

HISTOGRAM OF GRADIENTS (HOG)

Edges



HOG: Human Detection

Navneet Dalal and Bill Triggs "Histograms of Oriented Gradients for Human Detection" CVPR05

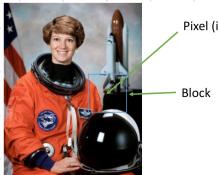


HOG: Oriented Gradients

- Given image I, Pixel location (i,j)
Compute the HOG feature for that pixel

steps

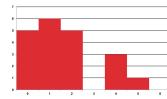
1. extract a block of some size



Revist Histogram

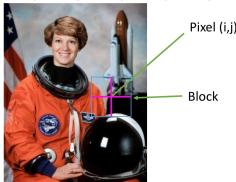
0	1	1	2	4
2	1	0	0	2
5	2	0	0	4
1	1	2	4	1

image



histogram

2. Divide block into sq. grid of sub blocks



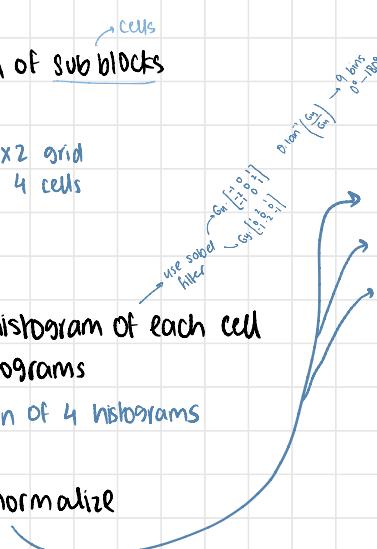
2x2 grid
so 4 cells

3. Compute orientation histogram of each cell
4. Concatenate the 4 histograms

V = Vector of concatenation of 4 histograms

5. Normalize V

There are 3 ways to normalize



1. Divide V by its Euclidean norm

2. Divide V by its L1 norm

3. Divide V by its Euclidean norm

↳ then clip any value over 0.2
in the resulting vector

↳ Then renormalize it by
dividing again by its Euclidean

sum of all absolute
values of v

Histogram of Oriented Gradients

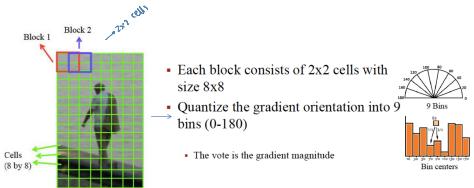
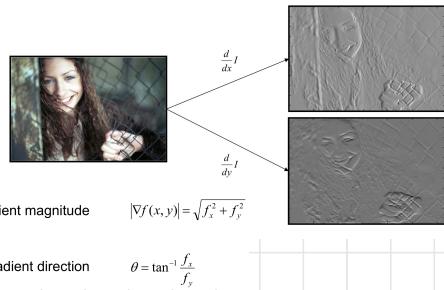


Image gradients



HOG



Histograms of Oriented Gradients

• Parameters and design options:

- Angles range from 0 to 180 or from 0 to 360 degrees?
 - In the Dalal & Triggs paper, a range of 0 to 180 degrees is used,
 - and HOGs are used for detection of pedestrians.
- Number of orientation bins.
 - Usually 9 bins, each bin covering 20 degrees.
- Cell size.
 - Cells of size 8x8 pixels are often used.
- Block size.
 - Blocks of size 2x2 cells (16x16 pixels) are often used.

• Usually a HOG feature has 36 dimensions.

- 4 cells * 9 orientation bins.

SCALE INVARIANT FEATURE TRANSFORM (SIFT)

SIFT

- ↳ Image content is transformed into local features coordinates
- ↳ Invariant local features

never changing

Invariant Local Features

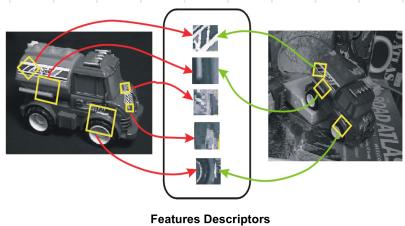
A good feature should be invariant to transformations:

Translation – Interest points should move consistently.

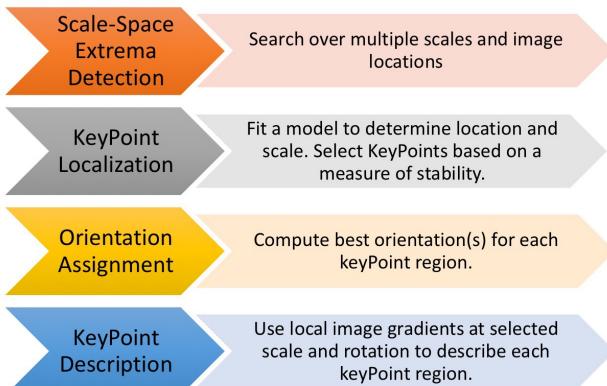
Rotation – Descriptors should stay the same even if the image rotates.

Scaling – Features should be detectable at multiple scales.

Lighting Changes – Interest points should be robust to brightness changes.



Overall Procedure at a High Level



AUTOMATIC SCALE SELECTION



$$f(I_{l_1 \dots l_n}(x, \sigma)) = f(I_{l_1 \dots l_n}(x', \sigma'))$$

How to find patch sizes at which f response is equal?

What is a good f ?

- Function responses for increasing scale (scale signature)



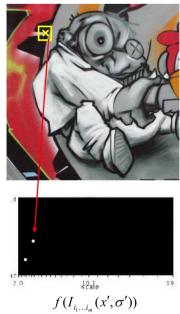
10/3/2024 $f(I_{l_1 \dots l_n}(x, \sigma))$ CAP5415 - Lecture 9



10/3/2024 $f(I_{l_1 \dots l_n}(x', \sigma'))$ CAP5415 - Lecture 9



10/3/2024 $f(I_{l_1 \dots l_n}(x, \sigma))$ CAP5415 - Lecture 9



10/3/2024 $f(I_{l_1 \dots l_n}(x', \sigma'))$ CAP5415 - Lecture 9

84

10/3/2024 $f(I_{l_1 \dots l_n}(x, \sigma))$ CAP5415 - Lecture 9



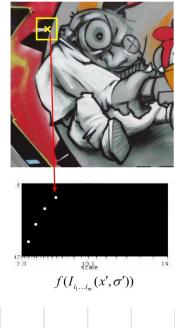
10/3/2024 $f(I_{l_1 \dots l_n}(x, \sigma))$ CAP5415 - Lecture 9



