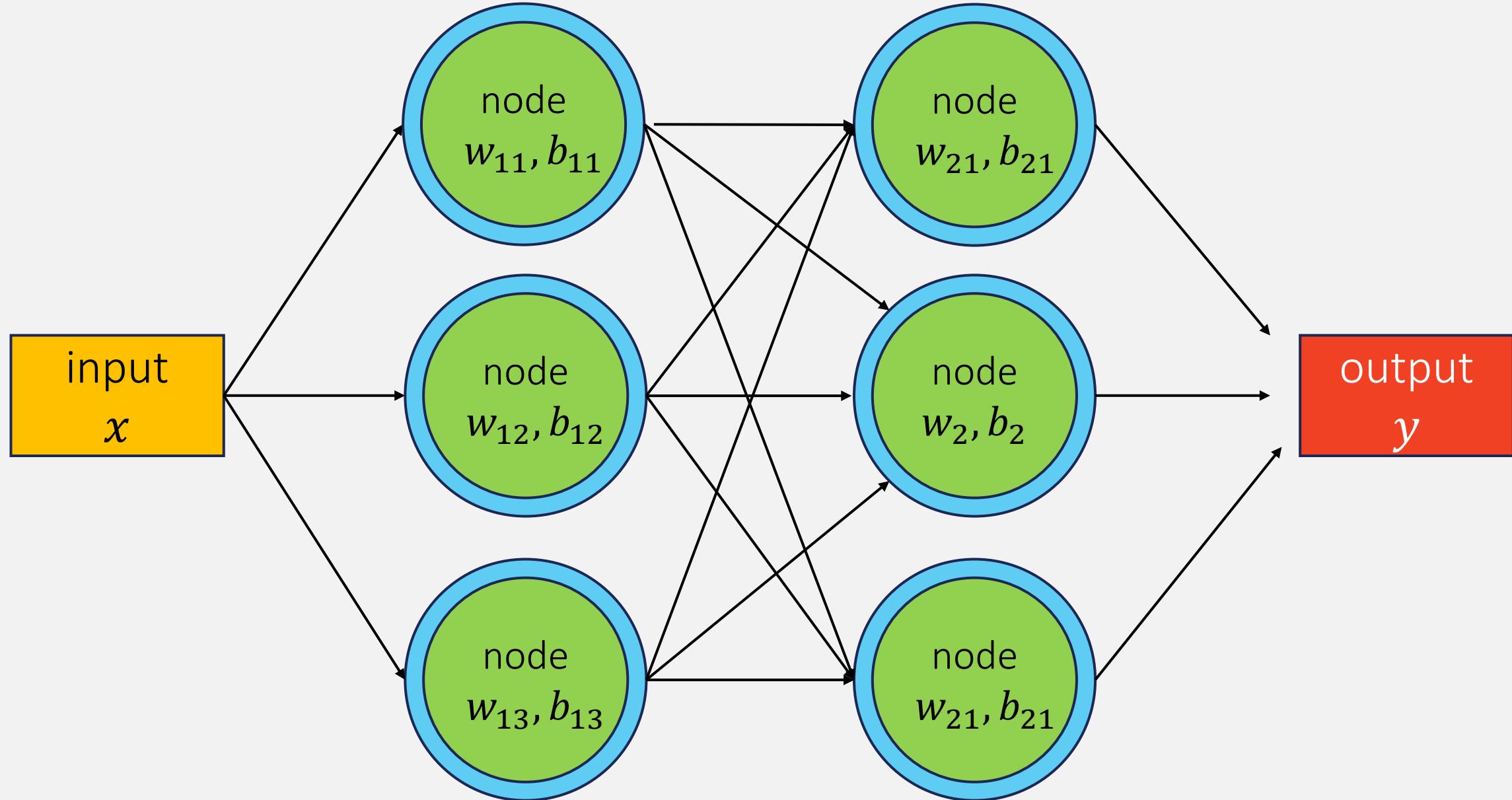


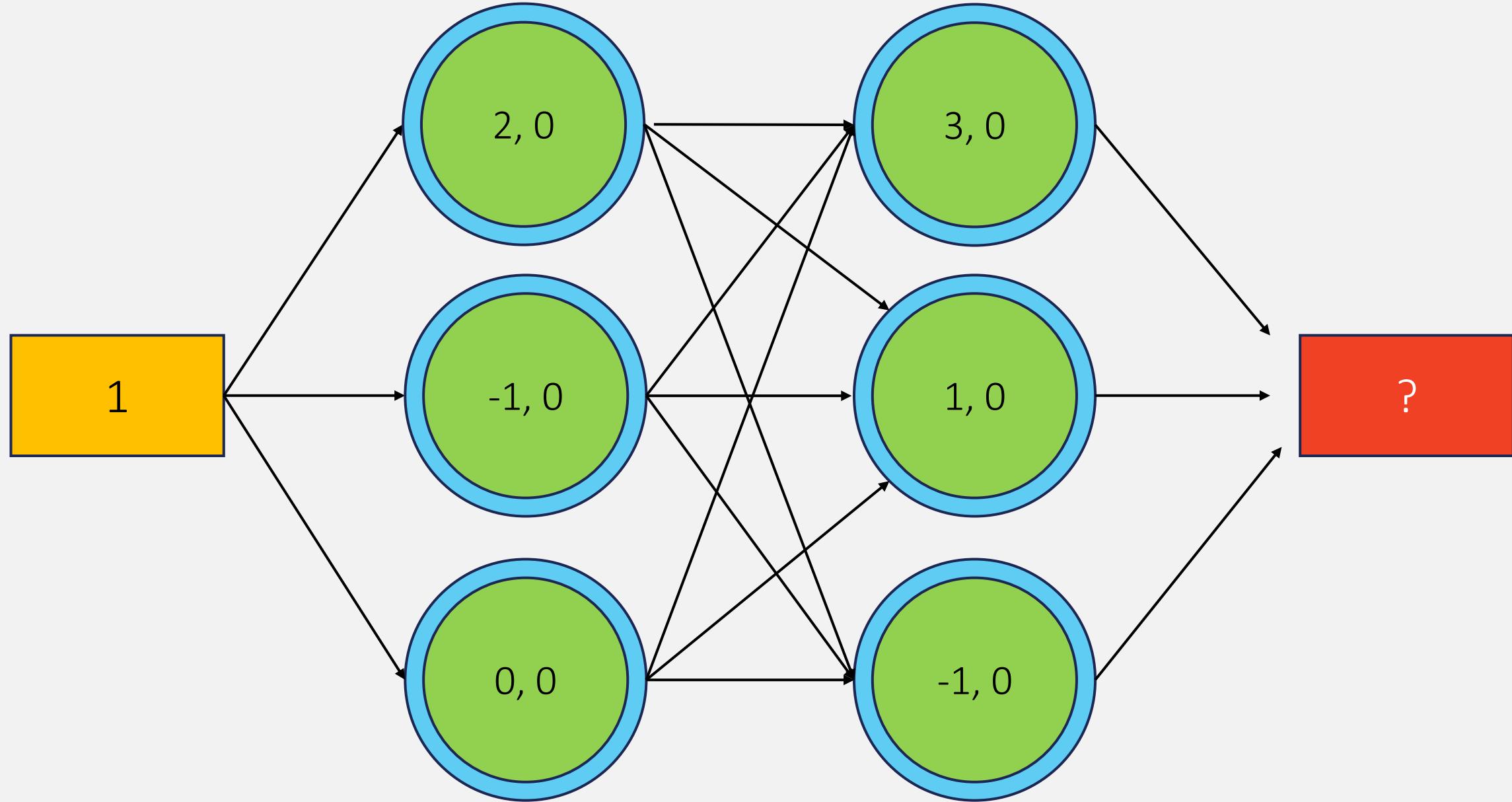
$$\begin{aligned}x'' &= f(w_2 x' + b_2) \\&= f(w_2 f(w_1 x + b_1) + b_2)\end{aligned}$$

