

# Machine Learning with python

MNIT, Jaipur

Falak Shah  
Lead Research Scientist  
Infocusp Innovations Pvt. Ltd.

# Contents

— — —

Machine Learning Basics

Supervised Learning

Unsupervised (Clustering)

Deep Learning – ANN/ CNN/ RNN/ LSTM

# Machine Learning

- Machine learning is the science of getting computers to act without being explicitly programmed.

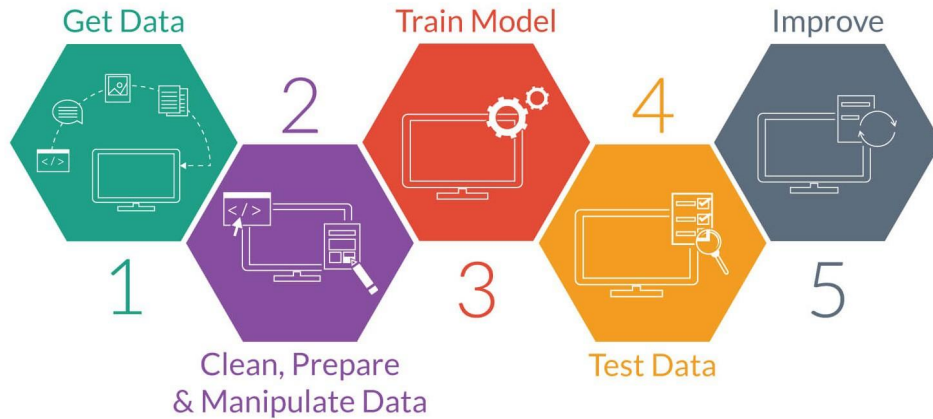


Image Courtesy: <https://upxacademy.com/artificial-intelligence-foundation/>

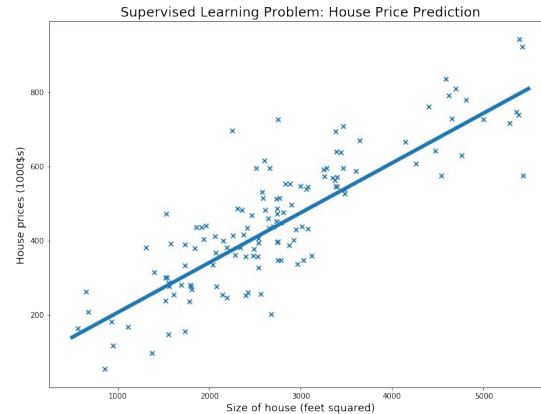
# Supervised Learning

— — —

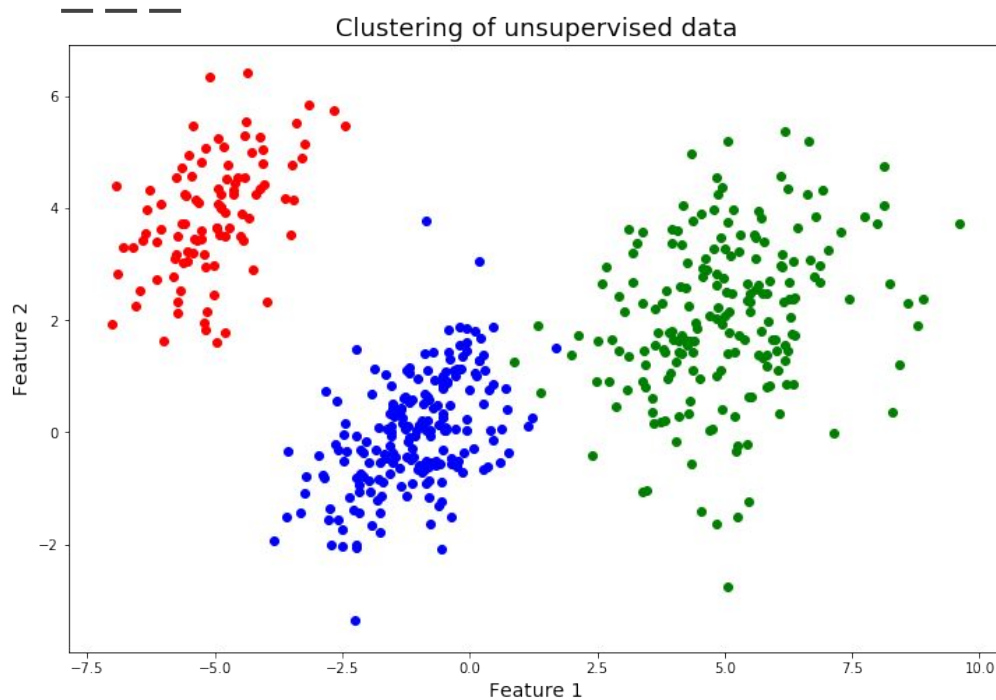
Input data with correct output labels which are used by machine in learning process.

**Example:** Housing prices

Area (feet <sup>2</sup> )	Price (1000\$s)
1210	234
1450	266
750	172
1200	232
1380	257
840	184



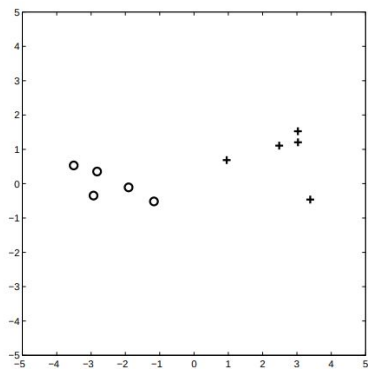
# Unsupervised Learning



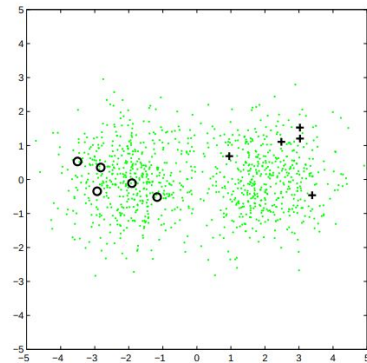
Actual labels are not provided in the data. Patterns and clusters are determined from the given data.

**Example** - Clustering news stories into categories

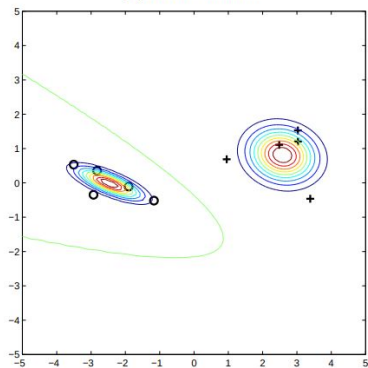
# Semi-supervised Learning



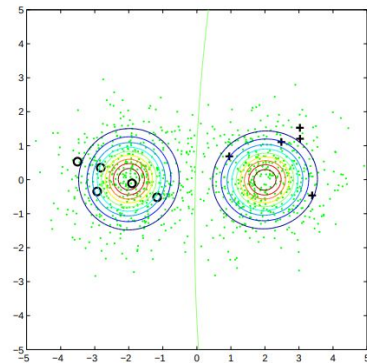
(a) labeled data



(b) labeled and unlabeled data (small dots)



(c) model learned from labeled data



(d) model learned from labeled and unlabeled data

Unsupervised Learning can be interpreted as:

- Supervised Learning + Additional Unlabelled Data
- Unsupervised Learning + Additional Labelled Data

Here, we are trying to predict Gaussian distributions that generated the given data.

