Recap

Standard **Bag-of-Words** Pipeline (for image classification)

Dictionary Learning: Learn Visual Words using clustering

Encode: build Bags-of-Words (BOW) vectors for each image

Classify: Train and test data using BOWs

K-Means Clustering

K-means clustering alternates between two steps:

- **1**. Assume the cluster centers are known (fixed). Assign each point to the closest cluster center.
- **2.** Assume the assignment of points to clusters is known (fixed). to the cluster.
- The algorithm is initialized by choosing K random cluster centers
- K-means converges to a local minimum of the objective function Results are initialization dependent

Compute the best center for each cluster, as the mean of the points assigned

Expectation Maximization

A simpler version

Given a model repeat Not exactly the hard assignments of K-Means

The K-Means centers 1. Create an "expectation" of the (log-)likelihood with the current hypothesis 2. Update the hypothesis to one that **maximizes** the expectation above







2. Encode: build Bag-of-Words (BOW) vectors for each image

2. Histogram: count the number of visual word occurrences







Support Vector Machines (SVM)

What's the best w?





Want a hyperplane that is far away from 'inner points'

Support Vector Machines (SVM)

Find hyperplane w such that ...



Distance to the border



Becomes 1 because it's the thing at the border (+1)

Support Vectors



• The active constraints are due to the data that define the classification boundary, these are called support vectors

> Final classifier can be written in terms of the support vectors:

$$\hat{y} = \operatorname{sign} \left(\hat{w}_0 + \sum_{\alpha_i > 0} \alpha_i y_i \mathbf{x}_i^T \mathbf{x} \right)$$

Non-Linear SVM

Replace inner product with kernel





- $\mathbf{x}_i^T \mathbf{x} \to \phi(\mathbf{x}_i)^T \phi(\mathbf{x}) \to k(\mathbf{x}_i, \mathbf{x})$
 - Data are (ideally) linearly separable in $\Phi(x)$ **0**^{0.7} But we don't need to know $\phi(x)$, we just specify k(x,y)Points with $\alpha > 0$ (circled) are support vectors Other data can be removed without affecting classifier

Bag-of-Words Representation

Algorithm:

Initialize an empty K -bin histogram, where K is the number of codewords Extract local descriptors (e.g. SIFT) from the image For each local descriptor **x**

Map (Quantize) **x** to its closest codeword \rightarrow **c**(**x**) Increment the histogram bin for c(x)Return histogram

vector machine or k-Nearest Neighbor classifier

We can then classify the histogram using a trained classifier, e.g. a support

- Won the Imagenet Large Scale Visual Recognition Challenge (ILSVRC) in 2012 by a large margin
- Some ingredients: Deep neural net (Alexnet), Large dataset

(Incompate) I and a factor (1 (DI) and a contraction -IM GENET Large Scale Visual Recognition Challenge



Alexnet

$[J.Johnson]_{12}$



THE UNIVERSITY OF BRITISH COLUMBIA

CPSC 425: Computer Vision



Lecture 19: Visual Classification 2, Linear Classification

Menu for Today

Topics:

- Nearest Neighbour, nearest mean
- Linear Classification

Readings:

- Today's Lecture: Szeliski 11.4, 12.3-12.4, 9.3, 5.1-5.2

Reminders:

Assignment 5: Stereo and Optical Flow due Apr 3rd

Bayesian classification



Visual Classification 2

Nearest neighbour, nearest mean Linear classification, CIFAR10 case study

- Bayesian Classification, Gaussian distributions, priors

CIFARIO Dataset

• 60,000 32x32 images in 10 classes (50k train, 10k test)

airplane	the state of the s
automobile	
bird	
cat	
deer	
dog	W. 1.
frog	
horse	- the set
ship	
truck	

Good test set for visual recognition problems

Hand labelled set of 10 categories from Tiny Images dataset



CIFARIO Classification

• Let's build an image classifier!











airplane automobile

bird

cat

deer

• Start by vectorizing the image data



• x = 3072 element vector of 0-255 • Note this throws away spatial structure, we'll bring it back later when we look at feature extraction and CNNs







dog

frog

horse

ship

truck

 $32 \times 32 \times RGB$ (8 bit) image \rightarrow x = [65 102 33 57 54 ...]