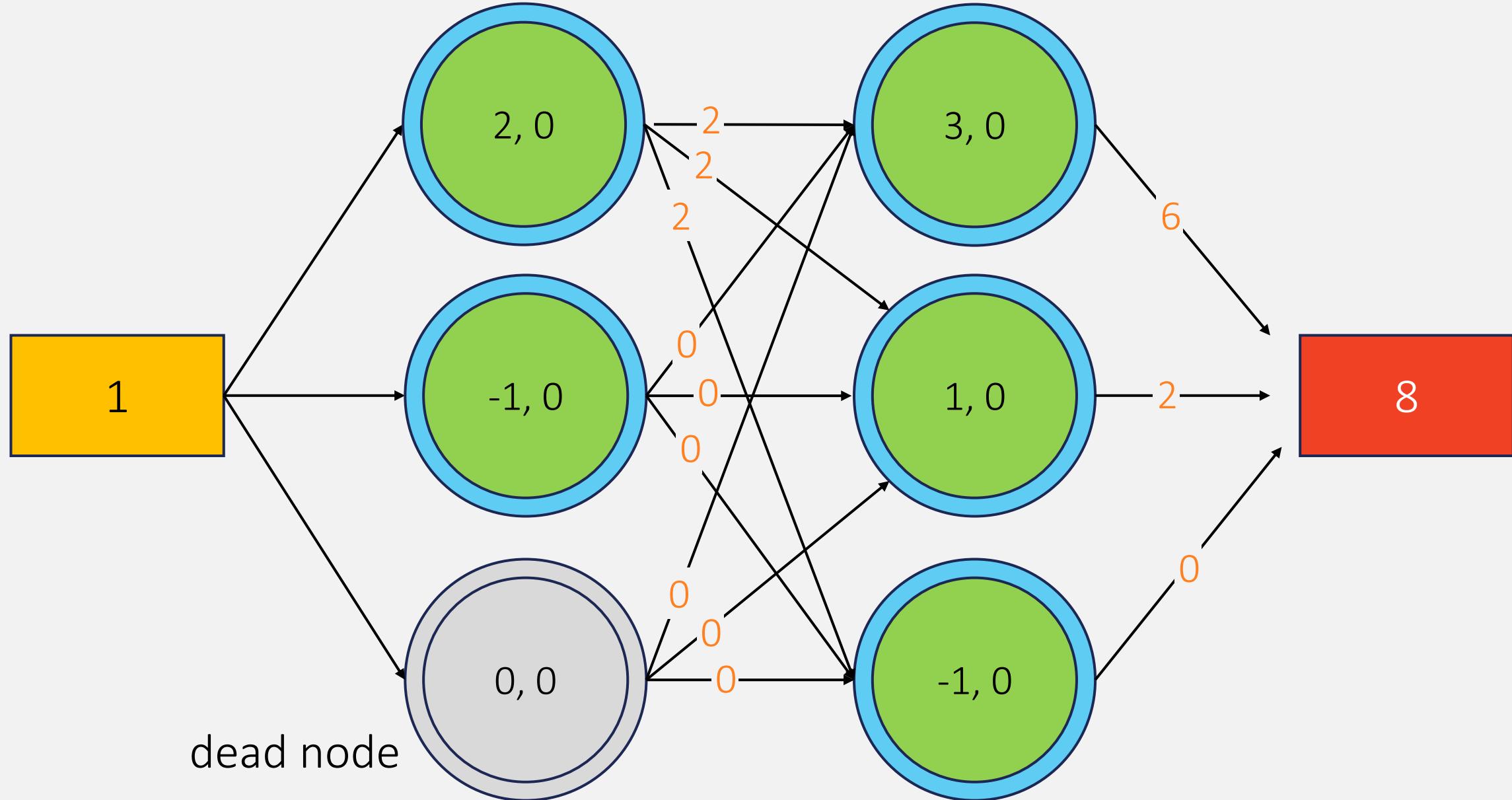


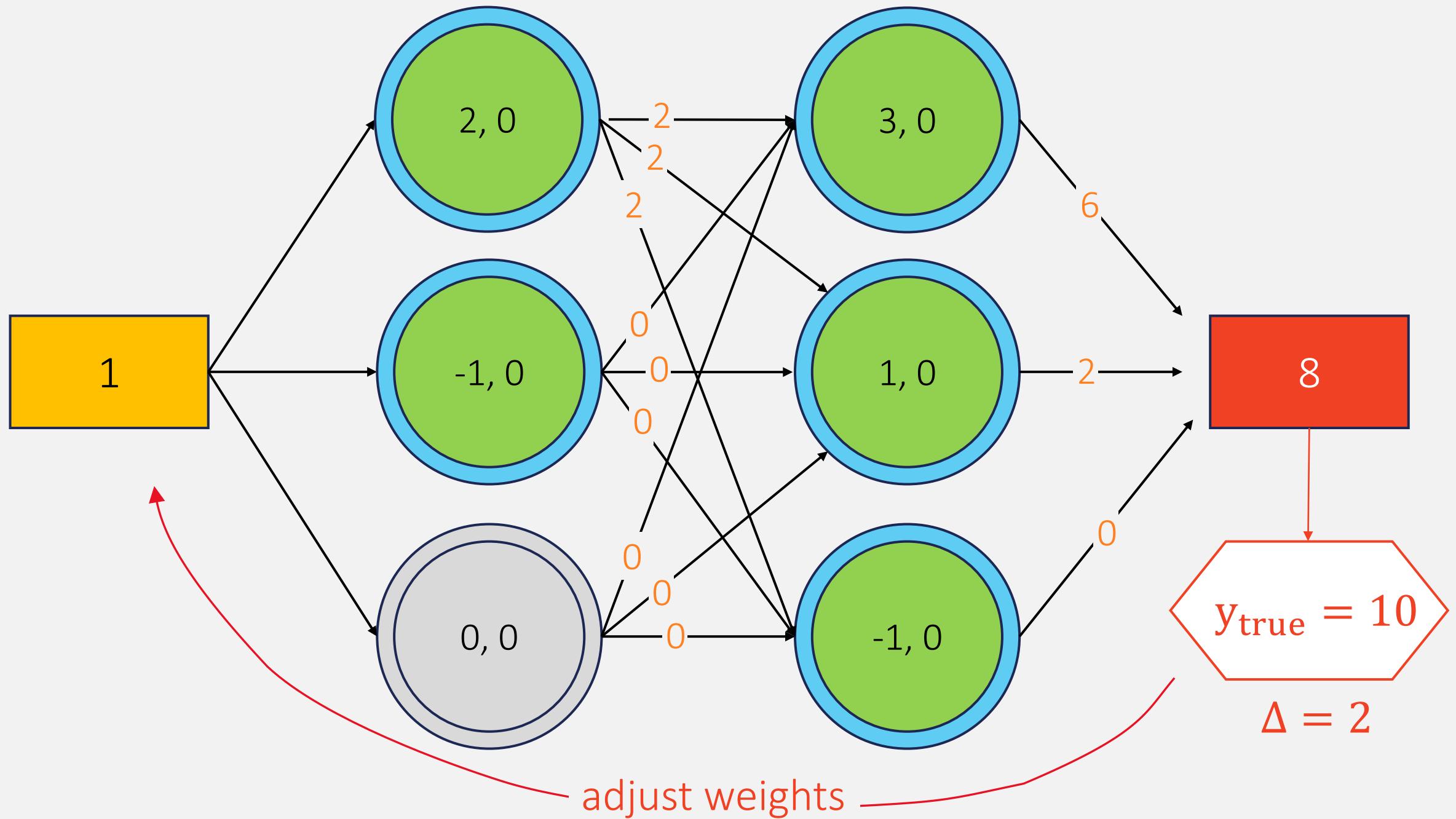


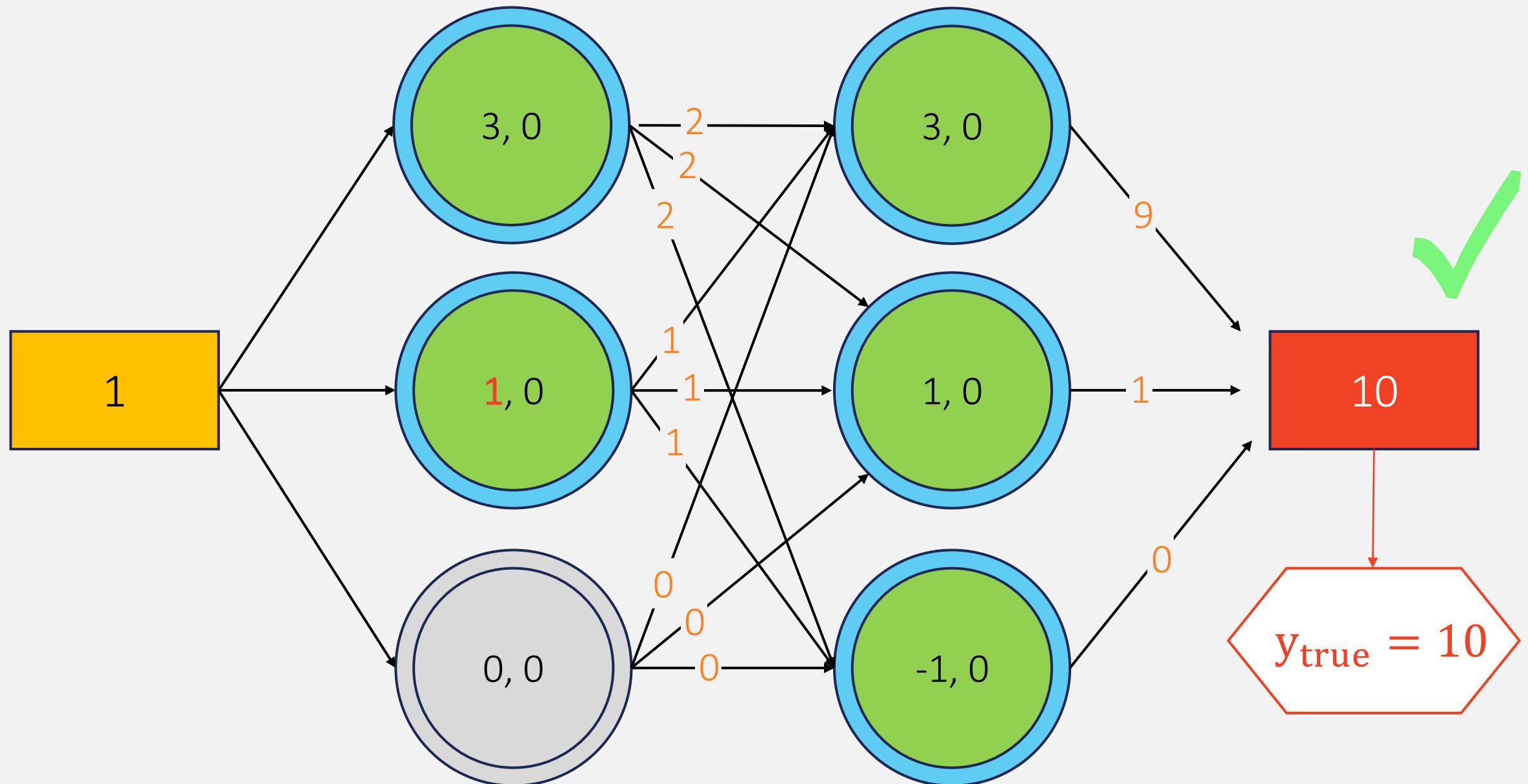
At each node:  
 $y = f(wx + b)$



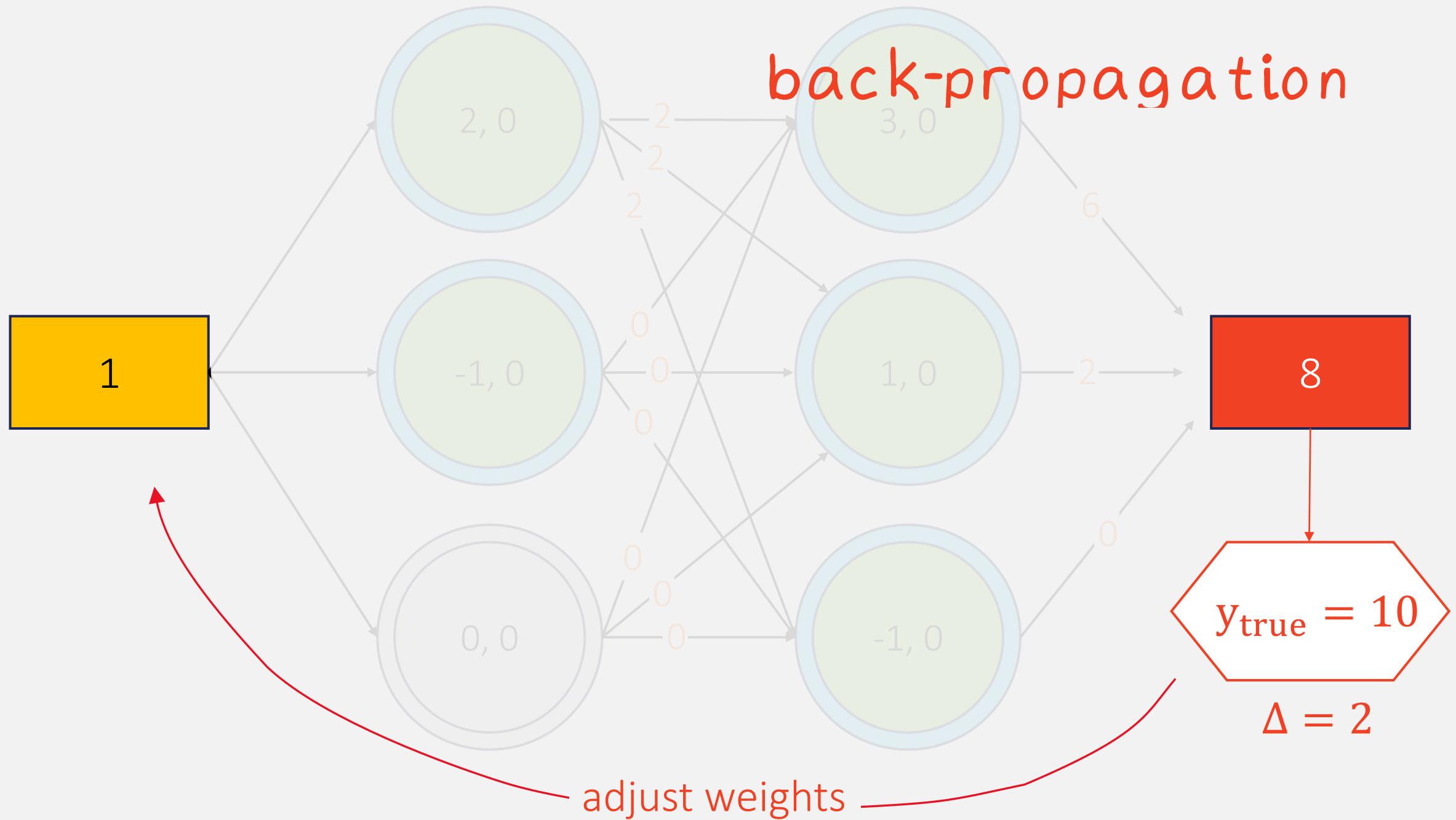




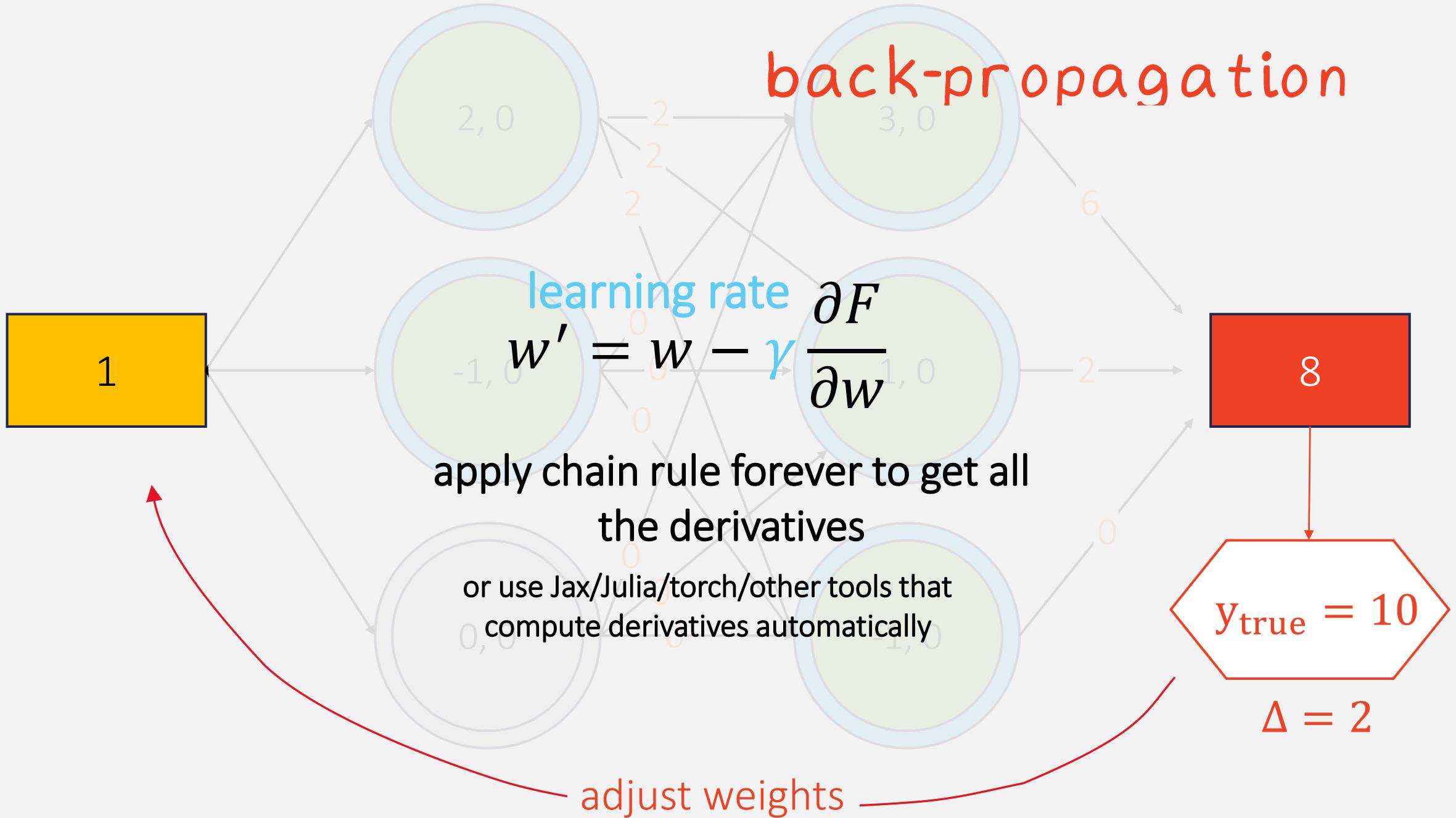




back-propagation



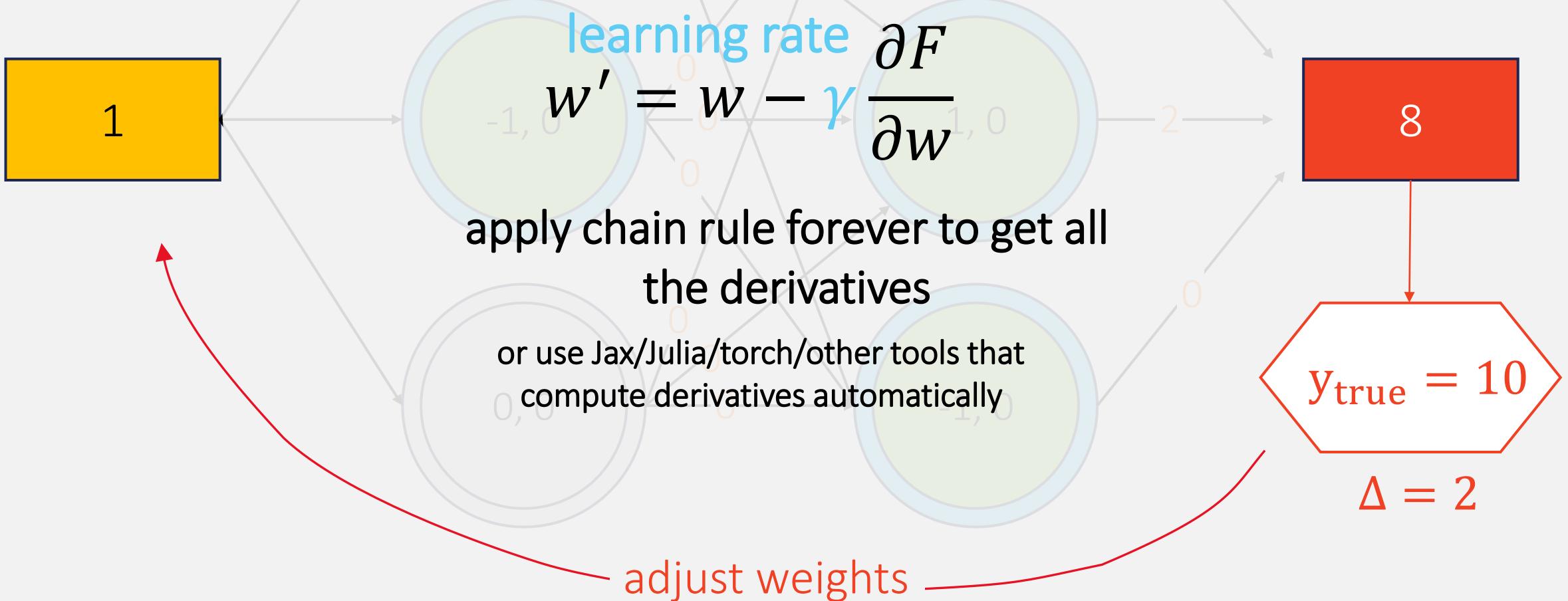
# back-propagation



Different optimization choices

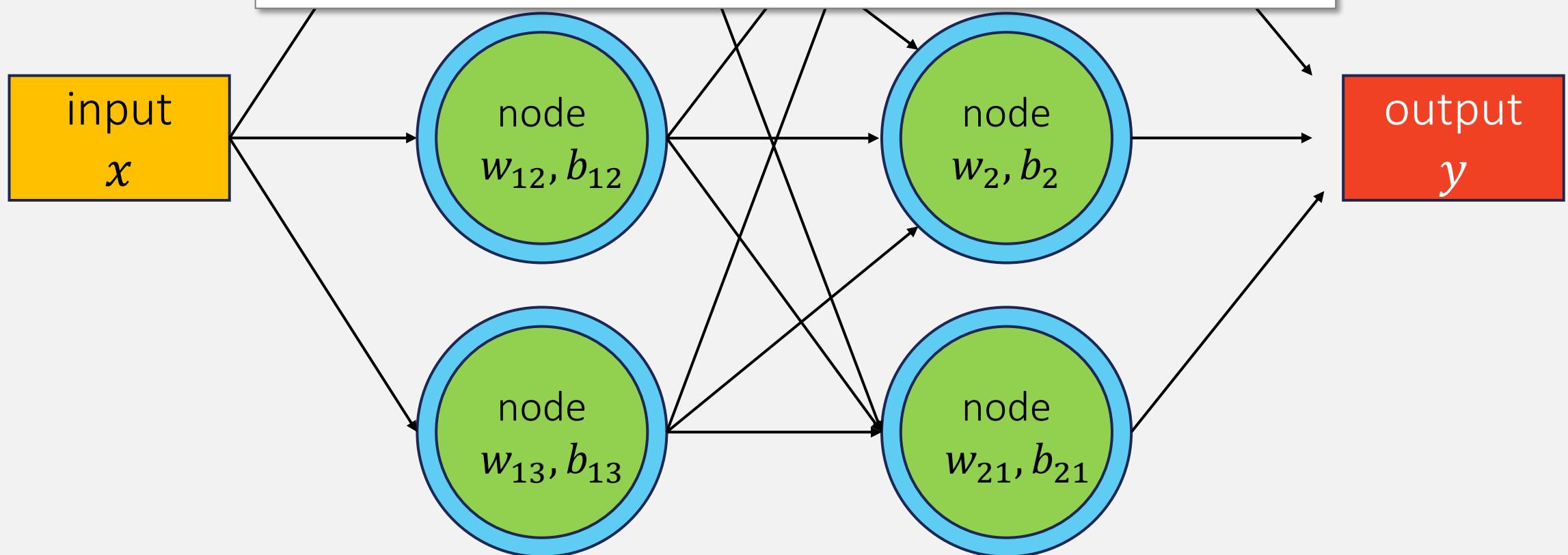
- Stochastic gradient descent
- Learning rate optimization (Adam, Adamax, ...)

# back-propagation



## fully-connected networks

these networks are basic but still great for most tasks:  
classification, regression with simple data  
(vectors, time series, small images)



`sklearn.neural_network.MLPClassifier`

`sklearn.neural_network.MLPRegressor`

## fully-connected networks

these networks are basic but still great for most tasks:  
