

# Primer on ML [2]

# Overview

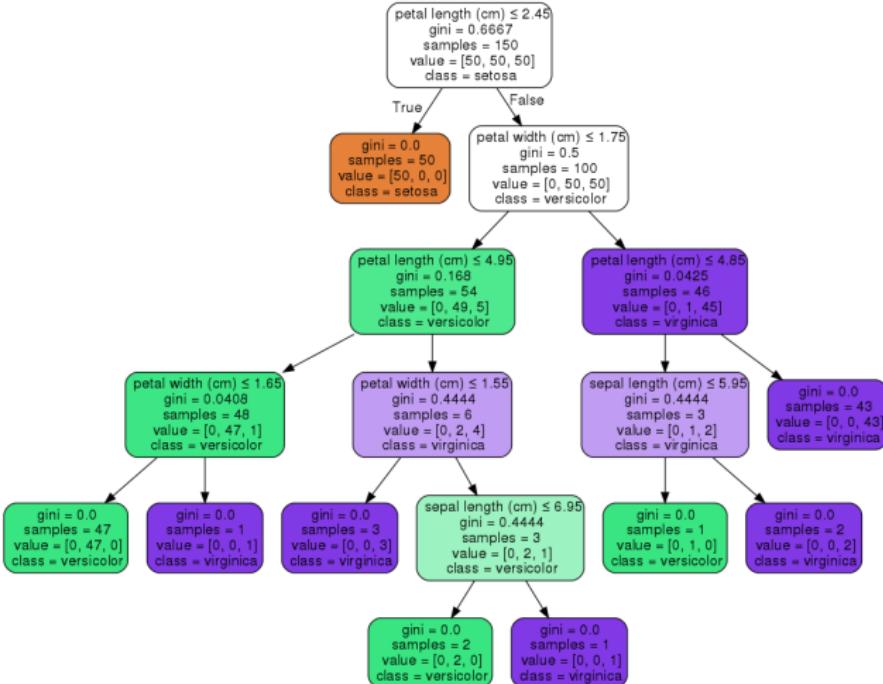
- Ensembles
    - Motivation
    - Bootstrap & Bagging → Random Forest
    - Boosting → BDT
  - Neural Nets
    - Motivation
    - Building
    - Training
    - Special tricks

next lecture

  - Thrilling applications in HEP
- \* Bonus
  - Convolutional NN
  - Recurrent NN

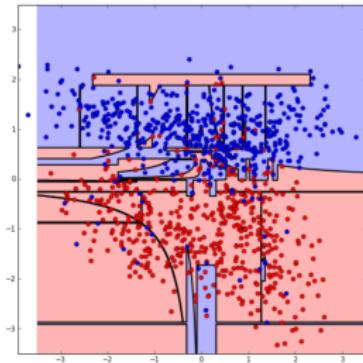
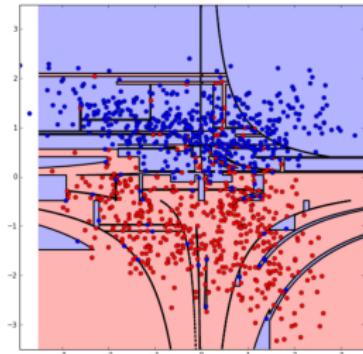
# Trees [recap]

# Trees [recap]



## Trees [recap]

- + interpretable
- + categorical features
- + don't care about feature linear dependence
- + no need to scale inputs
- weak learners
- can overfit
- large variance with sample variation

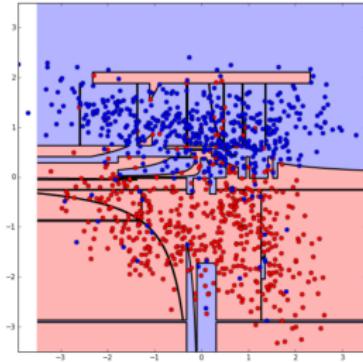
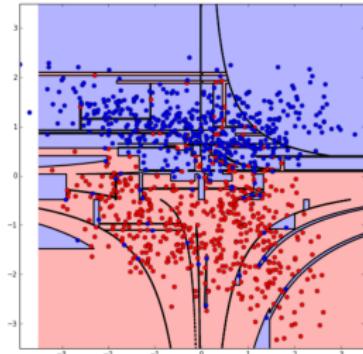


## Trees [recap]

- + interpretable
- + categorical features
- + don't care about feature linear dependence
- + no need to scale inputs
- weak learners
- can overfit
- large variance with sample variation

Maybe combine to produce a better algorithm?

Yes!



# Ensembles

more on bias-variance in ML [link]

## Ensembles

### Motivation

with data set variation:

$$\text{Error [ model]} = \text{Noise} + \text{Bias} + \text{Variance}$$

**Noise** - performance of ideal model (noise in data)

**Bias** - average deviation from ideal model's prediction

- ▷ smaller for complex models

**Variance** - average sensitivity to sample replacement

- ▷ smaller for simple models

with data set variation:

$$\text{Error [ model]} = \text{Noise} + \text{Bias} + \text{Variance}$$

$$\text{Error} \left[ \frac{1}{N} \sum_i \text{model}_i \right] = \text{Noise} + \text{Bias} + \frac{1}{N} \cdot \text{Variance} + \text{Correlation}$$

**Noise** - performance of ideal model (noise in data)

**Bias** - average deviation from ideal model's prediction

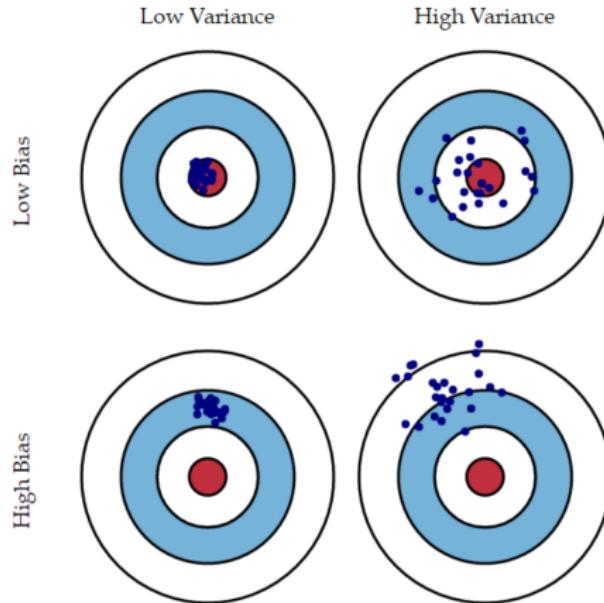
- ▷ smaller for complex models

**Variance** - average sensitivity to sample replacement

- ▷ smaller for simple models

# Ensembles

## Motivation



# Ensembles

## Motivation

$$\text{Error} \left[ \frac{1}{N} \sum_i \text{model}_i \right] = \text{Noise} + \text{Bias} + \frac{1}{N} \cdot \text{Variance} + \text{Correlation}$$

Can we reduce Correlation?

# Ensembles

## Motivation

$$\text{Error} \left[ \frac{1}{N} \sum_i \text{model}_i \right] = \text{Noise} + \text{Bias} + \frac{1}{N} \cdot \text{Variance} + \text{Correlation}$$

Can we reduce Correlation? Yes

- Bagging
  - train on subsample of events
- Random subspace method (feature bagging)
  - train on subsample of features

# Ensembles

## Motivation

$$\text{Error} \left[ \frac{1}{N} \sum_i \text{model}_i \right] = \text{Noise} + \text{Bias} + \frac{1}{N} \cdot \text{Variance} + \text{Correlation}$$

Can we reduce Correlation? Yes

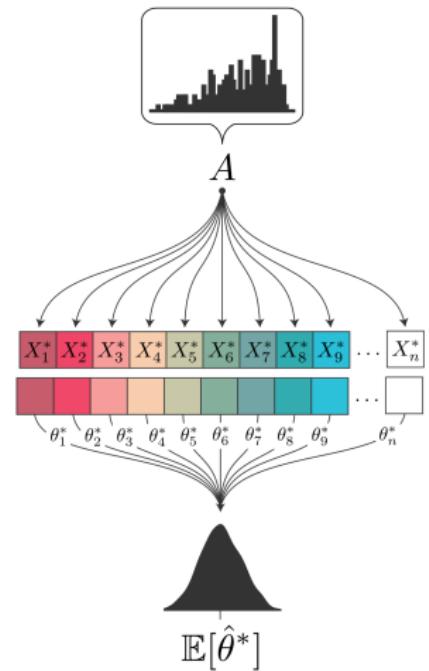
- Bagging
  - train on subsample of events ← hyperparameter :(
- Random subspace method (feature bagging)
  - train on subsample of features

more on that [link]

# Ensembles

## Bootstrap

- same-size sampling with replacement
- widely used in
  - confidence intervals
  - hypothesis tests



# Random Forest

$$\text{RF} = \frac{1}{N} \sum_{i=1}^N \text{tree}_i + \text{bootstrap each tree} + \text{bagging each split}$$

→ for  $i = 1, 2, \dots, N$ :

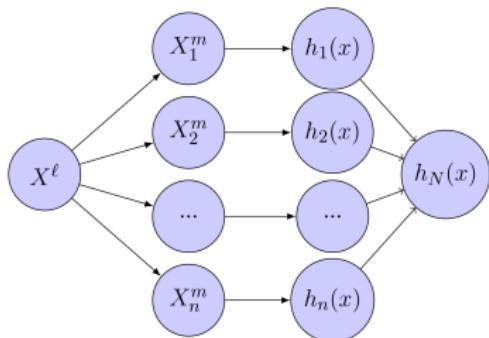
- bootstrap sample  $X_i$
- while !stopping criteria:
  - pick  $p$  random features out of  $d$
  - make split

→ prediction:

r  $\text{RF}(x) = \frac{1}{N} \sum_{i=1}^N \text{tree}_i(x)$

c  $\text{RF}(x) = \text{sign} \frac{1}{N} \sum_{i=1}^N \text{tree}_i(x)$

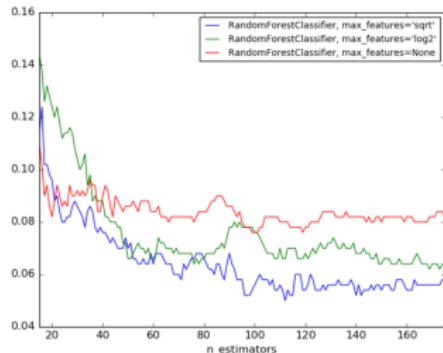
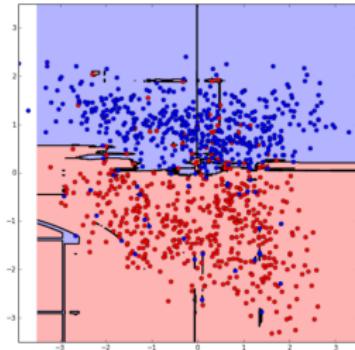
NB: proba average in sklearn



