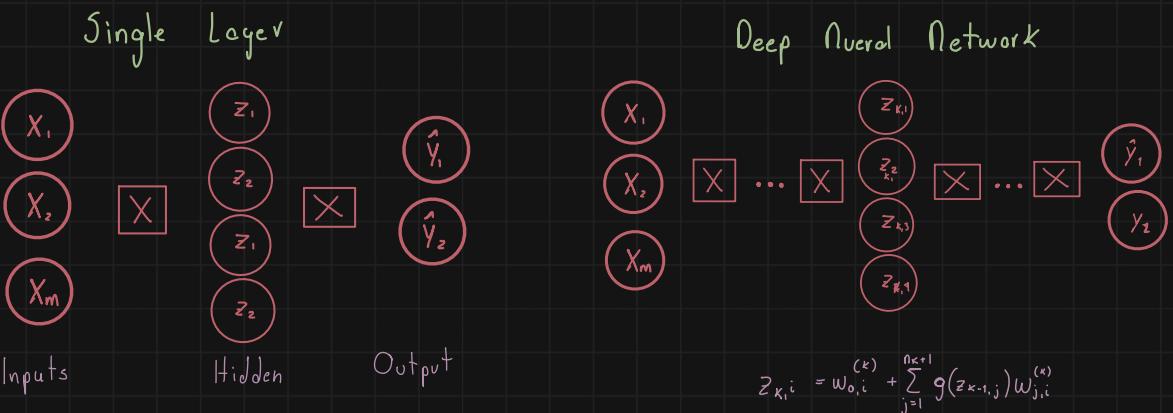
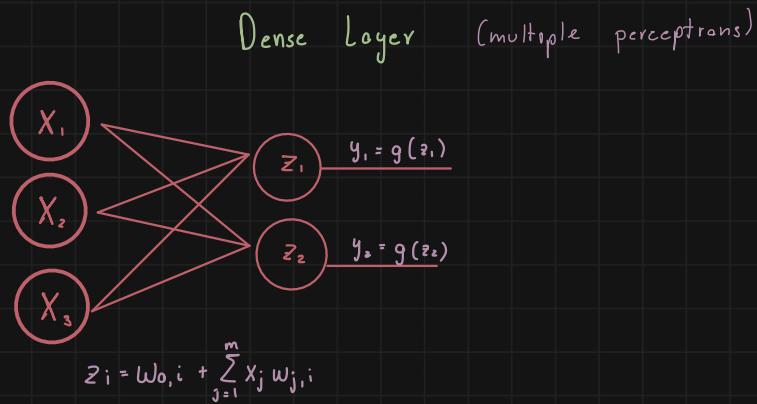
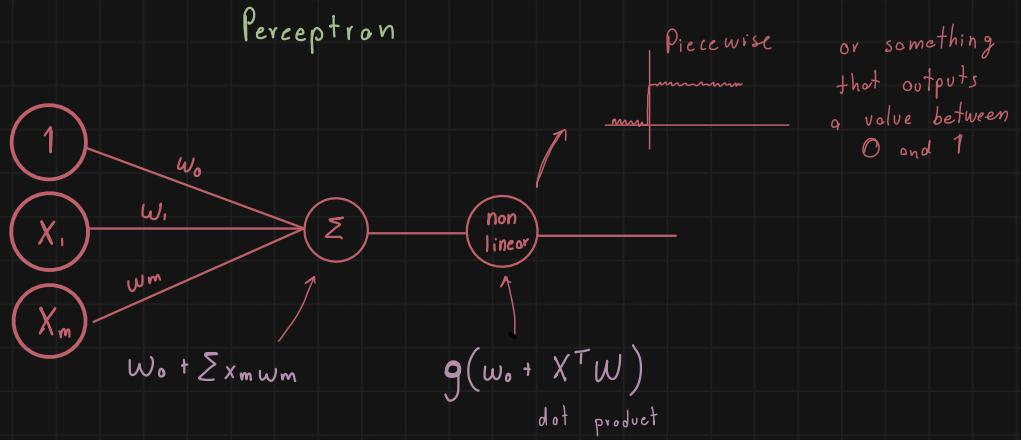


lecture One



Loss:

