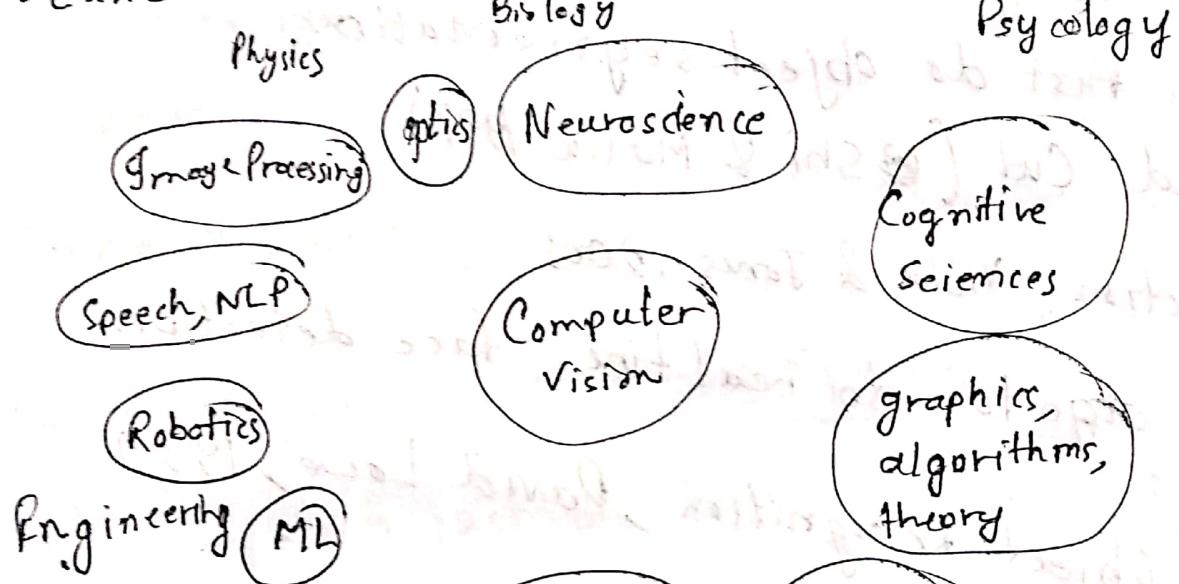


CS231n: Convolutional Nets for Visual Recognition

Fei-Fei Li; Justin Johnson; Serena Yueng

Computer vision is the study of visual data

• called the dark matter of internet



Hubel & Wiesel, 1959

Electrical signal from brain \rightarrow stimulus

\rightarrow mapping \leftarrow

Block world [Larry Roberts, 1963]

1966 → Summer Vision Project

Generalized Cylinder

1979

Pictorial Structure

1973

Nonna.

If object recognition is too hard, maybe we should first do object segmentation

Normalized Cut (shi & Malik 1997)

Face Detection, Viola & Jones, 2001

Adaboost algo to do real time face detection

'SIFT' & Object Recognition, David Lowe, 1990

HOG

Pascal Visual Object Challenge (10 categories)

Imagenet (Deng, Dong, Socher, Li, & Fei-Fei, 2009)
Stanford & Princeton

(22K categories & 14M images)

Took 3 years (took help from Amazon Mechanical Turk)

② Imagenet Classification Challenge

There is a number of visual recognition problems that are related to image classification, such as object detection, image captioning

Lecun et al (conv-net)

Alexnet

Image \rightarrow Captions

Lecture 2: Image Classification Pipeline

- An image is just a big grid of numbers between [0, 255]

Attempts: finding edges (not scalable)