# Random Forest

## pros & cons

- + trees' pros
- + less overfitting with increasing N
- + better capturing of non-linearities
- + parallelized
- + out-of-bag score
- + more robust (than single tree)
- not lightweight
- previous trees' improvements lost



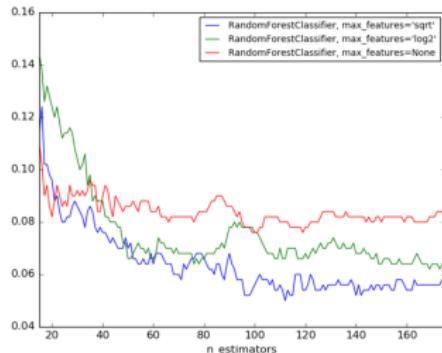
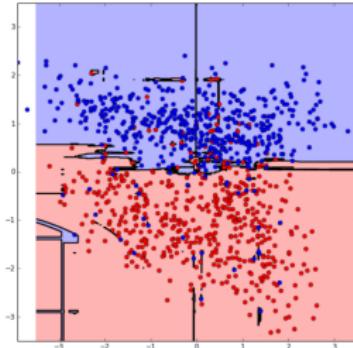
# Random Forest

## pros & cons

- + trees' pros
- + less overfitting with increasing N
- + better capturing of non-linearities
- + parallelized
- + out-of-bag score
- + more robust (than single tree)
- not lightweight
- previous trees' improvements lost



Boosting



# Gradient Boosting

# Gradient Boosting

Regression: simplified

Init. learner e.g.  $a_0 = h_0 = \frac{1}{N} \sum_i y_i$

Add learners greedily to improve results of previous ones:

① residuals:

$$s_i = y_i - \sum_{n=1}^{N-1} h_n(x) = y_i - a_{N-1}(x_i)$$

② train next algorithm on residuals:

$$h_N = \arg \min_{h \in H} \sum_i (h(x_i) - s_i)^2$$

③ add to composition:

$$a_N = a_{N-1} + h_N$$

# Gradient Boosting

Regression: simplified

Init. learner e.g.  $a_0 = h_0 = \frac{1}{N} \sum_i y_i$

Add learners greedily to improve results of previous ones:

① residuals:

$$s_i = y_i - \sum_{n=1}^{N-1} h_n(x) = y_i - a_{N-1}(x_i) \leftarrow \text{looks like L2 gradient}$$

② train next algorithm on residuals:

$$h_N = \arg \min_{h \in H} \sum_i (h(x_i) - s_i)^2$$

③ add to composition:

$$a_N = a_{N-1} + h_N$$

# Gradient Boosting

Regression: simplified

Init. learner e.g.  $a_0 = h_0 = \frac{1}{N} \sum_i y_i$

Add learners greedily to improve results of previous ones:

- ① residuals:

$$s_i = y_i - \sum_{n=1}^{N-1} h_n(x) = y_i - a_{N-1}(x_i) \leftarrow \text{looks like L2 gradient}$$

- ② train next algorithm on residuals:  $\leftarrow$  and here we learn this gradient

$$h_N = \arg \min_{h \in H} \sum_i (h(x_i) - s_i)^2$$

- ③ add to composition:

$$a_N = a_{N-1} + h_N$$

# Gradient Boosting

Regression: simplified

Init. learner e.g.  $a_0 = h_0 = \frac{1}{N} \sum_i y_i$

Add learners greedily to improve results of previous ones:

- ① residuals:

$$s_i = y_i - \sum_{n=1}^{N-1} h_n(x) = y_i - a_{N-1}(x_i) \leftarrow \text{looks like L2 gradient}$$

- ② train next algorithm on residuals:  $\leftarrow$  and here we learn this gradient

$$h_N = \arg \min_{h \in H} \sum_i (h(x_i) - s_i)^2$$

- ③ add to composition:  $\leftarrow$  and finally descend in algorithm space

$$a_N = a_{N-1} + h_N$$

# Gradient Boosting

GBM [Friedman, 2001]

- More general approach:

Define  $\sum_{i=1}^{\ell} L(y_i, a_{N-1}(x_i) + \gamma h(x_i)) \rightarrow \min_{\gamma, h}$

At each iteration:

$$① s_i = -\frac{\partial L(y_i, z)}{\partial z} \Big|_{z=a_{N-1}(x_i)}$$

- regression on  $s_i$ :

$$h_N = \arg \min_{h \in H} \sum_{i=1}^{\ell} (h(x_i) - s_i)^2$$

- optimal  $\gamma$  with plain grad. descent:

$$\gamma_N = \arg \min_{\gamma} \sum_{i=1}^{\ell} L(y_i, a_{N-1}(x_i) + \gamma h_N(x_i))$$

- $a_N \leftarrow a_{N-1} + \gamma_N h_N$

# Gradient Boosting

GBM [Friedman, 2001]

- More general approach:

Define  $\sum_{i=1}^{\ell} L(y_i, a_{N-1}(x_i) + \gamma h(x_i)) \rightarrow \min_{\gamma, h}$

At each iteration:

①  $s_i = -\frac{\partial L(y_i, z)}{\partial z} \Big|_{z=a_{N-1}(x_i)}$

- classification?

② regression on  $s_i$ :

$$h_N = \arg \min_{h \in H} \sum_{i=1}^{\ell} (h(x_i) - s_i)^2$$

③ optimal  $\gamma$  with plain grad. descent:

$$\gamma_N = \arg \min_{\gamma} \sum_{i=1}^{\ell} L(y_i, a_{N-1}(x_i) + \gamma h_N(x_i))$$

④  $a_N \leftarrow a_{N-1} + \gamma_N h_N$

# Gradient Boosting

GBM [Friedman, 2001]

- More general approach:

Define  $\sum_{i=1}^{\ell} L(y_i, a_{N-1}(x_i) + \gamma h(x_i)) \rightarrow \min_{\gamma, h}$

At each iteration:

$$① s_i = -\frac{\partial L(y_i, z)}{\partial z} \Big|_{z=a_{N-1}(x_i)}$$

- classification?
- BDT = GB with trees
- but you can boost whatever you want\*

- regression on  $s_i$ :

$$h_N = \arg \min_{h \in H} \sum_{i=1}^{\ell} (h(x_i) - s_i)^2$$

- optimal  $\gamma$  with plain grad. descent:

$$\gamma_N = \arg \min_{\gamma} \sum_{i=1}^{\ell} L(y_i, a_{N-1}(x_i) + \gamma h_N(x_i))$$

- $a_N \leftarrow a_{N-1} + \gamma_N h_N$

# Gradient Boosting

GBM [Friedman, 2001]

- More general approach:

Define  $\sum_{i=1}^{\ell} L(y_i, a_{N-1}(x_i) + \gamma h(x_i)) \rightarrow \min_{\gamma, h}$

At each iteration:

$$① s_i = -\frac{\partial L(y_i, z)}{\partial z} \Big|_{z=a_{N-1}(x_i)}$$

- classification?
- BDT = GB with trees
- but you can boost whatever you want\*
- modern libraries are slightly different

- regression on  $s_i$ :

$$h_N = \arg \min_{h \in H} \sum_{i=1}^{\ell} (h(x_i) - s_i)^2$$

- optimal  $\gamma$  with plain grad. descent:

$$\gamma_N = \arg \min_{\gamma} \sum_{i=1}^{\ell} L(y_i, a_{N-1}(x_i) + \gamma h_N(x_i))$$

- $a_N \leftarrow a_{N-1} + \gamma_N h_N$

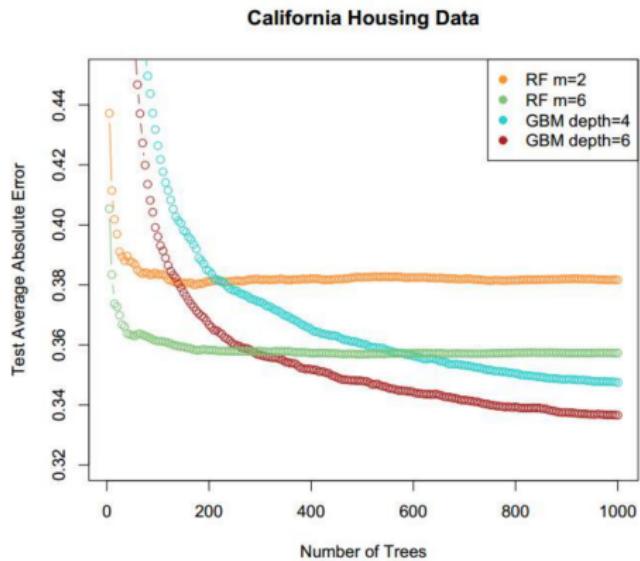
intuition and formulas [1], [2]  
and nice visuals [3], [4]

## BDT pros & cons

- + damn good performance
- not interpretable (directly)
- overfits

You might want try these ones:

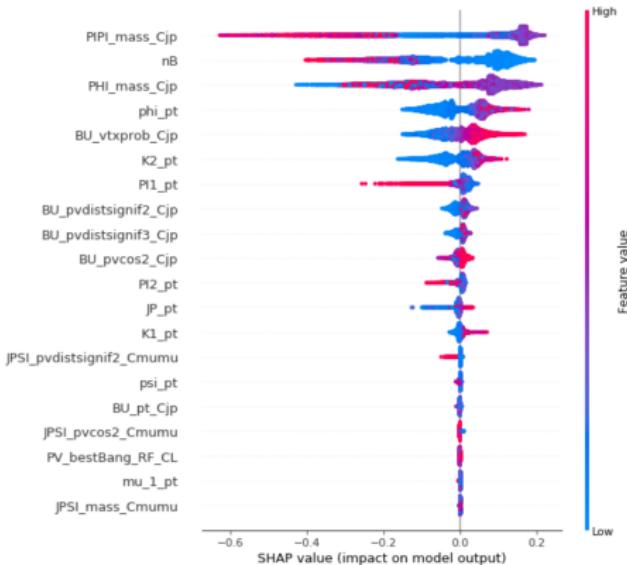
- LightGBM
- XGBoost
- CatBoost
- ? AdaBoost (used in TMVA)



# Sidenote

## Feature importance

- Linear models:
  - Weights' values
- Trees:
  - Gain
  - Frequency
- But you should probably use  
**SHAP values** [link]
  - Makes BDT interpretable!