10/3/2024 $f(I_{l_1 \dots l_n}(x', \sigma'))$ CAP5415 - Lecture 9



10/3/2024 $f(I_{l_1 \dots l_n}(x, \sigma))$ CAP5415 - Lecture 9



10/3/2024 $f(I_{l_1 \dots l_n}(x', \sigma'))$ CAP5415 - Lecture 9



10/3/2024 $f(I_{l_1 \dots l_n}(x, \sigma))$ CAP5415 - Lecture 9



10/3/2024 $f(I_{l_1 \dots l_n}(x', \sigma'))$ CAP5415 - Lecture 9



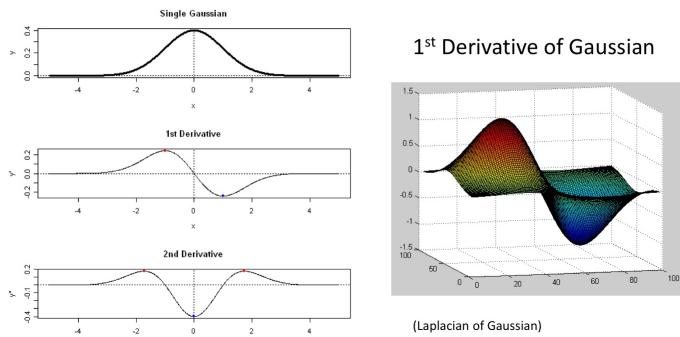
10/3/2024 $f(I_{l_1 \dots l_n}(x, \sigma))$ CAP5415 - Lecture 9



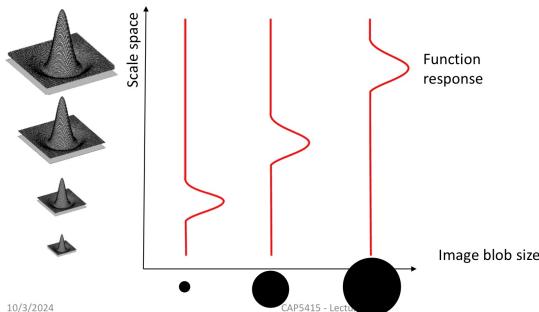
10/3/2024 $f(I_{l_1 \dots l_n}(x', \sigma'))$ CAP5415 - Lecture 9

88

What Is A Useful Signature Function f ?



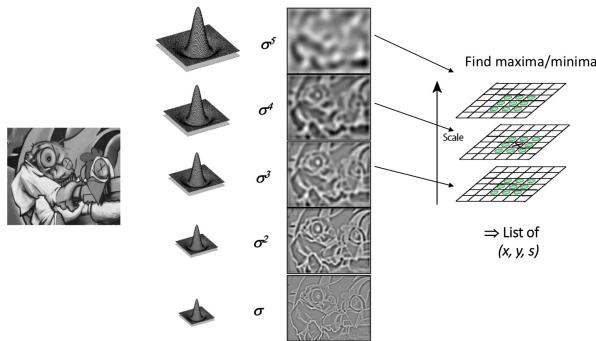
"Blob" detector is common for corners
• Laplacian (2nd derivative) of Gaussian (LoG)



10/3/2024

93

Find local maxima in position-scale space



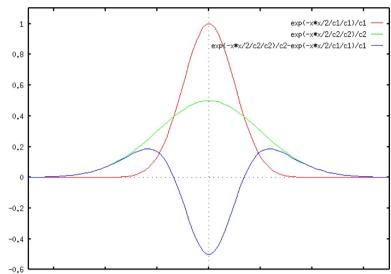
10/3/2024

CAPS415 - Lecture 9

94

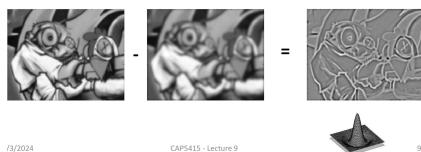
Alternative kernel

Approximate LoG with Difference-of-Gaussian (DoG).



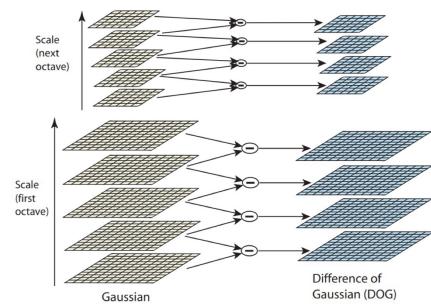
Approximate LoG with Difference-of-Gaussian (DoG).

1. Blur image with σ Gaussian kernel
2. Blur image with $k\sigma$ Gaussian kernel
3. Subtract 2. from 1.

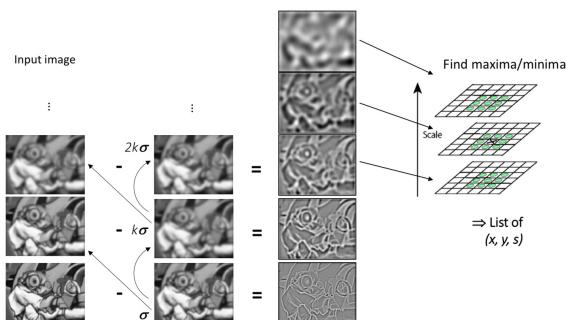


96

Scale-space

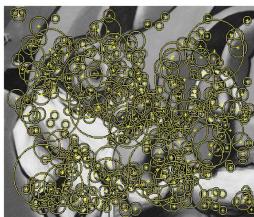


Find local maxima in position-scale space of DoG



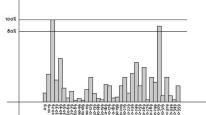
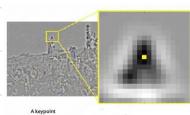
Results: Difference-of-Gaussian

- Larger circles = larger scale
- Descriptors with maximal scale response



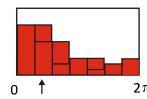
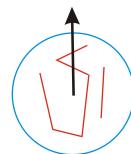
SIFT Orientation Estimation

- ↳ Compute gradient orientation histogram
- ↳ Select dominant orientation θ



SIFT Orientation Normalization

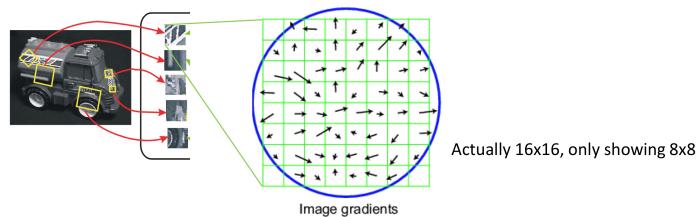
- ↳ Compute gradient orientation histogram
- ↳ Select dominant orientation θ
- ↳ Normalize/rotate to fixed orientation