Cost of incorrect predictions

$$\mathcal{L}(f(x^{(i)}; \mathbf{w}), y^{(i)})$$

[if the difference between predicted and real is very high, the loss will be very high]

Empirical Loss

Total Loss of entire data set

$$X = \begin{bmatrix} x_1 \\ x_2 \\ \dots \end{bmatrix} \quad \begin{bmatrix} z_1 \\ z_2 \\ \dots \end{bmatrix} \quad \begin{bmatrix} f(x) \\ y \end{bmatrix}$$

$$\mathcal{J}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(\text{predicted}, \text{actual})$$

Minimize loss

(just the average)

Mean Squared Error Loss
For continuous numbers as outputs

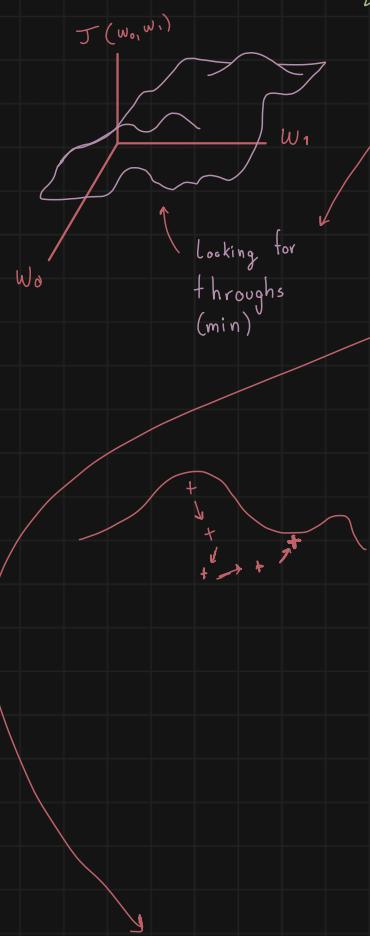
Binary Cross Entropy Loss
For binary (0,1) outputs

Loss Optimization

Find the weights for the lowest loss

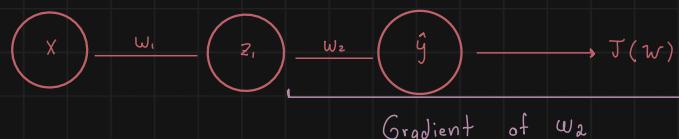
$$W^* = \underset{w}{\operatorname{argmin}} J(w)$$

$$W = \{W_0, W_1, \dots\}$$



- ① Take random point and take gradient $\frac{\partial J(w)}{\partial w}$ AKA which way is up/down
 - ② Take small step towards opposite of gradient AKA to the min
 - ③ Repeat until local (or global, we don't know) minimum AKA update weights by small amount
 - ④ Return weights
- Gradient Descent**

How to compute gradient: Back propagation



$$\frac{\partial J}{\partial w_2} = \frac{\partial J}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial w_2} \quad \left(\begin{array}{l} \text{small change in} \\ w_2 \text{ affects } J(w) \end{array} \right)$$

Gradient of w_1

$$\frac{\partial J}{\partial w_1} = \frac{\partial J}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_1} \frac{\partial z_1}{\partial w_1}$$

(w_1 affects w_2 which in turn affects \hat{y})

Back Propagation

(Repeat for every weight in the network)
[Apply chain rule over and over]

How to actually train NNs to get the precise weights (precise prediction)

Training in Practice: Optimization

Training is extremely hard

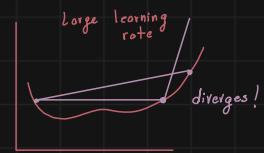
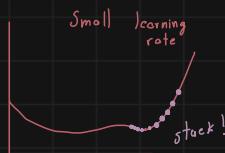


Optimization at gradient descent

$$w \leftarrow w - \eta \frac{\partial J}{\partial w}$$

learning rate

(How large every step should be)



We want stable learning rates

(converge smoothly but avoid local min)

We can

- Try a lot of learning rates and see which fits.