Nearest Mean Classification

• How about a single template per class



Nearest Mean Classification

• Find nearest mean and assign class

• CIFAR 10 class means



- $c_q = \arg\min_i |\mathbf{x}_q \mathbf{m}_i|^2$

Nearest Mean Classifier

• Suppose we have 2 classes linearly separable



• Suppose we have 2 classes of 2-dimensional data that are not

- A simple approach could be to assign to the class of the nearest mean
- Can we do better if we know about the data distribution?



observing the data given a class/parameters



Bayesian Classificaion

A probabilistic view of classification models the likelihood of

e.g., we might assume that the distribution of data given the class is Gaussian

Multi-dimensional Gaussian

• The Gaussian probability density is given by





These estimates maximise the probability of the data x given parameters m, Σ

$$\exp{-\frac{1}{2}(\mathbf{x}-\mathbf{m})^T \mathbf{\Sigma}^{-1}(\mathbf{x}-\mathbf{m})}$$

• To estimate from data (x)

$$\hat{\mathbf{m}} = \frac{1}{N} \sum_{i=1}^{N} \mathbf{x}_i$$
$$\hat{\mathbf{\Sigma}} = \frac{1}{N} \sum_{i=1}^{N} (\mathbf{x}_i - \hat{\mathbf{m}}) (\mathbf{x}_i - \hat{\mathbf{m}})^T$$

2-Class Gaussian Classifier

- Simple classification rule: choose class #1 if $p(\mathbf{x}|c_1) > p(\mathbf{x}|c_2)$
- taking -2 x ln of both sides (reverses sign) $-2\ln p(\mathbf{x}|c_1) < -2\ln p(\mathbf{x}|c_2)$
- negative log of Gaussian density $-2\ln p(\mathbf{x}) = -2\ln \frac{1}{|2\pi \mathbf{\Sigma}|^{\frac{1}{2}}}$
- decision rule becomes (class #1 if...)

$$\frac{1}{\frac{1}{2}} \exp{-\frac{1}{2}(\mathbf{x}-\mathbf{m})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\mathbf{m})}$$

 $= \ln(2\pi^d) + \ln|\mathbf{\Sigma}| + (\mathbf{x} - \mathbf{m}^T)\mathbf{\Sigma}^{-1}(\mathbf{x} - \mathbf{m})$

 $\ln \Sigma_1 + (\mathbf{x} - \mathbf{m}_1)^T \Sigma_1^{-1} (\mathbf{x} - \mathbf{m}_1) < \ln \Sigma_2 + (\mathbf{x} - \mathbf{m}_2)^T \Sigma_2^{-1} (\mathbf{x} - \mathbf{m}_2)$

2-Class Gaussian Classifier

• Suppose we've modelled our 2 classes with Gaussian distributions



$$p(\mathbf{x}|c_1) = N(\mathbf{x}; \mathbf{m}_1, \boldsymbol{\Sigma}_1)$$
$$p(\mathbf{x}|c_2) = N(\mathbf{x}; \mathbf{m}_2, \boldsymbol{\Sigma}_2)$$

- Our decision rule, class #1 if $p(\mathbf{x}|c_1) > p(\mathbf{x}|c_2)$
 - is called a maximum likelihood classifier

Incorporating Prior Knowledge

- What if red is more common than blue?
 - Weight each likelihood by prior probabilities $p(c_1), p(c_2)$
- Decision rule (MAP classifier) choose class #1 if:



 $p(\mathbf{x}|c_1)p(c_1) > p(\mathbf{x}|c_2)p(c_2)$

$$p(c_1) = 0.99$$

 $p(c_2) = 0.95$

Nearest Neighbor Classifier

space.