SIFT Descriptor formation

- ↳ Compute on local 16x16 window around detection
- ↳ Rotate and scale window according to discovered orientation θ and scale r
- ↳ Compute gradients weighted by a Gaussian of variance half the window

for smooth fallout

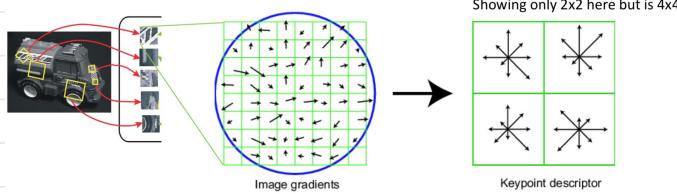


- ↳ 4x4 array of gradient orientation histogram

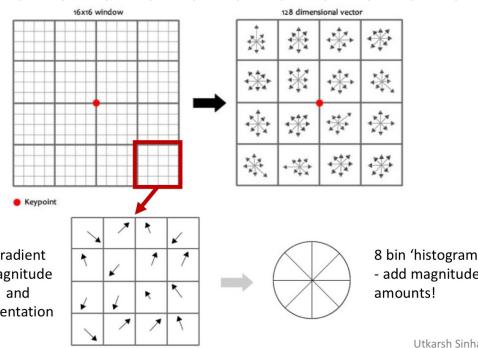
weighted by gradient magnitude

- ↳ Bin into 8 orientations

$$8 \times 4 \times 4 \text{ array} = 128 \text{ dimensions}$$

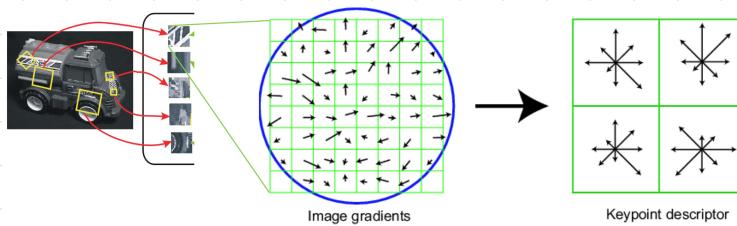


SIFT Descriptor Extraction



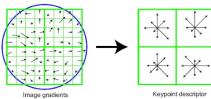
Reduce effect of illumination

- ↳ 128-dim vector normalize to 1
- ↳ Threshold gradient magnitudes
 - ↳ to avoid excessive influence of high gradients
 - ↳ After normalization, clamp gradients > 0.2
- ↳ Renormalize



Review: Local Descriptors

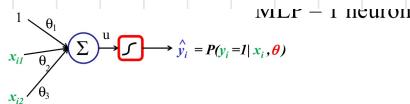
- Most features can be thought of as
 - templates,
 - histograms (counts),
 - or combinations
- The ideal descriptor should be
 - Robust and Distinctive
 - Compact and Efficient
- Most available descriptors focus on edge/gradient information
 - Capture texture information
 - Color rarely used



ishma hafeez
notes
represent

MULT LAYER PERCEPTRON

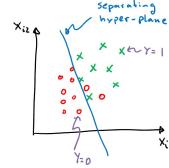
1. MLP - 1 layer 1 neuron



We are given the data $\{x_i, y_i\}_{i=1}^n$

e.g.

	x_{i1}	x_{i2}	y_i
i=1	0.2	6	0
i=2	0.3	22	1
i=3	0.6	-0.6	1
i=4	-0.4	58	0
:			



$$0 < \hat{y}_i < 1$$

$$u = \theta_0 + \theta_1 x_{i1} + \theta_2 x_{i2}$$

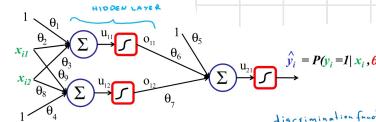
$$\hat{y}_i = \frac{1}{1+e^{-u}} = \frac{1}{1+e^{-\theta_0-\theta_1 x_{i1}-\theta_2 x_{i2}}} = P(Y_i = 1 | x_i, \theta)$$

$$P(Y_i | x_i, \theta) = \hat{y}_i^{Y_i} (1 - \hat{y}_i)^{1-Y_i} = \begin{cases} \hat{y}_i & \text{When } Y_i = 1 \\ 1 - \hat{y}_i & \text{Otherwise} \end{cases}$$

For n independent observations (Bernoulli):

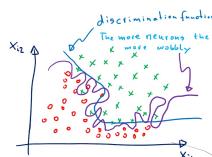
$$P(Y | x, \theta) = \prod_{i=1}^n P(Y_i | x_i, \theta)$$

2. MLP - 2 layers 3 neurons



Data:

	x_{i1}	x_{i2}	y_i
i=1	4	2	0
i=2	0.2	-5	1
i=3	-100	3.1	1
i=4	6	9	0
i=5	5	8	0



$$u_{11} = \theta_0 + \theta_1 x_{i1} + \theta_2 x_{i2}$$

$$u_{12} = \theta_0 + \theta_3 x_{i1} + \theta_4 x_{i2}$$

$$u_{21} = \theta_0 + \theta_5 x_{i1} + \theta_6 x_{i2}$$

$$\hat{y}_j = P(Y_j = 1 | x_j, \theta)$$

$$\hat{y}_{11} = \frac{1}{1+e^{-u_{11}}} \quad \hat{y}_{12} = \frac{1}{1+e^{-u_{12}}}$$

$$\hat{y}_{21} = \frac{1}{1+e^{-u_{21}}}$$

$$P(Y | x, \theta) = \prod_{i=1}^n \hat{y}_{ij}^{Y_{ij}} (1 - \hat{y}_{ij})^{1-Y_{ij}}$$

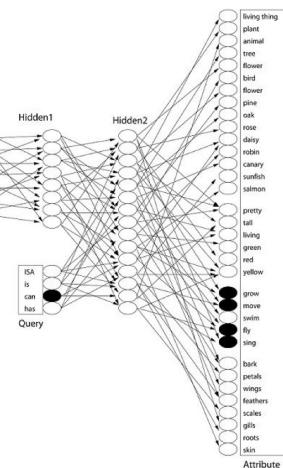
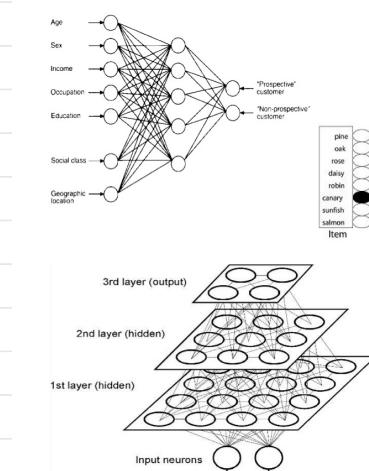
For n independent observations.

$$P(Y | x, \theta) = \prod_{i=1}^n \hat{y}_{ij}^{Y_{ij}} (1 - \hat{y}_{ij})^{1-Y_{ij}} = \prod_{i=1}^n P(Y_i | x_i, \theta)$$