Interested folks can read up more on Expectation Maximization (EM)

# Reinforcement Learning

— — —

The "cause and effect" idea can be translated into the following steps for an RL agent:

1. The agent observes an input state
2. An action is determined by a decision making function (policy)
3. The action is performed
4. The agent receives a scalar reward or reinforcement from the environment
5. Information about the reward given for that state / action pair is recorded

# Supervised Learning

---

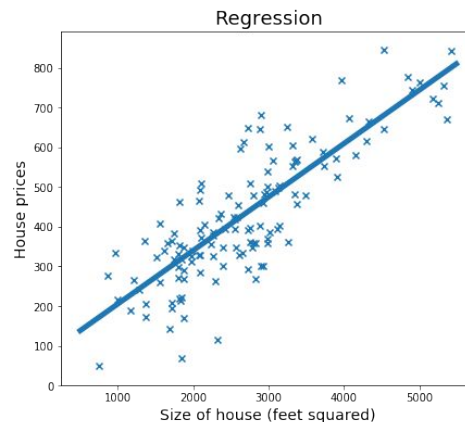
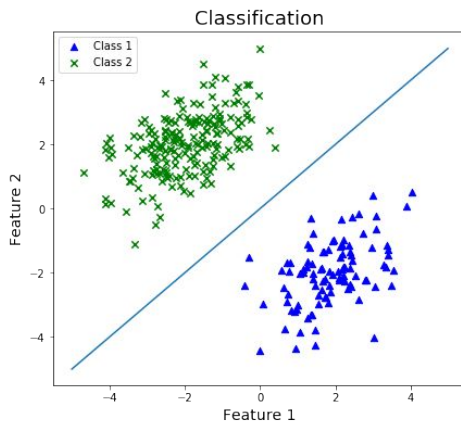
**Given:** Pair of input and output values.

**Goal:** Define a loss function and minimize that to predict accurate output values.



# Regression and Classification problems

— — —



Examples:

1. Classify images into different classes.
2. Housing price prediction based on features of house.

# Linear Regression

---

One of the simplest ML algorithms

Trying to find a linear mapping from input to output

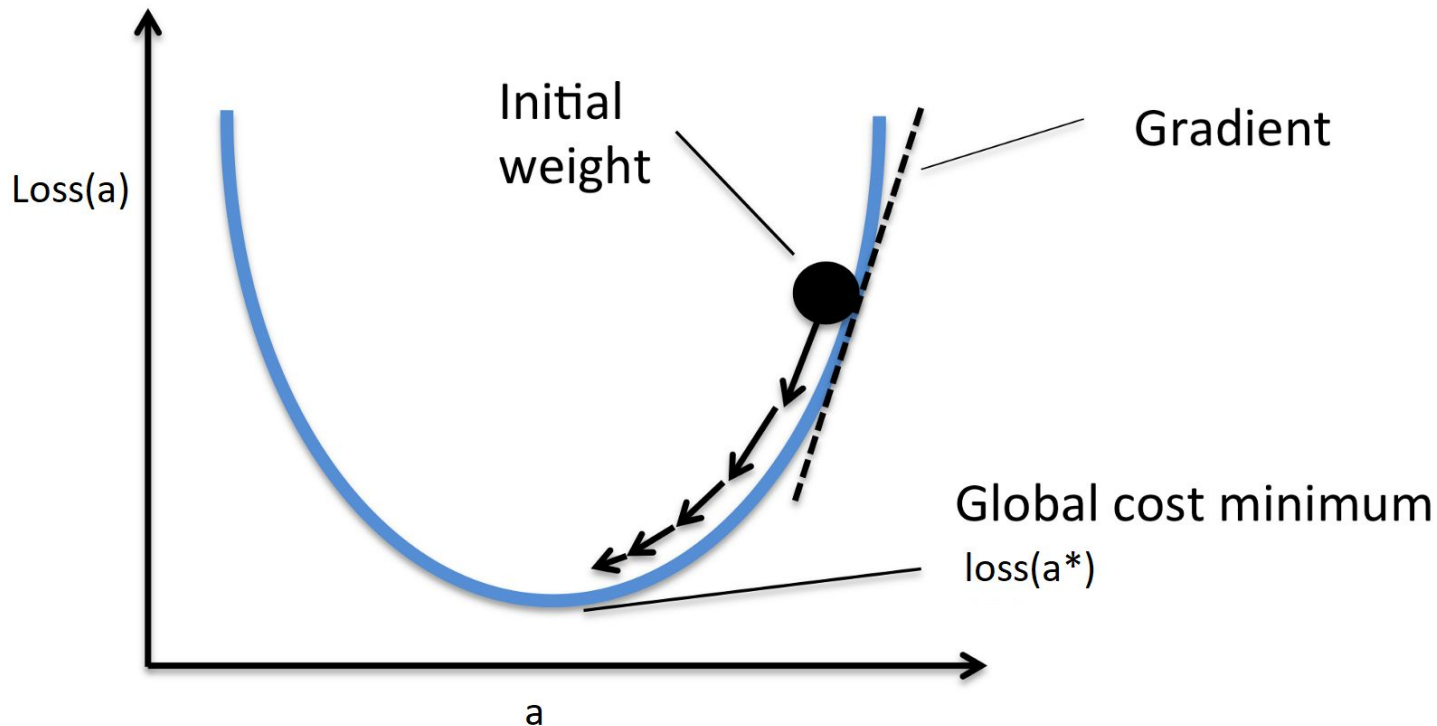
[Example](#)

$$y = ax + b$$

$$J = \frac{1}{2} \sum_{i=1}^n (y_i - (ax_i + b))^2$$

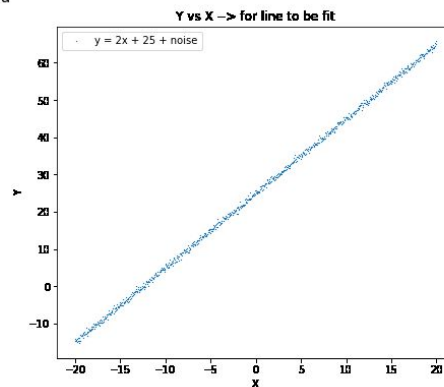
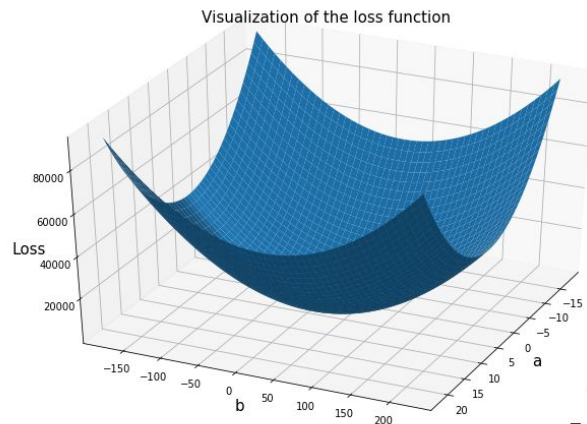
$$(a_{opt}, b_{opt}) = \underset{a, b}{\operatorname{argmin}}(J)$$