Ο O \mathbf{O} 0 0 OC 0 0 0

Given a new data point, assign the label of nearest training example in feature



Nearest Neighbor Classifier

space.

Given a new data point, assign the label of nearest training example in feature



Nearest Neighbour Classification

• Find nearest neighbour in training set

• Assign class to class of the nearest neighbour

 $\hat{y}(\mathbf{x}_q)$







Calculate $|\mathbf{x}_q - \mathbf{x}_i|$ for all training data

- $i_{NN} = \arg\min_{i} |\mathbf{x}_{q} \mathbf{x}_{i}|$

$$= y(\mathbf{x}_{i_N N})$$

Nearest Neighbour Classification

• We can view each image as a point in a high dimensional space



bird

Nearest Neighbour Classifier



• What is the decision boundary for a nearest-neighbour classifier?





k-Nearest Neighbor (kNN) Classifier

by majority vote.

various dimensions

minimizing probability of error

- We can gain some robustness to noise by voting over **multiple** neighbours.
- Given a **new** data point, find the k nearest training examples. Assign the label

Simple method that works well if the **distance measure** correctly weights the

For **large data sets**, as k increases kNN approaches optimality in terms of

- Identify k nearest neighbours of the query
- Assign class as most common class in set
- k-NN decision boundaries:



k = I

Good performance depends on suitable choice of k

k-NN Classifier

k-Nearest Neighbor (kNN) Classifier

1-Nearest Neighbor Classifier



15-Nearest Neighbor Classifier

kNN decision boundaries respond to local clusters where one class dominates

Figure credit: Hastie, Tibshirani & Friedman (2nd ed.)

What do nearest neighbours look like with 80 million images?











Tiny Image Recognition

• Recognition performance (categories vary in semantic level)



Nearest neighbour becomes increasingly accurate as N increases, but do we need to store a dataset of 80 million images?

yellow = 7900, red = 790,000, blue = 79,000,000
Linear Classification

Linear classification, 2-class, N-class Regularization, softmax, cross entropy • SGD, learning rate, momentum

Linear Classification

• Let's start by using 2 classes, e.g., bird and plane • Apply labels (y) to training set:





y = +1





• Use a linear model to regress y from x



Linear Classification



Let's start by using 2 classe
 Apply labels (y) to training



Use a linear model to regress y from x

s, e.g., bird and plane set:



2-class Linear Classification

• Separating hyperplane, projection to a line defined by w



 $\hat{y} = \operatorname{sign} h = \operatorname{sign} \mathbf{w}^T \mathbf{x}_q$