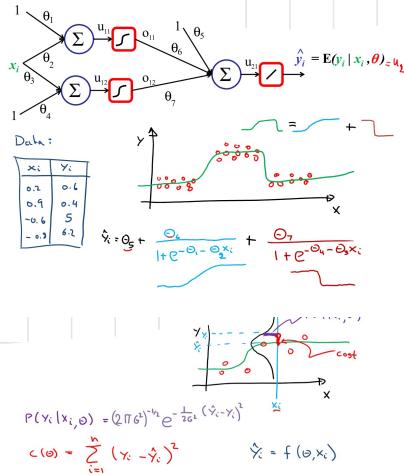
Cost:

$$C(\theta) = -\log P(Y | x, \theta) = -\sum_{i=1}^n Y_i \log \hat{y}_{ij} + (1 - Y_i) \log (1 - \hat{y}_{ij})$$

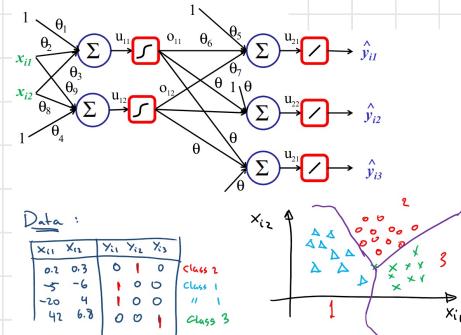
i.e. minimize the cross-entropy error.



3. MLP - Regression



4. MLP - Multiclass



To get a probabilistic model, define: SOFTMAX

$$P(Y_i = k | x_i, \theta) = P(Y_i = k | x_i, \theta) = \frac{e^{\hat{y}_k}}{e^{\hat{y}_1} + e^{\hat{y}_2} + e^{\hat{y}_3}}$$

$$\hat{y}_2(Y_i) \approx \begin{cases} 1 & Y_i = 2 \\ 0 & \text{else} \end{cases}$$

Then,

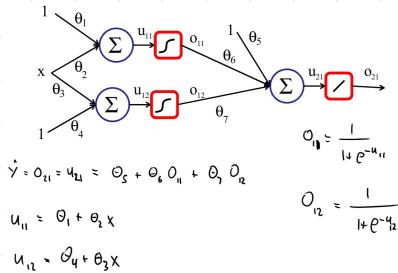
$$P(Y_i = k | x_i, \theta) = \left[\frac{e^{\hat{y}_{i1}}}{e^{\hat{y}_{i1}} + e^{\hat{y}_{i2}} + e^{\hat{y}_{i3}}} \right] \hat{y}_1(x_i) \left[\frac{e^{\hat{y}_{i2}}}{e^{\hat{y}_{i1}} + e^{\hat{y}_{i2}} + e^{\hat{y}_{i3}}} \right] \hat{y}_2(x_i) \left[\frac{e^{\hat{y}_{i3}}}{e^{\hat{y}_{i1}} + e^{\hat{y}_{i2}} + e^{\hat{y}_{i3}}} \right] \hat{y}_3(x_i)$$

$$= \begin{cases} e^{\hat{y}_1} / \sum_m e^{\hat{y}_m} & Y_i = 1 \\ e^{\hat{y}_2} / \sum_m e^{\hat{y}_m} & Y_i = 2 \\ e^{\hat{y}_3} / \sum_m e^{\hat{y}_m} & Y_i = 3 \end{cases}$$

Cost:

$$C(\theta) = -\log P(Y_i | x_i, \theta) = -\sum_{i=1}^n \sum_{j=1}^3 \hat{y}_j(Y_i) \log \frac{e^{\hat{y}_j}}{\sum_m e^{\hat{y}_m}}$$

BACKPROPAGATION



$$E(\theta) = (\hat{y}_i - y_i)^2$$

$$\frac{\partial E(\theta)}{\partial \theta_j} = -2 (\hat{y}_i - y_i) \frac{\partial \hat{y}_i}{\partial \theta_j}$$

$$\frac{\partial \hat{y}_i}{\partial \theta_5} = 1 \quad \frac{\partial \hat{y}_i}{\partial \theta_6} = o_1 \quad \frac{\partial \hat{y}_i}{\partial \theta_7} = o_2$$

$\hat{y} = o_4 = \theta_5 o_1 + \theta_6 o_2 + \theta_7 o_3$

$o_1 = \frac{1}{1 + e^{-u_1}}$

$u_1 = \theta_1 + \theta_2 x$

$\frac{\partial \hat{y}}{\partial \theta_3} = \frac{\partial \hat{y}}{\partial o_1} \frac{\partial o_1}{\partial u_1} \frac{\partial u_1}{\partial \theta_3}$

$= \theta_1 \theta_3 [1 - o_1] x$

$\frac{\partial \hat{y}}{\partial \theta_3} = \frac{\partial \hat{y}}{\partial o_1} \frac{\partial o_1}{\partial u_1} \frac{\partial u_1}{\partial \theta_3}$

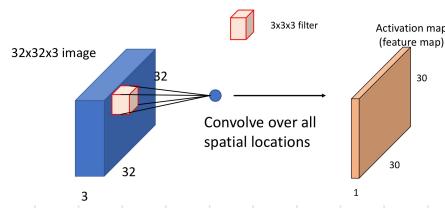
$= \theta_1 \theta_3 [1 - o_1] x$

$\text{Batch: } \hat{o}_j^{(k+1)} = \hat{o}_j^{(k)} + \sum_{i=1}^n (\hat{y}_i - y_i) \frac{\partial \hat{y}_i}{\partial \theta_j}$

$\text{Update: } \hat{o}_j^{(k+1)} = \hat{o}_j^{(k)} + (\hat{y}_i - y_i) \frac{\partial \hat{y}_i}{\partial \theta_j} x_i$

TRAINING NEURAL NETWORKS

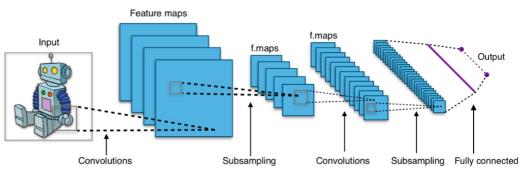
Network Parameters - Recap



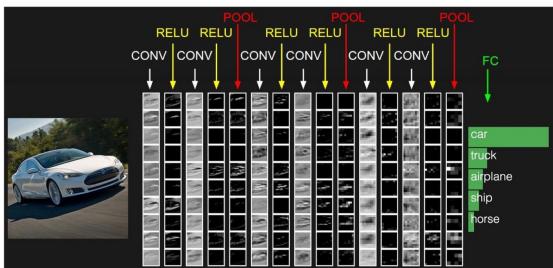
Convolution - Intuition



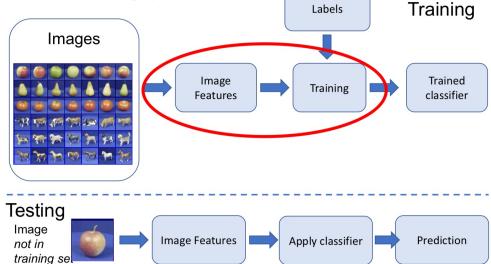
General CNN architecture - recap



Visualizing CNN



Learning phases



LOSS FUNCTION

↳ tells how good the network performs
in terms of Prediction

↳ Network training (optimization)

↳ find the best network parameters \rightarrow to minimize the loss

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

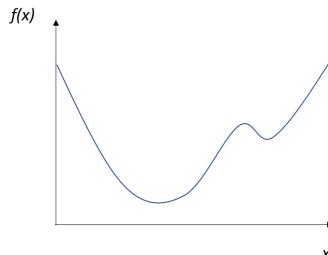
Loss function
network
input
Network parameters
Ground truth

Loss function measures how 'good' our network is at classifying the training examples wrt. the parameters of the model (the perceptron weights).