classification, regression with simple data  
(vectors, time series, small images)

### pros

- Few ( $\sigma(100s)$ ) free parameters
- Fast: no need for a GPU
- Easy to implement
- Still capture a lot of complexity
- Easy to adapt for different tasks
- Can be pre-trained

### cons

- Not as powerful as complex NNs
- Still prone to overfitting

`sklearn.neural_network.MLPClassifier`

`sklearn.neural_network.MLPRegressor`

# overfitting

NNs have 100s to  $\sigma(10^6)$  free parameters – overfitting can be a real problem

1

Always keep your train/test sets separate

Easy to mix them up if data becomes more complex or is generated on the fly

2

Neural networks train in batches

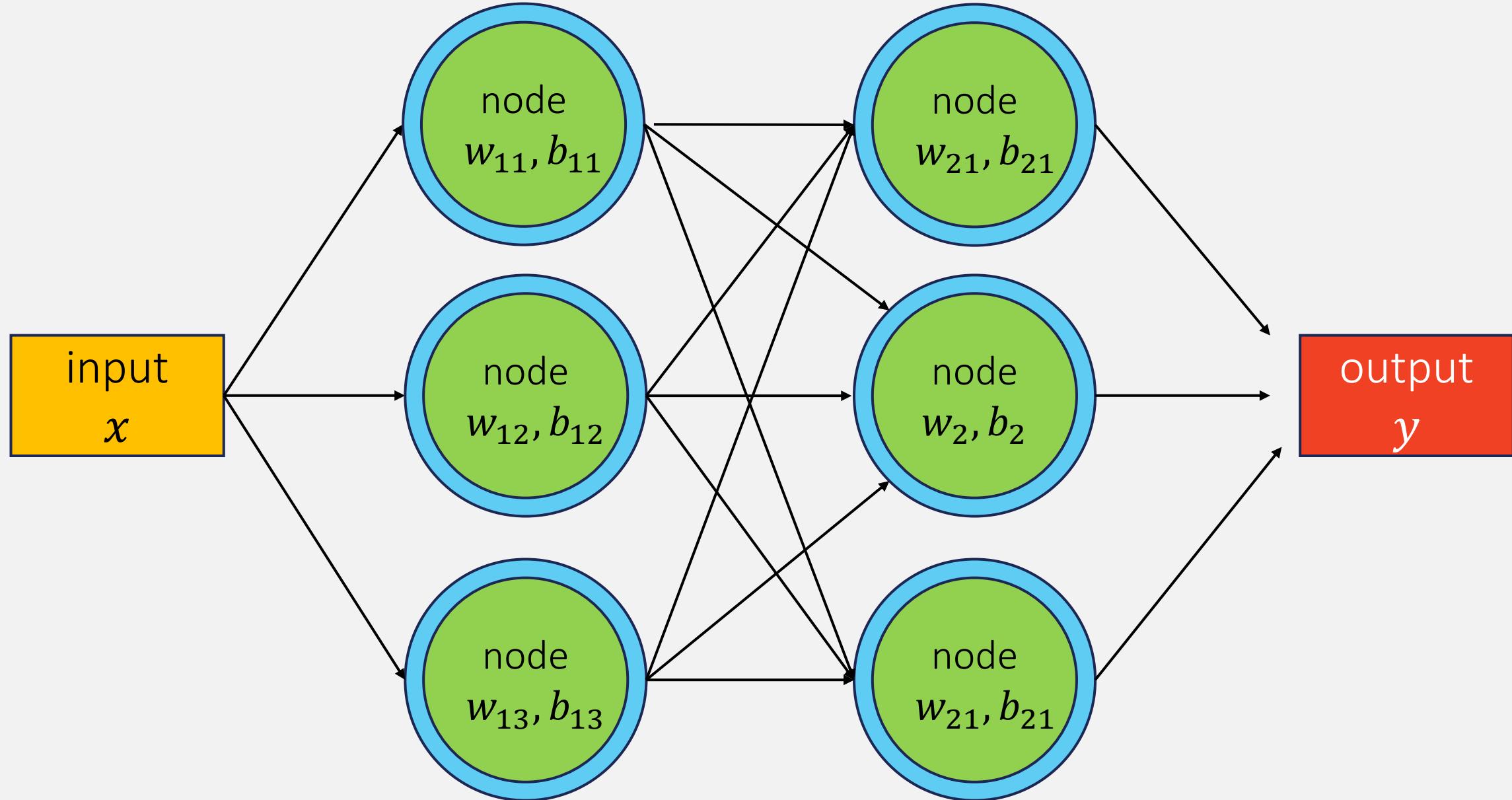
Take a subsample → run inference → calculate loss → adjust weights

Sometimes a subsample is just one input vector!

