# Machine Learning

# Computer Vision (UW EE/CSE 576)



Richard Szeliski Facebook & UW Lecture 6 – Apr 16, 2020



# Class calendar

Date	Lecture	Reading	Homework
March 31	Introduction	Szeliski, Chapter 1	
April 2	Human Vision, Color Spaces and Transforms	Brown, M. S. (2019). ICCV 2019 tutorial on understanding color and the in-camera image processing pipeline for computer vision.	HW1 assigned
April 7	Image coordinates, resizing	Szeliski, Chapter 2.1, 3.6	
April 9	Filters and convolutions	Szeliski, Chapter 3	HW1 due, HW2 assigned
April 14	Interpolation and Optimization	Szeliski, Chapter 4	
April 16	Machine Learning	Szeliski, Chapter 5.1-5.2	

# Readings

## Chapter 5

## **Deep Learning**

5.1	Superv	vised learning
	5.1.1	Nearest neighbors
	5.1.2	Bayesian classification
	5.1.3	Logistic regression
	5.1.4	Support vector machines
	5.1.5	Decision trees and forests
5.2	Unsup	ervised learning
	5.2.1	Clustering
	5.2.2	K-means and mixtures of Gaussians
	5.2.3	Principal component analysis
	5.2.4	Semi-supervised learning
5.3	Deep n	neural networks

# References







# Traditional machine learning

- Why learning?
- Nearest neighbors
- Bayesian classification
- Logistic regression
- Support vector machines
- Clustering
- Principal component analysis





# ... finish off previous lecture ...

# Why Bayesian (probabilistic) modeling?

- Merge uncertain measurement in an optimal way
- Estimate the uncertainty in the final answer
- Build in (quantitative) prior assumptions

## **Bayesian modeling and inference**

<b>B</b> .1	Estimation theory						
	B.1.1 Likelihood for multivariate Gaussian noise						
B.2	Maximum likelihood estimation and least squares						
<b>B</b> .3	Robust statistics						
<b>B</b> .4	Prior models and Bayesian inference						
B.5	Markov random fields						
	B.5.1 Gradient descent and simulated annealing						
	B.5.2 Dynamic programming						
	B.5.3 Belief propagation						
	B.5.4 Graph cuts						
	B.5.5 Linear programming						
<b>B</b> .6	Uncertainty estimation (error analysis)						

Appendix B

# Example: merge GPS readings

- *Measurement 1*: 10.1 ± 0.1667 (1/6) Gaussian noise
- *Measurement 2*: 10.2 ± 0.1250 (1/8) Gaussian noise
- What is the optimal combined measurement?



Exercise for next lecture (send me your answer on Slack)

## Answer: use inverse variances (Appendix B)

$$\begin{aligned} \mathbf{y}_{i} &= \mathbf{f}_{i}(\mathbf{x}) + \mathbf{n}_{i} \\ L &= p(\mathbf{y}|\mathbf{x}) = \prod_{i} p(\mathbf{y}_{i}|\mathbf{x}) = \prod_{i} p(\mathbf{y}_{i}|\mathbf{f}_{i}(\mathbf{x})) \\ &= \prod_{i} |2\pi \mathbf{\Sigma}_{i}|^{-1/2} \exp\left(-\frac{1}{2} ||\mathbf{y}_{i} - \mathbf{f}_{i}(\mathbf{x})||_{\mathbf{\Sigma}_{i}^{-1}}^{2}\right) \\ E &= -\log L = \frac{1}{2} \sum_{i} (\mathbf{y}_{i} - \mathbf{f}_{i}(\mathbf{x}))^{T} \mathbf{\Sigma}_{i}^{-1} (\mathbf{y}_{i} - \mathbf{f}_{i}(\mathbf{x})) + k \\ \mathbf{f}_{i}(\mathbf{x}) &= \mathbf{H}_{i} \mathbf{x}. \end{aligned}$$