NETWORK TRAINING

↳ Gradient Descent

↳ way to minimize a cost function



LOSS FUNCTION TYPES

1. CROSS ENTROPY

for multiclass classification

↳ measures diff b/w 2 probability distributions

True Distribution (actual tasks)
Predicted Distribution

$$L = -\frac{1}{N} \sum_{i=1}^N [y_i \log(\hat{y}_i) + (1-y_i) \log(1-\hat{y}_i)]$$

Actual
Predicted

2. MEAN SQUARED ERROR (MSE)

↳ used in regression tasks

↳ measures the avg of squared diff b/w actual and predicted values

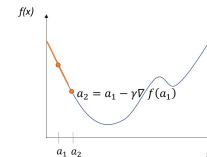
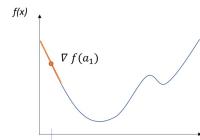
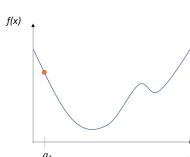
$$L = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

General APPROACH

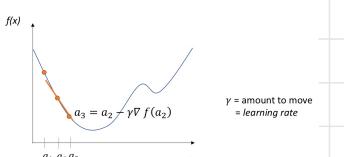
Pick random starting point.

Compute gradient at point (analytically or by finite differences)

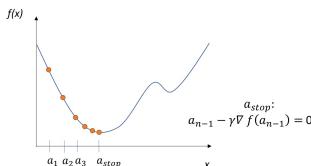
Move along parameter space in direction of negative gradient



Move along parameter space in direction of negative gradient.



Stop when we don't move any more.

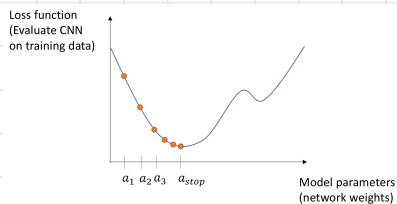


Gradient Descent

- ↳ Optimizer for functions
- ↳ Guaranteed to find optimum for convex functions
 - ↳ non-convex → find local optimum
 - ↳ most vision problems aren't convex
- ↳ Works for multi-variate functions
 - ↳ need to compute matrix of partial derivatives

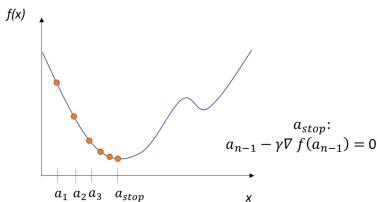
Train CNN with Gradient Descent

- $x^i, y^i = n$ training examples
- $f(x) = \text{feed forward network}$
- $L(x, y; \theta) = \text{some loss function}$



General approach - recap

Stop when we don't move any more.

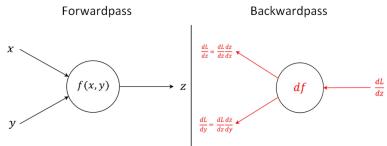


Differentiability

- Loss function
- Activation function
- Convolution
- Pooling
- ...

Backpropagation - Chain Rule

- Chain rule $\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial x}$



- Add gate:
 - gradient distributor
- Max gate:
 - gradient router
- Mul gate:
 - gradient switcher

Stochastic Gradient Descent

↳ dataset can be too large

↳ can not apply gradient descent w.r.t all data points

↳ randomly sample a data point

↳ perform gradient descent per sample and iterate

↳ picking a subset of points: "mini batch" → batchsize

↳ randomly initialize starting W and pick learning rate γ

while not at minimum :

↳ shuffle training set

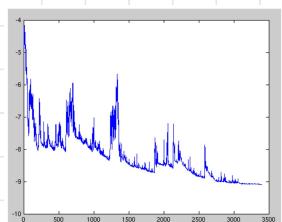
↳ for each data point $i=1 \dots n \rightarrow$ maybe as mini-batch] → EPatch

↳ Gradient Descent

↳ loss will not always decrease TOCALLY

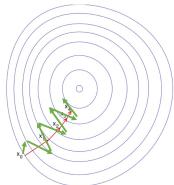
↳ as training data point is random

↳ still converges over time



Gradient Descent Oscillations

Slow to converge to the (local) optimum



9/18/2024

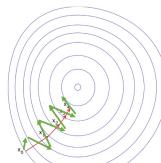
CAPS415 - Lecture 7

Momentum

- ↳ Adjust the gradient by a weighted sum of the prev amount plus the current amount.
- ↳ without momentum: $\theta_{t+1} = \theta_t - \gamma \frac{\partial L}{\partial \theta}$
- ↳ with momentum (new & popular)
$$\theta_{t+1} = \theta_t - \gamma \left(\alpha \left[\frac{\partial L}{\partial \theta} \right]_{t-1} + \left[\frac{\partial L}{\partial \theta} \right]_t \right)$$

Lowering the learning rate

-Takes longer to get to the optimum

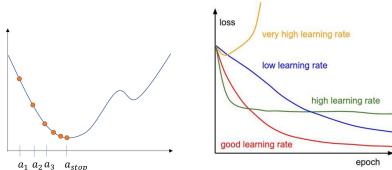


9/18/2024

CAPS415 - Lecture 7

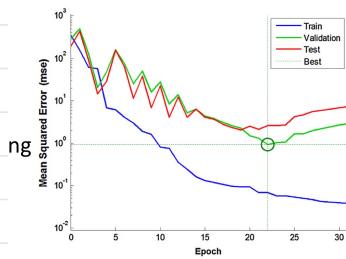
40

Learning rate

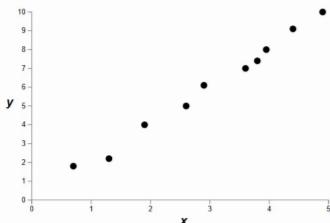


Problem of fitting

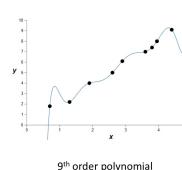
- ↳ Too many parameters = overfitting
- ↳ Not enough parameters = underfitting
- ↳ more data = less chance to overfit
- ↳ How do we know what is required



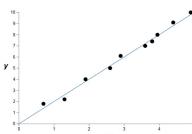
Data fitting problem



Which is better?