Data Driven approach

- 1) Collect data
- 2) Use ML to train
- 3) Evaluate the classifier

Distance Metric

$$L_1 \text{ Distance} \Rightarrow d_1(I_1, I_2) = \sum_p |I_1^p - I_2^p|$$

Nearest Neighbor $O(1)$

Memoize data

$$L_2 \text{ dist} \Rightarrow d_2(I_1, I_2) = \sqrt{\sum_p (I_1^p - I_2^p)^2}$$

for each test image

find closest train image $O(N)$

Predict label of nearest image

This is bad: we want classifiers that are fast at prediction; slow for training is ok

k-nearest neighbor (to get rid of noisy answer)

Hyperparameters: choices about the algo that we set rather than learn

Setting hyper params

K-NN on images never used

- very slow at testing

- distance metric not informative (not good to check similarity)

- curse of dimensionality

Linear Classification

Parametric Approach

$f(x, w)$ \rightarrow 10 numbers giving class scores
 ↑
 weights

$$f(x, w) = \underbrace{w^T x}_{10 \times 1} + b \quad \begin{matrix} \uparrow & \uparrow \\ 10 \times 3072 & 3072 \times 1 \end{matrix}$$

\rightarrow Linear classifier is learning only one template

Lecture 3 - Loss Function & Optimization

Challenges of recognition

i) Viewpoint

ii) Illumination

iii) Deformation

iv) Occlusion

v) Clutter

vi) Intraclass variation

- Define a loss function that quantifies our unhappiness with the scores across the training set.
- Come up with a way of efficiently finding parameters that minimize the loss function

$$L = \frac{1}{N} \sum_i L_i(f(x_i, w), y_i)$$

Multiclass SVM loss

$$L_i = \sum_{j \neq y_i} \begin{cases} 0 & \text{if } s_j > s_{y_i} + 1 \\ s_j - s_{y_i} + 1 & \text{otherwise} \end{cases}$$

true value

$$= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1) \quad s = f(x_i, w)$$

$$\text{Then, } L = \frac{1}{N} \sum_{i=1}^N L_i$$

What happens to loss if a score is changed a bit?

→ Nothing happens because it still returns zero loss.

What is the min/max possible loss for SVM

→ min → 0
max → ∞

Q3] At initialization ω is small so all $s_i > 0$.

What is the loss?

(number of classes - L)

Q4] What if sum as ~~is~~ over all classes?

The loss increases by L

This is just for convention we omit the correct class so that our minimum loss is zero.

Q5] What if we used mean instead of sum?

→ answer would be same because we don't care about true scores

Q5] What if $\sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)^2$

→ This would end up giving a different loss function that's not linear

Q6] If we find a ω that gives $L(\omega)$ is this ω unique?
No, there are many other ω s. Like 2ω

Regularization: Model should be simple so that it works on test data

$$L(\omega) = \frac{1}{N} \sum_{i=1}^N l_i(f(x_i, \omega), y_i) + \lambda R(\omega)$$

Pecan's Razor:

Among competing hypotheses, the simplest is the best

L2 regularization $R(\omega) = \sum_k \sum_e \omega_{k,e}^2$

L1 regularization $R(\omega) = \sum_k \sum_e |\omega_{k,e}|$

Elastic net ($L1 + L2$) $R(\omega) = \sum_k \sum_e \beta \omega_{k,e}^2 + \gamma |\omega_{k,e}|$

Max norm regularization

Dropout

Batchnorm

Stochastic depth

If you're Bayesian: L2 regularization also corresponds

MAP inference using a Gaussian prior on ω

Softmax classifier (Multinomial Logistic Regression)

$$P(Y=k | X=x_i) = \frac{e^{s_i}}{\sum_j e^{s_j}} \quad \text{where } s = f(x_i; \omega)$$

$$L_i = -\log P(Y=y_i | X=x_i) \quad L_j = -\log \left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}} \right)$$

Want to maximize log likelihood, or to minimize the negative log likelihood of the ~~the~~ correct class

Q1: What is the min. and max loss?

$$\min = 0 \\ \max = \infty$$

Q2: Usually at initialization w is small so all $s \approx 0$ what is the loss?

$$= -\log\left(\frac{1}{c}\right)$$

Softmax vs SVM

Q: Suppose I take a datapoint and jiggles it (changing its score slightly). What happens to the loss in both cases?