# Gradient Descent Intuition



[https://github.com/falakth/eoptimist/gradient\\_descent\\_optimizers/blob/master/Optimizers%20Linear%20Regression%20example.ipynb](https://github.com/falakth/eoptimist/gradient_descent_optimizers/blob/master/Optimizers%20Linear%20Regression%20example.ipynb)

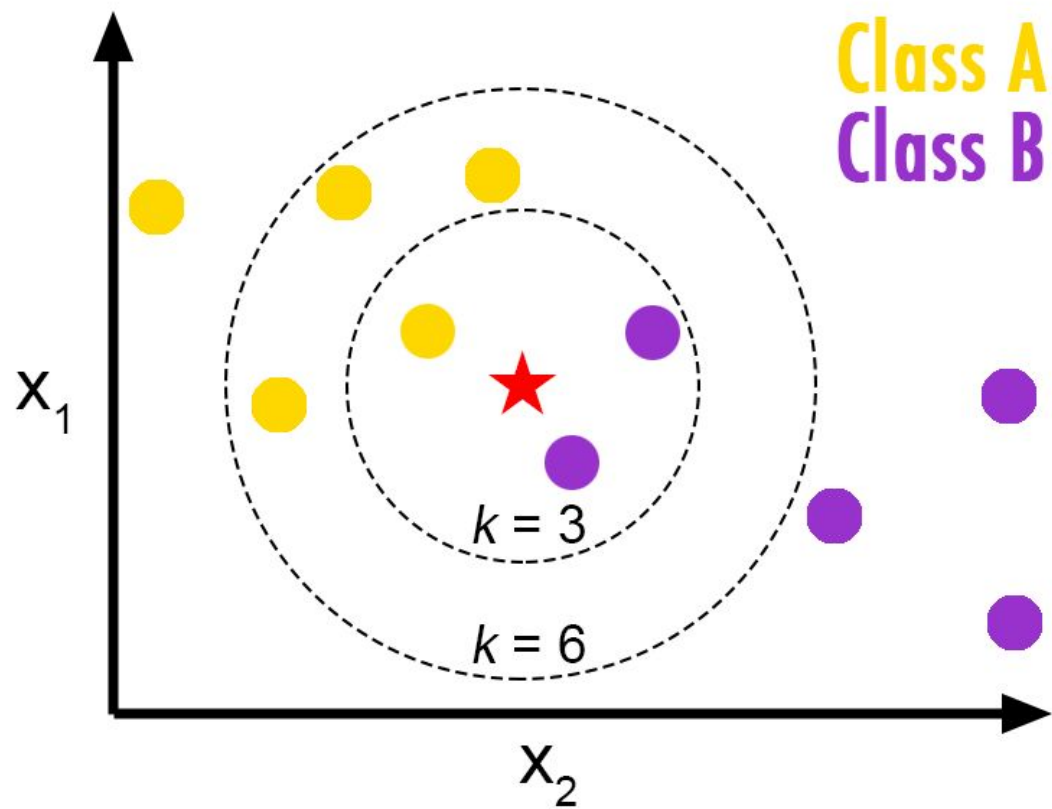
# Data and loss visualization



# Classification: Simple approach

---

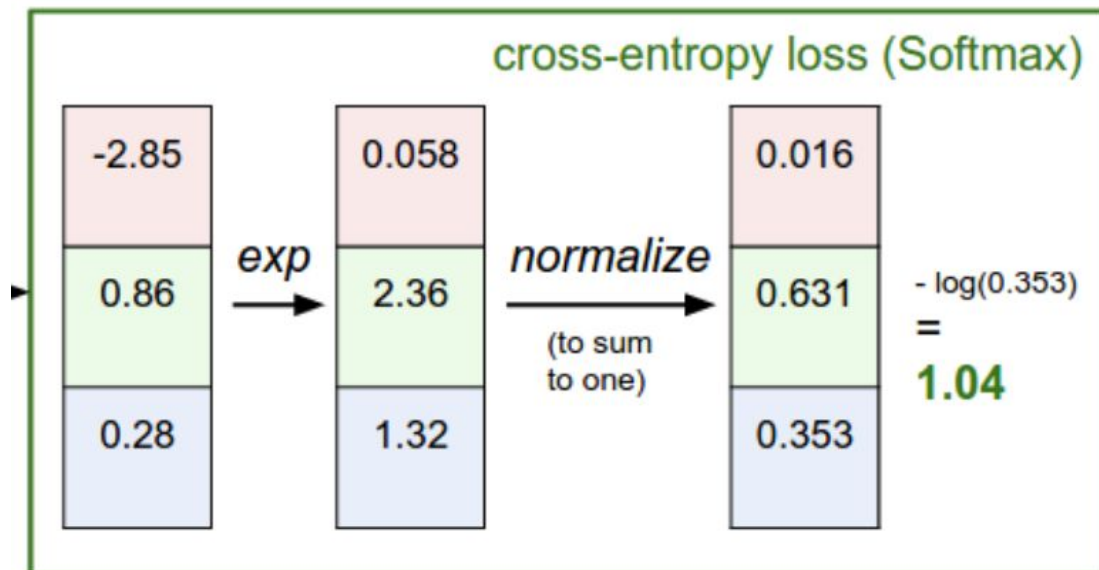
- Find  $k$  nearest neighbours
- This is why ML is intuitive
- No training needed
- Vectorized approach possible (50 seconds with for loops to 0.5 seconds with matrices)
- Important to learn linear algebra
- CS231n - Stanford - Begins with this
- How to decide  $k$ ?  $K$  folds holdout accuracy



# Classification problem

---

Softmax Loss



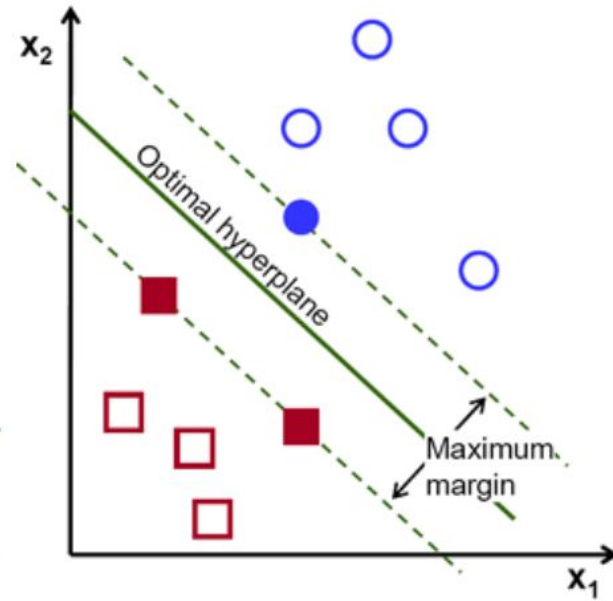
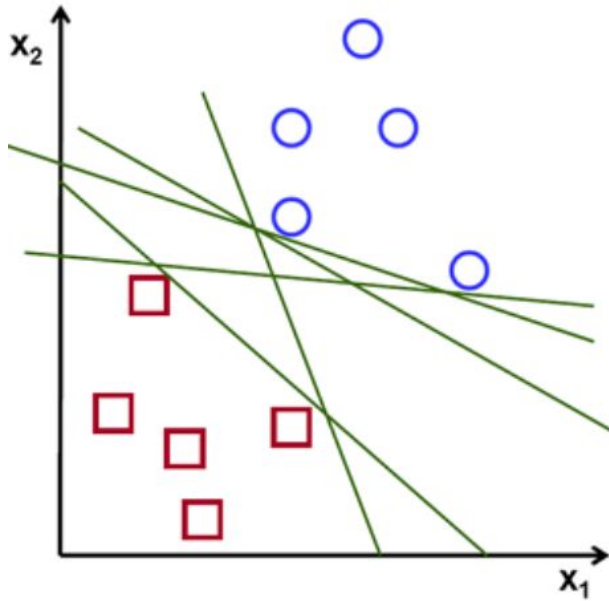
$$t = \begin{bmatrix} t1 \\ t2 \\ t3 \end{bmatrix}$$

$$loss = - \sum_{i=1}^3 y_i \log(P_i)$$

$$\text{where } P_i = \frac{e^{t_i}}{\sum_{j=1}^3 e^{t_j}}$$

# Support Vector Machines

---





# Kernels: to handle non linearly separable data

---

Take advantage of computing just the distance metric in higher dimension

Can map upto infinite dimensional hyperspace (RBF kernel)

