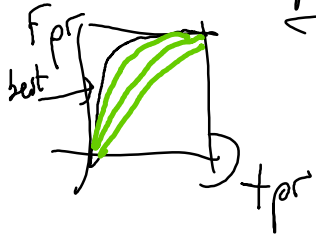


Thursday, May 27, 2021 7:19 AM

## Neyman-Pearson Lemma

"Optimal" binary classifier is the likelihood ratio

best possible  
ROC curve



$$R(x) = \frac{P(x|S)}{P(x|B)}$$

$R(x) > R_c$  defines  $tpr$  &  $fpr$

In lieu of a proof, let's check that  $R(x)$  comes out of minimizing the binary cross entropy.

$$L = - \left( \sum_{i \in S} \log f(x_i) + \sum_{i \in B} \log (1 - f(x_i)) \right)$$

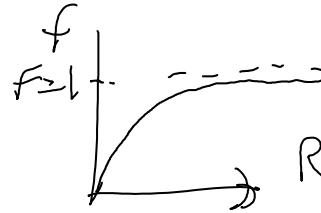
$f \rightarrow f + \delta f$

$$= - \int dx \left( P(x|S) \log f(x) + P(x|B) \log (1 - f(x)) \right)$$

$$0 = \frac{\partial L}{\partial f} = \frac{P(x|S)}{f(x)} - \frac{P(x|B)}{1 - f(x)}$$

Thursday, May 27, 2021 7:50 AM

$$F = \frac{P(x|S)}{P(x|S) + P(x|B)} = \frac{R}{R+1} \quad (\text{monotonic wrt } R)$$



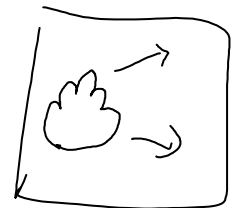
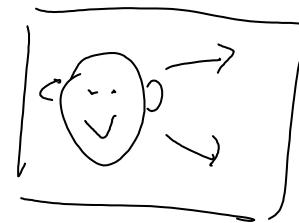
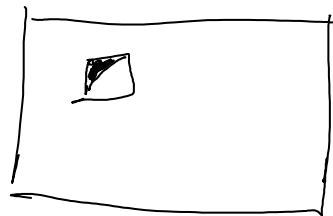
## Convolutional Neural Network (CNN)

DNN is the most general feed forward NN, but has many weights

Can get better performance by reducing # of indep weights using symmetry.

For images, CNN architecture was a big breakthrough (2012)

- leverages <sup>2d</sup> translation symmetry in images
- finds 2d "features" in images



Thursday, May 27, 2021 8:18 AM

To develop "feature detectors" or "feature maps" CNN will drag  
 $n$  filter over the image  
 ↓  
 "convolution"



e.g. 2x2 filter  $\begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix} \rightarrow$  trainable weights, will learn to detect specific feature in image  
 (eg noses, faces, edges, ...)

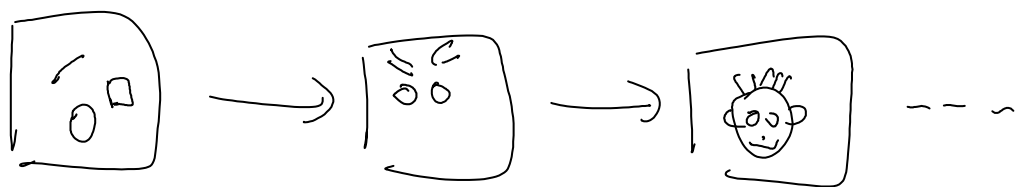
multiply pixels elementwise

$$\begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \rightarrow \square = A(w_{11}a_{11} + w_{12}a_{12} + w_{21}a_{21} + w_{22}a_{22} + b) = a'_{11}$$

$$a'_{ij} = A \left( \sum_{r,s=1}^2 a_{i+r, j+s} w_{rs} + b \right)$$

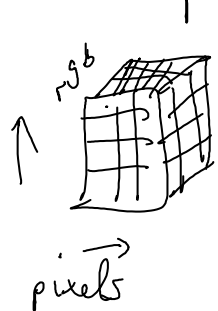
Thursday, May 27, 2021 8:24 AM

Can stack 2d convolutional layers like we do in DNN



At every successive layer, CNN is learning higher level features in images  
What about color?

Can operate 2d convolution on multiple dimensions "channels"



(r, g, b)

multiply a filter map w/ 3 channels

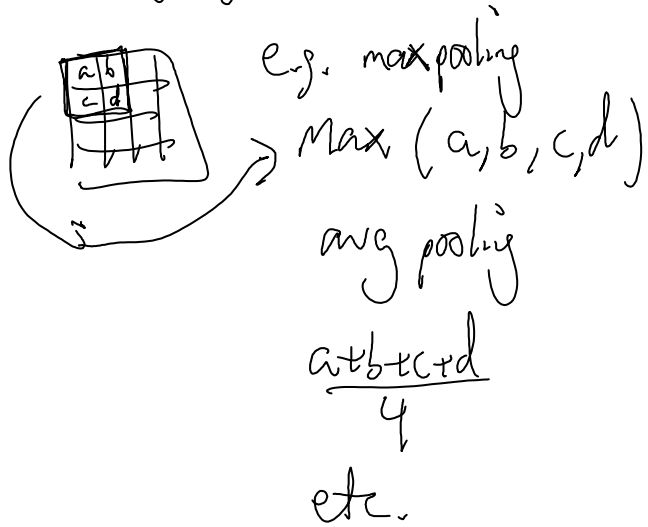
sum over channels, output image w/ single channel.