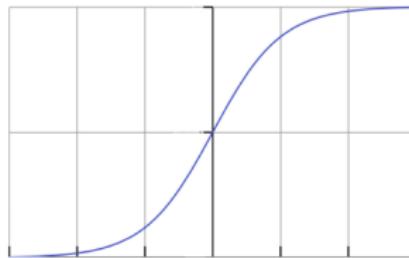


the following are the slides from [1], [2]  
huge credit to them!

# Principles of linear vs. nonlinear models

# Recap: linear models

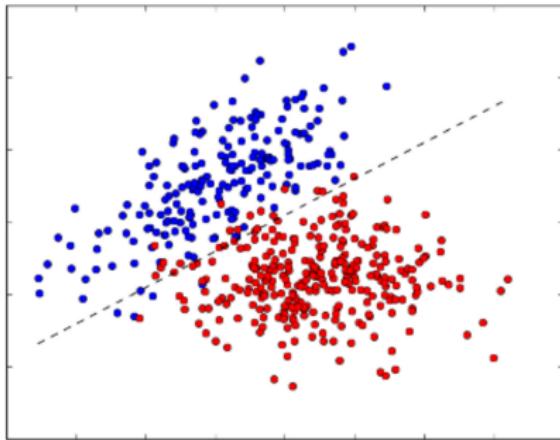
$W \cdot x + b \longrightarrow$



$\longrightarrow P(y = +1|x)$

- $x$ : features vector
- $W, b$ : model slope and intercept

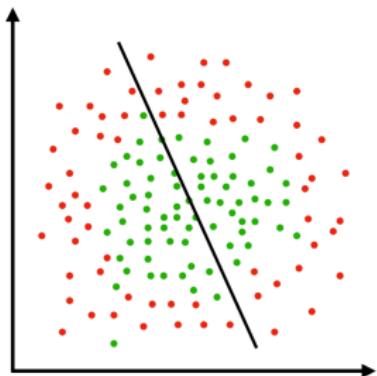
# Linear dependency



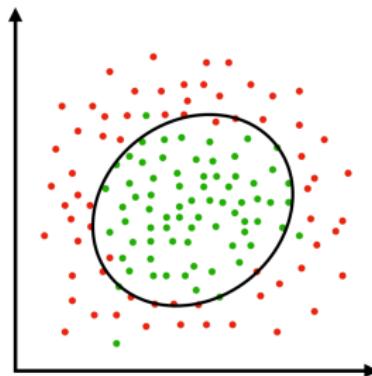
5

# Nonlinear dependencies

What we have

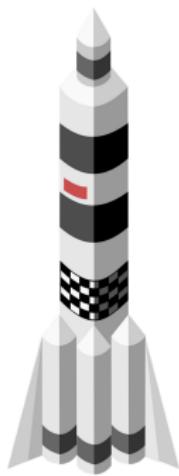


What we want



# Extremely nonlinear dependencies

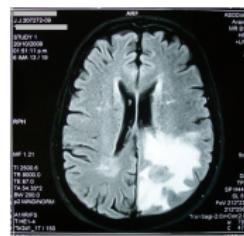
Most of the dependencies in this world!



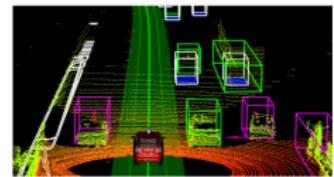
rocket



cat

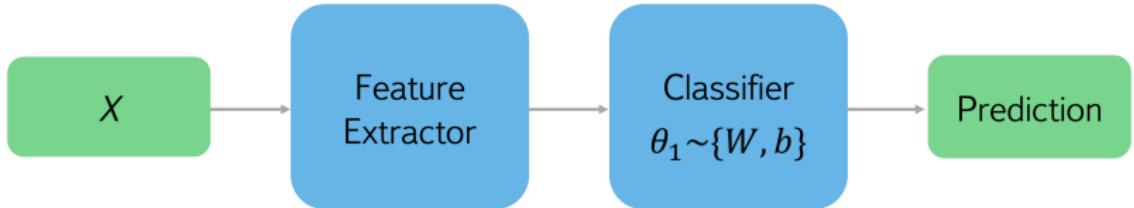


Brain tumors  
as seen on MRI



Self-driving LiDAR

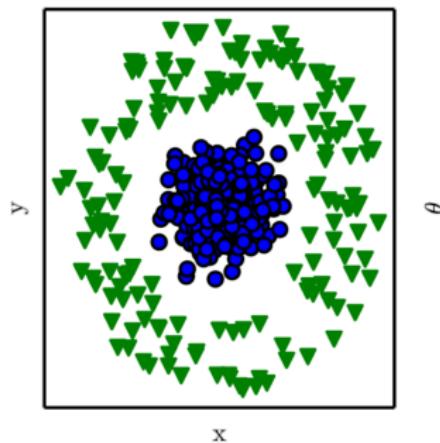
# Extremely nonlinear dependencies



- Decouple **feature extractor** from the **classifier**
- Training and inference really can have multiple stages!

# Feature extraction?

Cartesian coordinates



Polar coordinates

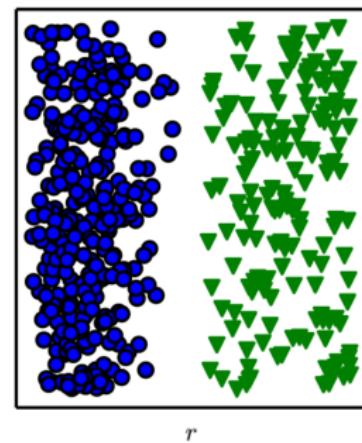
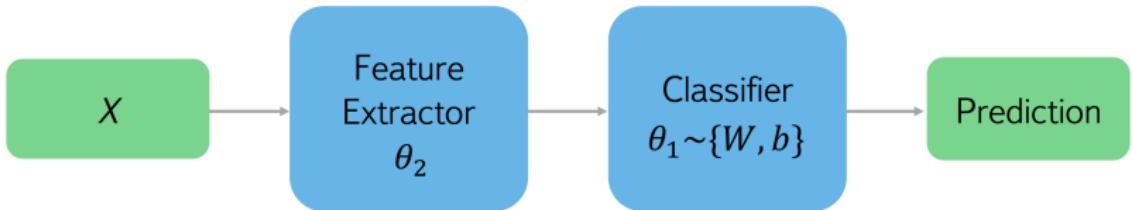


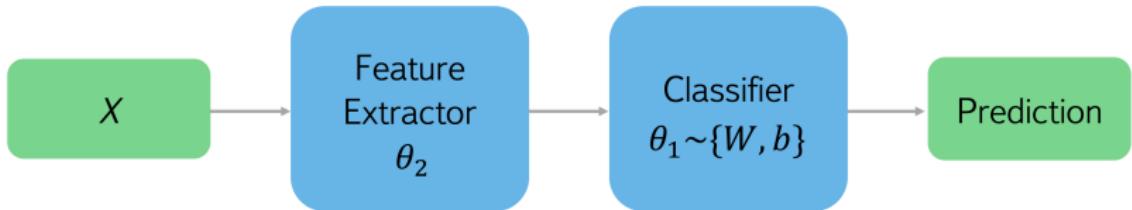
Image: Ian Goodfellow et al.

# Feature extraction



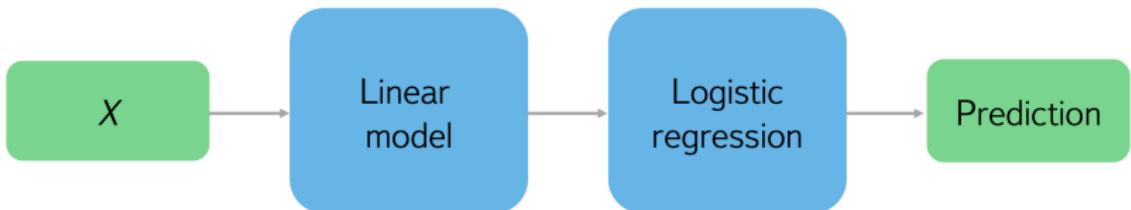
- *Manually* extracted features
- Training is left with finding  $\underset{\theta_1}{\operatorname{argmin}} L(y, P(y | x))$

# Can it be done automatically?



- *Automatically* extracted features
- Training still needs to find  $\underset{\theta_1}{\operatorname{argmin}} L(y, P(y | x))$
- Yet, we face a different challenge of finding  $\theta_2$

# Try stacked linear models



- Compute features 
$$h_j = \sum_i w_{ij}^h x_i + b_j^h \quad j \in \{1, 2, \dots, n\}$$
- Eventual output of the model 
$$y_{pred} = \sigma \left( \sum_j w_j^o h_j + b^o \right)$$
- Train by jointly minimizing loss to search for 
$$\operatorname{argmin}_{w^h, w^o, b^h, b^o} L(y, P(y | x))$$

# A question

- Will stacking linear functions improve quality?

# Answer: no

- Why?
- A combination of linear models is a linear model:

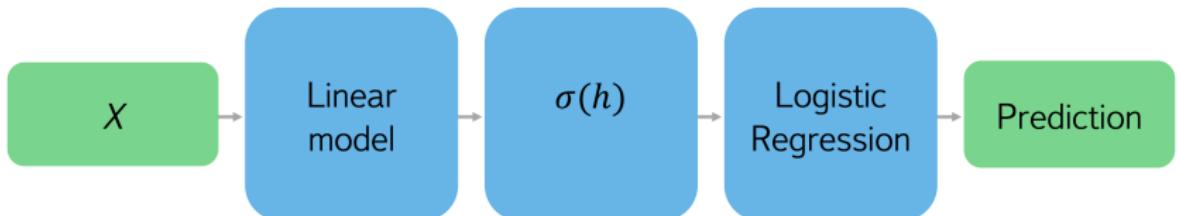
$$P(y | x) = \sigma \left( \sum_j w_j^o \left( \sum_i w_{ij}^h x_i + b_j^h \right) + b^o \right)$$

$$w'_i = \sum_j w_j^o w_{ij}^h \quad b' = \sum_j w_j^o b_j^h + b^o$$

$$P(y | x) = \sigma \left( \sum_i w'_i x_i + b' \right)$$

16

# The nonlinearity



- Compute features  $h_j = \sigma\left(\sum_i w_{ij}^h x_i + b_j^h\right) \quad j \in \{1, 2, \dots, n\}$
- Eventual output of the model  $y_{pred} = \sigma\left(\sum_j w_j^o h_j + b^o\right)$
- Compositionality:  $P(y | x) = \sigma\left(\sum_j w_j^o \sigma\left(\sum_i w_{ij}^h x_i + b_j^h\right) + b^o\right)$