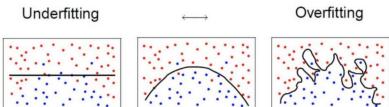
9th order polynomial



1st order polynomial

[Nielson]

Data fitting problem



- Early stopping
- Regularization
- Dropout
- ...

2. Regularization

- ↳ Attempt to guide solution to not overfit
- ↳ But still give freedom with many parameters
- ↳ So basically

Penalise the use of Parameters
to Prefer small weights

- ↳ Add a cost to having high weights

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

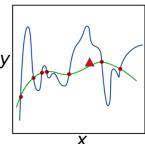
Loss Weight Decay Regularization

- In common use,

$$\text{L1 Norm} - R(W) = \sum_i \sum_j |W_{ij}|$$

$$\text{L2 Norm} - R(W) = \sum_i \sum_j W_{ij}^2$$

$$\text{Elastic net} - R(W) = \sum_i \sum_j \beta W_{ij}^2 + (1-\beta) |W_{ij}|$$



9/18/2024

CAPS415 - Lecture 7

51

3. Dropout

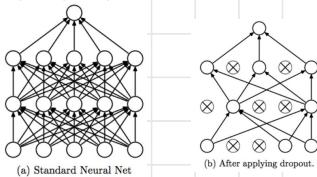
- ↳ Stochastically switch neurons off
- ↳ Each neuron is set to 0 $\xrightarrow{\text{with probability } p}$

↳ hidden units cannot co-adapt
to each other

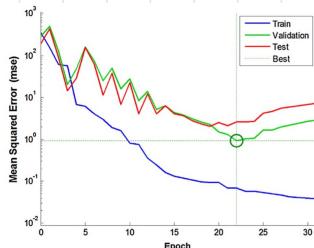
↳ Units are useful independently

Hyperparameter

↳ p is usually set to 0.5

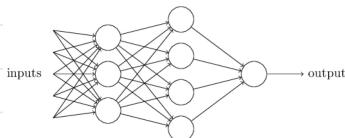


1. Early Stopping



Regularization: Ensemble

- ↳ Networks start with random weights
- ↳ at every train = slightly different outcome
- ↳ Random weights are used cuz
if weights are all equal,
response across filters will
be equivalent
- ↳ so network closest train



$$w \cdot x \equiv \sum_j w_j x_j$$

Why not train 5 different networks with random starts and vote on their outcome?

- Works fine!
- Helps generalization because error due to overfitting is averaged; reduces variance.

Training Steps

- Define network
- Loss function
- Initialize network parameters
- Get training data
 - Prepare batches
- Feedforward one batch
 - Compute loss
 - Backpropagate gradients
 - Update network parameters
 - Repeat

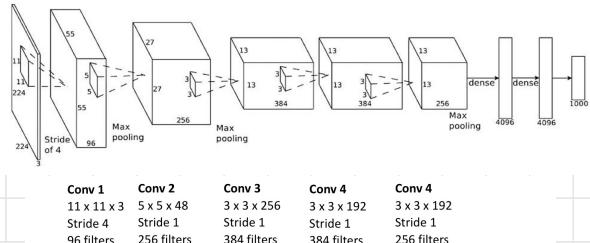
AlexNet - Training

Parameters:

- First use of ReLU
- Dropout 0.5
- Batch size 128
- Optimizer SGD
- Momentum 0.9
- Learning rate 1e-2
- Decay – lr reduced by 10 manually when val accuracy plateaus
- L2 weight decay 5e-4

AlexNet for ImageNet

	AlexNet	FLOPs
params	4M	4M
FC 1000	4M	16M
FC 4096 / ReLU	16M	37M
FC 4096 / ReLU	37M	
Max Pool 3x3x2	442K	
Conv 3x3x1, 256 / ReLU	1.3M	74M
Conv 3x3x1, 384 / ReLU	884K	112M
Conv 3x3x1, 384 / ReLU	1.3M	149M
Max Pool 3x3x2	884K	
Local Response Norm	307K	
Conv 5x5x1, 256 / ReLU	Conv 5x5x1, 256 / ReLU	223M
Max Pool 3x3x2	307K	
Local Response Norm	35K	
Conv 11x11x4, 96 / ReLU	35K	105M



Residual Networks

- Deep networks performs worse
 - As we add more layers
- Problem
 - Vanishing gradients
- It models
 - $H(x) = F(x) + x$
- Skip connections
 - Help in backpropagation

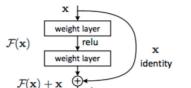
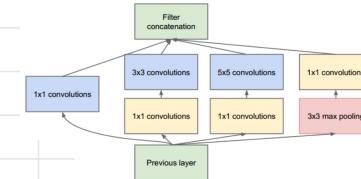


Figure 2. Residual learning: a building block.

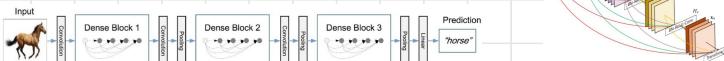
GOOGLENET - Inception

- ResNet is about going deeper
- Inception is about going wider
- Focused on computational efficiency
- The network learns
 - Which features are useful



DenseNet

- Densely Connected Convolutional Network
 - Similar to ResNet
 - Concat features instead of summation
 - A layer is passed all previous maps
 - Sequence of dense blocks



Huang, Gao, et al. "Densely connected convolutional networks." (2016).

CNN AND LEARNING RATE

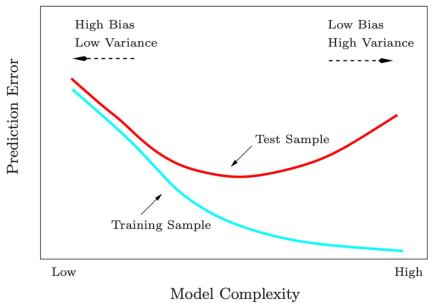
Bias

error caused because the model lacks the ability to represent the (complex) concept

Variance

error caused because the learning algorithm overreacts to small changes (noise) in the training data

Total Loss = Bias + Variance (+ noise)



Universality Theorem

↳ Any continuous function f

$$f: \mathbb{R}^N \rightarrow \mathbb{R}^M$$

can be realized by a network
with one hidden layer given enough hidden neurons

Universality is Not Enough

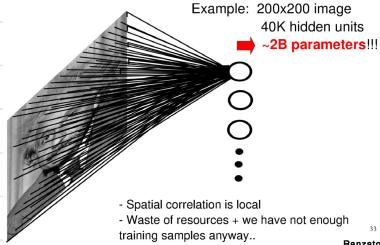
- Neural network has very high capacity (millions of parameters)
- By our basic knowledge of bias-variance tradeoff, so many parameters should imply very low bias, and very high variance. The test loss may not be small.
- Many efforts of deep learning are about mitigating overfitting!