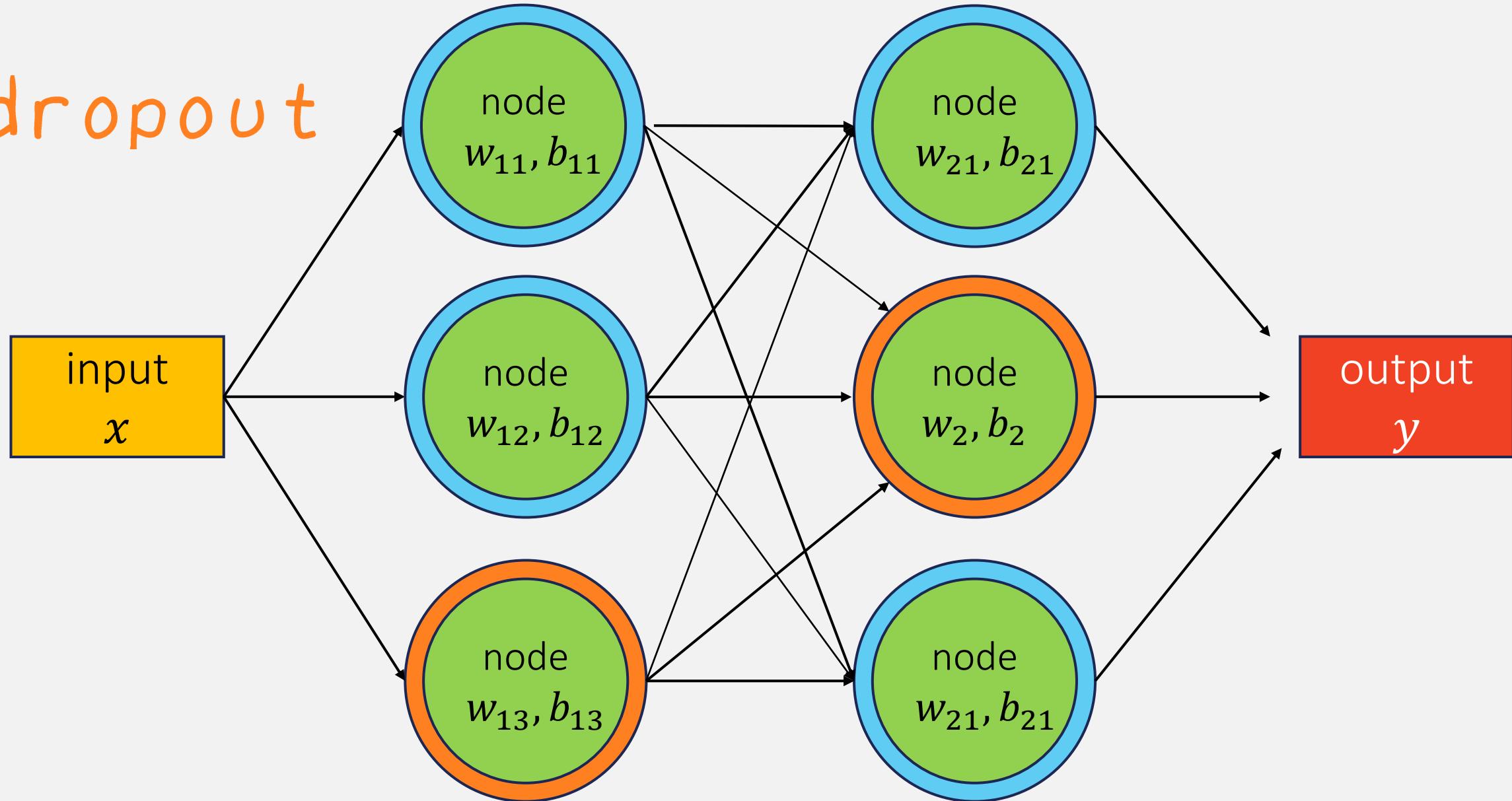
3

Dropout

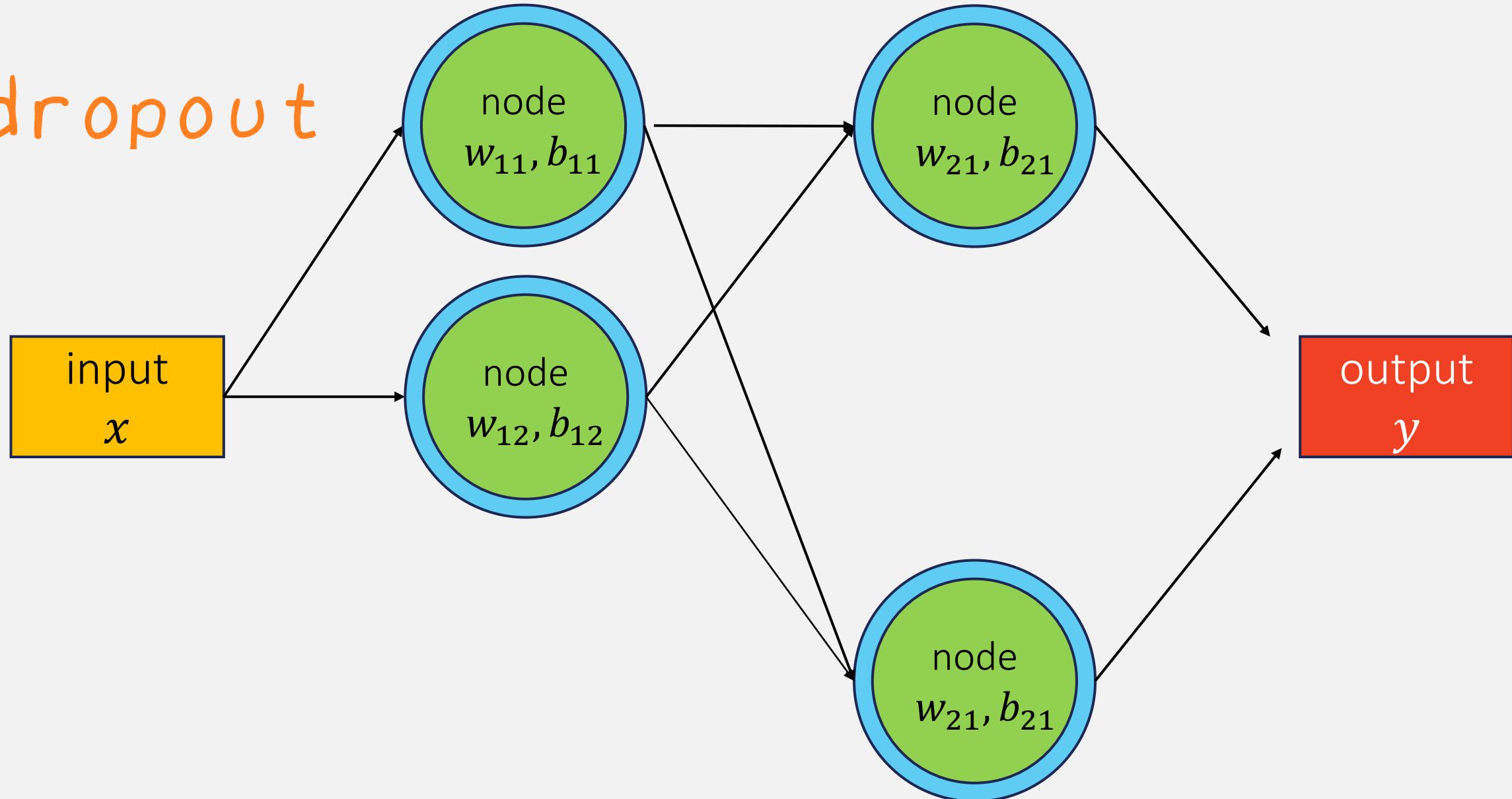
Adds stochasticity to the network



dropout



dropout



# dropout

1

Assign probability of dropout,  $p(\text{drop}) \in (0,1)$  to the nodes you want  
Usually this will range between 0.1 and 0.5. Higher values = more generalized, longer to train

2

At each training step (each batch), randomly drop nodes  
Different nodes will be dropped at different iterations

3

The model should generalize better  
As different nodes won't learn to overcompensate each other

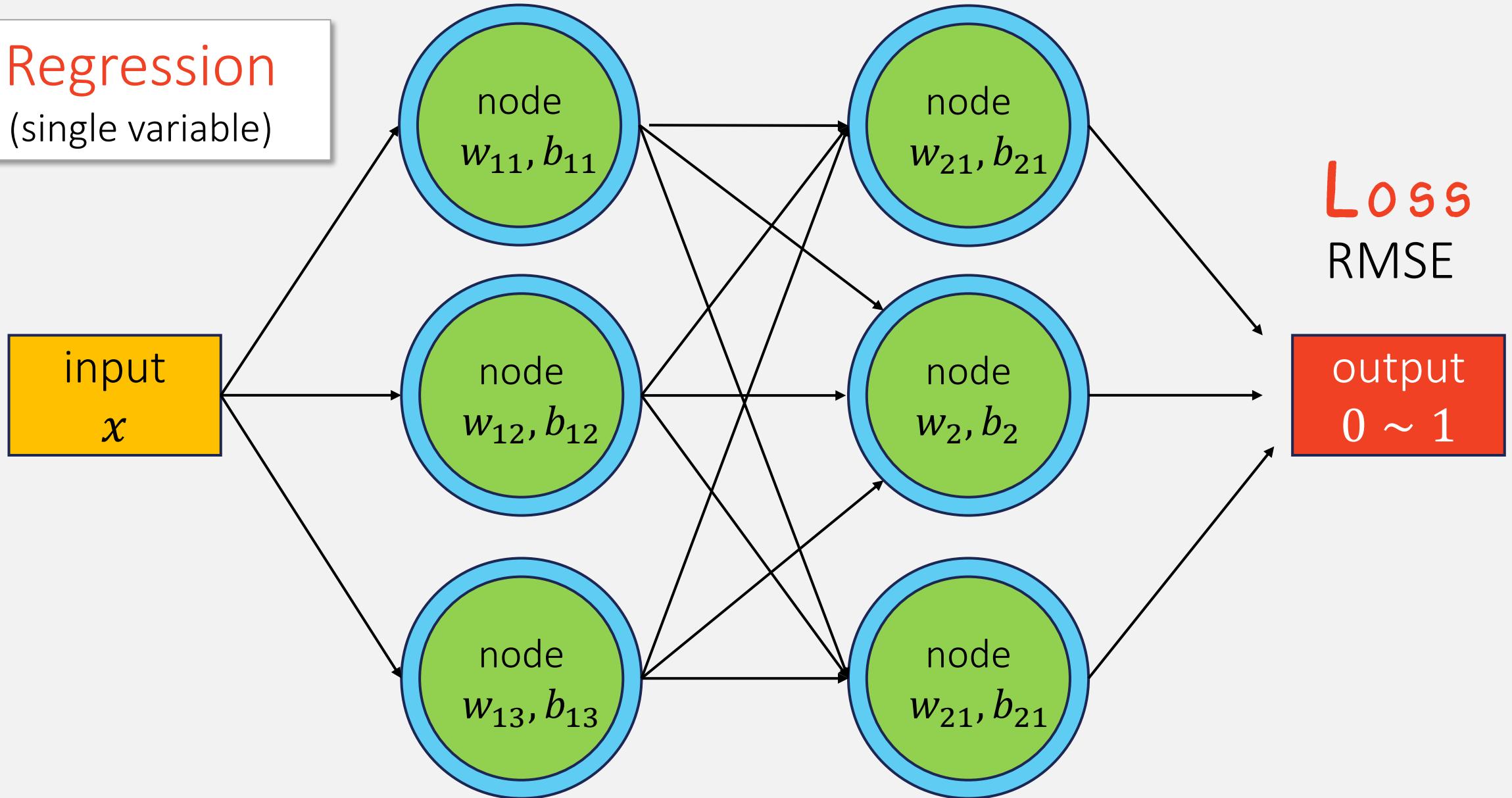
4

Turn dropout off when running inference  
Produces consistent results

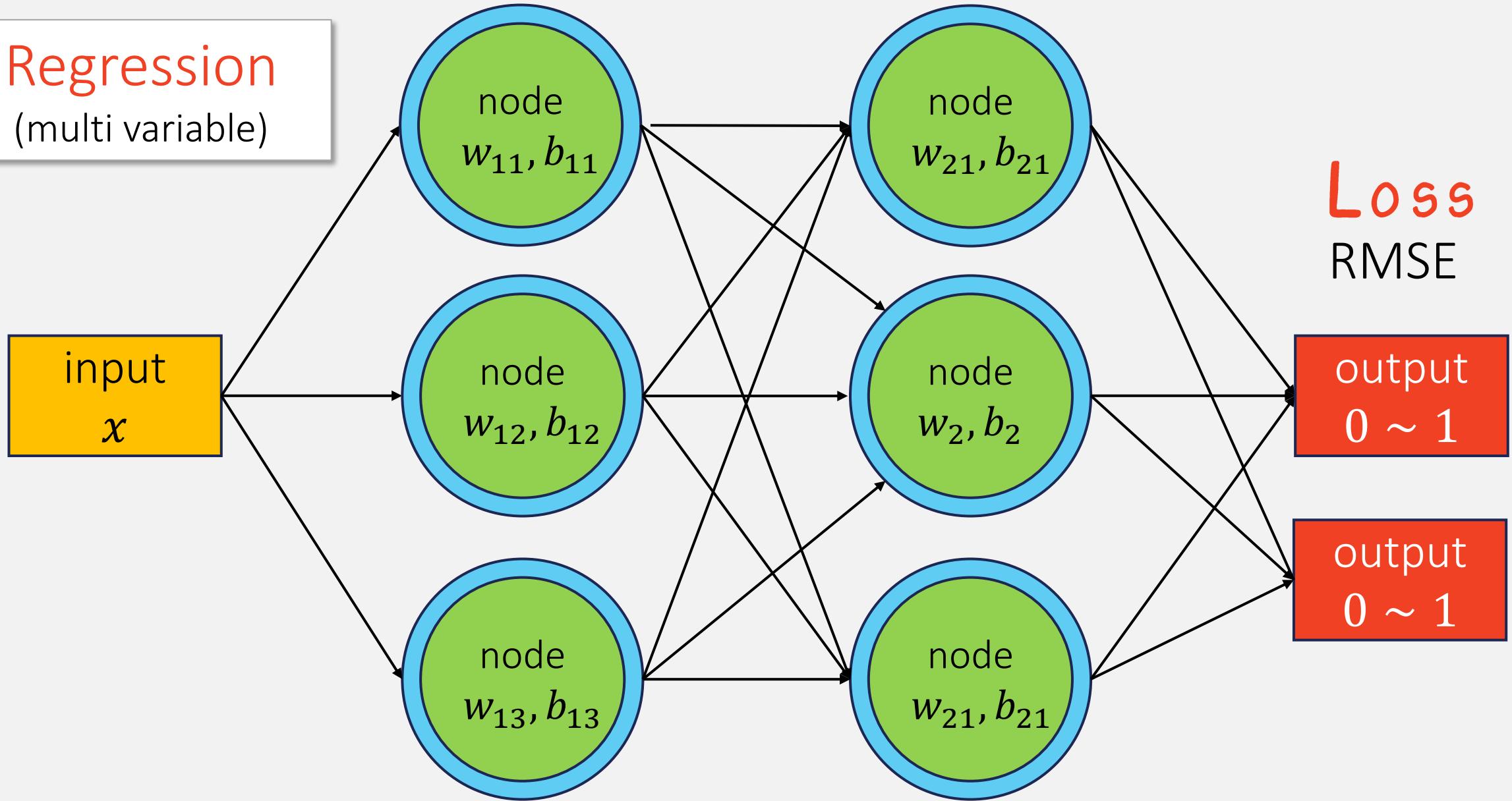
## NN use cases

what the network learns to do depends mostly on the labels, the output layer, and the loss function

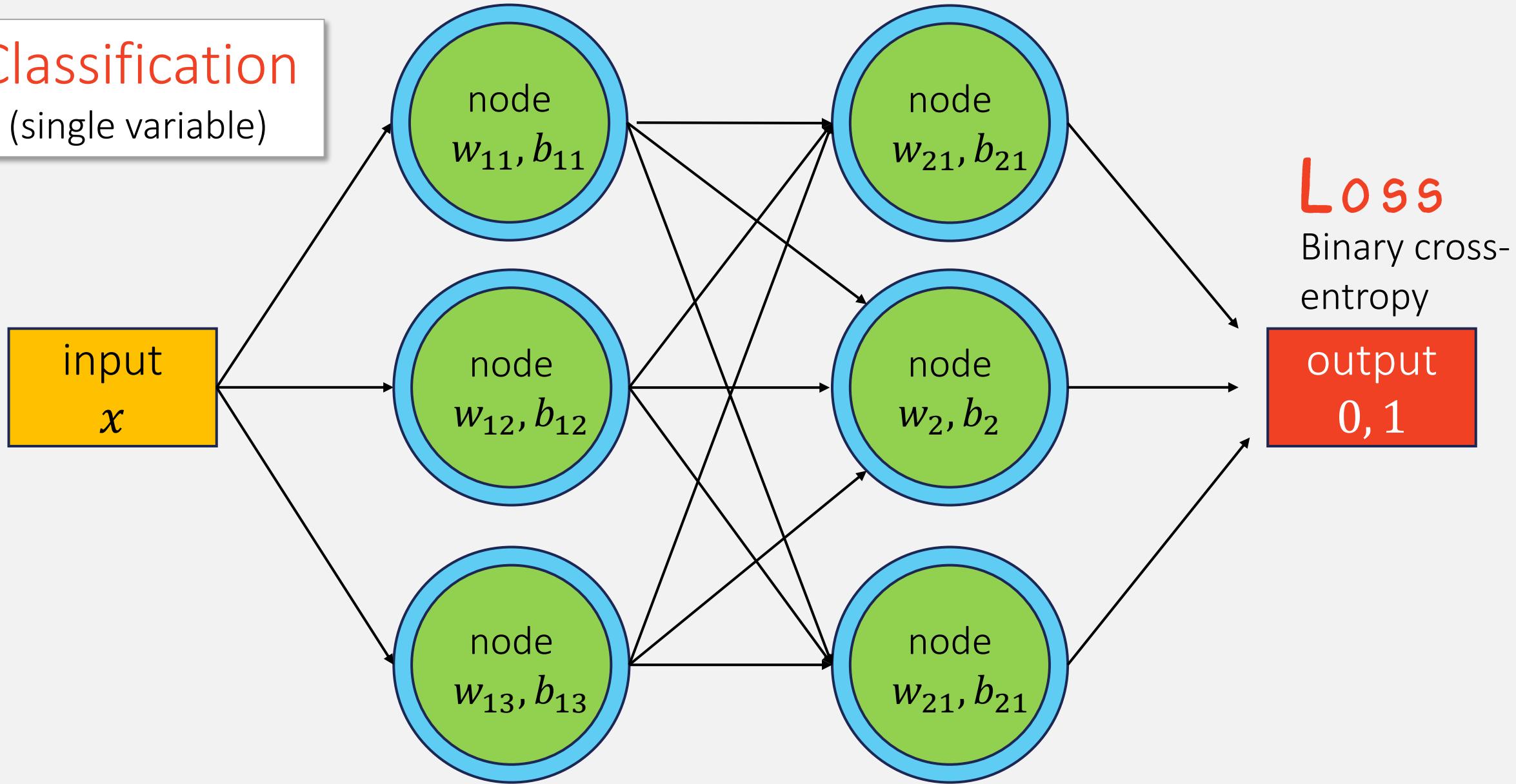
## Regression (single variable)