- Use adaptive learning rate

(Optimizers, for example: SGD, Adam, Adadelta, etc)

Training in Practice: Mini Batching

Gradient descent

A derivative is computationally intensive

$$\frac{d J(w)}{d w}$$

All weights

(slow, precise)

Stochastic Gradient descent

We can compute a single point, to be faster

$$\frac{d J(w)}{w}$$

Single point

(fast, very noisy)

There is a middle ground

Batching

We can compute a number B of points

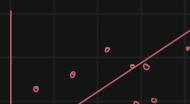
$$\frac{1}{B} \sum_{k=1}^B \frac{\partial J_k(w)}{\partial w}$$

Multi point

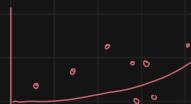
(fast, kinda precise)

Mini-batches allow parallelization !

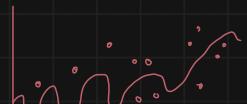
Training in Practice: Overfitting



underfitting



Ideal fit



overfitting

Regularization

Technique to discourage overfitting

- Dropout

Randomly set activations to 0 (lol)

- Forces to rely on 1 node



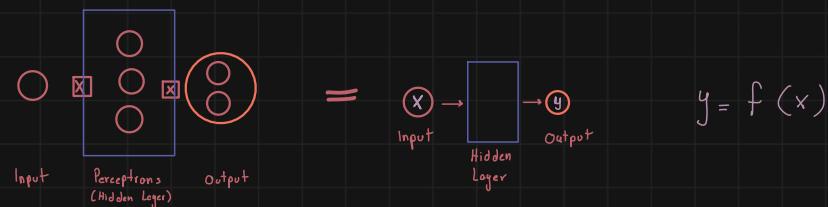
Lecture 2

Sequential Modeling
history

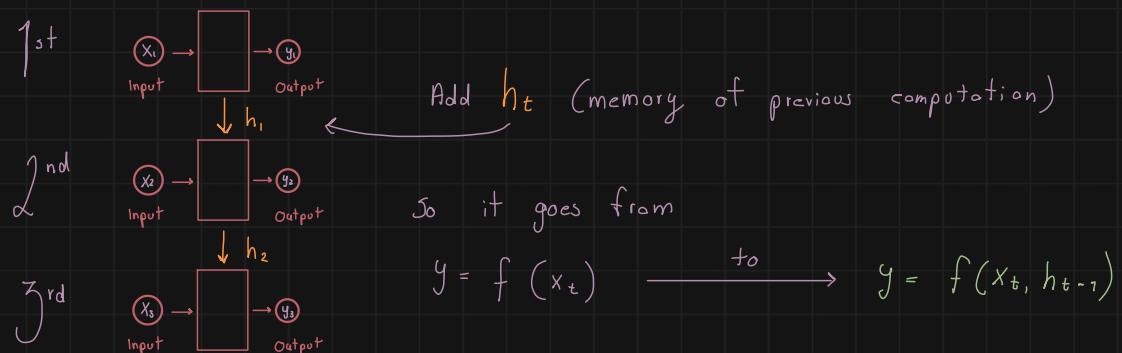
Sequence Modelling Applications

1 to 1	Many to 1	1 to Many	Many to Many
binary (will I pass?)	sentiment (text sentiment?)	image caption (text of this image)	Machine Translation (Übersetzer)

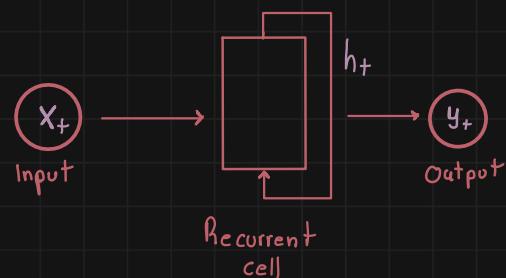
Neurons with Recurrence



How to make this time dependent?



Now the output depends on the memory (previous calculation).