17

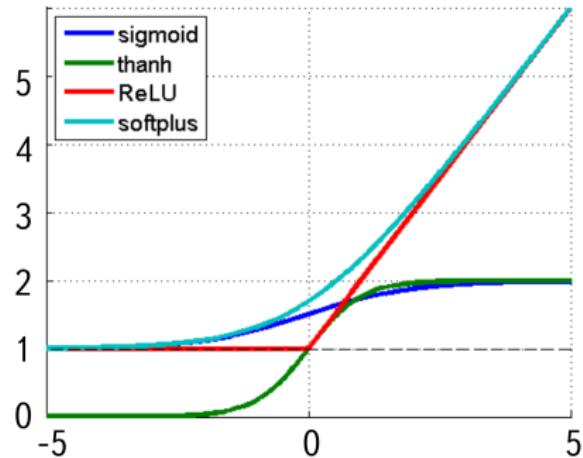
# Types of nonlinearity

$$f(a) = \frac{1}{1 + e^{-a}}$$

$$f(a) = \tanh(a)$$

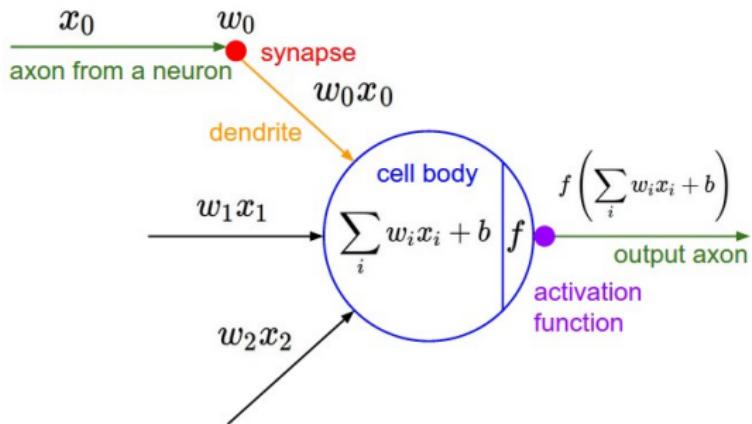
$$f(a) = \max(0, a)$$

$$f(a) = \log(1 + e^a)$$

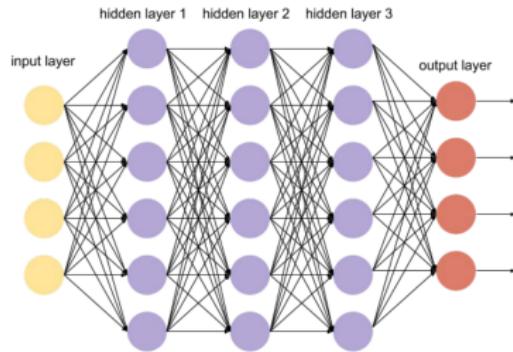


Wei Di, <https://imiloainf.wordpress.com/2013/11/06/rectifier-nonlinearity/>

# Neuron



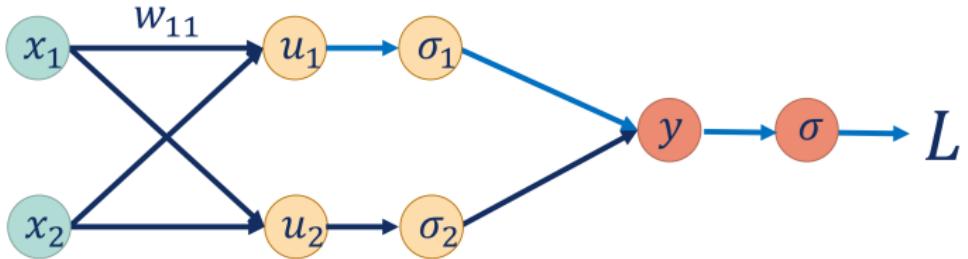
# Neural network



$$h_t = f(Wh_{t-1} + b)$$

$$a(x) = f_3(W_3f_2(W_2f_1(W_1f_0(W_0x + b_0) + b_1) + b_2) + b_3)$$

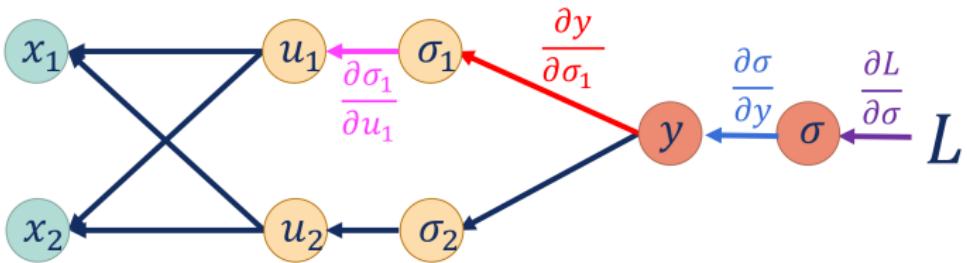
# The backpropagation algorithm



$$w_{11}^{(t+1)} = w_{11}^{(t)} - \gamma_t \frac{\partial L}{\partial w_{11}}$$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial u_1} \frac{\partial u_1}{\partial w_{11}} = \frac{\partial L}{\partial u_1} x_1$$

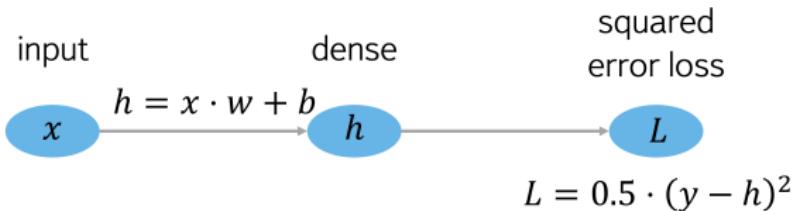
need to know gradients at each graph's node!



$$\frac{\partial L}{\partial u_1} = \frac{\partial L}{\partial \sigma} \frac{\partial \sigma}{\partial y} \frac{\partial y}{\partial \sigma_1} \frac{\partial \sigma_1}{\partial u_1} \leftarrow \text{chain rule}$$

**Backprop in a nutshell:** using the chain rule  
 propagate loss derivatives back going over each node step by step  
 starting from the NN's end

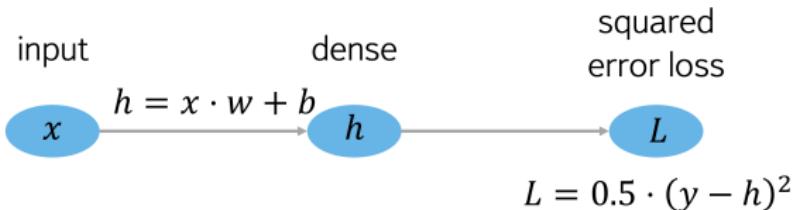
# The simplest NN ever



- Parameters:
  - Weight  $w$  and bias  $b$
- Input:  $x$
- Target:  $y$

Also known as  
least squares linear regression

# The simplest NN ever

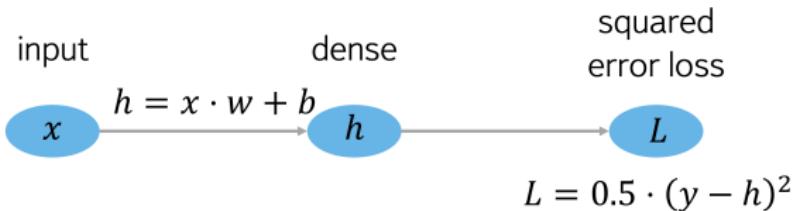


- Parameters:
  - Weight  $w$  and bias  $b$
- Input:  $x$
- Target:  $y$

$L$  is just a function of parameters, features and target:

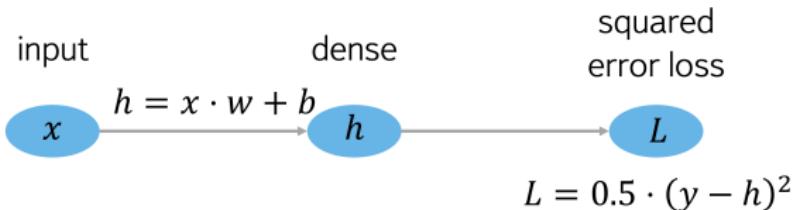
$$L = f(y, g(x, w, b))$$

# The simplest NN ever



- Gradient?
- $\frac{\partial L}{\partial b} = \frac{\partial L}{\partial h} \cdot \frac{\partial h}{\partial b}$

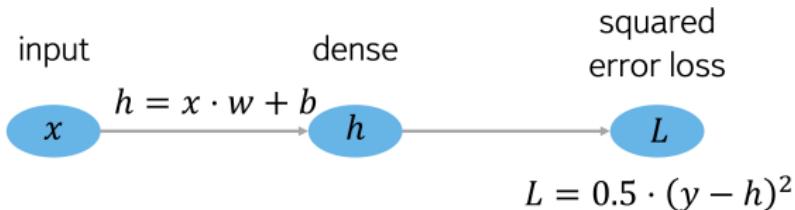
# The simplest NN ever



- Let's fit
  - $y = 3, x = 1$
- Initial
  - $w = 0.1, b = 1$

$h$	$L$	$\frac{\partial L}{\partial h}$	$\frac{\partial L}{\partial w}$	$\frac{\partial L}{\partial b}$	$w$	$b$

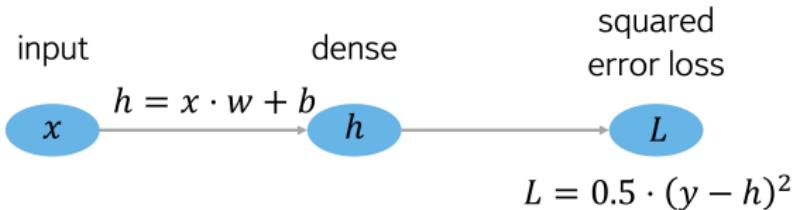
# Forward pass



- Let's fit
  - $y = 3, x = 1$
- Initial
  - $w = 0.1, b = 1$

$h$	$L$	$\frac{\partial L}{\partial h}$	$\frac{\partial L}{\partial w}$	$\frac{\partial L}{\partial b}$	$w$	$b$
1.1	1.80					