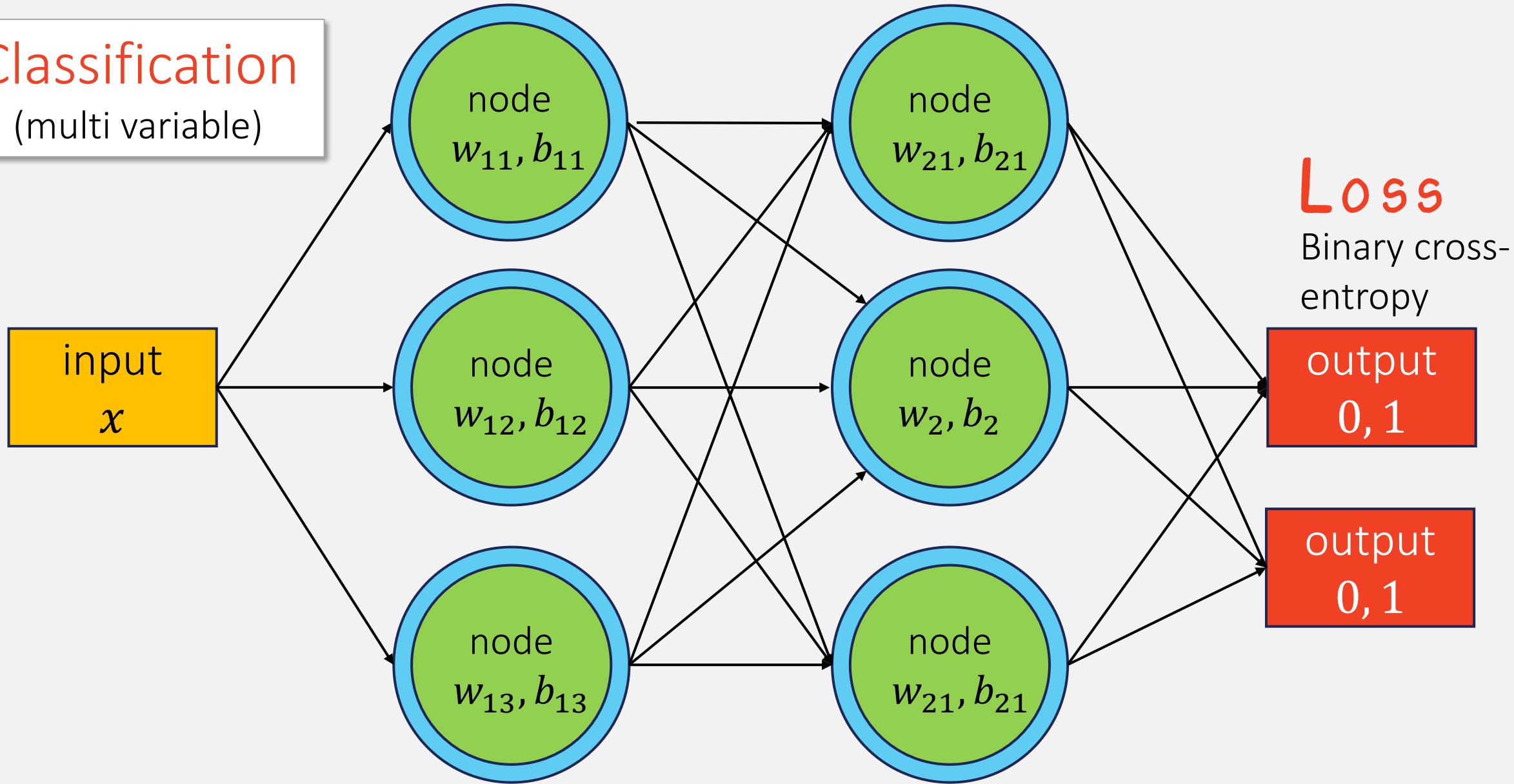
## Regression (multi variable)



**Classification**  
(single variable)



# Classification (multi variable)



## NN use cases

what the network learns to do depends mostly on the labels, the output layer, and the loss function

- + Easy to adapt the same model for different tasks
- + No need to re-train much if the task is similar (fine-tuning)

# Other things to know about

1

## Normalization

Your input data should be normalized!!! Networks perform **best** if the data is in (0,1) range

Everything is additive: if some values cap at 1 and some at  $10^5$ , the  $10^5$  ones will dominate

Learning rate: the same rate used for all nodes, so it can't change based on data amplitude

2

## Balanced training sets

Still important – networks will not learn about data it hasn't seen, or seen little of

3

## Uncertainty estimation

How to get the error on the prediction?

# Other things to know about

1

## Normalization

Your input data should be normalized!!! Networks perform **best** if the data is in (0,1) range

Everything is additive: if some values cap at 1 and some at  $10^5$ , the  $10^5$  ones will dominate

Learning rate: the same rate used for all nodes, so it can't change based on data amplitude

2

## Balanced training sets

Still important – networks will not learn about data it hasn't seen, or seen little of

3

## Uncertainty estimation

How to get the error on the prediction?

**Uncertain input data:** resample within uncertainties

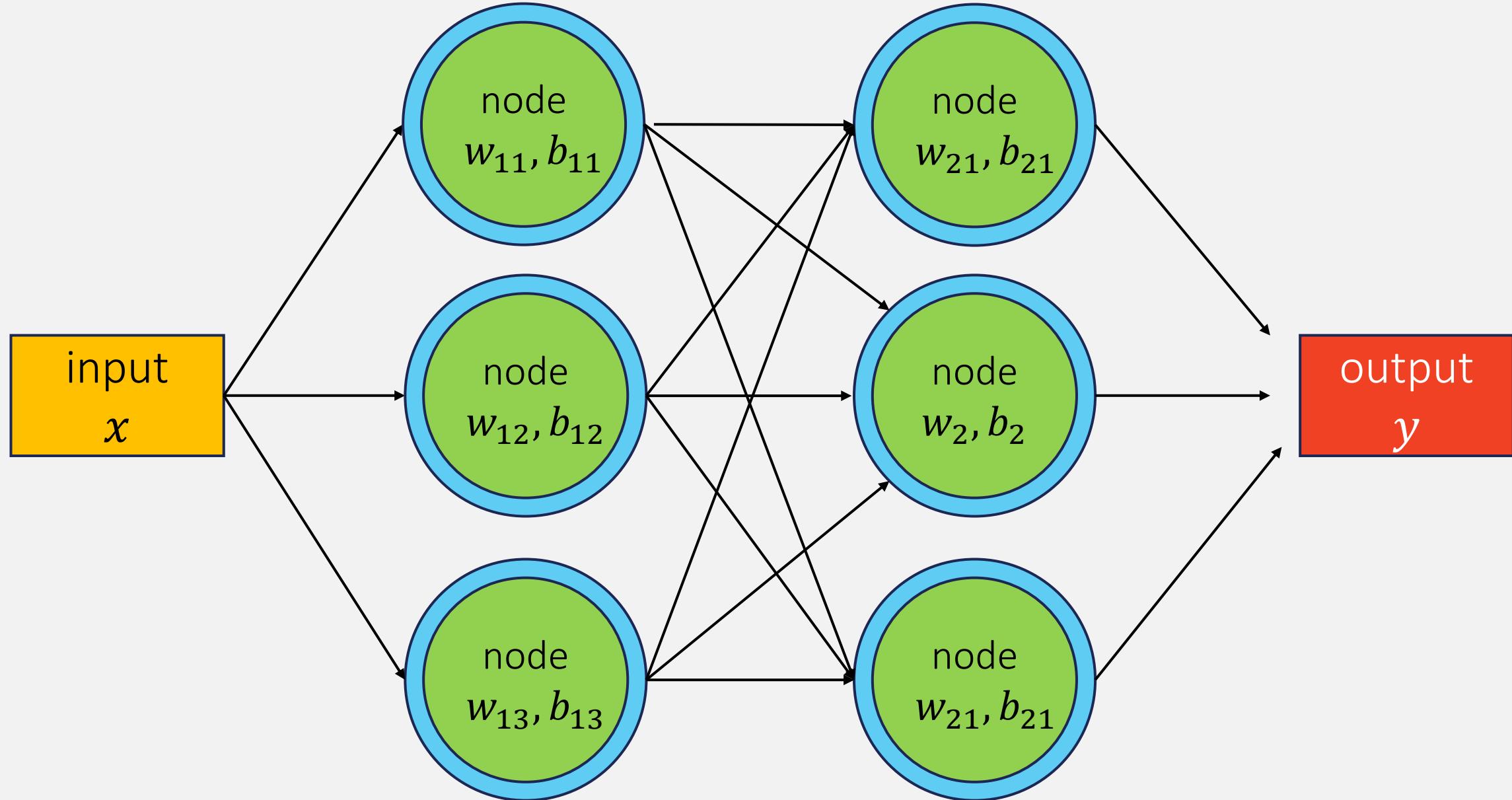
**Dropout:** turn the dropout on during inference!