SVM doesn't care about the score. Softmax continuously try to make the datapoints like pushing the correct to ∞ and pushing the incorrect to $-\infty$

$$\text{Softmax}, t_i = \log\left(\frac{e^{s_i}}{\sum_j e^{s_j}}\right)$$

$$\text{SVM}, L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

$$\text{Full Loss } L = \frac{1}{N} \sum_{i=1}^N L_i + R(w)$$

Optimization:

Strategy #1 Random search

Strategy #2 follow the slope

In 1 dimension, the derivative of a function

$$\frac{\delta f(x)}{\delta x} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} \quad [\text{slope}]$$

In multiple dimensions, the gradient is the vector of partial derivatives along each direction

- The slope in any direction is the dot product of the direction with the gradient
- The direction of steepest descent is the negative gradient.
- This is super slow and super bad.

We can use compute analytic gradient

Numerical gradient: approximate, slow, easy to write

Analytic gradients: exact, fast, error-prone

Gradient check: Using analytic gradient to find

grads but checking with numerical gradient to
• This is an interesting debugging tool.

Gradient descent

- i) find grads
- ii) weights $\leftarrow \frac{\text{stepsize}}{\uparrow \text{hyperparameter}} \times \cancel{\text{grads}}$
learning rate

Stochastic Gradient Descent:

Full sum is expensive when it's large

Approx sum using minibatch of examples

32/64/128 common

Image Features

Color Histogram

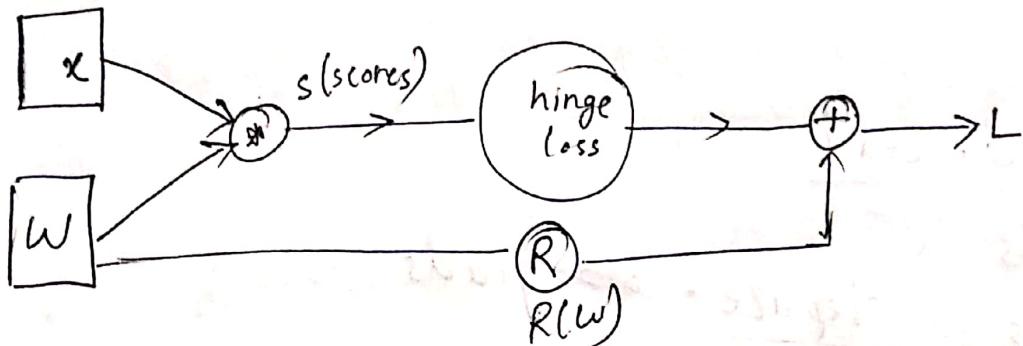
Histogram of Oriented Gradients

Bag of Words (Build Codebook, Encode Images)

Image Features vs ConvNets

ConvNets (Krizhevsky 2012) AlexNet

Lecture 4: Intro to Neural Nets



Leverage chain rule

$$\delta(x) = \frac{1}{1+e^{-x}}$$

$$\frac{d}{dx}(e^x) = e^x$$

$$\frac{d}{dx}(\alpha x) = \alpha$$

$$\frac{d}{dx}\left(\frac{1}{x}\right) = -\frac{1}{x^2}$$

$$f(x) = c+x = \alpha$$

$$\frac{d}{dx}(c+x) = 1$$

$$\begin{aligned} \frac{d \delta(x)}{dx} &= \frac{e^{-x}}{(1+e^{-x})^2} \\ &= \left(\frac{1+e^{-x}-1}{1+e^{-x}} \right) \left(\frac{1}{1+e^{-x}} \right) \\ &= [1-\delta(x)] \delta'(x) \end{aligned}$$