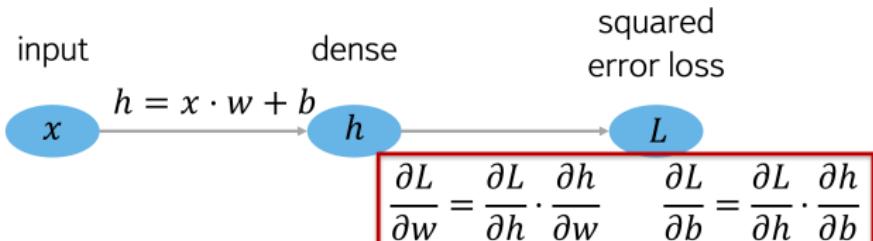
# Backward pass



- Let's fit
  - $y = 3, x = 1$
- Initial
  - $w = 0.1, b = 1$

$h$	$L$	$\frac{\partial L}{\partial h}$	$\frac{\partial L}{\partial w}$	$\frac{\partial L}{\partial b}$	$w$	$b$
1.1	1.80	-1.9				

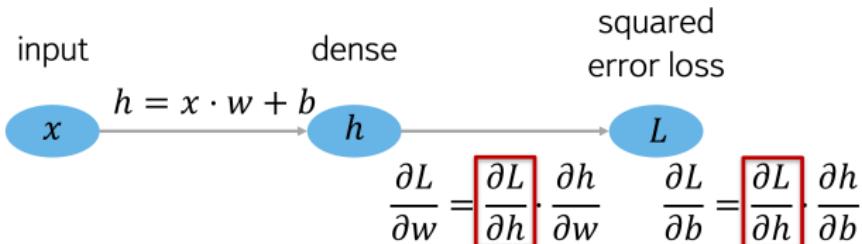
# Backward pass



- Let's fit
  - $y = 3, x = 1$
- Initial
  - $w = 0.1, b = 1$

$h$	$L$	$\frac{\partial L}{\partial h}$	$\frac{\partial L}{\partial w}$	$\frac{\partial L}{\partial b}$	$w$	$b$
1.1	1.80	-1.9				

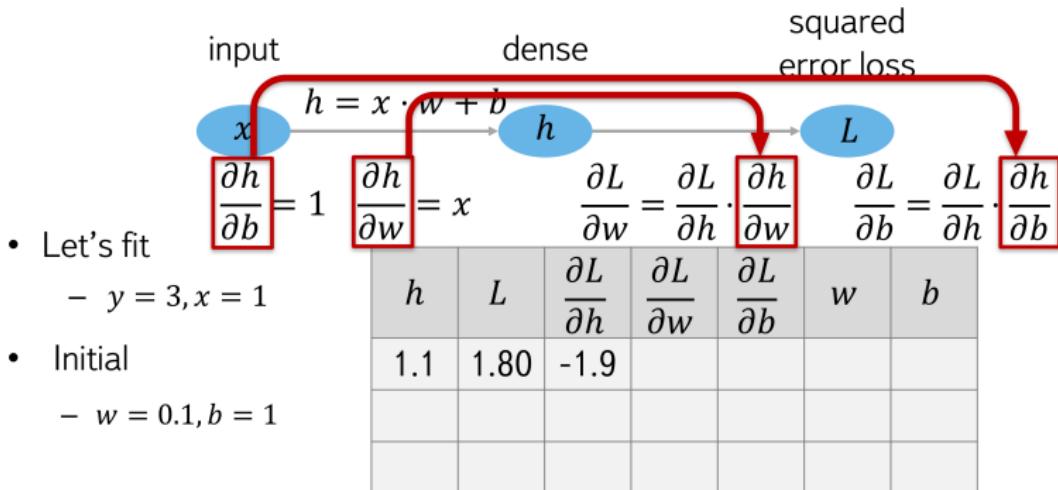
# Backward pass



- Let's fit
  - $y = 3, x = 1$
- Initial
  - $w = 0.1, b = 1$

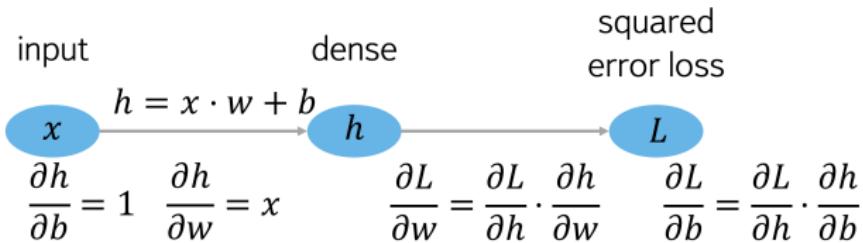
$h$	$L$	$\frac{\partial L}{\partial h}$	$\frac{\partial L}{\partial w}$	$\frac{\partial L}{\partial b}$	$w$	$b$
1.1	1.80	-1.9				

# Backward pass



31

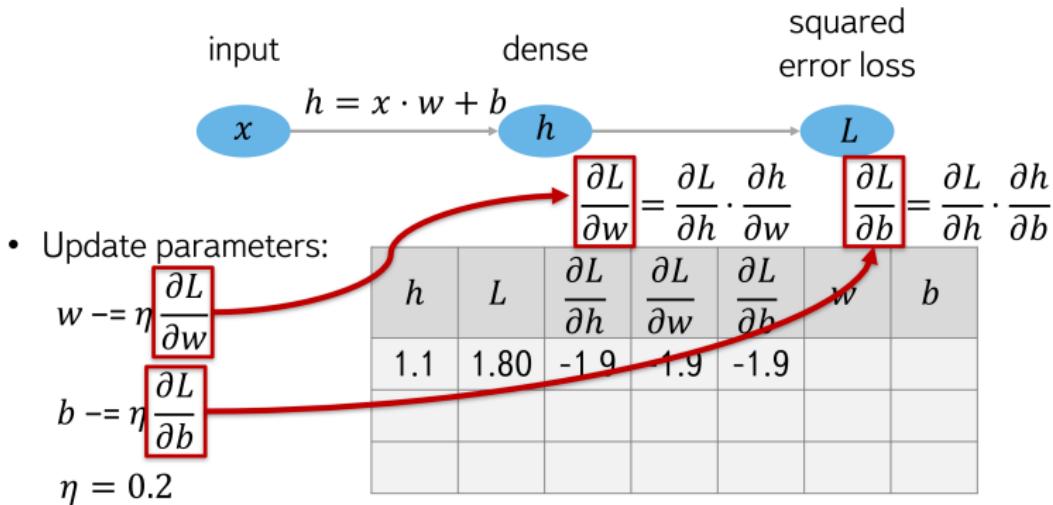
# Backward pass



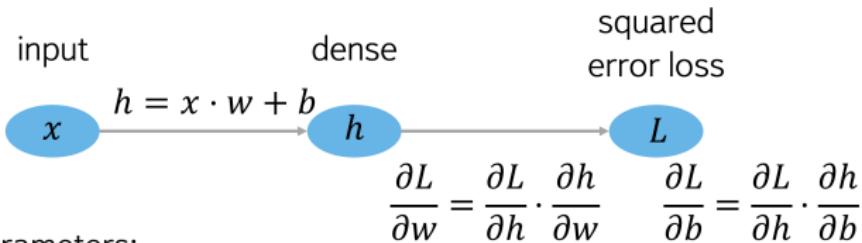
- Let's fit
  - $y = 3, x = 1$
- Initial
  - $w = 0.1, b = 1$

$h$	$L$	$\frac{\partial L}{\partial h}$	$\frac{\partial L}{\partial w}$	$\frac{\partial L}{\partial b}$	$w$	$b$
1.1	1.80	-1.9	-1.9	-1.9		

# Backward pass



# Backward pass



- Update parameters:

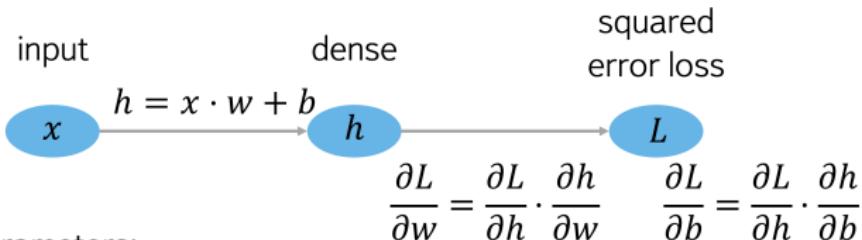
$$w \leftarrow \eta \frac{\partial L}{\partial w}$$

$$b \leftarrow \eta \frac{\partial L}{\partial b}$$

$$\eta = 0.2$$

$h$	$L$	$\frac{\partial L}{\partial h}$	$\frac{\partial L}{\partial w}$	$\frac{\partial L}{\partial b}$	$w$	$b$
1.1	1.80	-1.9	-1.9	-1.9	0.48	1.38

## After a few more updates...



- Update parameters:

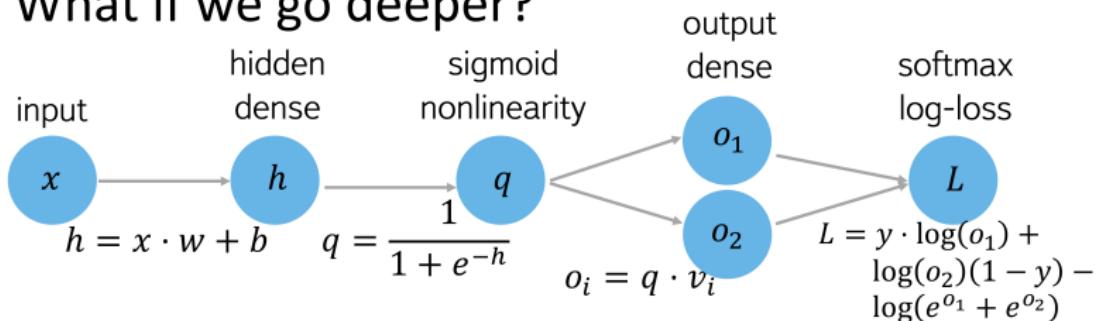
$$w \leftarrow \eta \frac{\partial L}{\partial w}$$

$$b \leftarrow \eta \frac{\partial L}{\partial b}$$

$$\eta = 0.2$$

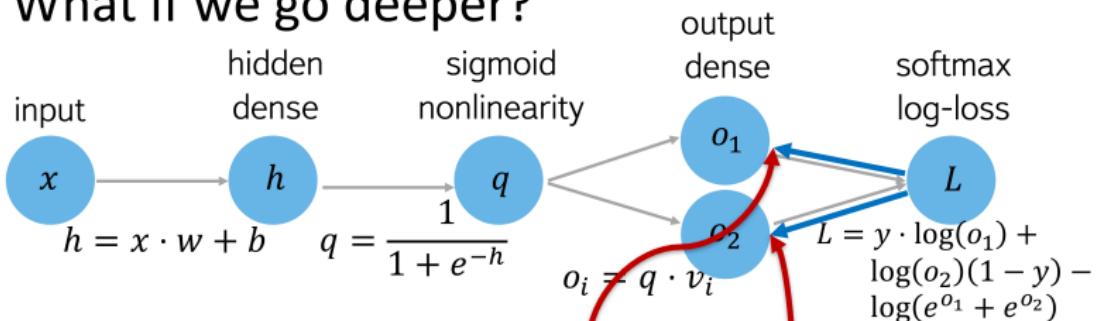
$h$	$L$	$\frac{\partial L}{\partial h}$	$\frac{\partial L}{\partial w}$	$\frac{\partial L}{\partial b}$	$w$	$b$
1.1	1.80	-1.9	-1.9	-1.9	0.48	1.38
1.86	0.65	-1.14	-1.14	-1.14	0.71	1.61
2.32	0.23	-0.68	-0.68	-0.68	0.84	1.75