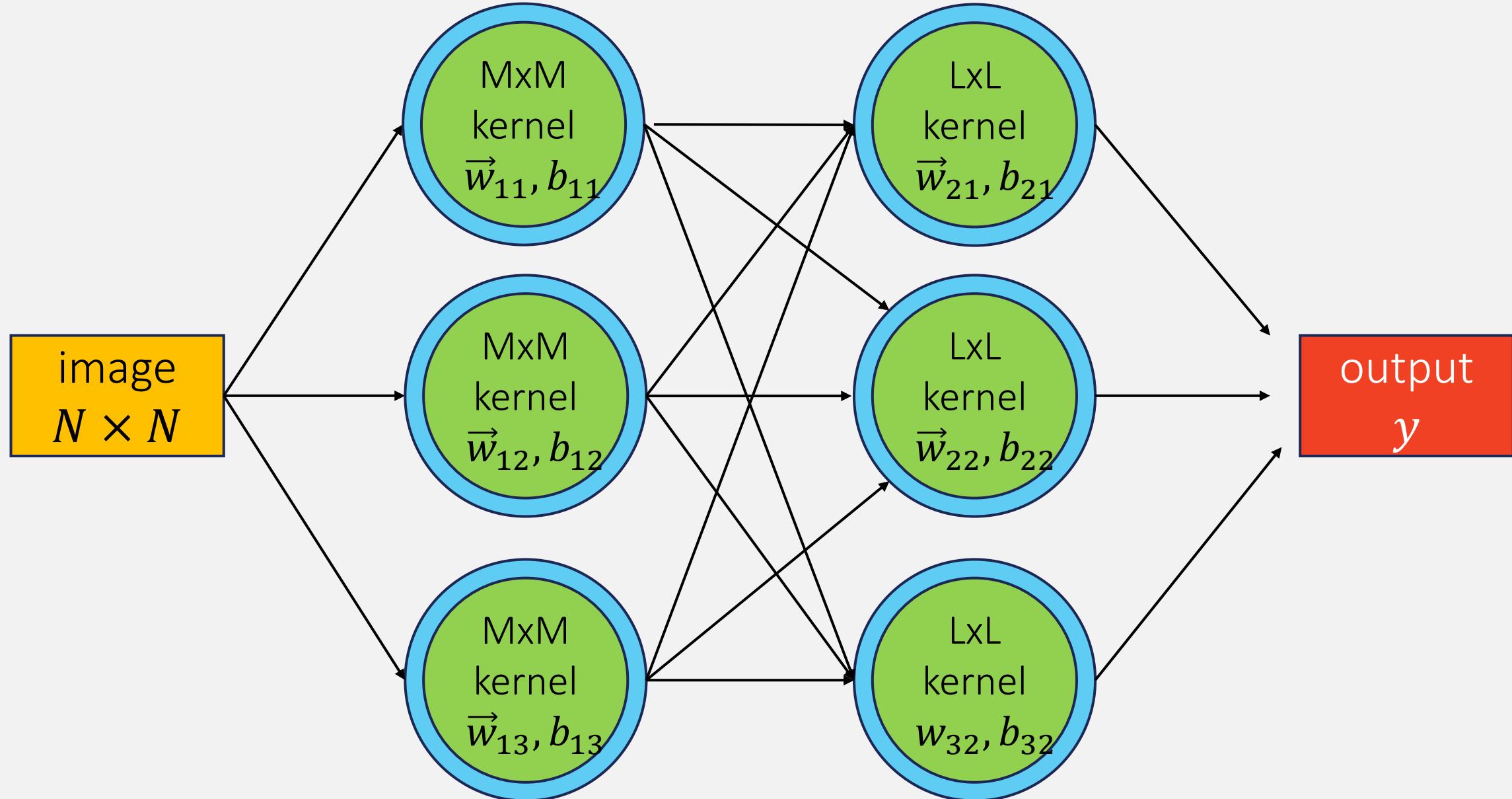


re-run the prediction N times  
to get uncertainties

# Machine Learning 201

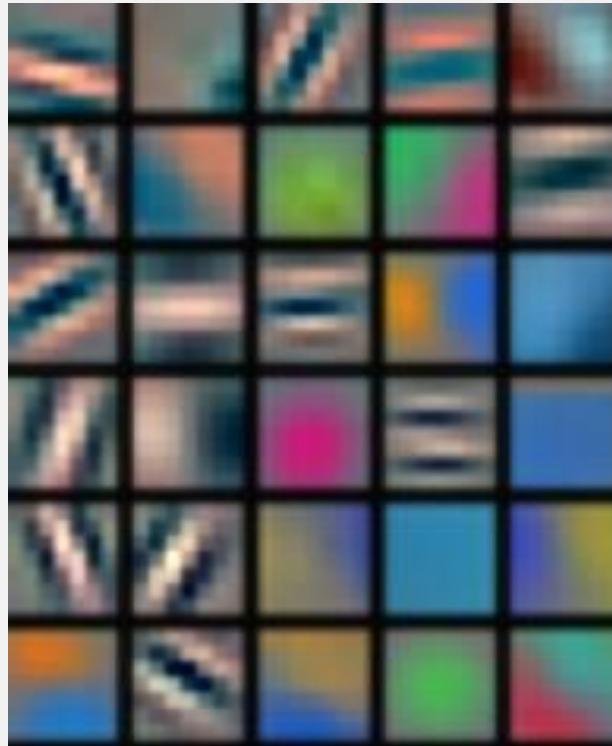
## Convolutional neural nets





# Convolution kernels

Simpler patterns

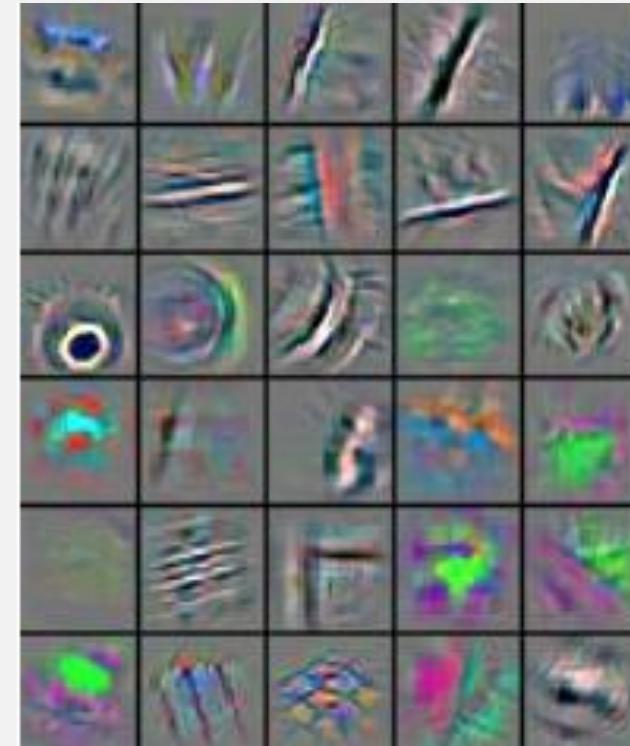


Layer 1

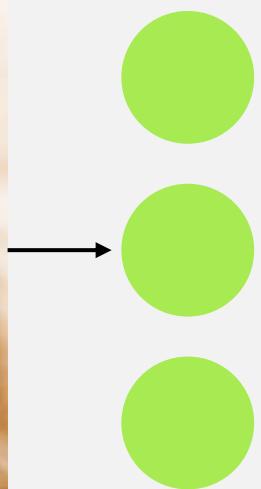
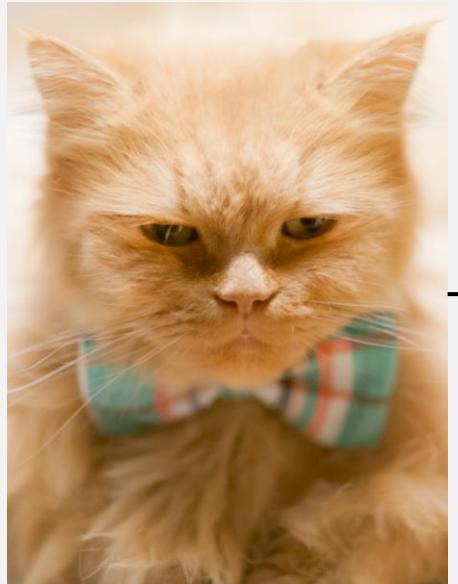
Complex patterns



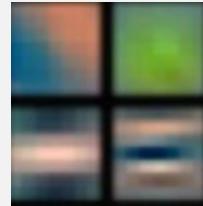
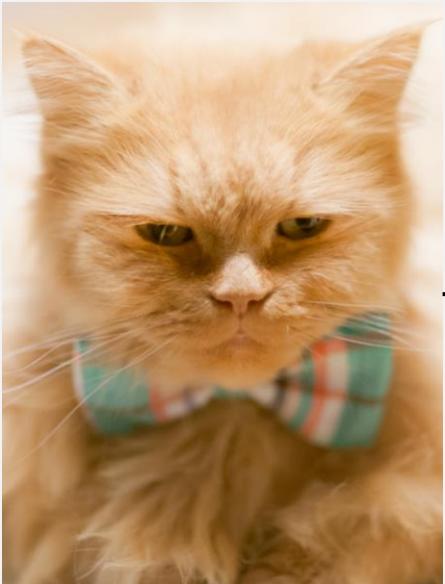
Layer 3



Layer 2



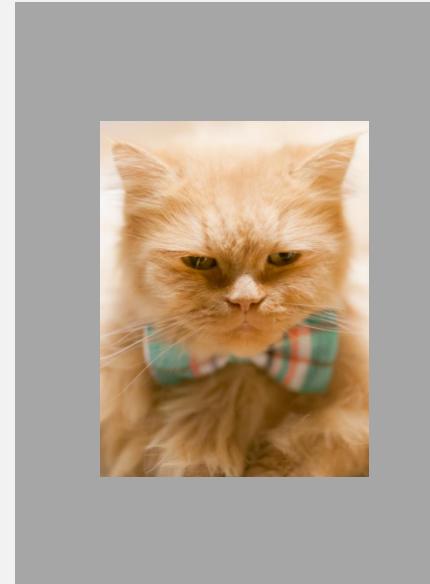
Layer 1  
5x5 kernels



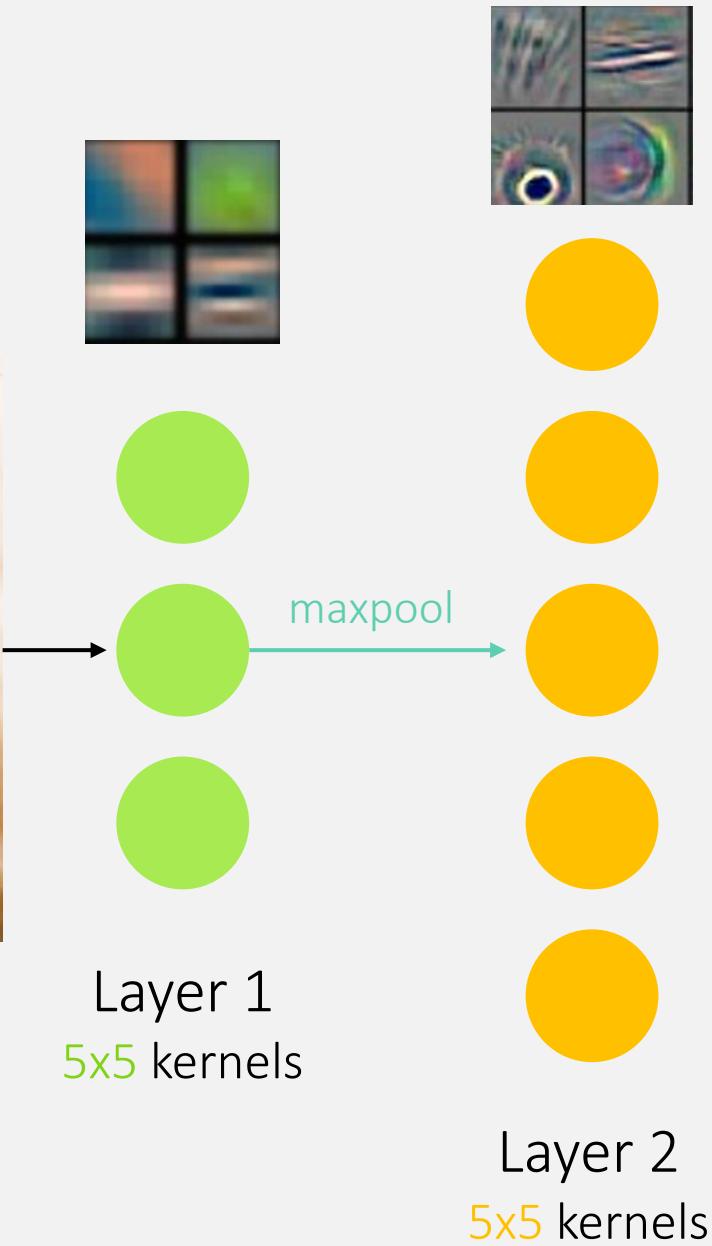
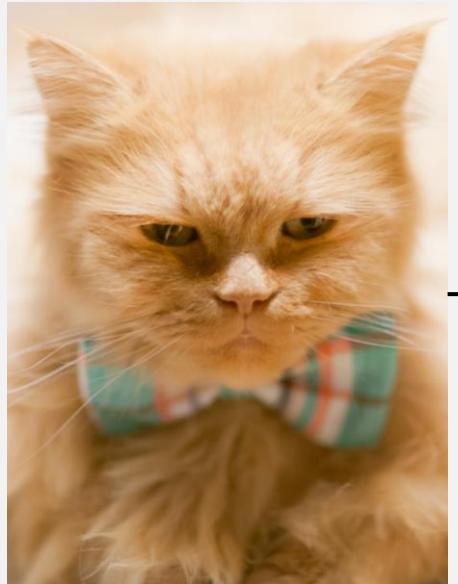
Layer 1  
5x5 kernels

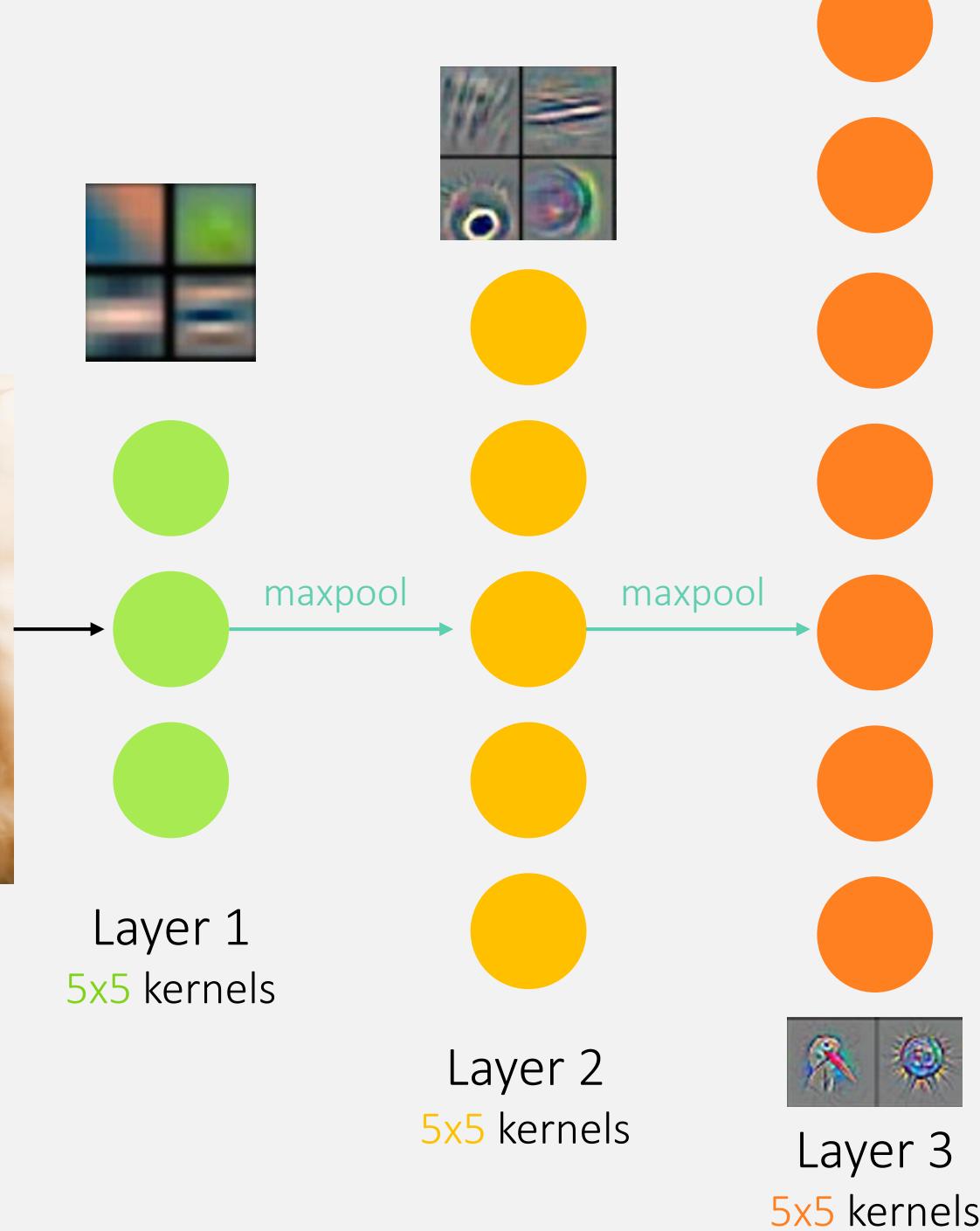
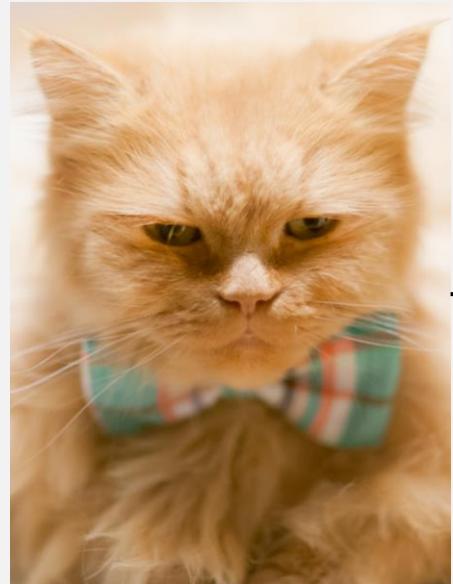
maxpooling  
or [other]pooling