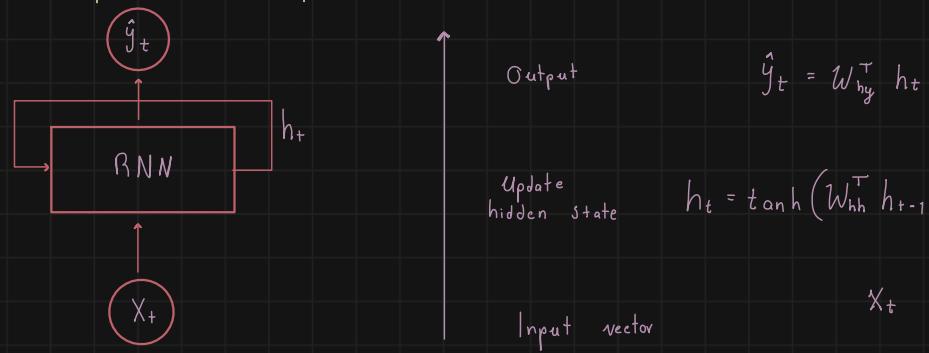
h_t : State
is updated at each step

$$\underline{h_t} = \underline{f}_{W^h} (\underline{x_t}, \underline{h_{t-1}})$$

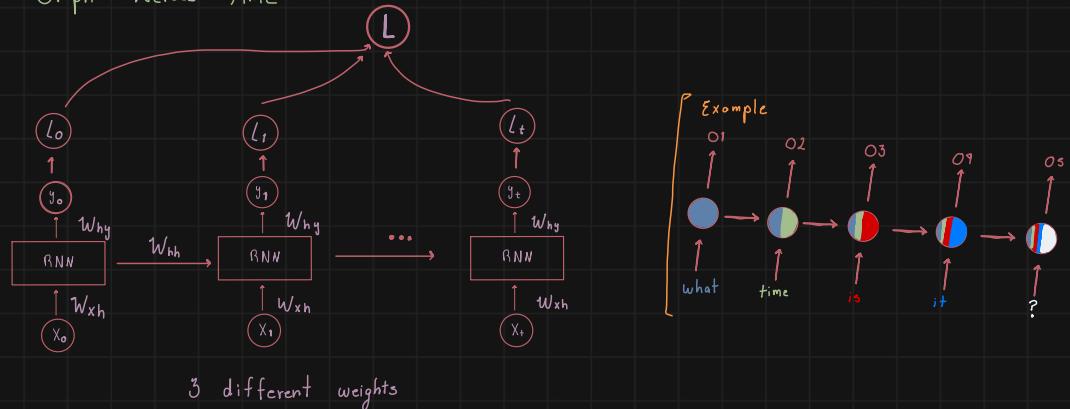
cell state function weights input old state

Some function \circ

State update and Output



RNN: Graph Across Time



Sequence Modelling: Design Criteria

- Variable length
- share states (memory)
- Maintain order
- Track long-term dependencies

Models work on numerical inputs, of course (math functions)

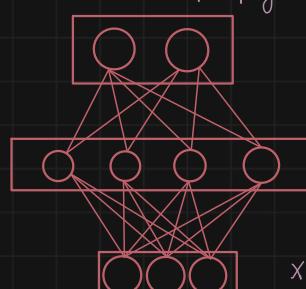
Embedding: Transform indexes into a vector

Example

1. Vocabulary	2. Indexing	3. Embedding
"o" "dog" "took" "my" "I"	Value to every unique word $\begin{array}{ccc} o & \longrightarrow & 1 \\ dog & \longrightarrow & 2 \\ \dots & & \end{array}$	<u>One hot embedding</u> "dog" = $[0, 1, 0, \dots]$ (Binary index) <u>Learn embedding</u> Feed index map to model. Similar words will have similar embeddings

Back Propagation Through Time

Normal backpropagation:



How to train

First

Forward pass
makes a prediction

Then

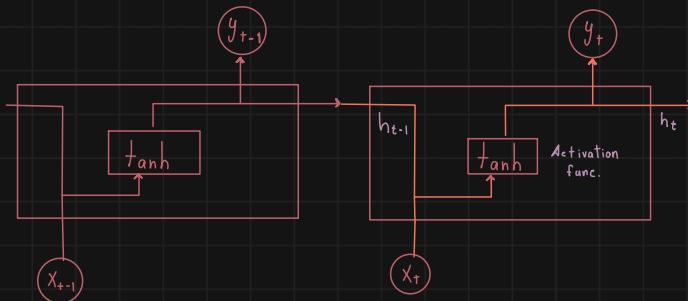
Computes Loss (pred, truth)

To train, we backpropagate:

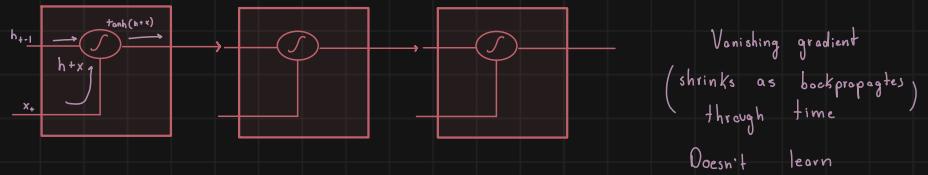
1. Derivate loss with respect to each parameter
2. Shift parameters to minimize loss

Gated Cell (LSTM, GRU, etc)

Controls which information goes through.