# What if we go deeper?



- Parameters:
  - Weight  $w$  and bias  $b$
  - Weights  $v_1, v_2$

# What if we go deeper?

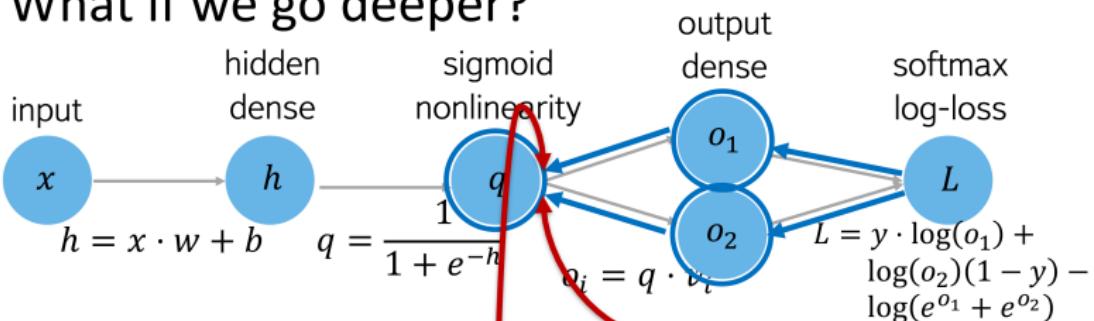


- Parameters:
  - Weight  $w$  and bias  $b$
  - Weights  $v_1, v_2$

$$\boxed{\frac{dL}{do_1}} = \frac{y}{o_1} - \frac{e^{o_1}}{e^{o_2} + e^{o_1}}$$

$$\boxed{\frac{dL}{do_2}} = \frac{1 - y}{o_2} - \frac{e^{o_2}}{e^{o_2} + e^{o_1}}$$

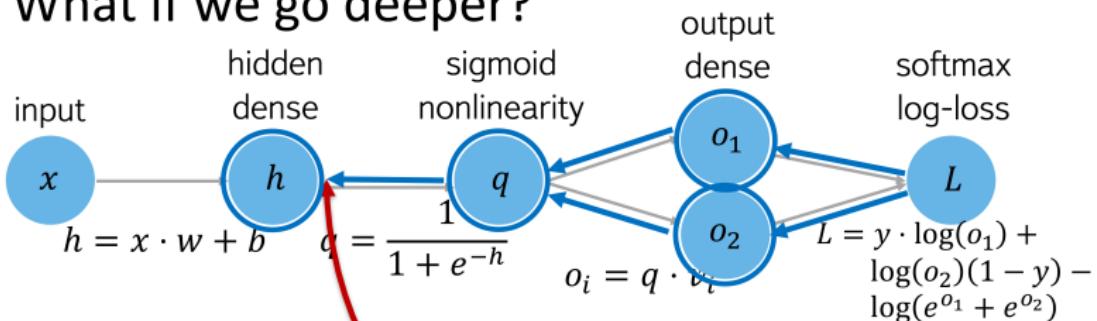
# What if we go deeper?



- Parameters:
  - Weight  $w$  and bias  $b$
  - Weights  $v_1, v_2$

38

# What if we go deeper?



- Parameters:

- Weight  $w$  and bias  $b$
- Weights  $v_1, v_2$

$$\boxed{\frac{\partial L}{\partial h}} = \frac{\partial L}{\partial q} \frac{e^{-q}}{(1 + e^{-q})^2}$$

[more maths on NN](#) [link]

## Backpropagation: the algorithm

- Chain rule can be evaluated numerically!
- Compute the network output and the loss value
- Compute "dLoss" / "dActivation\_of\_output\_layer"
- For each layer, starting from the last:
  - Compute "dActivation" / "dLayer\_parameters",  
"dActivation" / "dLayer\_input"
  - Multiply it by "dLoss" / "dActivation", get "dLoss" /...
- Make optimization step for the parameters

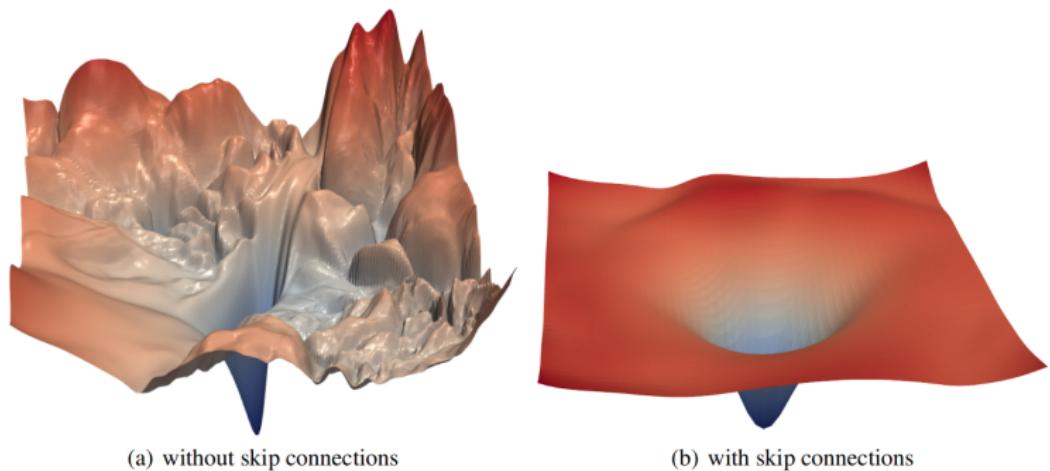


Figure 1: The loss surfaces of ResNet-56 with/without skip connections. The proposed filter normalization scheme is used to enable comparisons of sharpness/flatness between the two figures.

# Practical issues

- bad convergence
  - lots of local minima, can get stuck there
  - overfitting as a piece of cake
- gradient flow problems
  - explosion
  - vanishing

# Practical issues

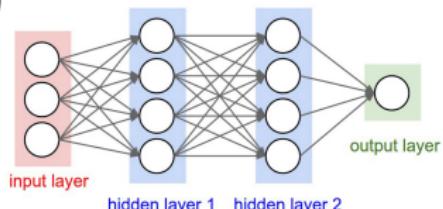
- bad convergence
  - lots of local minima, can get stuck there
  - overfitting as a piece of cake
- gradient flow problems ⇒ need addit. heuristics
  - explosion
  - vanishing

- Optimization

- grad. descent modifications (e.g., SGD, Adam, Nesterov momentum)
- learning rate schedules (e.g., LR decay)
- early stopping

- Architecture

- activations (e.g., SELU, Leaky RELU, Swish)
- dropout
- batch normalization
- weight initialization (why 0-init is bad?)
- penalty loss terms (L1, L2)



TF vs PyTorch comparison [link]

## Frameworks

you should know about



- TensorFlow
- Keras (TF & Theano back-end)
- PyTorch
- ? Theano



# Baseline

you should start from

- normalize data
- train in batches
- RELU activation
- He et al. weight init
- Adam optimizer
- hyperparams with random search (not grid) or go bayesian
- try BatchNorm
- try Dropout ( $p=0.5$ ) and L1/L2 regularization

# Neural Nets

## pros & cons

- + damn very good performance
- + general to adapt to various problem domains
- + modular and easy to construct
- + automatic feature engineering
- + breaks No Free Lunch theorem (almost) [1],[2]
- overfits
- requires careful baby-sitting
- you will probably need to get GPU
- "black box" (no)

# Here goes fancy

# Convolutional NN



# Convolutional NN

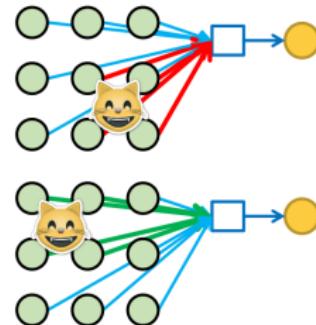
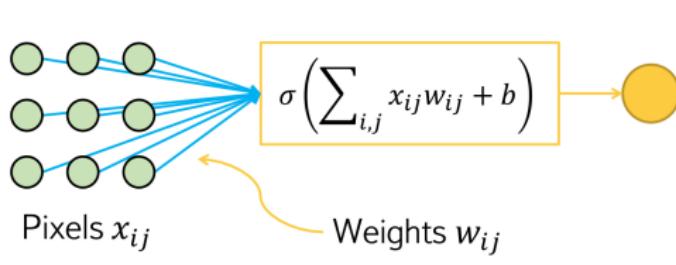


		Blue channel						
		Green channel						
		Red channel						
1	120	67	89	107	...	13	39	
2	12	216	145	26	...	181	18	
3	0	16	4	45	...	44	8	
4	0	78	90	167	...	25	81	
...	...	...	...	...	...	...	71	
64	12	67	82	141	...	12	56	
	1	2	3	4	...	64	...	

Image array: {64 x 64 x 3}

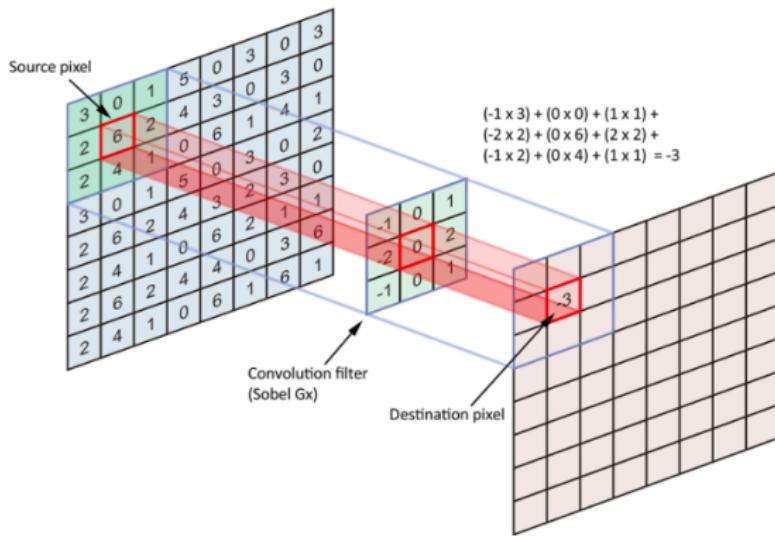
# Convolutional NN

why not Dense?