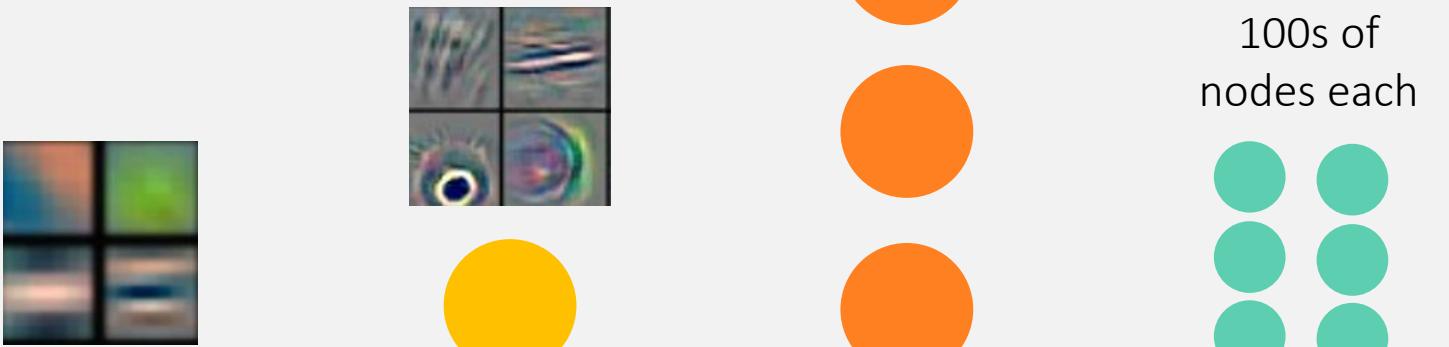
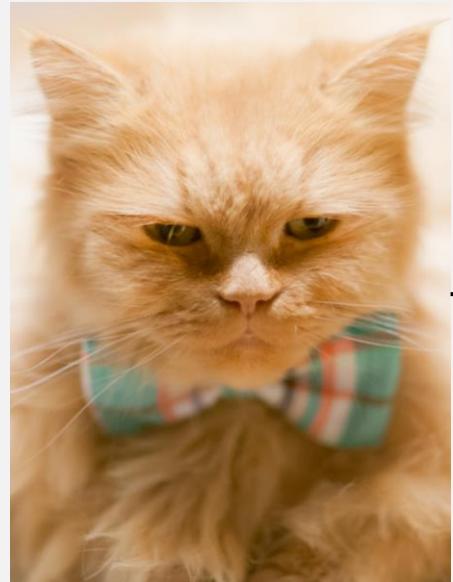
re-sample the image to lower resolution by some factor (2)  
maxpooling: take the max. of each 2x2 slice  
avgpooling: take the mean



then pad it with 0s to keep the shape the same







Layer 1  
5x5 kernels

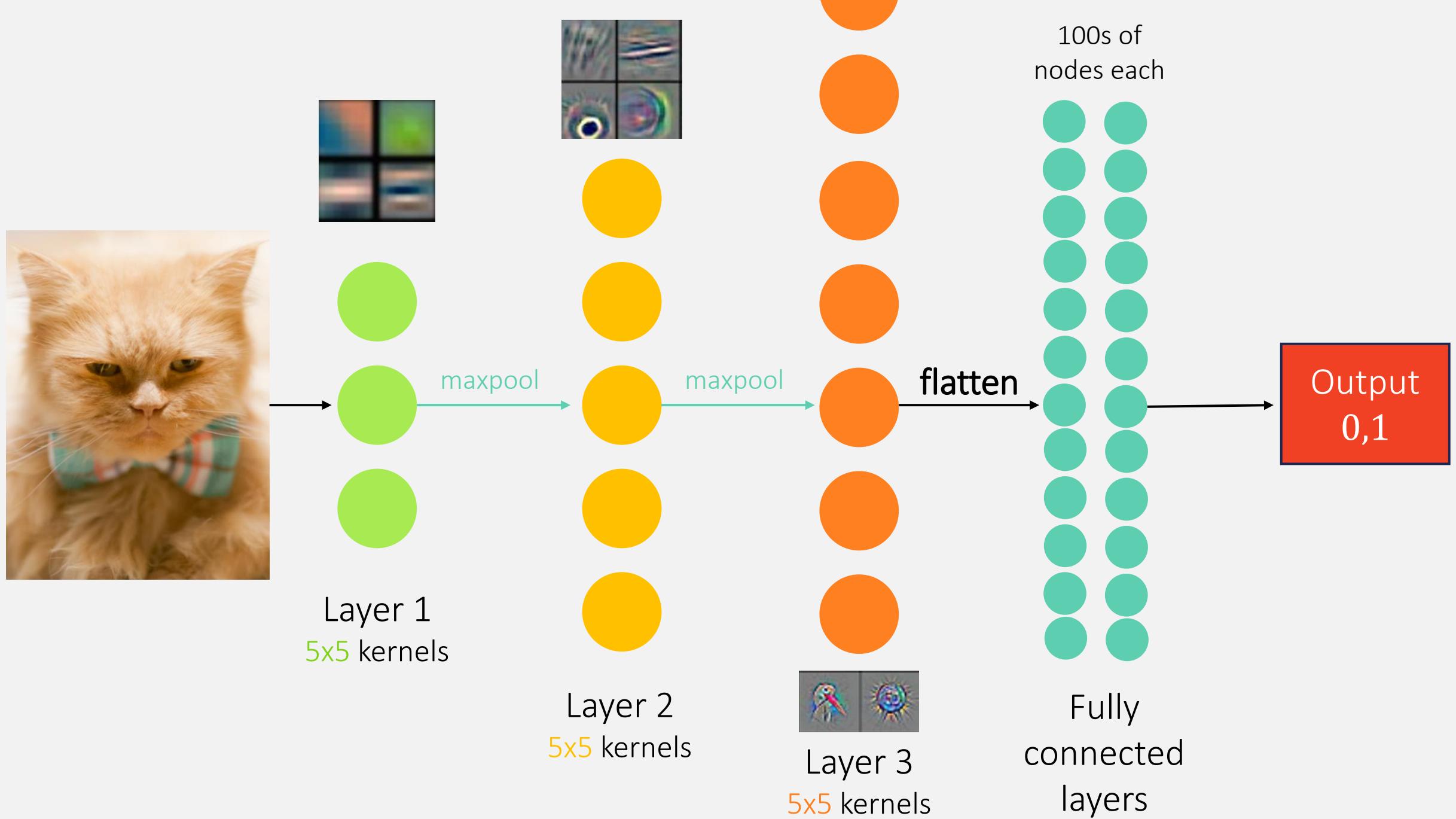
Layer 2  
5x5 kernels

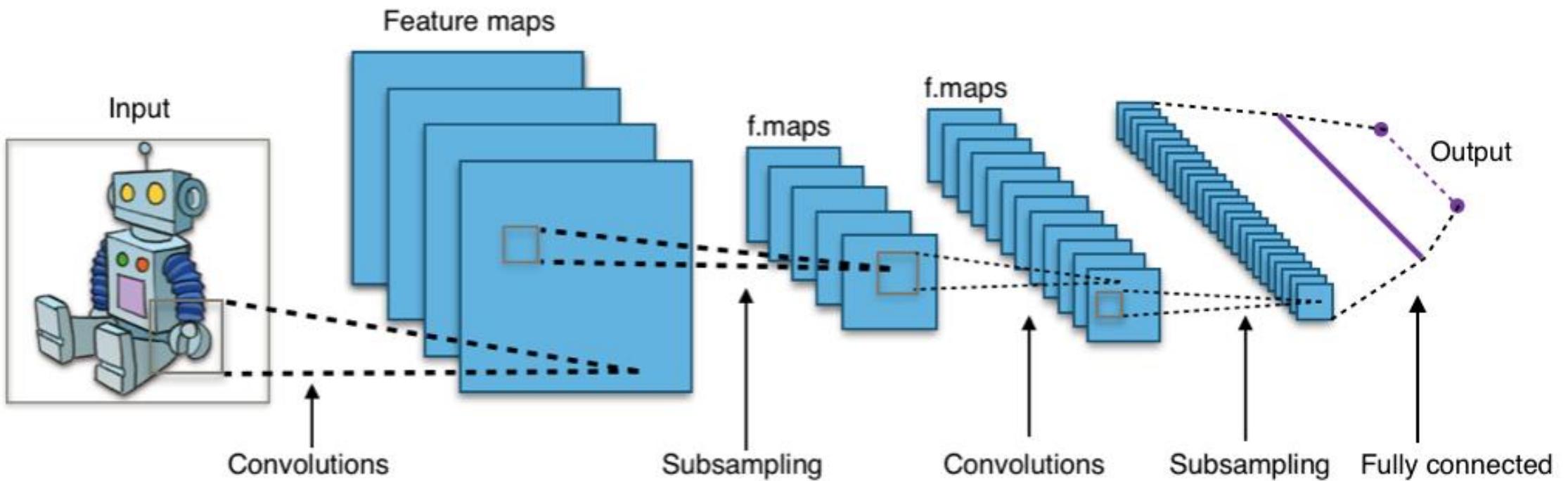
Layer 3  
5x5 kernels

Fully  
connected  
layers

100s of  
nodes each

usually dropout  
is applied here  
But you can add  
dropout to the conv  
layers too





# Galaxy Zoo

20 citizen scientists classify each galaxy,  
giving  $p(\text{smooth})$

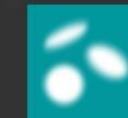


⌚ ⓘ You should sign in!

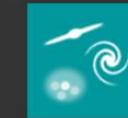
## TASK

## TUTORIAL

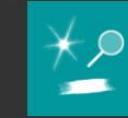
Is the galaxy simply smooth and rounded, with no sign of a disk?



Smooth



Features or Disk



Star, Artifact, or Bad Zoom

Is the galaxy simply smooth and rounded, with no sign of a disk?



Smooth



Features or Disk

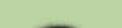


Star or Artifact

How rounded is it?



Round



In Between



Cigar Shaped

Could this be a disk viewed edge-on?



Yes - Edge On Disk



No - Something Else

Does the galaxy have a bulge at its centre? If so, what shape?



Rounded



Boxy



No bulge

How many spiral arms are there?



1



2



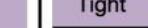
3



4



More than 4



Can't tell

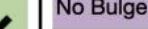
Is the galaxy merging or disturbed?



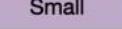
Merging



Major Disturbance



Minor Disturbance



None

Do you see any of these rare features?



Ring



Lens or arc



Irregular



Dust lane



Overlapping



Something Else



Nothing Unusual

END

# Galaxy Zoo

20 citizen scientists classify each galaxy,  
giving  $p(\text{smooth})$



➊ ⓘ You should sign in!

**TASK**

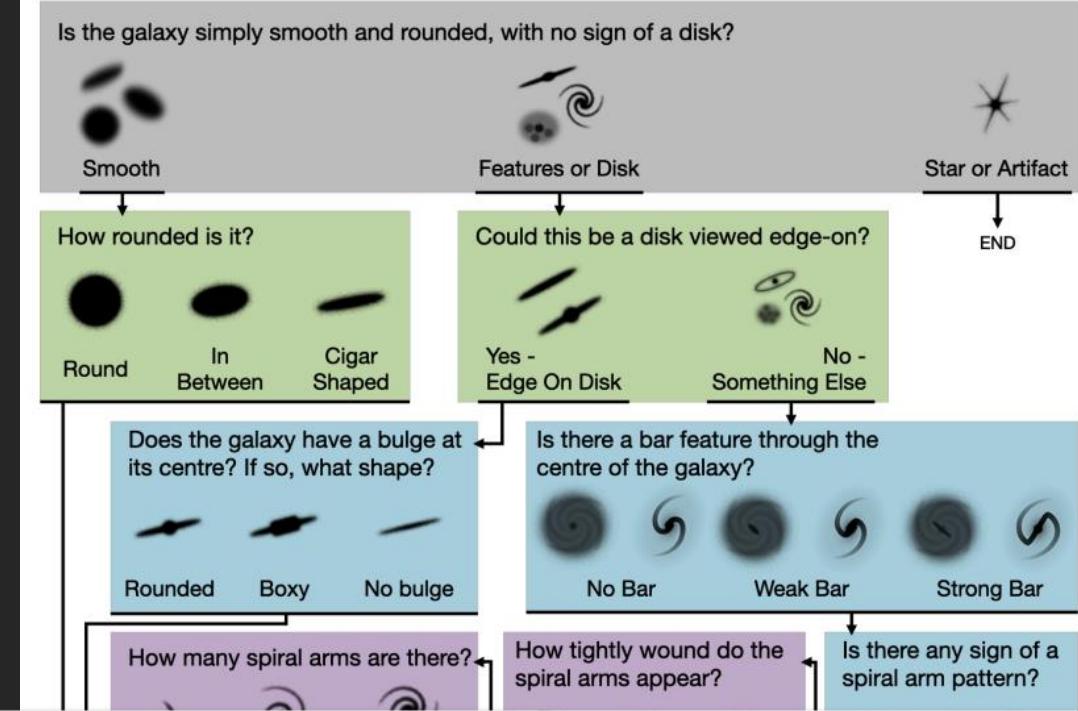
Is the galaxy simply smooth and rounded, with no sign of a disk?