N-class Linear Classification



N-class Linear Classification



bird

• We could regress directly to integer class id, y = {0,1,2,3...9}

classifiers



• A better solution is to regress to one-hot targets = I vs all

Stack into matrix form







• Transpose





Solve regression problem by Least Squares

$\begin{bmatrix} x_{11} & x_{12} & x_{13} & \dots \\ x_{21} & x_{22} & x_{23} & \dots \\ x_{31} & x_{32} & x_{33} & \dots \\ & & & & & & \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} & x_{13} & \dots \\ 0 & 0 & 0 & \dots \\ 0 & 0 & 0 & 1 & \dots \\ 0 & 0 & 0 & 1 & \dots \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} & x_{13} & \dots \\ 0 & 0 & 0 & 1 & \dots \\ 0 & 0 & 0 & 1 & \dots \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} & x_{13} & \dots \\ 0 & 0 & 0 & 1 & \dots \\ 0 & 0 & 0 & 1 & \dots \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} & x_{13} & \dots \\ 0 & 0 & 0 & 1 & \dots \\ 0 & 0 & 0 & 1 & \dots \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} & x_{13} & \dots \\ 0 & 0 & 0 & 1 & \dots \\ 0 & 0 & 0 & 1 & \dots \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} & x_{13} & \dots \\ 0 & 0 & 0 & 1 & \dots \\ 0 & 0 & 0 & 1 & \dots \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} & x_{13} & \dots \\ 0 & 0 & 0 & 1 & \dots \\ 0 & 0 & 0 & 1 & \dots \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} & x_{13} & \dots \\ 0 & 0 & 0 & 1 & \dots \\ 0 & 0 & 0 & 1 & \dots \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} & x_{13} & \dots \\ 0 & 0 & 0 & 1 & \dots \\ 0 & 0 & 0 & 1 & \dots \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} & x_{13} & \dots \\ 0 & 0 & 0 & 1 & \dots \\ 0 & 0 & 0 & 1 & \dots \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} & x_{13} & \dots \\ 0 & 0 & 0 & 1 & \dots \\ 0 & 0 & 0 & 1 & \dots \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} & x_{13} & \dots \\ 0 & 0 & 0 & 1 & \dots \\ 0 & 0 & 0 & 1 & \dots \\ 0 & 0 & 0 & 1 & \dots \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} & x_{13} & \dots \\ 0 & 0 & 0 & 1 & \dots \\ 0 & 0 & 0 & 1 & \dots \\ 0 & 0 & 0 & 1 & \dots \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} & x_{13} & \dots \\ 0 & 0 & 0 & 1 & \dots \\ 0 & 0 & 0 & 1 & \dots \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} & x_{13} & \dots \\ 0 & 0 & 0 & 1 & \dots \\ 0 & 0 & 0 & 1 & \dots \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} & x_{13} & \dots \\ 0 & 0 & 0 & 1 & \dots \\ 0 & 0 & 0 & 1 & \dots \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} & x_{13} & \dots \\ 0 & 0 & 0 & 1 & \dots \\ 0 & 0 & 0 & 1 & \dots \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} & x_{13} & \dots \\ 0 & 0 & 0 & 1 & \dots \\ 0 & 0 & 0 & 1 & \dots \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} & x_{13} & \dots \\ 0 & 0 & 0 & 1 & \dots \\ 0 & 0 & 0 & 0 & 1 & \dots \\ 0 & 0 & 0 & 0 & 1 & \dots \\ 0 & 0 & 0 & 0 & 0 & \dots \\ 0 & 0 & 0 & 0 & 0 & \dots \\ 0 & 0 & 0 & 0 & 0 & \dots \\ 0 & 0 & 0 & 0 & 0 & \dots \\ 0 & 0 & 0 & 0 & 0 & \dots \\ 0 & 0 & 0 & 0 & 0 & \dots \\ 0 & 0 & 0 & 0 & 0 & \dots \\ 0 & 0 & 0 & 0 & 0 & \dots \\ 0 & 0 & 0 & 0 & \dots \\ 0 & 0 & 0 & 0 & 0 & \dots \\ 0 & 0 & 0 & 0 & 0 & \dots \\ 0 & 0 & 0 & 0 & \dots \\ 0 & 0 & 0 & 0 & \dots \\ 0 & 0 & 0 & 0 & \dots \\ 0 & 0 & 0 & 0 & \dots \\ 0 & 0 & 0 & 0 & \dots \\ 0 & 0 & 0 & 0 & \dots \\ 0 & 0 & 0 & 0 & \dots \\ 0 & 0 & 0 & 0 & \dots \\ 0 & 0 & 0 & 0 & \dots \\ 0 & 0 & 0 & 0 & \dots \\ 0 &$

 $\mathbf{XW} = \mathbf{T}$

N-class Linear Classification

• One hot regression = I vs all classifiers



bird

cat

• Visualise class templates for the least squares solution











airplane automobile

bird

cat

deer



What is happening here?