$$\mathbf{C}_{i} = \mathbf{\Sigma}_{i}^{-1}$$



$$E = \sum_{i} \|\tilde{\mathbf{y}}_{i} - \mathbf{H}_{i}\mathbf{x}\|_{\boldsymbol{\Sigma}_{i}^{-1}} = \sum_{i} (\tilde{\mathbf{y}}_{i} - \mathbf{H}_{i}\mathbf{x})^{T} \mathbf{C}_{i} (\tilde{\mathbf{y}}_{i} - \mathbf{H}_{i}\mathbf{x}), \qquad (B.16) \qquad \mathbf{H}_{i} = \mathbf{I}$$
$$\mathbf{x} = \left(\sum_{i} \mathbf{C}_{i}\right)^{-1} \left(\sum_{i} \mathbf{C}_{i} \tilde{\mathbf{y}}_{i}\right)$$

# Answer: merge GPS readings

- *Measurement 1*: 10.1 ± 0.1667 (1/6) Gaussian noise
- Measurement 2: 10.2 ± 0.1250 (1/8) Gaussian noise
- What is the optimal combined measurement?



					wght.
<u>Measurement</u>	<u>value</u>	<u>std.dev.</u>	<u>variance</u>	<u>inv.var.</u>	<u>contr.</u>
1	10.1	0.166667	0.02778	36	363.6
2	10.2	0.125	0.01563	64	652.8
summed	10.16	0.1	0.01	100	1016



Excellent work! What's your academic background?

Nathan Hatch 4:32 PM I think the answer you're looking for is 10.164 \pm 0.1



Rick Szeliski 4:37 PM

Nathan Hatch 4:54 PM

Robotics and machine learning.

# Why a Bayesian formulation?

- Wider range of probability distributions
- Learn distributions from data
- Wider range of inference algorithms [Kappes, Andres, *et al.*, IJCV 2015]





# **Energy minimization**



(a) initial labeling (b) standard move (c)  $\alpha$ - $\beta$ -swap (d)  $\alpha$ -expansion

**Figure 4.15** *Multi-level graph optimization from (Boykov, Veksler, and Zabih 2001)*  $\bigcirc$  2001 *IEEE: (a) initial problem configuration; (b) the standard move only changes one pixel; (c) the*  $\alpha$ - $\beta$ -swap optimally exchanges all  $\alpha$  and  $\beta$ -labeled pixels; (d) the  $\alpha$ -expansion move optimally selects among current pixel values and the  $\alpha$  label.

**Richard Szeliski** 

# **Digital Photomontage**



**Figure 4.17** An unordered label MRF (Agarwala, Dontcheva, Agrawala et al. 2004) © 2004 ACM: Strokes in each of the source images on the left are used as constraints on an MRF optimization, which is solved using graph cuts. The resulting multi-valued label field is shown as a color overlay in the middle image, and the final composite is shown on the right.

## Markov Random Fields

According to Bayes' Rule (Appendix B.4), the *posterior* distribution for a given set of measurements  $\mathbf{y}$ ,  $p(\mathbf{y}|\mathbf{x})$ , combined with a prior  $p(\mathbf{x})$  over the unknowns  $\mathbf{x}$ , is given by

$$p(\mathbf{x}|\mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{x})p(\mathbf{x})}{p(\mathbf{y})},$$
(4.33)

where  $p(\mathbf{y}) = \int_{\mathbf{x}} p(\mathbf{y}|\mathbf{x})p(\mathbf{x})$  is a normalizing constant used to make the  $p(\mathbf{x}|\mathbf{y})$  distribution *proper* (integrate to 1). Taking the negative logarithm of both sides of (4.33), we get

$$-\log p(\mathbf{x}|\mathbf{y}) = -\log p(\mathbf{y}|\mathbf{x}) - \log p(\mathbf{x}) + C, \qquad (4.34)$$

which is the *negative posterior log likelihood*.