**TUTORIAL**

Smooth

Features or Disk

Star or Artifact



## Galaxy Zoo: Probabilistic Morphology through Bayesian CNNs and Active Learning

Mike Walmsley<sup>1\*</sup>, Lewis Smith<sup>2</sup>, Chris Lintott<sup>1</sup>, Yarin Gal<sup>2</sup>, Steven Bamford<sup>3</sup>, Hugh Dickinson<sup>4,8</sup>, Lucy Fortson<sup>4,8</sup>, Sandor Kruk<sup>5</sup>, Karen Masters<sup>6,7</sup>, Claudia Scarlata<sup>4,8</sup>, Brooke Simmons<sup>9,10</sup>, Rebecca Smethurst<sup>1</sup>, Darryl Wright<sup>4,8</sup>

<sup>1</sup>Oxford Astrophysics, Department of Physics, University of Oxford, Denys Wilkinson Building, Keble Road, Oxford, OX1 3RH, UK

<sup>2</sup>Oxford Computer Science, University of Oxford, 15 Parks Rd, Oxford, OX1 3QD, UK

<sup>3</sup>School of Physics and Astronomy, University of Nottingham, University Park, Nottingham NG7 2RD, UK

<sup>4</sup>School of Physics and Astronomy, University of Minnesota, 116 Church St SE, Minneapolis, MN 55455, USA

<sup>5</sup>European Space Agency, ESTEC, Keplerlaan 1, NL-2201 AZ, Noordwijk, The Netherlands

<sup>6</sup>Haverford College, Department of Physics and Astronomy, 370 Lancaster Avenue, Haverford, Pennsylvania 19041, USA

<sup>7</sup>Institute for Cosmology and Gravitation, University of Portsmouth, Dennis Sciama Building, Burnaby Road, Portsmouth, PO1 3FX, UK

<sup>8</sup>Minnesota Institute for Astrophysics, University of Minnesota, Minneapolis, MN 55455, USA

<sup>9</sup>Physics Department, Lancaster University, Lancaster, LA1 4YB, UK

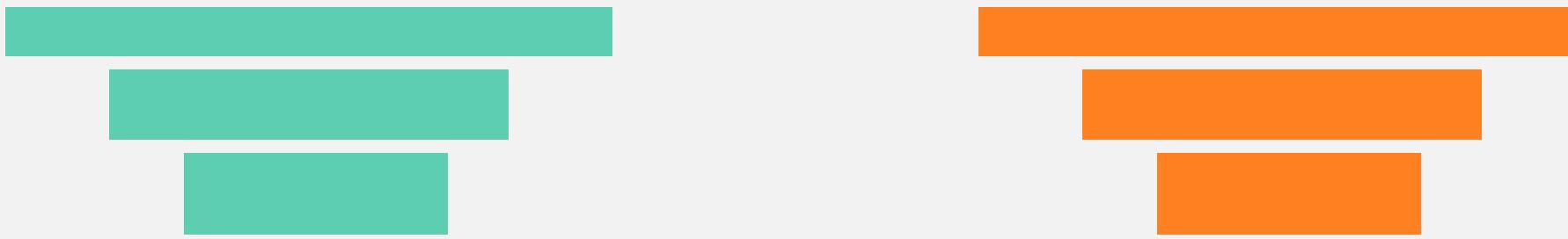
<sup>10</sup>Center for Astrophysics and Space Sciences (CASS), Department of Physics, University of California, San Diego, CA 92093, USA

Zoobot: uses dropout at inference  
to get 20 network votes on each class

# Machine Learning 301

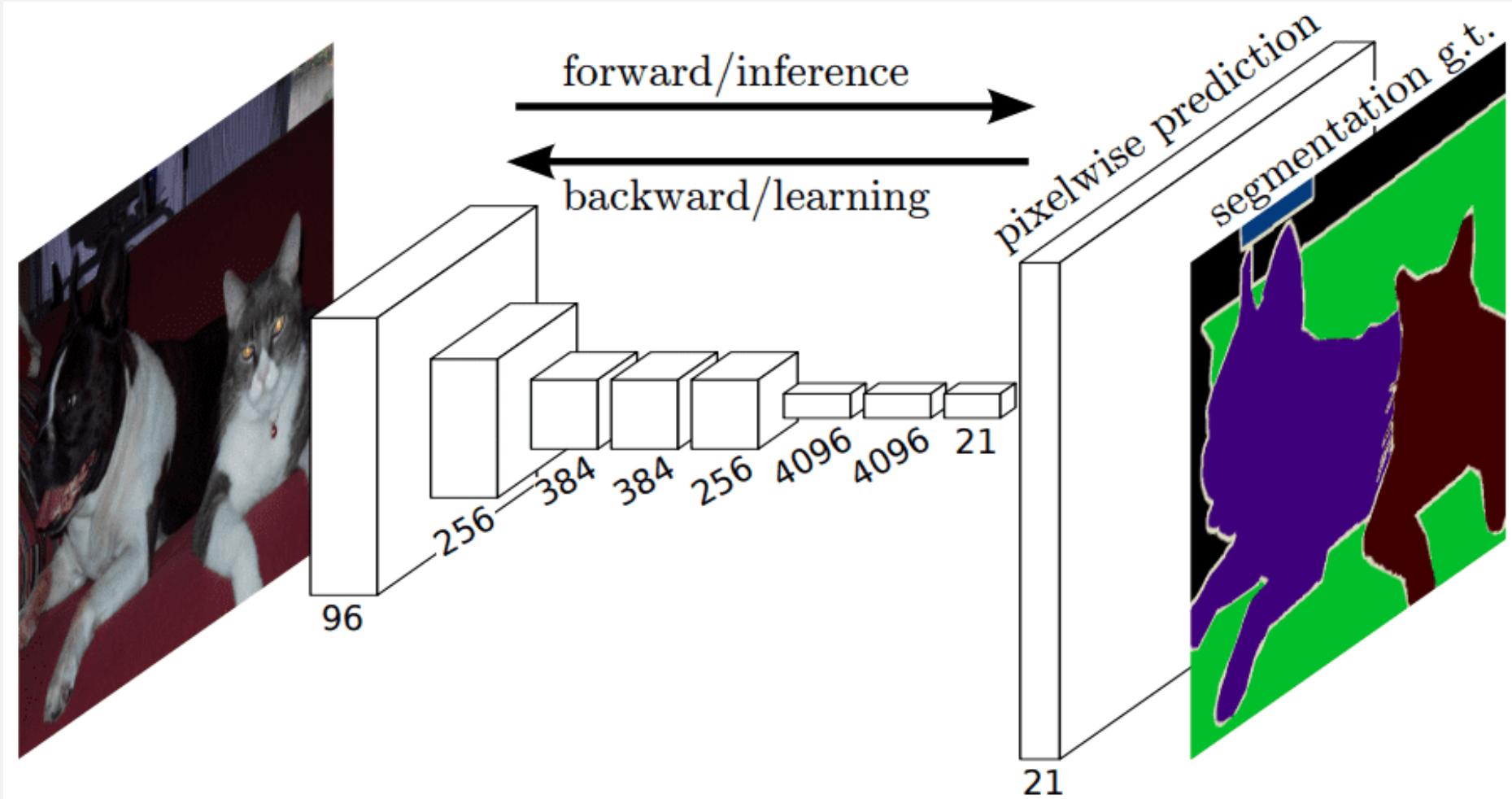
## GANs, contrastive nets, U-nets

# encoder-decoder Or U-net



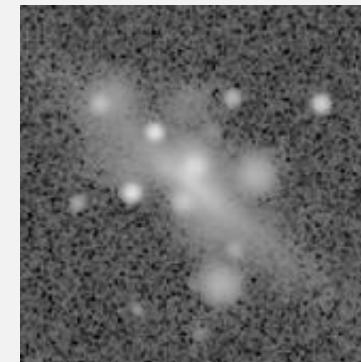
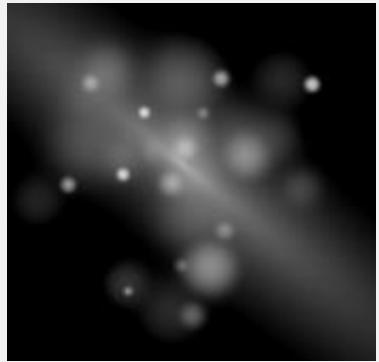
latent space

## U-nets: image segmentation



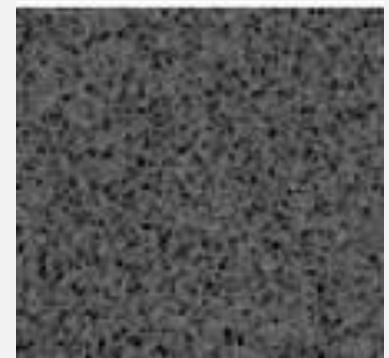
Loss: still similar to binary cross-entropy

# Contrastive Learning



latent space

Loss: similarity score

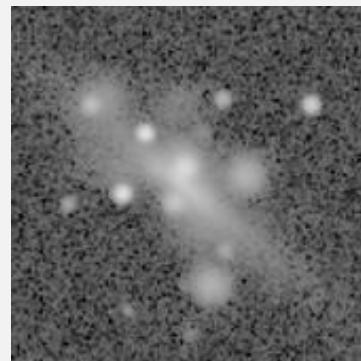


update weights

## Adversarial Networks

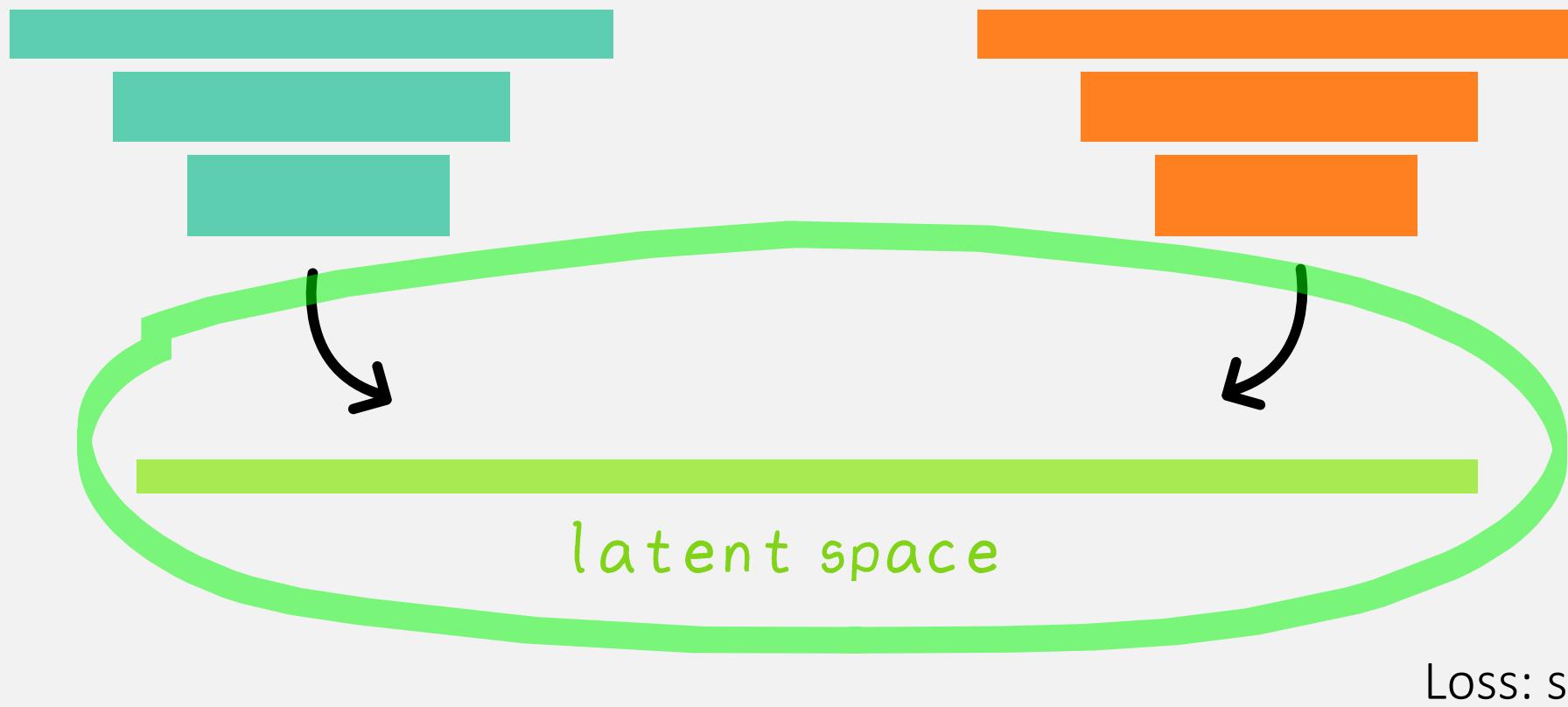
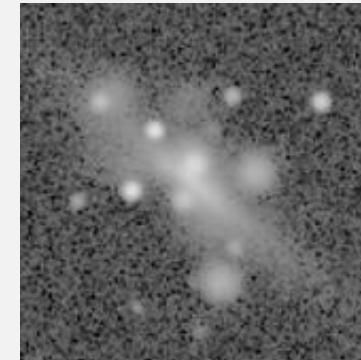
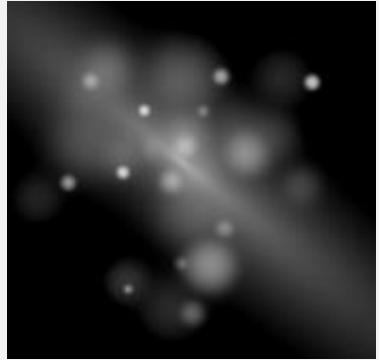
Is the image  
real?

update  
weights



Loss: similarity score

# Contrastive Learning



## Latent space

1024-dim vector that  
represents *something*  
about the data

Optimized to learn  
the real properties  
and be invariant to  
noise, etc.

## Latent space

1024-dim vector that represents *something* about the data

Optimized to learn the real properties and be invariant to noise, etc.