- Classifier accuracy = 35% (not bad, c.f., nearest mean = 27%)





• Consider fitting a polynomial to some data by linear regression

• Multiple data points (y_i, x_i)

In matrix form



Least Squares (if overconstrained)

 $y_1 = a_0 + a_1 x_1 + a_2 x_1^2 + a_3 x_1^3$ $y_2 = a_0 + a_1 x_2 + a_2 x_2^2 + a_3 x_2^3$ $y_3 = a_0 + a_1 x_3 + a_2 x_3^2 + a_3 x_3^3$

 $\mathbf{v} = \mathbf{M}\mathbf{a}$

• • •

• Solve linear system by Gaussian elimination (if square) or

• Fit Nth order polynomial by least squares



• Overfitting

Validation

• Fit the model to a subset of data, and evaluate the fit on a held out validation set



Validation

minimum for the best model order



• Training error always decreases, but validation error has a

• For large N, coefficients become HUGE!

	N=1	N=2	N=4	N=10
a_0	0.90	2.03	-2.88	48.50
a_1		-1.54	29.76	-1294.90
a_2			-57.43	14891.41
a_3			31.86	-95161.10
a_4				367736.84
a_5				-885436.68
a_6				1331063.41
a_7				-1212056.89
a_8				610930.32
a_9				-130727.39

Regularization

• L2 penalty on polynomial coefficients



Regularized Linear Regression

10th order polynomial, prior on the coefficients weight $|\lambda|$



• Over-smoothing...

• Test error vs lambda



- Training error always decreases as lambda is reduced



• Test error reaches a minimum, then increases \Rightarrow overfitting

Regularized Classification

• Add regularization to CIFAR10 linear classifier



• Row I = overfitting, Row 3 = oversmoothing?

Non-Linear Optimisation

- With a linear predictor and L2 loss, we have a closed form solution for model weights W
- How about this (non-linear) function

• Previously (e.g., bundle adjustment), we locally linearised the error function and iteratively solved linear problems

$$e = \sum_{i} |\mathbf{h}_{i} - \mathbf{t}_{i}|^{2} \approx |\mathbf{J}\Delta\mathbf{W} + \mathbf{r}|^{2}$$
$$\Delta\mathbf{W} = -(\mathbf{J}^{T}\mathbf{J})^{-1}\mathbf{J}^{T}\mathbf{r}$$

Does this look like a promising approach?



 $\mathbf{h} = \mathbf{W}_2 \max(0, \mathbf{W}_1 \mathbf{x})$



Vanilla Gradient Descent

Vanilla Gradient Descent

while True:

weights_grad = evaluate_gradient(loss_fun, data, weights) weights += - step_size * weights_grad # perform parameter update







Problem with vanilla GD

What if loss changes quickly in one direction and slowly in another? What does gradient descent do? Very slow progress along shallow dimension, jitter along steep direction



singular value of the Hessian matrix is large

Based on slides for Stanford cs231n by Li, Jonson, and Young. Modified and reused with permission

Loss function has high **condition number**: ratio of largest to smallest





Problem with vanilla GD

What if loss changes quickly in one direction and slowly in another? What does gradient descent do? Very slow progress along shallow dimension, jitter along steep direction



singular value of the Hessian matrix is large

Based on slides for Stanford cs231n by Li, Jonson, and Young. Modified and reused with permission

Loss function has high **condition number**: ratio of largest to smallest





Problem with vanilla GD

What if loss changes quickly in one direction and slowly in another? What does gradient descent do? Very slow progress along shallow dimension, jitter along steep direction



singular value of the Hessian matrix is large

Based on slides for Stanford cs231n by Li, Jonson, and Young. Modified and reused with permission

Loss function has high **condition number**: ratio of largest to smallest





Optimization: problem with SGD

What if the loss function has a local minima or saddle point?







Optimization: problem with SGD

What if the loss function has a local minima or saddle point?









Optimization: problem with SGD

What if the loss function has a local minima or saddle point?









Optimization: problem with SGD

What if the loss function has a local minima or saddle point?

Zero gradient, gradient descent gets stuck

Based on slides for Stanford cs231n by Li, Jonson, and Young. Modified and reused with permission









Optimization: problem with SGD

What if the loss function has a **local minima** or **saddle point**?