To find the most likely (*maximum a posteriori* or MAP) solution x given some measurements y, we simply minimize this negative log likelihood, which can also be thought of as an *energy*,

$$E(\mathbf{x}, \mathbf{y}) = E_{\mathrm{D}}(\mathbf{x}, \mathbf{y}) + E_{\mathrm{P}}(\mathbf{x}).$$
(4.35)

(We drop the constant C because its value does not matter during energy minimization.) The first term  $E_D(\mathbf{x}, \mathbf{y})$  is the *data energy* or *data penalty*; it measures the negative log likelihood that the data were observed given the unknown state  $\mathbf{x}$ . The second term  $E_P(\mathbf{x})$  is the *prior energy*; it plays a role analogous to the smoothness energy in regularization. Note Richard Szeliski UW CSE 576 - Machine Learning



# **Conditional Random Fields**

 Interaction potentials depend on guide image (sound familiar?)

Since the smoothness term now depends on the data, Bayes' Rule (4.44) no longer applies. Instead, we use a direct model for the posterior distribution  $p(\mathbf{x}|\mathbf{y})$ , whose negative log likelihood can be written as

$$E(\mathbf{x}|\mathbf{y}) = E_{\mathrm{D}}(\mathbf{x}, \mathbf{y}) + E_{\mathrm{S}}(\mathbf{x}, \mathbf{y})$$
$$= \sum_{p} V_{p}(x_{p}, \mathbf{y}) + \sum_{(p,q) \in \mathcal{N}} V_{p,q}(x_{p}, x_{q}, \mathbf{y}), \qquad (4.45)$$

UW CSE 576 - Machine Learning

using the notation introduced in (4.43). The resulting probability distribution is called a *conditional random field* (CRF) and was first introduced to the computer vision field by Kumar and Hebert (2003), based on earlier work in text modeling by Lafferty, McCallum, and Pereira (2001).

Figure 4.18 shows a graphical model where the smoothness terms depend on the data values. In this particular model, each smoothness term depends only on its adjacent pair of data values, i.e., terms are of the form  $V_{p,q}(x_p, x_q, y_p, y_q)$  in (4.45).

f(i,j+1) f(i+1,j+1) f(i,j) f(i,j) f(i+1,j)

## Dense CRF [Krähenbühl and Koltun 2011]

• Use a wide color-based support region (like bilateral solver)



# Traditional machine learning (today's lecture)

- Why learning?
- Nearest neighbors
- Bayesian classification
- Logistic regression
- Support vector machines
- Clustering
- Principal component analysis





# I'll be borrowing slides from

UNIVERSITY OF MICHIGAN

### EECS 498-007 / 598-005 Deep Learning for Computer Vision Fall 2019

#### Course Description

Computer Vision has become ubiquitous in our society, with applications in search, image understanding, apps, mapping, medicine, drones, and self-driving cars. Core to many of these applications are visual recognition tasks such as image classification and object detection. Recent developments in neural network approaches have greatly advanced the performance of these state-of-the-art visual recognition systems. This course is a deep dive into details of neural-network based deep learning methods for computer vision. During this course, students will learn to implement, train and debug their own neural networks and gain a detailed understanding of cutting-edge research in computer vision. We will cover learning algorithms, neural network architectures, and practical engineering tricks for training and fine-tuning networks for visual recognition tasks.