Vanilla Recurrent NN (standard)

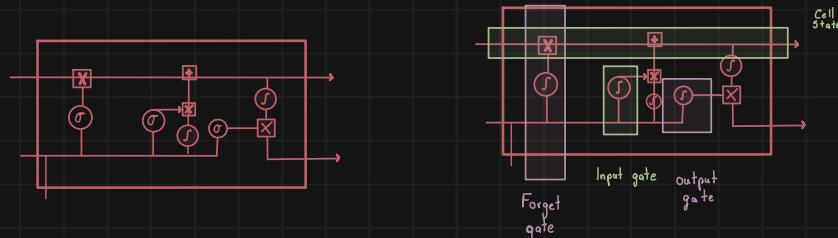


LSTM (Long Short Term Memory)

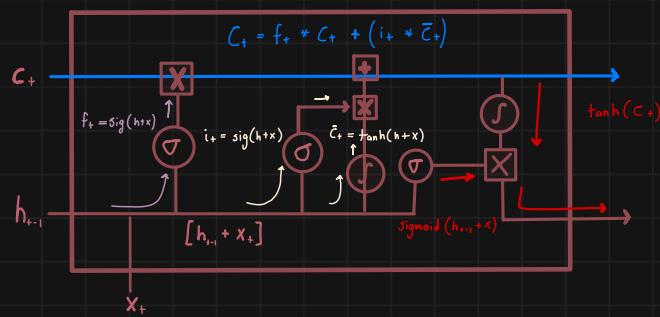
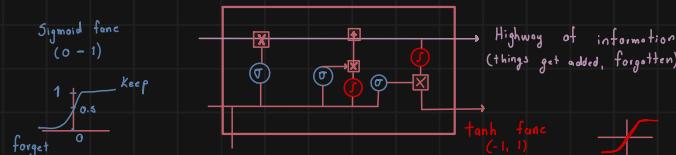
Mitigate short term memory (vanishing gradient)

Rely on gated cell to track info through steps

Gate: Regulate flow of information



Description of Components



What to keep ($\text{sigmoid}(h+x) \approx 1$) or to forget ($\text{sigmoid}(h+x) \approx 0$)

Input gate

Updates the state. The sigmoid determines which value is important (1: important, 0: Not important) and then gets multiplied by tanh to regulate. Sig decides which info to keep from tanh.

New Cell State

Previous state multiplied by forget vector, added to the input gate selection (update values)

Output gate

Previous hidden state / input into sigmoid and then multiply with tanh of new cell state to produce a new hidden state.

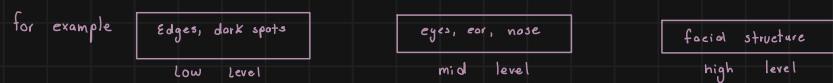
Lecture 3

Convolutional Neural Networks

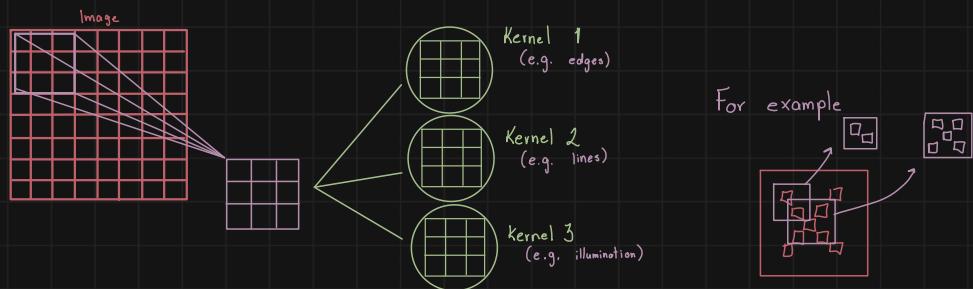
Computer vision (context)

Learning Feature Representations

Learn from a hierarchy of features directly from data.