If any problem with gradients,
break down to computational
graph

add gate: gradient distributor

Q: What is a max gate?

the highest one back is gonna get the max value. other will be zeroed out

multip gate: gradient switcher

local gradient is the value of the other variable

Neural Nets:

2-layer Network $f = w_2 \max(0, w_1 x)$

Biological Neurons are far more complicated

Activations

Sigmoid

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

Leaky ReLU

$$\max(0.01x, x)$$

tanh

$$\tanh(x)$$

Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

ReLU

$$\max(0, x)$$

ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$

$$\nabla_w f = 2q \cdot w^T$$

$$\nabla_w f = 2w^T q$$

Lecture 5: Convolution Neural Networks

- Try to maintain spatial structure

History of NN and CNN
→ Mark I perception machine was the first implementation of the perceptron algorithm.
* no backprop

$$f(x) = \begin{cases} 1 & w \cdot x + b \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

Frank Rosenblatt (1957)

→ Widrow and Hoff (1960) : Adaline/Madaline
Multilayer perception \hookrightarrow no backprop

Rumelhart (1986) \rightarrow Backprop

Hinton and Salakhutdinov 2006] \rightarrow Reinigorated research in Deep Learning

First strong results [2012] [Krizhevsky, 2012]
- ImageNet classification with deep CNN
AlexNet

Neurocognition [Fukushima, 1980]

[Lecun, Bottou, Bengio, Haffner, 1998] -
Gradient-based learning applied to
document recognition

ConvNets for classification, retrieval, detection,
segmentation, self driving car, pose estimation, game play
Image Captioning

NVIDIA Tesla line [GPU] powers Neural Style Transfer

Convolutional Neural Networks

Convolutional layers → preserve spatial structure

→ We call the layer convolutional because it is related to convolution of two signals

$$f[x, y] * g[x, y] = \sum_{n_1=-\infty}^{\infty} \sum_{n_2=-\infty}^{\infty} f[n_1, n_2] \cdot g[x-n_1, y-n_2]$$

elementwise multiplication
and sum of a filter and
the signed (image)

Conv Layer

- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires four hyperparams:
 - $K \rightarrow \#$ of filters
 - $F \rightarrow \#$ their spatial extent
 - $S \rightarrow \#$ stride
 - $P \rightarrow$ Amount of zero padding
- * Produces a volume of size $W_2 \times H_2 \times D_2$, where
 - $W_2 = (W_1 - F + 2P)/S + 1$
 - $H_2 = (H_1 - F + 2P)/S + 1$
 - $D_2 = K$
- with parameter sharing. It introduces $F \cdot P \cdot D_1$ weights per filter for a total of $(P \cdot F \cdot D_1) \cdot K$ weights and K biases.
- * The output volume, the d -th depth slice (of size $W_2 \times H_2$) is the result of performing a valid convolution of the d -th filter over the input volume with a stride of S , and then offset by d -th bias.

Pooling layer

- makes the representation smaller and more manageable
- operates over each activation map independently

Pooling layer:

Accepts a volume of size $w_1 \times h_1 \times d_1$

- Requires three hyperparameters

→ their spatial extent f

→ the stride s

- Produces a volume of size $w_2 \times h_2 \times d_2$ where:

$$\rightarrow w_2 = (w_1 - f) / s + 1$$

$$\rightarrow h_2 = (h_1 - f) / s + 1$$

$$\rightarrow d_2 = d_1$$

- Introduces zero parameters since it computes a fixed function of the input
- Note that it is not common to use zero-padding for pooling layers

Fully connected Layer (FC-Layer)

→ Contains neurons that connect to the entire input volume, as in ordinary Neural Networks

Typical architectures

$$[(\text{Conv-ReLU})^*N - \text{POOL?}] * M - (\text{FC-ReLU}) * K, \text{ Softmax}$$

- ResNet/~~GoogleNet~~ challenges this paradigm