Learning from Data by Prof Yaser of Caltech

Excellent explanations on the kernel trick

# Depending on the complexity of function you fit

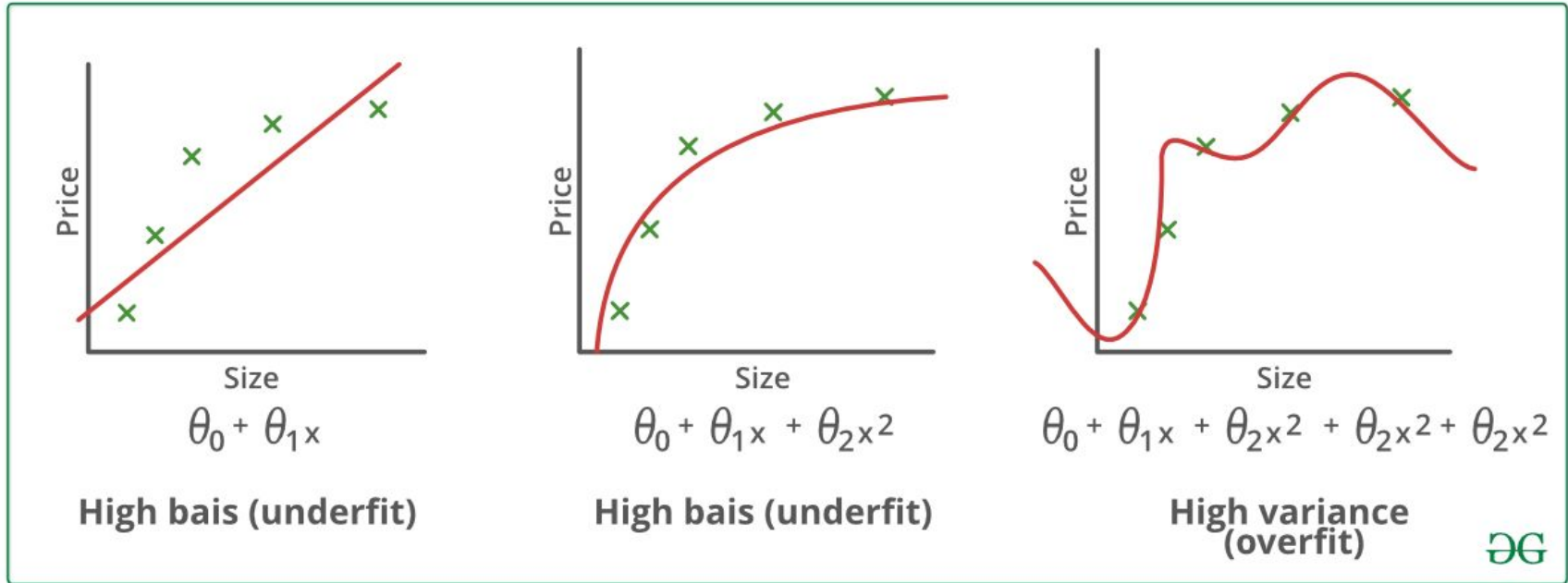


Image courtesy: Geeksforgeeks

# Let's try some hands on

— — —

<http://localhost:8889/notebooks/tutorial.ipynb>

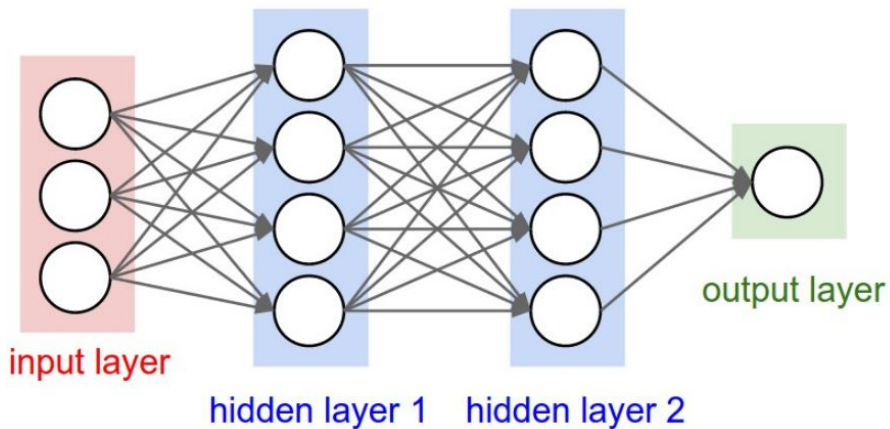
<https://colab.research.google.com/drive/1buDdba1UmgXYsXnE02ZZJSVJSjripGqE>

**Break?**  
**Back to deep learning**

# Artificial Neural Networks

— — —

- ANNs are used for modelling Non-linear hypothesis



$$H = \text{activation}(WX + b)$$
$$Y_{\text{pred}} = \text{activation}(WH + b)$$