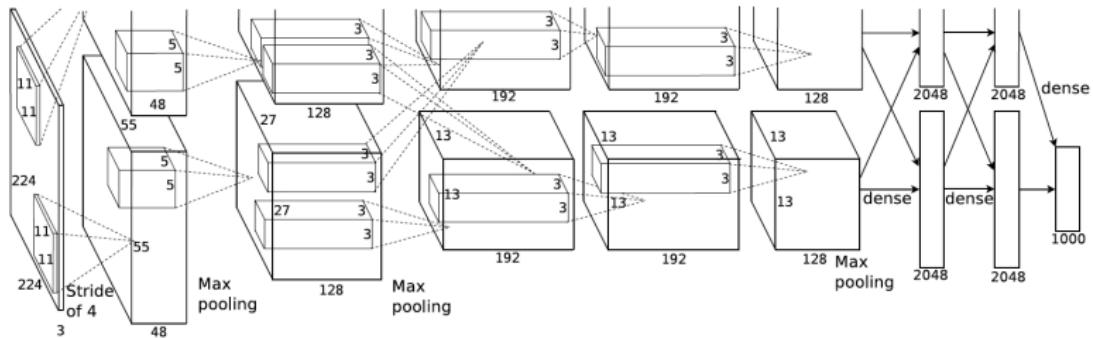
# Convolutional NN

## conv. layer



# Convolutional NN

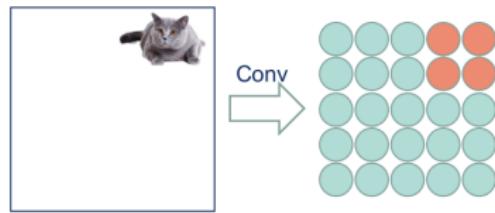
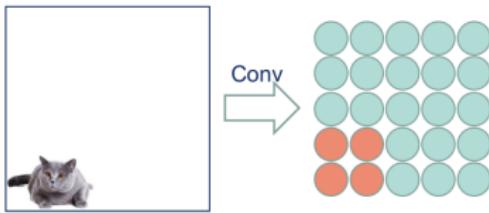
## AlexNet



# Convolutional NN

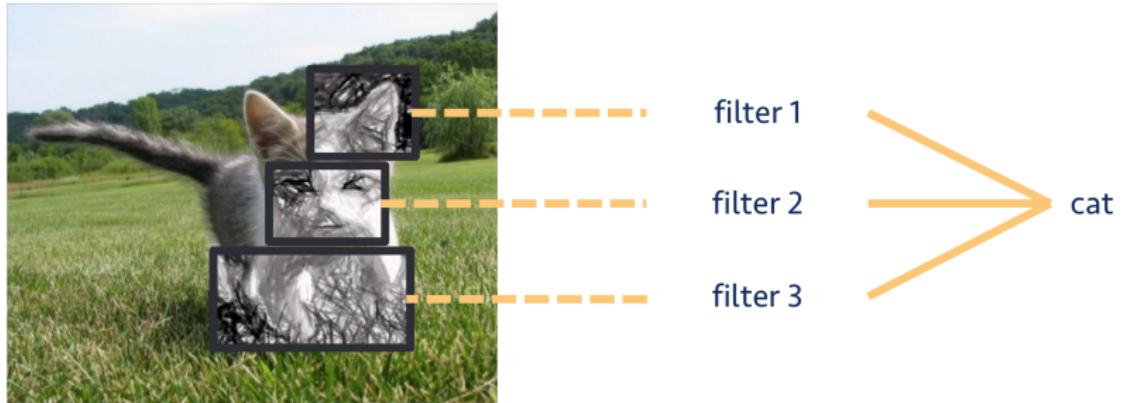
## equivariance

$$f(T(x)) = T(f(x))$$

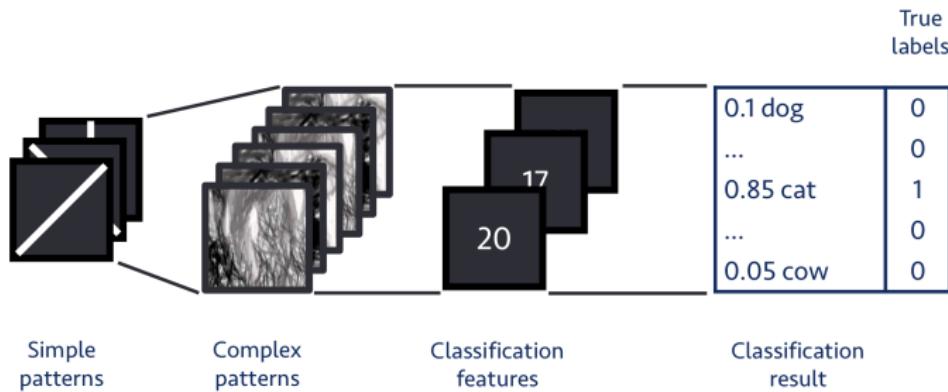


# Convolutional NN

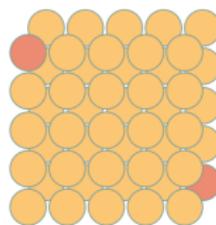
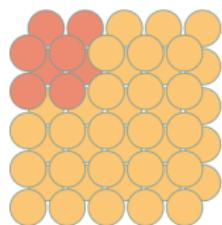
patterns



# Convolutional NN patterns

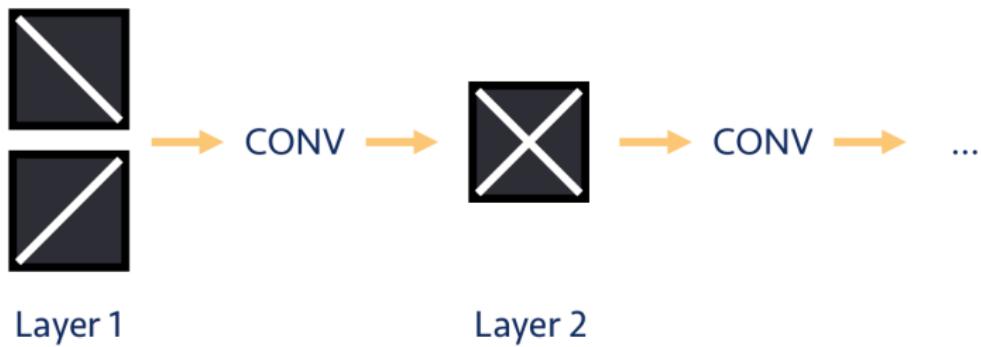


# Convolutional NN patterns

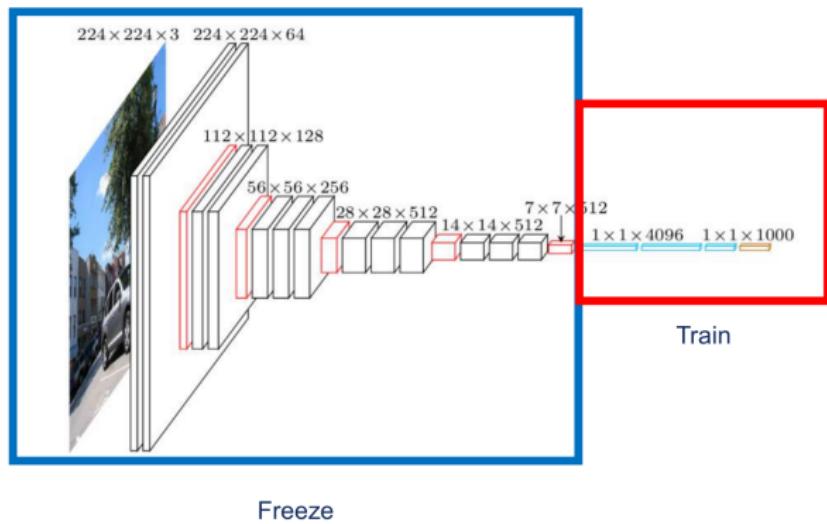


# Convolutional NN

patterns



# Convolutional NN modularity



# Recurrent NN

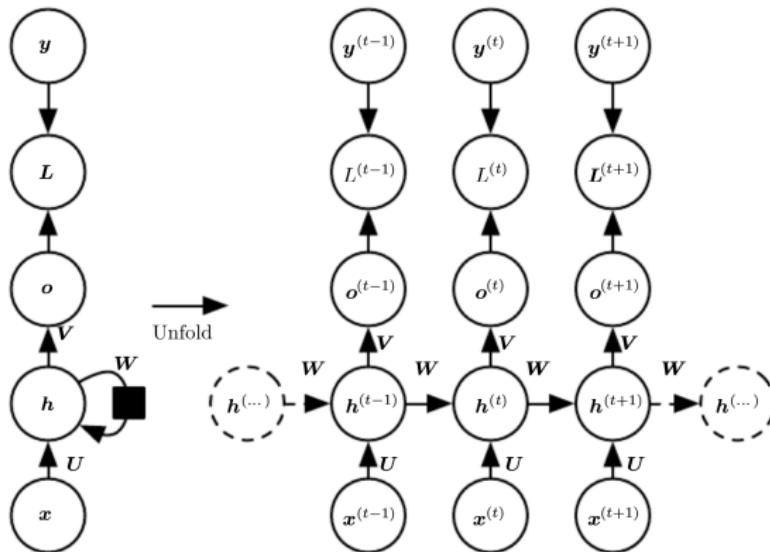
## motivation

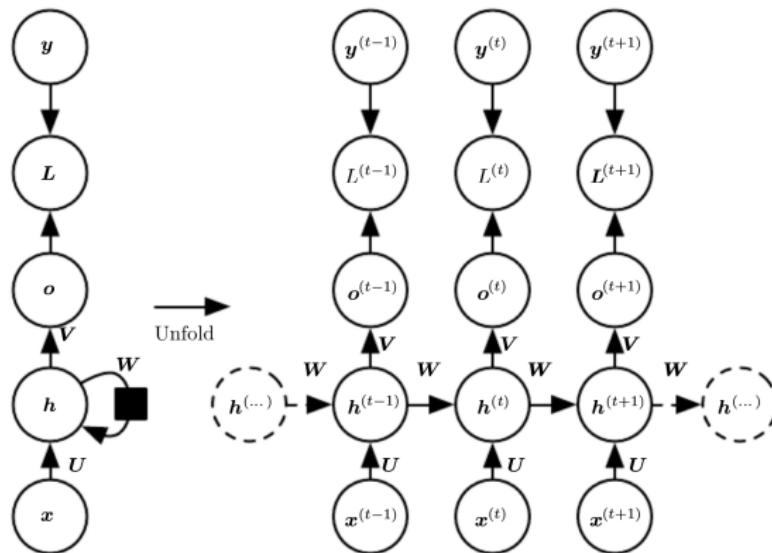
- sequence processing
  - sentence
  - song
  - jet
  - market stock
- problems
  - translation
  - speech recognition
  - jet grooming
  - price prediction
- why not Dense NN?