Address Overfitting for NN

- Use larger training data set
- Design better network architecture

CONVOLUTIONAL NN (CNN)

Images as input to neural networks



CNN = a multi-layer neural network with

- Local connectivity:
 - Neurons in a layer are only connected to a small region of the layer before it
- Share weight parameters across spatial positions:
 - Learning shift-invariant filter kernels

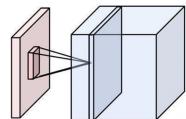
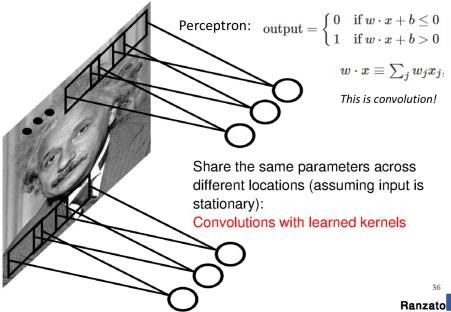


Image credit: A. Karpathy

Convolutional Layer



This is convolution!

This is convolution!

Share the same parameters across different locations (assuming input is stationary):

Convolutions with learned kernels

36

Recap: Image filtering

$f[\cdot, \cdot]$

$$h[.,.]$$

$$I[.,.]$$

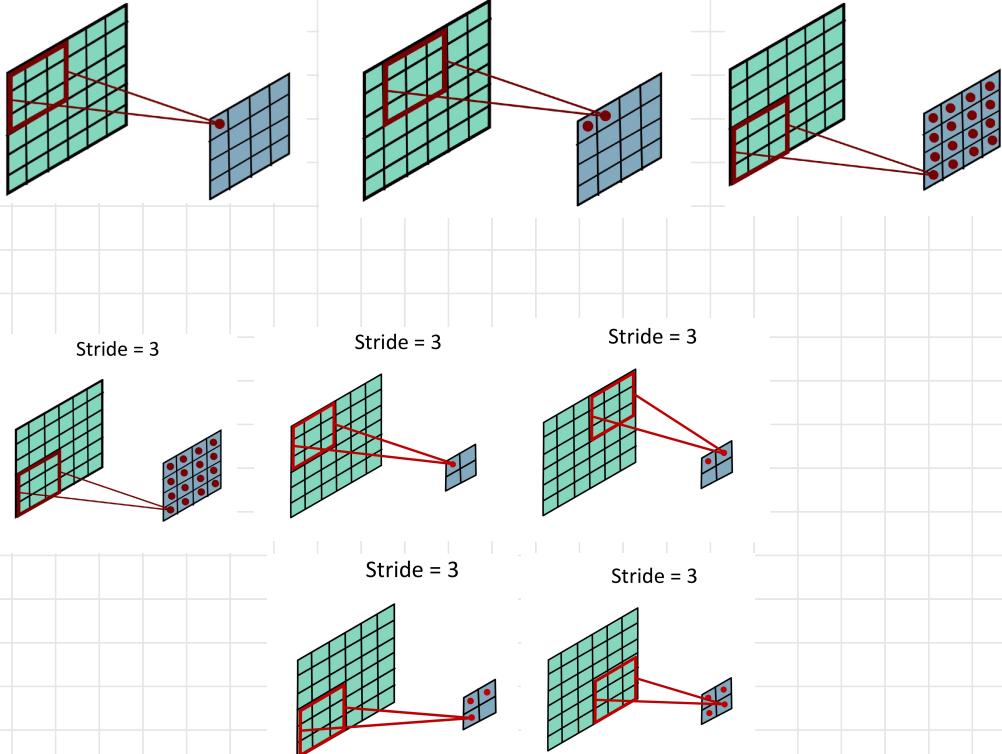
The diagram illustrates a 2D convolution process. The input layer (bottom) consists of a 10x10 grid with values 0, 1, and 2. The kernel (middle) is a 3x3 matrix with all elements set to 90. The stride is 2. The output layer (top) shows the result of the convolution, where the value at position (1,1) is 90, indicating the result of the first convolution step.

0	10	20	30	30	30	30	20	10
0	20	40	60	60	60	40	20	0
0	30	60	80	80	80	60	30	0
0	30	50	80	80	80	80	30	0
0	20	50	80	80	80	60	30	0
0	20	30	50	50	60	40	20	0
0	20	30	50	30	30	20	10	0
0	10	10	10	0	0	0	0	0

$$h[m, n] = \sum_{k,l} f[k, l] I[m + k, n + l]$$

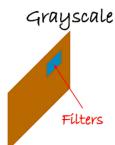
Credit: S. Sei

Convolutional Layer



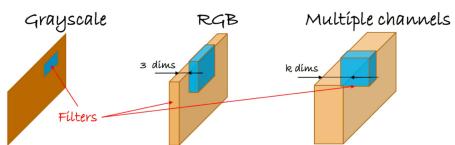
2D spatial features

- If images are 2-D, parameters should also be organized in 2-D
 - That way they can learn the local correlations between input variables
 - That way they can "exploit" the spatial nature of images

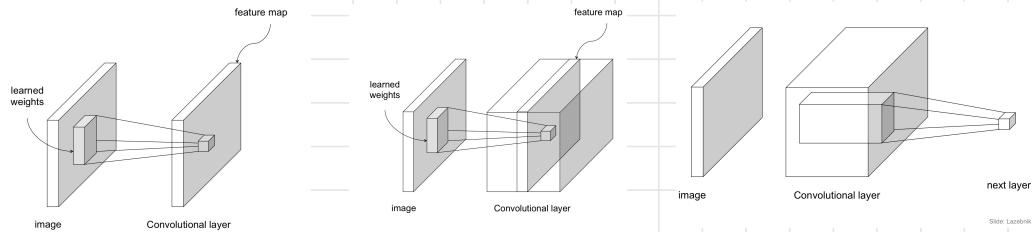


K-D spatial features

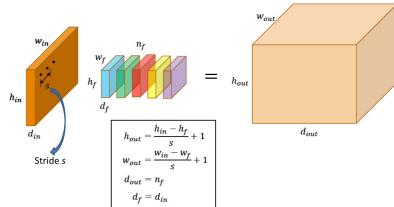
- Similarly, if images are k-D, parameters should also be k-D



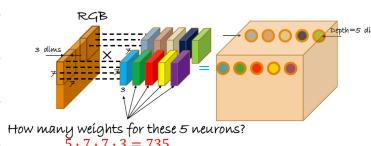
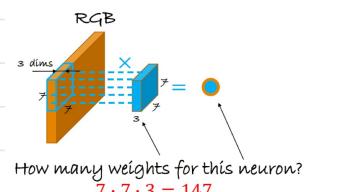
Dimensions of convolution



Dimensions of convolution



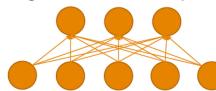
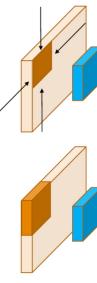
Number of weights