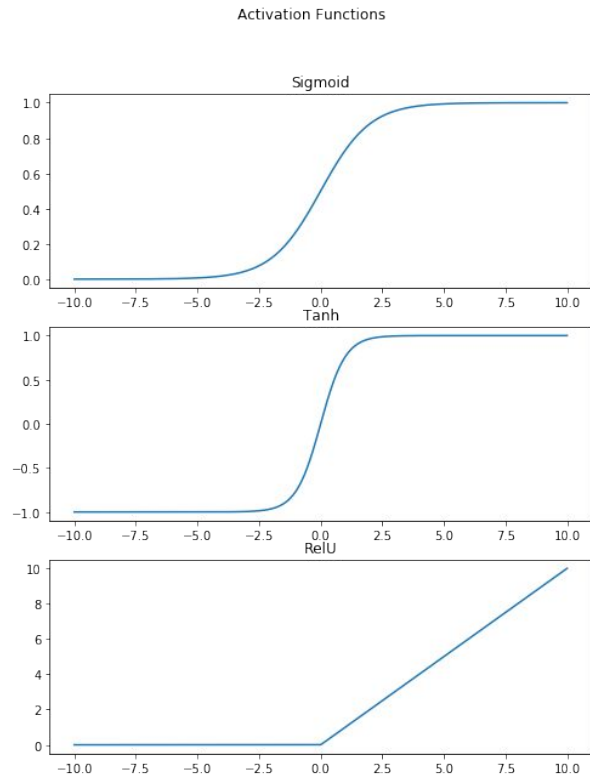
$$\text{Loss} = f(y, y_{\text{pred}})$$

$$\text{Softmax loss} = -\log(p_{y_{\text{true}}})$$

Image Courtesy: <http://cs231n.github.io/neural-networks-1/>

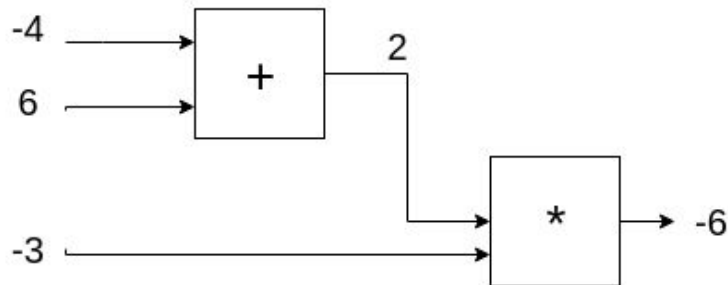
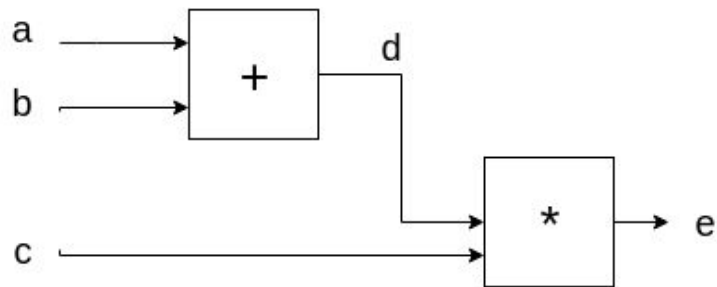
# Activation Functions

1. Sigmoid/ Logistic
2. Tanh - Hyperbolic tangent
3. ReLU - Rectified Linear Units



# Backpropagation Intuition

— — —



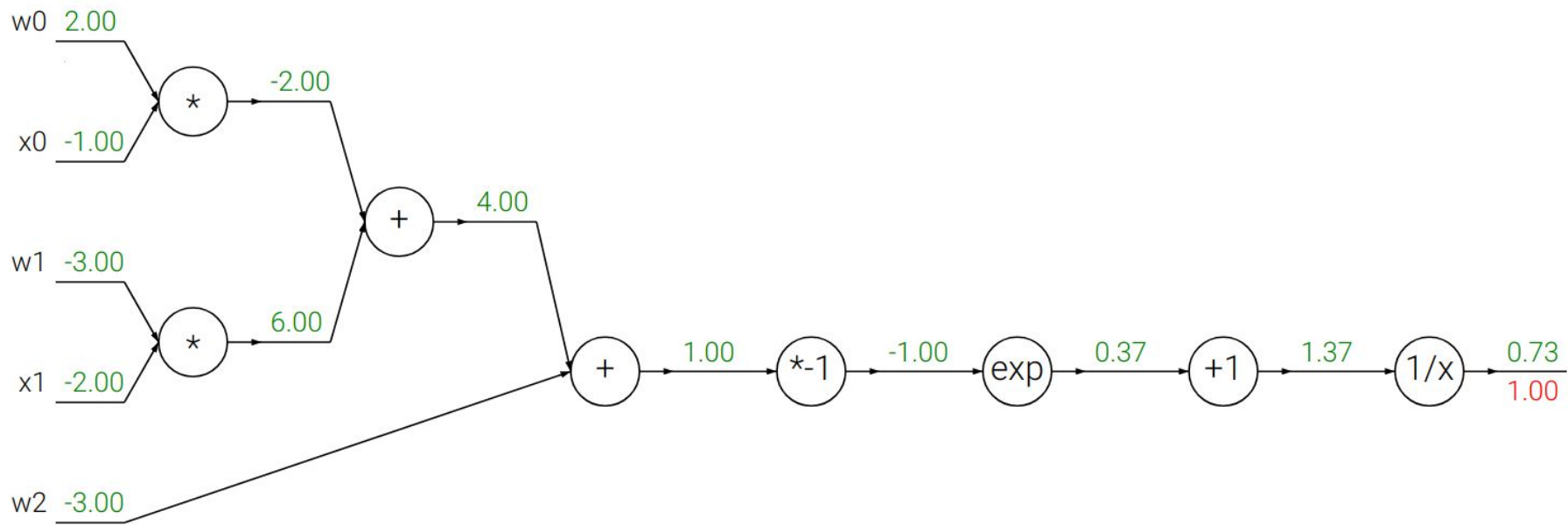
Once we have the loss at output layer, we want to back propagate that till the input layer such that we penalise the appropriate cells according to their contribution in the loss.

