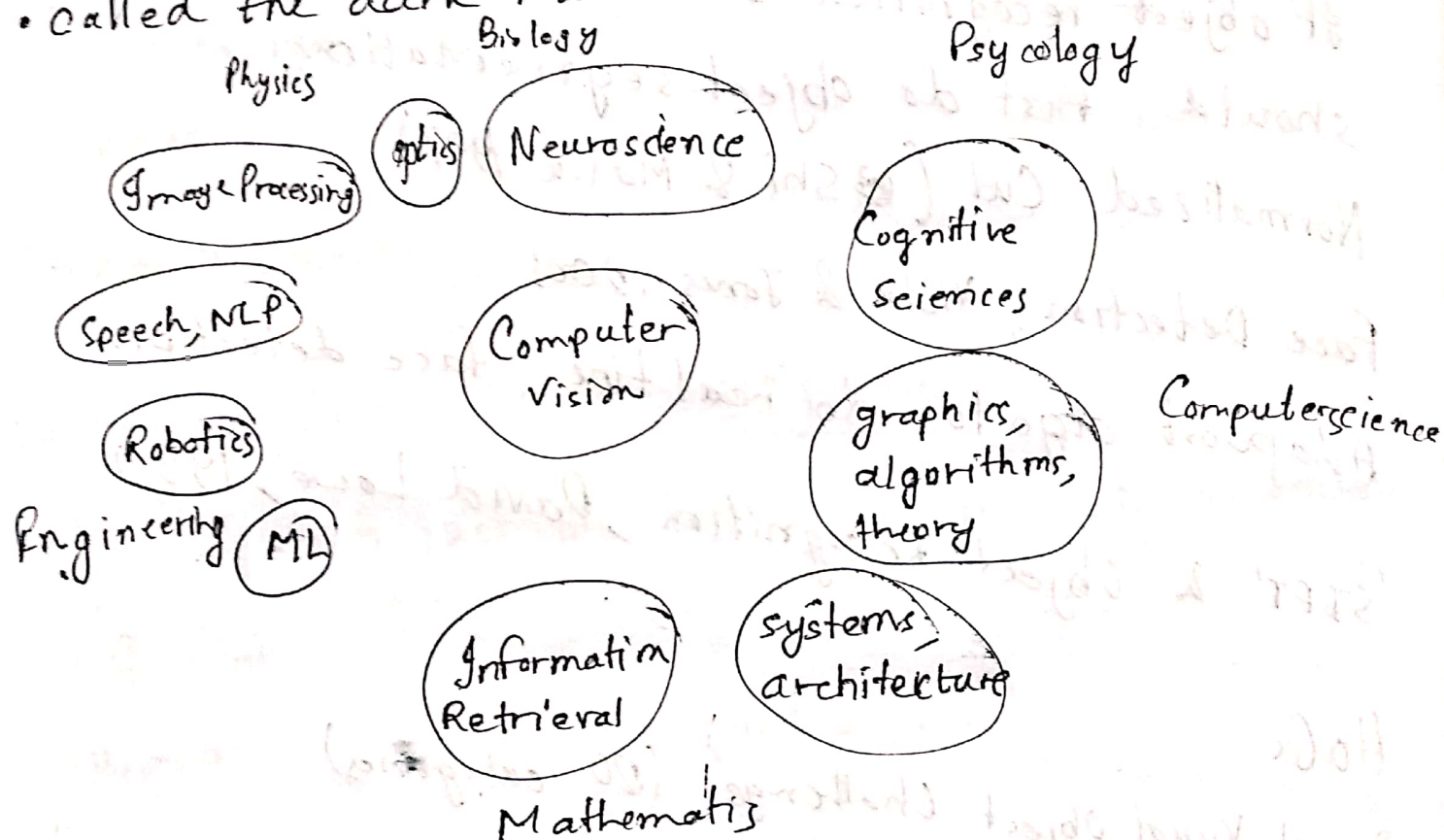


CS231n: Convolutional Nets for Visual Recognition

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... stimulus
└─ mapping ─┘

Block world [Larry Roberts, 1963]

1966 → Summer Vision Project

Generalized Cylinder

1979

Pictorial Structure

1973

Norm

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 (20 categories)

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

(22K categories & 14M images)

Took 3 years (took help from Amazon Mechanical Turk)