Saddle points much more common in high dimension

Dauphin et al, "Identifying and attacking the saddle point problem in high-dimensional non-convex optimization", NIPS 2014 67 Based on slides for <u>Stanford cs231n</u> by Li, Jonson, and Young. Modified and reused with permission





Optimization: problem with SGD

What if the loss function has a local minima or saddle point?

Or not?

"We show that gradient descent converges to a local minimizer, almost surely with random initialization. This is proved by applying the Stable Manifold Theorem from dynamical systems theory."

Lee et al, "Gradient Descent Only Converges to Minimizers", JLMR Workshop and Conference Proceedings, 2016 Dauphin et al, "Identifying and attacking the saddle point problem in high-dimensional non-convex optimization", NIPS 2014 68 Based on slides for Stanford cs231n by Li, Jonson, and Young. Modified and reused with permission





THE UNIVERSITY

OF BRITISH COLUMBIA

Stochastic gradient descent

Our gradients come from minibatches so they can be noisy!

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

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

Q: How would you remove the noise?







SGD + Momentum

SGD

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

while True: $dx = compute_gradient(x)$ x += learning_rate * dx

Rho gives "friction"; typically rho=0.9 or 0.99

Based on slides for Stanford cs231n by Li, Jonson, and Young. Modified and reused with permission



Build up "velocity" as a running mean of gradients



SGD + Momentum



71 Based on slides for <u>Stanford cs231n</u> by Li, Jonson, and Young. Modified and reused with permission

Gradient Noise





SGD + Momentum

Momentum update:



Nesterov, "A method of solving a convex programming problem with convergence rate O(1/k²)", 1983 Nesterov, "Introductory lectures on convex optimization: a basic course", 2004 Sutskever et al, "On the importance of initialization and momentum in deel learning", ICML 2013

72 Based on slides for <u>Stanford cs231n</u> by Li, Jonson, and Young. Modified and reused with permission


SGD + Momentum

Momentum update:



Nesterov, "A method of solving a convex programming problem with convergence rate O(1/k^2)", 1983 Nesterov, "Introductory lectures on convex optimization: a basic course", 2004 Sutskever et al, "On the importance of initialization and momentum in deel learning", ICML 2013

 $\frac{73}{1000}$ Based on slides for $\frac{\text{Stanford cs}231n}{1000}$ by Li, Jonson, and Young. Modified and reused with permission

Nesterov Momentum



ng", ICML 2013 73 ermission



Nesterov Momentum



74 Based on slides for <u>Stanford cs231n</u> by Li, Jonson, and Young. Modified and reused with permission



RMSProp





Q: What happens with RMSProp?

Tieleman and Hinton, 2012

Based on slides for Stanford cs231n by Li, Jonson, and Young. Modified and reused with permission

grad_squared = decay_rate * grad_squared + (1 - decay_rate) * dx * dx







RMSProp

76 Based on slides for Stanford cs231n by Li, Jonson, and Young. Modified and reused with permission

SGD SGD+Momentum RMSProp





Adam (almost)

first_moment = 0 second_moment = 0 while True: $dx = compute_gradient(x)$ first_moment = beta1 * first_moment + (1 - beta1) * dx second_moment = beta2 * second_moment + (1 - beta2) * dx * dx

Q: What happens at first the timestep?

Kingma and Ba, "Adam: A method for stochastic optimization", ICLR 2015

Based on slides for <u>Stanford cs231n</u> by Li, Jonson, and Young. Modified and reused with permission



RMSProp with momentum





Adam (full form)



Bias correction for the fact that first and second moment estimates start at zero

Kingma and Ba, "Adam: A method for stochastic optimization", ICLR 2015

Based on slides for Stanford cs231n by Li, Jonson, and Young. Modified and reused with permission

Adam with beta 1 = 0.9, beta2 = 0.999, and learning_rate = 1e-4is a great starting point for many models!









Learning rate: hyperparameter