$$d(a,b) = a + b$$

$$e(a,b,c) = (a+b)*c$$

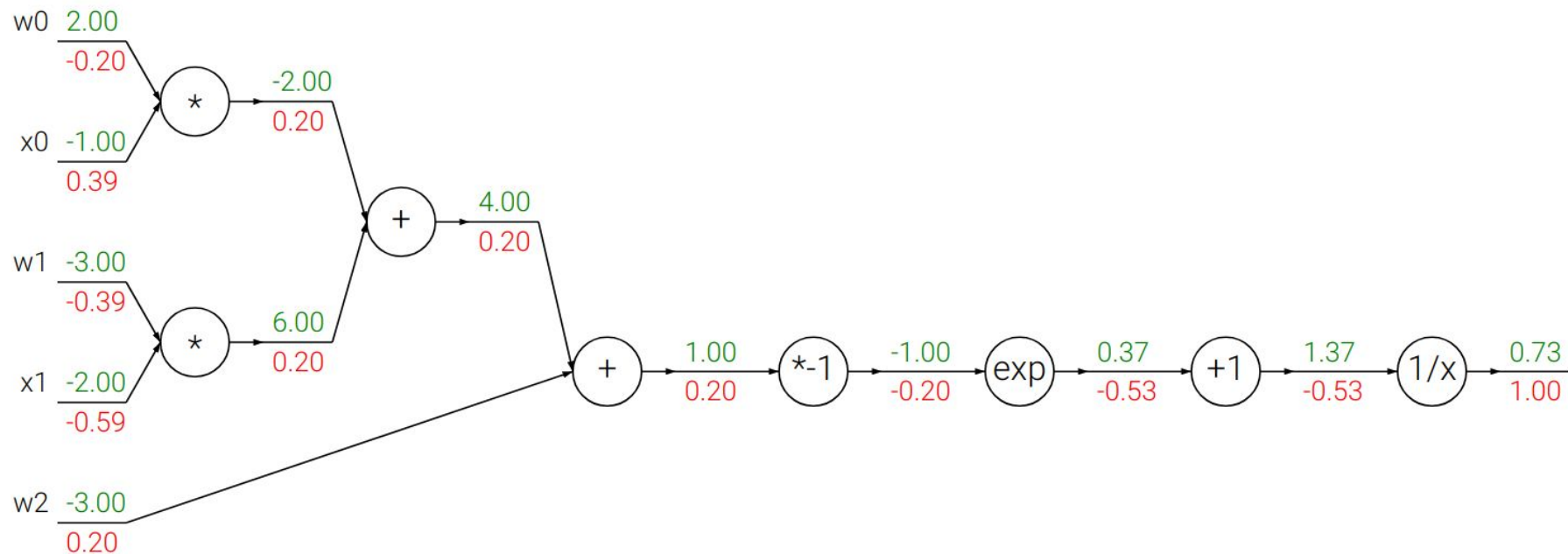
# Back-prop Example

— — —





# Back-prop Example

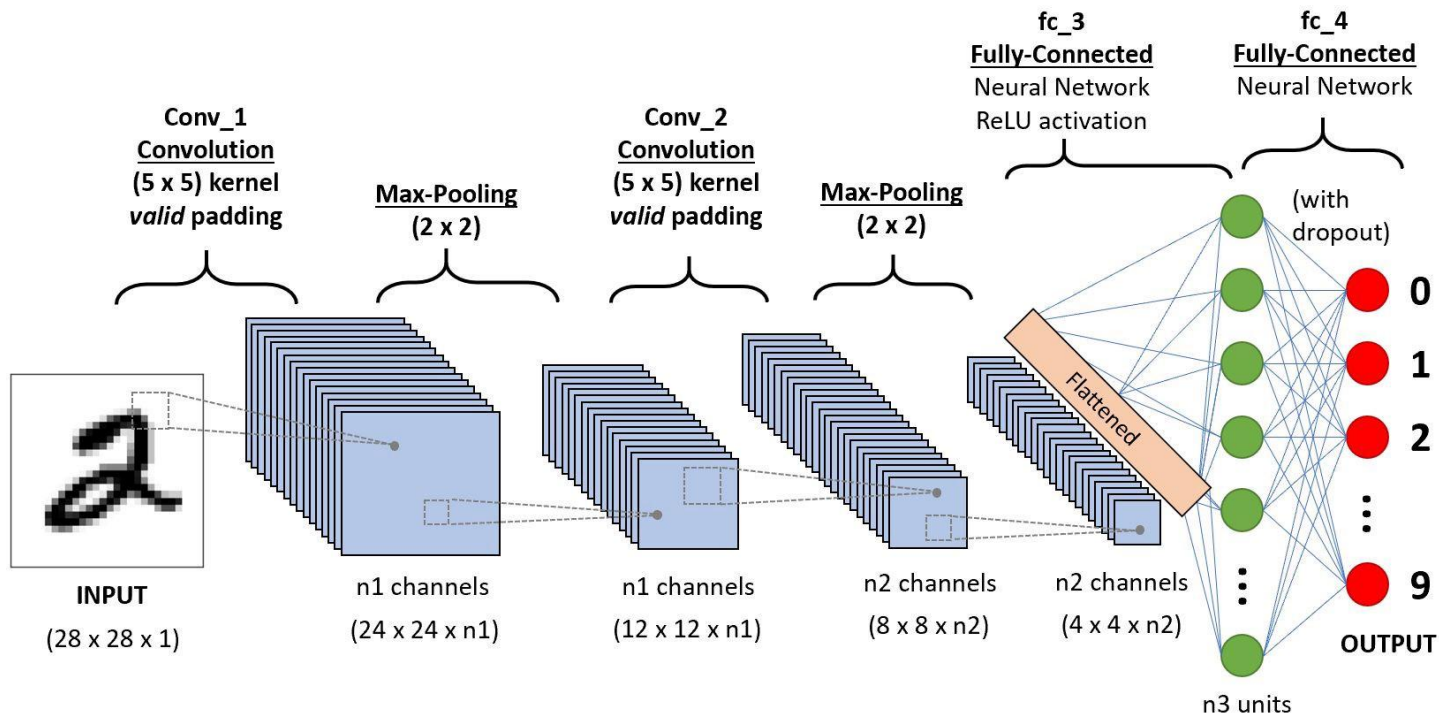


# Optimizers

— — —

- Different methods of taking steps when descending to the minimum of loss
- [Hands on](#)

# Convolutional Neural Network



# CNN used in

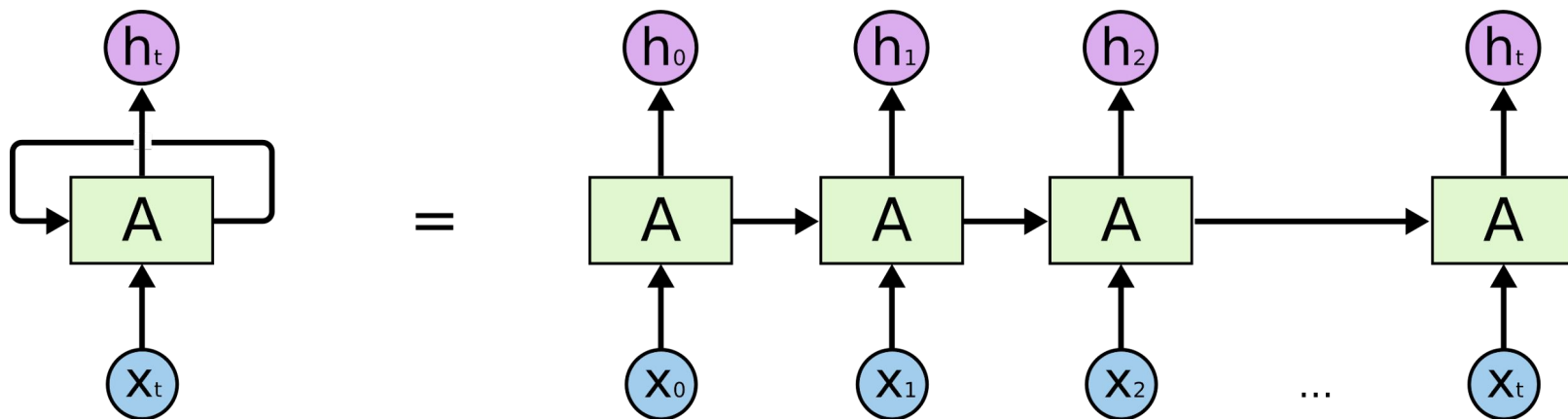
---

- Image classification
- Image segmentation (Bunch of contests on kaggle)
- Extracting features for caption generation
- Demo : caption generation
- Train CNN on MNIST in under 2 minutes: [Colab](#)
- Colabs and GPUs freely available - so are kaggle kernels
- Object detection from images : fast RCNN

# Recurrent Neural networks/ LSTMs

---

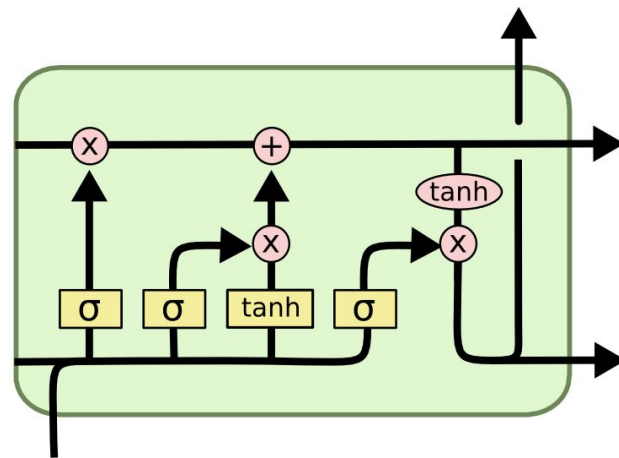
- Variable length sequences with dependencies across time
- Cannot be captured by static ANN/ CNN
- Try out on a small problem for understanding and then see its application: [Link1](#) and [Link 2](#)