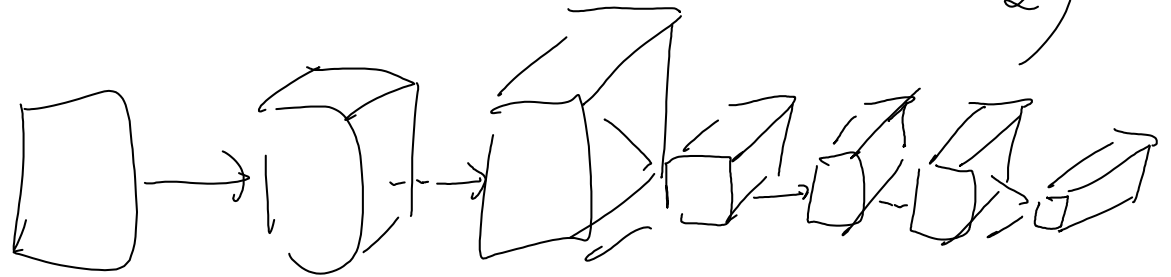
At each conv. layer, typically learn  $\approx N \sim 100$  filters  
 each conv layer then produces image w/  $N$  channels.

Thursday, May 27, 2021 8:33 AM

- Typically coarse grain image after some # of conv. layers  
"pooling layer"



→ reduce image size by some factor (e.g. factor of 2)



- Finally flatten reduced image to column vector and feed to a DNN  
MLF's that CNN has learned      classifier for MLF's.
- This is structure of the most basic CNN ("LeNet")

Thursday, May 27, 2021 9:06 AM

## Generative Modeling & LLM applications

- Given some data  $\{x_i\}$ , can I learn  $P(x)$  that data was drawn from and then sample from it to produce new examples of the data-
  - 3 main approaches
    - Generative Adversarial Networks (Goodfellow et al 2014) \*
    - Variational Autoencoders
    - Normalizing Flows. } also relevant for next topic, anomaly detection
- these give the best performance to date on image generation tasks

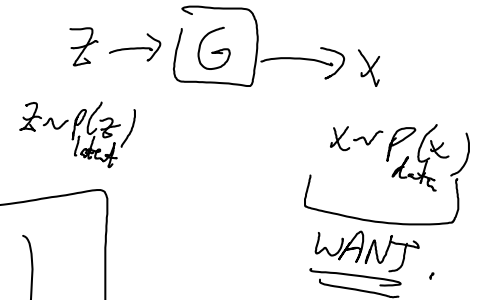
Thursday, May 27, 2021 9:16 AM

Idea: Adversarial training

- GAN has 2 NN's that are trying to defeat each other ("playing a game")
  1. G generator: generates fake data
  2. D discriminator: classifier real vs. fake.

- GAN loss: BCE btw generated & real data

$$L = \sum_{i \in \text{data}} \log D(x_i) + \sum_{z \in \text{gen}} \log (1 - D(G(z_i)))$$



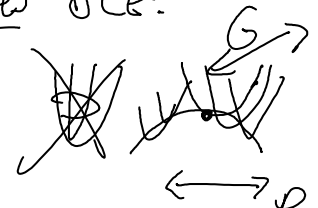
-  $z \in$  latent space w/ simple probability dist'n e.g.  $P(z) =$  uniform or gaussian  
 $\mathbb{R}^d$   $d \ll$  dim of data usually

-  $L$  is negative of usual classifier loss, good clf. maximizes BCE.

Objective:

$$\min_G \max_D L$$

- $\max_D$ : discriminator is as good as possible
- $\min_G$ : generator fools discriminator as much as possible.



Thursday, May 27, 2021 9:24 AM

- Saddle pt optimizer - tricky, unstable, how to guarantee convergence?



{ to train, typically alternate btw  
training G & D.  
- train G for some epochs training D  
- vice versa

↓  
you don't  
GANs don't often converge  
Select "best" epoch often based  
on subjective criteria.



Thursday, May 27, 2021 9:28 AM

# Optimality guarantee of GANs

$$\min_G \max_D L \rightarrow \max_D L | G \rightarrow \hat{D}_G(x) = \text{optimal classifier btw real images \& images generated by } G.$$

$$\downarrow \text{plug back into } L$$

$$= \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_G(x)}$$

"NP Lemma"  
"Likelihood trick"  
"Likelihood free learning"

$$L |_{D=\hat{D}_G} = \sum_{x \in \text{data}} \log \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_G(x)} + \sum_{x \in \text{gen}} \log \frac{p_G(x)}{p_{\text{data}}(x) + p_G(x)}$$

$$= \int dx p_{\text{data}}(x) \log p_{\text{data}}(x) + \int p_G(x) \log p_G(x) - \int (p_{\text{data}}(x) + p_G(x)) \log (p_{\text{data}}(x) + p_G(x))$$

$$= 2 \text{ JSD} (p_{\text{data}}, p_G) + \text{const.}$$

"Jensen-Shannon Distance" Measure of similarity btw 2 prob. distns.

$$\left. \begin{aligned}
 & \text{JSD} = 0 \text{ iff } p_G = p_{\text{data}} \\
 & 1 \geq \text{JSD} \geq 0 \text{ for all } p_G, p_{\text{data}} \\
 & \text{JSD} = 1 \text{ iff } p_G \text{ \& } p_{\text{data}} \text{ are disjoint}
 \end{aligned} \right\} \rightarrow \begin{array}{c} \checkmark \quad \checkmark \\ p_G \quad p_{\text{data}} \\ \text{[Two disjoint bell curves on a horizontal axis]} \end{array}$$

$$\boxed{\min_G L \mid \hat{p}_G \Rightarrow \text{JSD} = 0 \text{ i.e. } p_G = p_{\text{data}}}$$