Learning Visual Features



Convolution Operation

$$\text{ex. 1} \quad \begin{array}{|c|c|} \hline 3 & 2 \\ \hline 2 & 3 \\ \hline \end{array} \otimes \begin{array}{|c|c|} \hline 1 & -1 \\ \hline -1 & 1 \\ \hline \end{array} = \begin{array}{|c|c|} \hline 3 & -2 \\ \hline -2 & 3 \\ \hline \end{array} \xrightarrow{\text{sum}} [2]$$

$$\text{ex. 2.} \quad \begin{array}{|c|c|c|} \hline 1 & 2 & 3 \\ \hline 1 & 2 & 3 \\ \hline 1 & 2 & 3 \\ \hline \end{array} \otimes \begin{array}{|c|c|} \hline 2 & 2 \\ \hline 1 & 1 \\ \hline \end{array} = \begin{array}{|c|c|} \hline 1 & 9 & 9 & 6 \\ \hline 1 & 2 & 2 & 3 \\ \hline 1 & 9 & 9 & 6 \\ \hline 1 & 2 & 2 & 3 \\ \hline \end{array} \xrightarrow{\text{sum}} \begin{bmatrix} 8 & 15 \\ 8 & 15 \end{bmatrix}$$

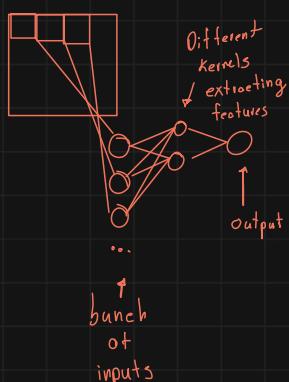
So, as coded on github, there are several filters

like

$\begin{bmatrix} 0 & 1 & 0 \\ -1 & 5 & -1 \\ 0 & -10 \end{bmatrix}$	$\begin{bmatrix} 0 & 1 & 0 \\ 1 & 9 & 1 \\ 0 & 1 & 0 \end{bmatrix}$
sharpen	edge detect

but CNN get these values themselves !

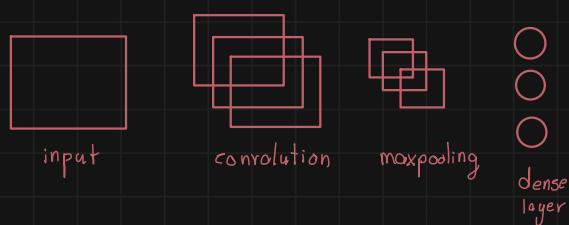
Simply put



So :

- Apply weights (filter / kernel) to extract local features
- Use multiple filters to get multiple features
- Spatially share parameters of each filter

CNN for classification



1. Convolution
2. Non linearity
3. Pooling (down-sampling)



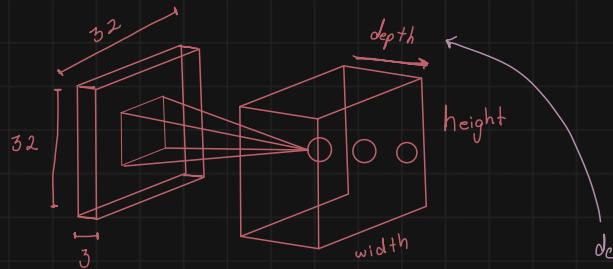
Each neuron sees a patch of its inputs

$$X \otimes \text{Kernel}$$

9×9
 W_{ij}

$$\sum_{i=1}^q \sum_{j=1}^q W_{ij} X_{i+p, j+q} + b$$

$$\text{non linear} \left(\sum_{i=1}^q \sum_{j=1}^q W_{ij} X_{i+p, j+q} + b \right)$$