- o Question: Spatial structure?
- Answer: Convolutional filters
- o Question: Huge input dimensionalities?
- Answer: Parameters are shared between filters
- o Question: Local variances?
- Answer: Pooling

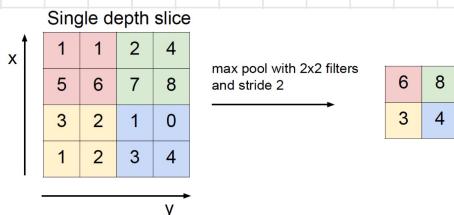
Local Connectivity

- The weight connections are surface-wise local
 - Local connectivity
- The weights connections are depth-wise global
- For standard neurons no local connectivity
 - Everything is connected to everything

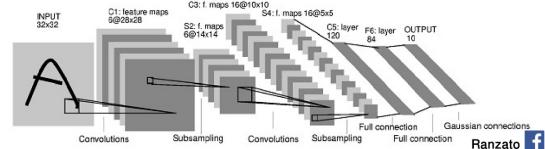


Pooling Operations

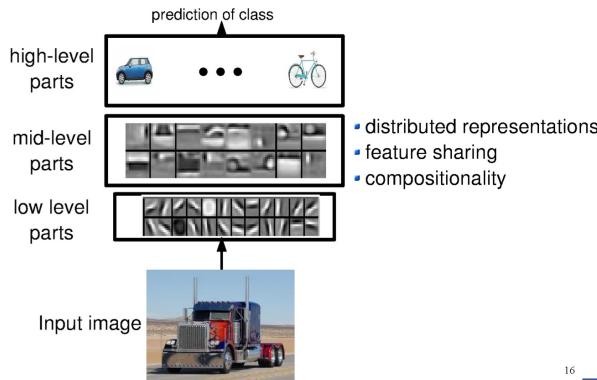
- Aggregate multiple values into a single value
- Invariance to small transformations
 - Keep only most important information for next layer
- Reduces the size of the next layer
 - Fewer parameters, faster computations
- Observe larger receptive field in next layer
 - Hierarchically extract more abstract features



Yann LeCun's MNIST CNN architecture



Interpretation



Lee et al. "Convolutional DBN's ..." ICML 2009

LEARNING NN

1. Setting Hyperparameters

Idea #1: Choose hyperparameters that work best on the data

BAD: big network always works perfectly on training data

Your Dataset
Idea #2: Split data into **train** and **test**, choose hyperparameters that work best on test data

BAD: No idea how algorithm will perform on new data

train

test

Idea #3: Split data into **train**, **val**, and **test**; choose hyperparameters on val and evaluate on test

Better!

train

validation

test

2. Select Optimizer

L Stochastic gradient descent

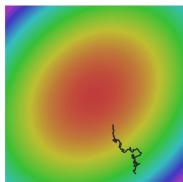
Gradient from entire training set:

$$\nabla C = \frac{1}{N} \sum_n \nabla c_n$$

- For large training data, gradient computation takes a long time
 - Leads to "slow learning"
- Instead, consider a mini-batch with m samples
- If sample size is large enough, properties approximate the dataset

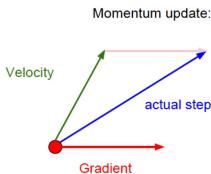
$$\frac{\sum_{j=1}^m \nabla c_{n_j}}{m} \approx \frac{\sum_n \nabla c_n}{N} \cdot \nabla C$$

Our gradients come from minibatches so they can be noisy!

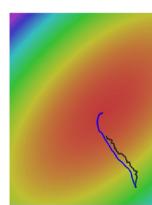


$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(x_i, y_i, W)$$

$$\nabla_W L(W) = \frac{1}{N} \sum_{i=1}^N \nabla_W L_i(x_i, y_i, W)$$



$$x_{t+1} = x_t - \alpha \nabla f(x_t)$$



$$v_{t+1} = \rho v_t + \nabla f(x_t)$$

$$x_{t+1} = x_t - \alpha v_{t+1}$$

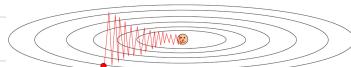
Build up velocity as a running mean of grad

Q) What if the loss function has a local minima or saddle point?
A) Zero gradient, gradient descent gets stuck Stochastic gradient descent

Q) What if loss changes quickly in one direction and slowly in another?

Q) What does gradient descent do?

A) Very slow progress along shallow dimension, jitter along steep direction



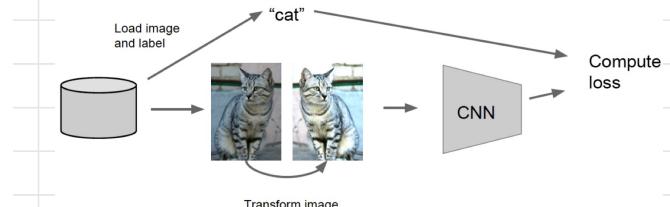
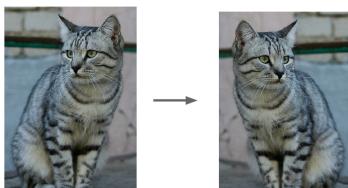
Loss function has high condition number: ratio of largest to smallest singular value of the Hessian matrix is large

Many variations of using momentum

- In PyTorch, you can manually specify the momentum of SGD
- Or, you can use other optimization algorithms with "adaptive" momentum, e.g., ADAM
- ADAM: Adaptive Moment Estimation
- Empirically, ADAM usually converges faster, but SGD gives local minima with better generalizability

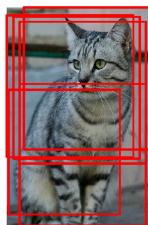
3. Data Augmentation

Horizontal flips



Random crops and scales

1. Pick random L in range [256, 480]
2. Resize training image, short side = L
3. Sample random 224 x 224 patch



Color jitter

Simple: Randomize contrast and brightness



More Complex:

1. Apply PCA to all [R, G, B] pixels in training set
2. Sample a "color offset" along principal component directions
3. Add offset to all pixels of a training image

Can do a lot more: rotation, shear, non-rigid, motion blur, lens distortions,

Exam

- Linear algebra, such as
 - rank, null space, range, invertible, eigen decomposition, SVD, pseudo inverse, basic matrix calculus
- Optimization:
 - Least square, low-rank approximation, statistical interpretation of PCA
- Image formation
 - diffuse/specular reflection, Lambertian lighting equation
- Filtering
 - Linear filter, filter vs convolution, properties of filters, filterbank, usage of filters, median filter
- Statistics:
 - Bias, variance, bias-variance tradeoff, overfitting, underfitting
- Neural network
 - Linear classifier, softmax, why linear classifier is insufficient, activation function, feed-forward pass, universality theorem, what does back-propagation do, stochastic gradient descent, concepts in neural networks, why CNN, concepts in CNN, how to set hyperparameter, moment in SGD, data augmentation

ishma hafeez
notes
represent