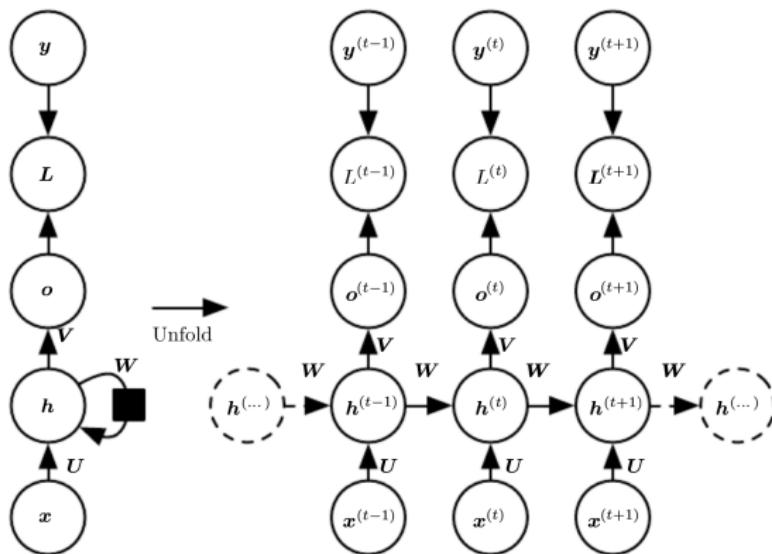
# Recurrent NN

## motivation

- sequence processing
  - sentence
  - song
  - jet
  - market stock
- problems
  - translation
  - speech recognition
  - jet grooming
  - price prediction
- why not Dense NN?
  - no time invariance
  - fixed sequence size

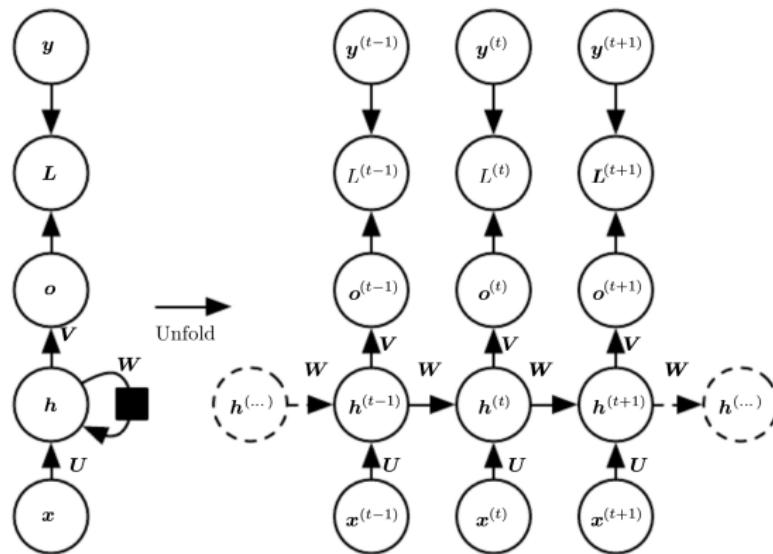


what is  $x$ ?

what is  $x$ ?

embeddings

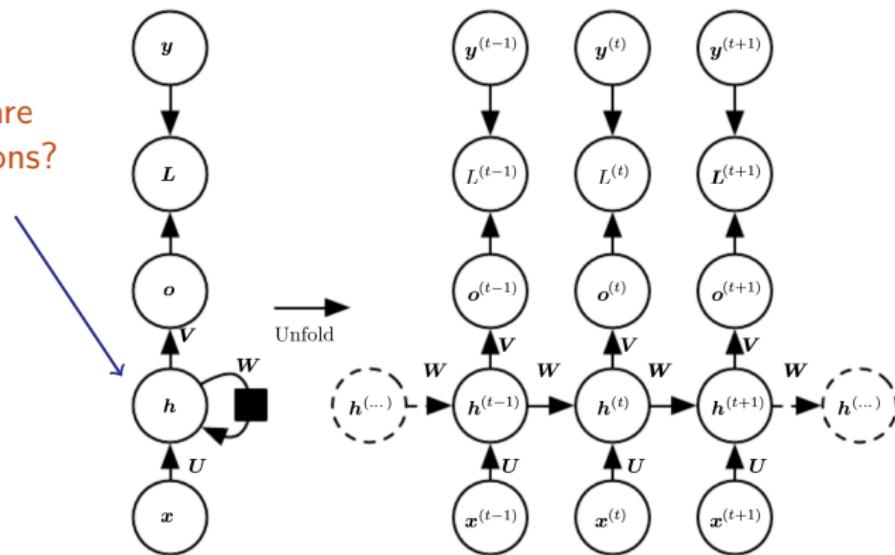
where are  
activations?



what is  $x$ ?

embeddings

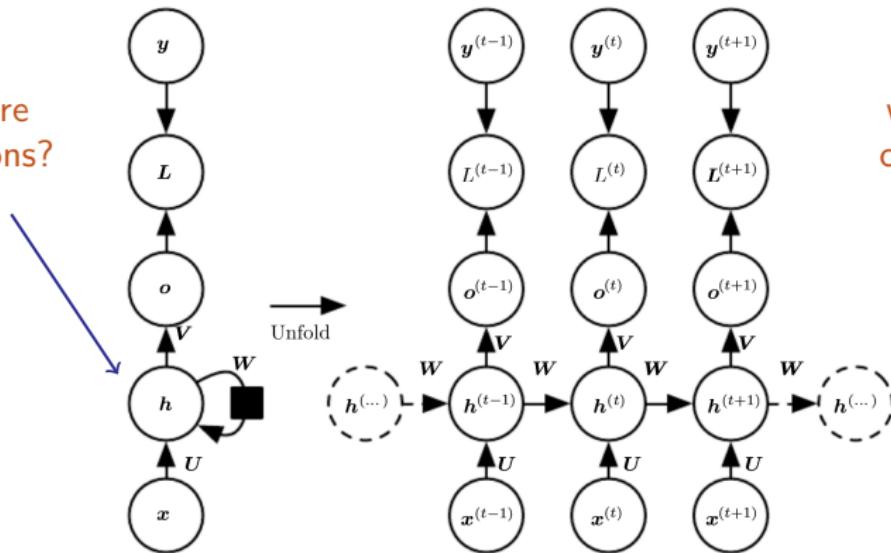
where are  
activations?



what is  $x$ ?

embeddings

where are  
activations?

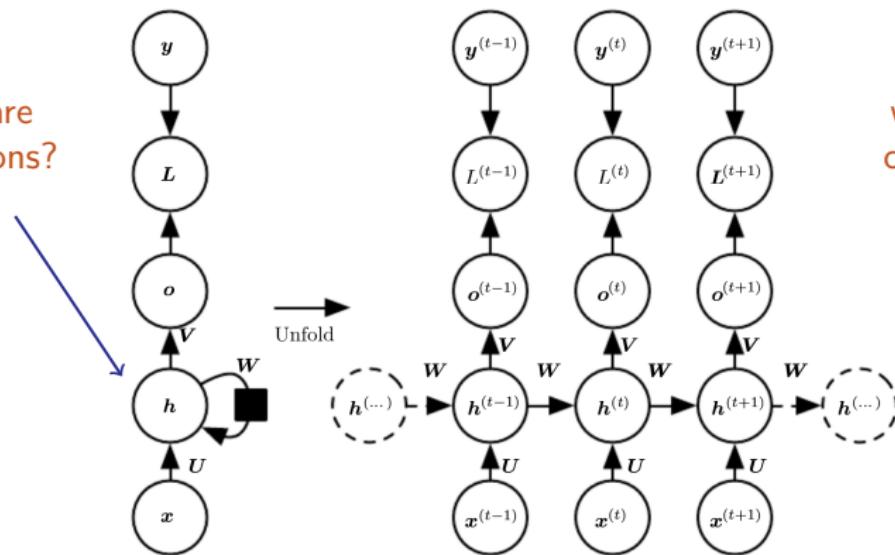


where this one  
can be useful?

what is  $x$ ?

embeddings

where are  
activations?



where this one  
can be useful?

## seq labeling

what is x?

## embeddings

# Summary

- Ensembles
    - Motivation
    - Bootstrap & Bagging → Random Forest
    - Boosting → BDT
  - Neural Nets
    - Motivation
    - Building
    - Training
    - Special tricks
  - \* Bonus
    - Convolutional NN
    - Recurrent NN
- next lecture
- Thrilling applications in HEP



Крафтер Колян  
28 июня 2018



- Здравствуйте, Николай. Поздравляем с успешным выступлением в контесте "Обучение скорости обучения". Нам хотелось бы обсудить ваше решение.
- Хорошо, без проблем.

-----  
Через 20 минут в офисе Яндекса (Отдел исследований)

- 
- Итак, по заданию надо было предсказать learning rate для разных выборок, чтобы алгоритм CatBoost давал максимальную точность. Как вы это делали?
- Ну, там данных немного, всего 80 примеров, я закинул их в Excel и решил посмотреть, а что собственно есть.
- Так, что дальше?
- Оставил только примеры с большим N, заметил, что у меня в обучающей выборке есть почти то же самое, что и для тестовой. Для них ответ сразу выписал.
- Вы нашли зависимость от N?
- Нет, этого же не было в задании.
- Вы использовали gradient boost или регрессию?
- Если честно, я ни слова не понимаю.
- Вы использовали какие-нибудь алгоритмы, математику?
- Нет, я просто искал для каждого теста похожий пример в обучающей выборке и давал такой же ответ.
- Вы хотя бы калькулятором пользовались?
- Нет.
- И у Вас лучшее решение?
- Да, и у него точность выше на 0.15%, чем у вашего решения, причём у всех остальных максимум на 0.04%.

- ...

[link]