# Some interesting LSTM applications

— — —

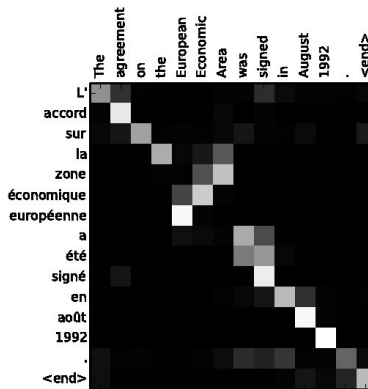
- Handwriting generation demo (Uses attention mechanism)
- Generated music using LSTMs: [Link](#)
- Stock trends prediction
- Copying Shakespeare



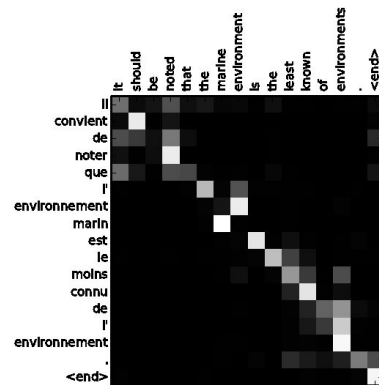
# Attention based networks

— — —

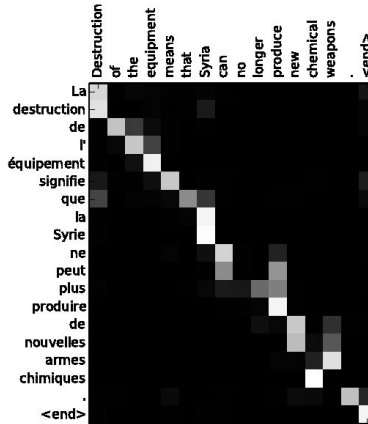
Has proved effective  
for machine  
translation tasks



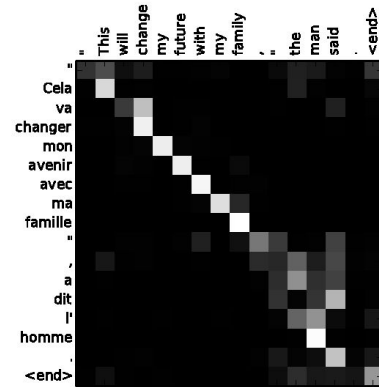
(a)



(b)



(c)



(d)

## Tips for data scientists contd.- ML specific



# Validation

— — —

- Check, double-check, triple-check if the model is learning generic concepts, and not just regurgitating the information contained in the training data set.
- Use advanced cross-validation techniques suitable for your data set and domain
- Keep various test data sets, untouched/unrelated from training process. Ideally mimicking the real-life situation where you want the models to be used.
- Using a pool of models, and selecting ones that have good results on validation set is often a good ploy
- External validation is a must.

# Be algorithm agnostic

— — —

- Don't make a particular algorithm (or a particular class of algorithm) your favourite.
- Best algorithm is the one that solves your problem most effectively.

# Employ good programming practices

— — —

- Many ML practitioners write sloppy code. Because their work is exploratory. (Roughly <5 out of 10/15/100 models they create go into final production).
- Good code means better re-using of code!
- Use of comments
- Good code architecting.
- Generic bits at generic places. Specific at specific places
- Object oriented when possible
- Refactoring when possible
- Module / Package structuring
- All this even if you are 1-person team and only user of your code!

# Enhance knowledge of ML systems

— — —

Specifically

- Databases
- Distributed systems
- Cloud computing

# Keep an eye on stability

— — —

- For your problem, you often don't know how quickly the data will grow!
- Some algorithms / approaches are prohibitively slow if the data becomes big. Make sure you know that – and keep some breathing space.

# My pathway

---

Machine Learning course by Andrew NG

Some bits from Learning from data by Caltech

Statistical Learning Theory by Tibshirani (Stanford)

CS231n

7 MOOCs

Kaggle contests

Most importantly - learn linear algebra/ probability/ stats

# Some good practices for ML/ DL projects

---

- **Make sure it can be replicated**
- Inputs/ outputs for all steps in the pipeline must be clear
- Make the pretrained models available online
- Share the project work on open source whenever possible
- Have full knowledge of the black boxes you use

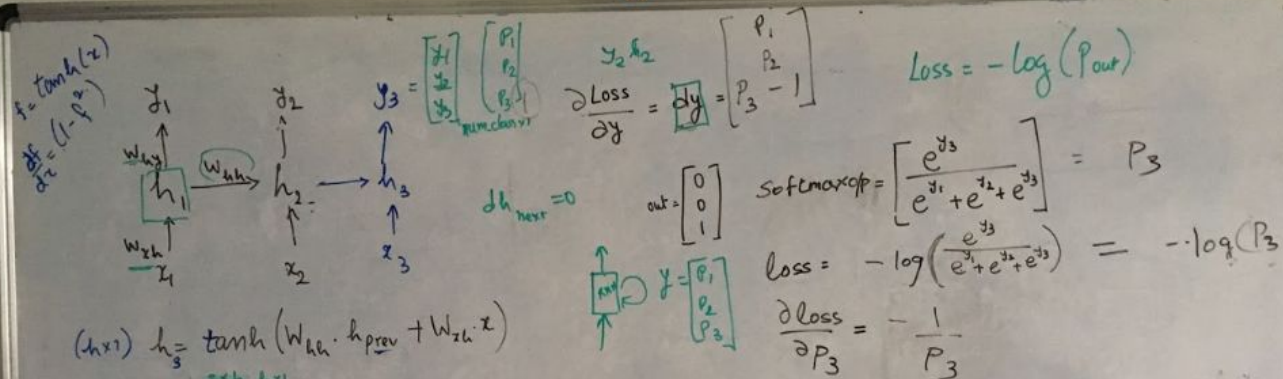
# Keeping up to pace with a growing field

---

- Read papers- required if you want to stay up to date
- Discuss - reading groups help keep up the pace
- Start simple and build up (almost upto 120 papers)
- We just whiteboard it and then dive into the code as inputs outputs are clearer there
- ICLR, ICML, ICCV, CVPR, NeurIPS, PAMI - some good conferences out there



# RNN backprop



$$(h \times 1) \quad h_s = \tanh(W_{hh} \cdot h_{prev} + W_{zh} \cdot x)$$

$$(3 \times 1) \quad y_3 = \overset{2 \times h}{W_{avg}} \cdot \overset{h \times 1}{h} + b_y$$

$$\textcircled{2} \frac{\partial \text{loss}}{\partial w_{xy}} = \frac{\partial \text{loss}}{\partial y} \cdot \frac{\partial y}{\partial w_{xy}} = \overset{(x_i)}{dy} \cdot \overset{(x_h)}{h}$$

$$\textcircled{3} \frac{\partial \text{Loss}}{\partial h} = \left( \frac{\partial \text{Loss}}{\partial y} \cdot \frac{\partial y}{\partial h} \right) + \frac{\partial \text{Loss}_{\text{next}}}{\partial h} \rightarrow \frac{\partial \text{Loss}_{\text{next}}}{\partial h_{\text{next}}} \cdot \frac{\partial h_{\text{next}}}{\partial h} = dy \cdot W_{hy} + dh_{\text{next}} \cdot (1 - h_{\text{next}})^2 \cdot W_{hh}$$