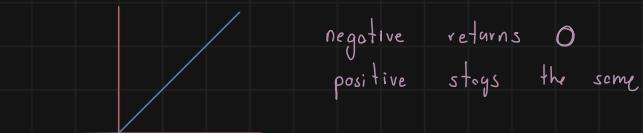
Apply Kernel

Compute linear combination

Finally, activate with
non-linear function

depth is the number
of filters

Usually non-linear operation is ReLU



Max pooling reduces dimensionality

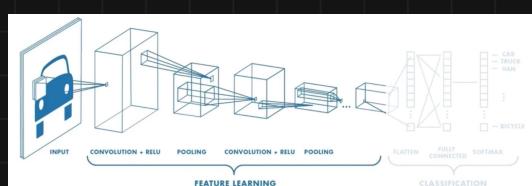
Divide into 2×2 patches

1	1	3	9
2	8	3	10
5	5	20	1
1	19	1	1

Take the
max value

8	10
19	20

- Mean pooling takes the average
- Min takes the minimum ...



```
import tensorflow as tf

def generate_model():
    model = tf.keras.Sequential([
        # first convolutional layer
        tf.keras.layers.Conv2D(32, filter_size=3, activation='relu'),
        tf.keras.layers.MaxPool2D(pool_size=2, strides=2),

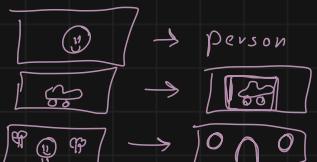
        # second convolutional layer
        tf.keras.layers.Conv2D(64, filter_size=3, activation='relu'),
        tf.keras.layers.MaxPool2D(pool_size=2, strides=2),

        # fully connected classifier
        tf.keras.layers.Flatten(),
        tf.keras.layers.Dense(1024, activation='relu'),
        tf.keras.layers.Dense(10, activation='softmax') # 10 outputs
    ])
    return model
```

Cool Applications

Classification

What is this image?



Object detection

Where are the objects?



Semantic segmentation

What class this belongs?



Lecture 9

Deep Generative Models

Supervised Learning

Data (x, y) Label

Function to map $x \rightarrow y$

VS

Unsupervised Learning

Data (x) No labels?

Learn the structure of the data

Generative Modelling

Input training for some distribution and learn a model to represent that distribution.

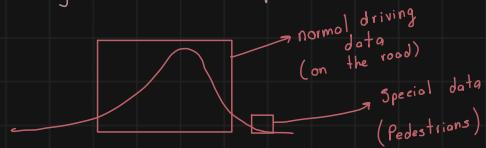
Use cases:

Debiasing

Uncovering underlying features to make the data more distributed

Outlier Detection

Detect very rare and unique cases in training data

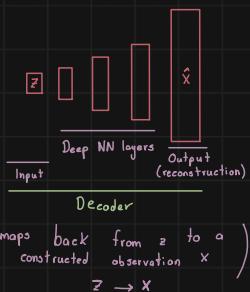
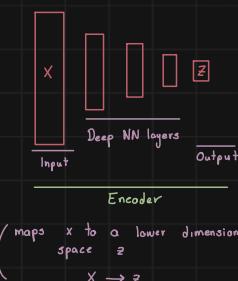
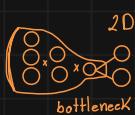


Latent Variable Models
true explanatory factors



Autoencoders

Learn a lower dimensional feature of raw data



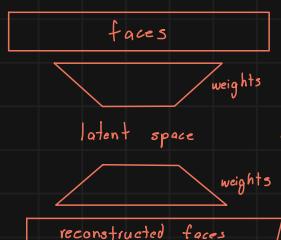
To train, we compare

x with \hat{x}

Loss doesn't use any labels?

$$\text{e.g. } L(x, \hat{x}) = |x - \hat{x}|^2 \quad (\text{mean squared error})$$

You could train a model



when trained just input latent vector and get a entirely new face

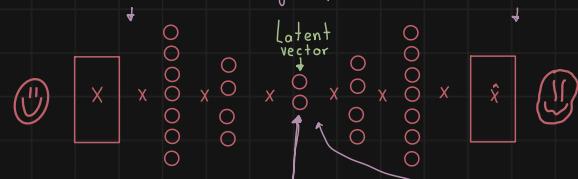
But what does each dimension on the vector mean??:
 ↗ inclination?
 ↗ skin tone
 ↗ eye color

We don't know!

Pool too wide



What comes through input is reconstructed in output



Forces the network to learn a representation of data very small

like really, really good compression

i.e. all info to reconstruct is in 2 neurons