Imagenet Classification Challenge

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

Lecanet. ol (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

L1 Distance $\Rightarrow d_1(I_1, I_2) = \sum_p |I_1^p - I_2^p|$

Nearest Neighbor

$O(1)$

L2 dist $\Rightarrow d_2(I_1, I_2) = \sqrt{\sum_p (I_1^p - I_2^p)^2}$

Memorize data

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 test time
- distance metric not ~~to~~ informative (not good to check similarity)
- curse of dimensionality

Linear Classification

Parametric Approach

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

$$\underbrace{f(x, w)}_{10 \times 1} = \underbrace{w}_{10 \times 3072} \underbrace{x}_{3072 \times 1} + \underbrace{b}_{10 \times 1}$$

\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) Intra-class variation

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

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

Multi-class SVM loss

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

true value

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

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

What happens to loss if a score ^{that's high} 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 W is small so all $s \approx 0$.

What is the loss?

$(\text{number of classes} - 1)$

Q4) What if sum as ~~as~~ over all classes?

The loss increases by 1

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 W that gives $L=0$. Is this W unique?

No, there are many other W s. Like $2W$

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

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

Occam's Razor:

Among competing hypotheses the simplest is the best

L2 Regularization $R(w) = \sum_k \sum_l w_{k,l}^2$

L1 Regularization $R(w) = \sum_k \sum_l |w_{k,l}|$

Elastic net (L1+L2) $R(w) = \sum_k \sum_l \beta w_{k,l}^2 + |w_{k,l}|$

Max norm regularization

Dropout

Batchnorm

Stochastic depth

If you're Bayesian: L2 regularization also corresponds

MAP inference using a Gaussian prior on w

Softmax Classifier (Multinomial Logistic Regression)

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

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

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

Q1: What is the min and max loss?

$$\begin{aligned} \min &= 0 \\ \max &= \text{infinity} \end{aligned}$$

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 jiggle a bit (changing its score slightly). What happens to the loss in both cases?

SVM doesn't care about the score
Softmax continuously try to ^{improve} make the datapoints
like pushing the ~~ma~~ correct to inf and pushing the incorrect to $-\text{inf}$

$$\text{Softmax}, L_i = \log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

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

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

Optimization:

Strategy #1 Random search

Strategy #2 follow the slope

In 1-dimension, the derivative of a function

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

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

- 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 gradient: exact, fast, error-prone

Gradient check Using analytic gradient to find grads but checking with numerical gradient to see if they match. This is an interesting debugging tool.

Gradient descent

i) find grads

ii) $\text{weights} -= \frac{\text{stepsize} \times \text{grads}}{\text{hyperparameter learning rate}}$

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

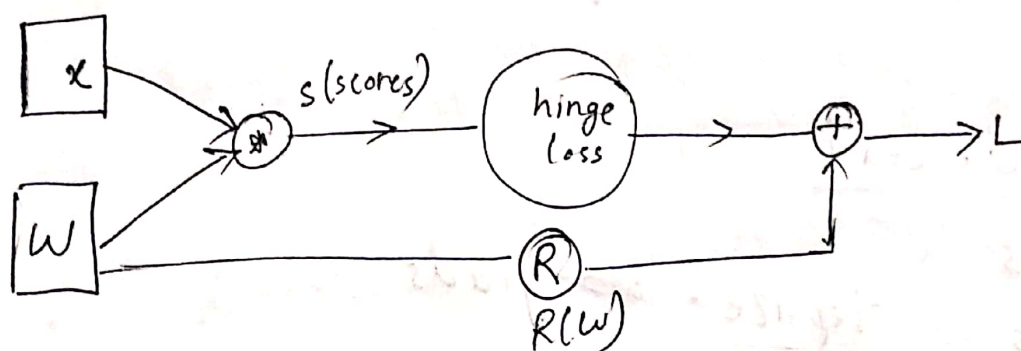
Histogram of Oriented Gradients

Bag of Words (Build Codebook, Encode Images)

Image Features vs ConvNets:

ConvNets (Krizhevsky 2012) AlexNet

Lecture 4: Intro to Neural Nets:



Leverage chain rule

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

$$\frac{d}{dx}(ax) = a$$

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

$$f(x) = c + x = a$$

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

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

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

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

mul 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

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

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

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