$$\textcircled{A} \frac{\partial \text{Loss}}{\partial w_{ih}} = \frac{\partial \text{Loss}}{\partial h} \cdot \frac{\partial h}{\partial w_{ih}} \\ = dh \cdot (1 - h^2) \cdot h_{prev}$$

$$\textcircled{5} \frac{\partial \text{loss}}{\partial w_{2n}} = \frac{\partial \text{loss}}{\partial n} \cdot \frac{\partial n}{\partial w_{2n}} \\ = \text{dn} \cdot (1 - n^2) \cdot x$$

$$\frac{\partial \text{Loss}_4}{\partial h_3} = 0$$

$$\frac{\partial \text{Loss}_3}{\partial h_2} = \frac{\partial \text{Loss}_{\text{next}}}{\partial h} \cdot \frac{\partial h}{\partial h_2} \cdot W_{hh}$$

$$\text{softmax}_p = \left[ \frac{e^{j_3}}{e^{j_1} + e^{j_2} + e^{j_3}} \right] = p_3$$

$$\text{loss} = -\log\left(\frac{e^{y_3}}{e^{y_1} + e^{y_2} + e^{y_3}}\right) = -\log(P_3)$$

$$\frac{\partial \text{loss}}{\partial P_3} = -\frac{1}{P_3}$$

$$P_3 = \frac{e^{y_3}}{e^{y_1} + e^{y_2} + e^{y_3}}$$

$$\frac{\partial P_3}{\partial y_i} = \frac{(e^{y_1} + e^{y_2} + e^{y_3}) e^{y_3} - e^{y_3} e^{y_i}}{(e^{y_1} + e^{y_2} + e^{y_3})^2}$$

$$\frac{e^{y_1} + e^{y_2}}{e^{y_1} + e^{y_2} + e^{y_3}} \times \frac{e^{y_3}}{e^{y_1} + e^{y_2} + e^{y_3}}$$

$$= P_3(1 - P_3)$$

$$\frac{\partial \text{loss}}{\partial y_i^3} = p_3 - 1$$

$$\frac{\partial \text{Loss}}{\partial y_1} = \frac{\partial \text{Loss}}{\partial p_3} \cdot \frac{\partial p_3}{\partial y_1} = -\frac{1}{p_3} \cdot (-p_3 \cdot p_1) = p_1$$

$$\frac{\partial \log}{\partial y_2} = p_1$$

# Attention models

input:  $T_x$  length sentence:  $x = (x_1, x_2, \dots, x_{T_x})$   $x_i \in K_x$   
 output:  $T_y$  length sentence:  $y = (y_1, y_2, \dots, y_{T_y})$   $y_i \in K_y \rightarrow$  vocab of o/p language

$$\vec{h}_i = \begin{cases} \vec{z}_i \circ \vec{h}_1 + (1 - \vec{z}_i) \circ \vec{h}_{i-1} & i > 0 \\ 0 & i = 0 \end{cases}$$

$i/p$  is  $K_x$  dimensional  
 intermediate state is  $n$  dimensional.

$$\vec{h}_i = \tanh(\vec{w}(\vec{E}x_i) + \vec{U}[\vec{x}_i \circ \vec{h}_{i-1}])$$

$(n \times 1)$   $(n \times 1)$   $(n \times 1)$   $(n \times 1)$   $(n \times 1)$   $(n \times 1)$

$$\vec{z}_i = \sigma(\vec{w}_z(\vec{E}x_i) + \vec{U}_z \vec{h}_{i-1})$$

$(n \times 1)$   $(n \times 1)$   $(n \times 1)$   $(n \times 1)$   $(n \times 1)$

$$\vec{r}_i = \sigma(\vec{w}_r(\vec{E}x_i) + \vec{U}_r \vec{h}_{i-1})$$

$(n \times 1)$   $(n \times 1)$   $(n \times 1)$   $(n \times 1)$   $(n \times 1)$

$\vec{h}_i$  computed same as above

$$s_i = (1 - z_i) \circ s_{i-1} + z_i \tilde{s}_i$$

$$\tilde{s}_i = \tanh(\vec{w}(\vec{E}y_{i-1}) + \vec{U}(\vec{r}_i \circ s_{i-1}) + \vec{C}_i)$$

$$z_i = \sigma(\vec{w}_z(\vec{E}y_{i-1}) + \vec{U}_z s_{i-1} + \vec{C}_z \vec{C}_i)$$

$$\vec{r}_i = \sigma(\vec{w}_r(\vec{E}y_{i-1}) + \vec{U}_r s_{i-1} + \vec{C}_r \vec{C}_i)$$

$(n \times 1)$   $(n \times 1)$   $(n \times 1)$   $(n \times 1)$   $(n \times 1)$   $(n \times 1)$   $(n \times 1)$   $(n \times 1)$

We can get  $T_y$  values of  $s_i, z_i, \vec{r}_i$

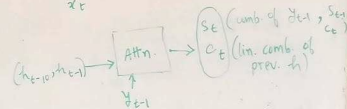
We already have  $y_i$

$$\tilde{t}_i = \vec{V}_0 s_{i-1} + \vec{V}_0(\vec{E}y_{i-1}) + \vec{C}_0 \vec{C}_i$$

$(2 \times 1)$   $(2 \times 1)$   $(2 \times 1)$   $(2 \times 1)$   $(2 \times 1)$   $(2 \times 1)$

$$p(y_i | s_i, y_{i-1}, \vec{C}_i) \propto \exp(\sum_j \vec{W}_0 t_i)$$

$(1 \times 1)$   $(1 \times 1)$   $(1 \times 1)$   $(1 \times 1)$   $(1 \times 1)$   $(1 \times 1)$



$\vec{h}_i = \begin{bmatrix} \vec{h}_i \\ \vec{h}_i \end{bmatrix}$   
 Flow B/W states for  $T_x$  time steps,  $h_0 = 0$

We compute  $C_i$  as

$$C_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j$$

$(n \times 1)$   $(n \times 1)$

$S_0 = \tanh(\vec{W}_S \vec{h}_1)$

$C_i \rightarrow f^n$  of  $s_{i-1}, h_j$   
 $f^n$  of  $y_{i-2}, s_{i-2}$

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})}$$

$$e_{ij} = \vec{V}_0^T \tanh(\vec{W}_0 s_{i-1} + \vec{U}_0 h_j)$$

$(1 \times 1)$   $(1 \times 1)$   $(1 \times 1)$   $(1 \times 1)$   $(1 \times 1)$   $(1 \times 1)$   $(1 \times 1)$   $(1 \times 1)$

$$t_i = \begin{bmatrix} \max_j(\tilde{t}_{i,j-1}, \tilde{t}_{i,j}) \end{bmatrix}_{j=1,2,\dots,l}^T$$

$(1 \times 1)$

# Conclusion

— — —

- ML \ DM landscape changes rapidly, with (smarter, yay!) people getting more and more familiar\comfortable with the techniques.
- Avoid pitfalls of overfitting.
- To use ML in industrial settings: just make sure you keep in mind the advisory nuggets (to certain extent). Special mention to
  - Know your data
  - Know your settings (business, customers, products etc.)
  - Know your maths
  - Know your tools
- The problem is not always well-defined. You have to define\design\transform a messy real-life problem into an ML problem.
- ML is possible in **any** field! – given that we have adequate data.