#### Instructor Graduate Student Instructors





EECS 498-007 / 598-005 Deep Learning for Computer Vision Fall 2019

# Lecture 2: Image Classification

Justin Johnson



## Image Classification

### **Input**: image



This image by Nikita is licensed under CC-BY 2.0

**Output**: Assign image to one of a fixed set of categories

cat bird deer dog truck

### Justin Johnson







#### Lecture 2 - 21

### Justin Johnson

## Challenges: Viewpoint Variation



<u>This image</u> by <u>Nikita</u> is licensed under <u>CC-BY 2.0</u>

### Justin Johnson



## Challenges: Intraclass Variation



This image is CC0 1.0 public domain

### Justin Johnson





## **Challenges**: Fine-Grained Categories

### Maine Coon

### Ragdoll

### American Shorthair



This image is free for for use under the Pixabay License

This image is CC0 public domain

<u>This image is CC0 public domain</u>

### Justin Johnson



## Challenges: Background Clutter



This image is CC0 1.0 public domain

This image is CC0 1.0 public domain

### Justin Johnson



## **Challenges**: Illumination Changes



This image is CC0 1.0 public domain

### Justin Johnson



## **Challenges**: Deformation



This image by Umberto Salvagnin is licensed under CC-BY 2.0

licensed under CC-BY 2.0

This image by Tom Thai is licensed under CC-BY 2.0

### Justin Johnson

### Lecture 2 - 27

under <u>CC-BY 2.0</u>



## Challenges: Occlusion



This image is CC0 1.0 public domain

This image is CC0 1.0 public domain

This image by jonsson is licensed under <u>CC-BY 2.0</u>

### Justin Johnson



### **Example: Object Detection**





This image is free to use under the Pexels license

### Justin Johnson



**Example: Object Detection** 



This image is free to use under the Pexels license

## Background

Horse

Person

Car Truck

### Justin Johnson



**Example: Object Detection** 



Background Horse Person Car

Truck

This image is free to use under the Pexels license

### Justin Johnson



**Example: Image Captioning** 



This mage is nee to use under the Pexels lic

Justin Johnson



Example: Image Captioning



riding What word cat to say next? horse man when Caption: Man riding

This image is free to use under the Pexels license

### Justin Johnson



Example: Image Captioning



This image is free to use under the Pexels license

riding What word to say next?
 horse
 man What word to say next?
 Man riding horse

### Justin Johnson



**Example: Image Captioning** 



This image is free to use under the Pexels license

riding cat horse man when . . .

## <STOP>

## What word to say next?

Caption: Man riding horse

### Justin Johnson



## An Image Classifier

def classify\_image(image):
 # Some magic here?
 return class\_label

Unlike e.g. sorting a list of numbers,

**no obvious way** to hard-code the algorithm for recognizing a cat, or other classes.

Justin Johnson


## You could try ...



#### John Canny, "A Computational Approach to Edge Detection", IEEE TPAMI 1986





### Machine Learning: Data-Driven Approach

- 1. Collect a dataset of images and labels
- 2. Use Machine Learning to train a classifier
- 3. Evaluate the classifier on new images

def train(images, labels): # Machine learning! return model

def predict(model, test\_images): # Use model to predict labels return test\_labels

airplane automobile bird cat deer

### **Example training set**

Justin Johnson



# Supervised learning

• Goal: provide best output predictions for novel inputs



# Supervised learning



- Goal: provide best output predictions for novel inputs
- How can we predict future performance?
- Placeholder answer: minimize *empirical risk*

$$E_{\text{Risk}}(\mathbf{w}) = \frac{1}{N} \sum L(\mathbf{y}_i, \mathbf{f}(\mathbf{x}_i; \mathbf{w})).$$
(5.1)

The loss function L measures the "cost" of predicting an output  $f(x_i; w)$  for input  $x_i$  and model parameters w when the corresponding target is  $y_i$ .

- Where have we seen this before?
- What are potential *models* and *loss functions*?

# Traditional and deep learning



# Traditional machine learning

- Why learning?
- Nearest neighbors
- Bayesian classification
- Logistic regression
- Support vector machines
- Clustering
- Principal component analysis







EECS 498-007 / 598-005 Deep Learning for Computer Vision Fall 2019

# Nearest neighbors



Justin Johnson



### First classifier: Nearest Neighbor

def train(images, labels):
 # Machine learning!
 return model

Memorize all data and labels

def predict(model, test\_images):
 # Use model to predict labels
 return test\_labels

→ Predict the label of
 → the most similar
 training image

### Distance Metric to compare images

L1 distance:

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



```
import numpy as np
```

```
class NearestNeighbor:
 def init (self):
   pass
 def train(self, X, y):
    """ X is N x D where each row is an example. Y is 1-dimension of size N """
   # the nearest neighbor classifier simply remembers all the training data
    self.Xtr = X
    self.ytr = y
 def predict(self, X):
    """ X is N x D where each row is an example we wish to predict label for """
    num test = X.shape[0]
   # lets make sure that the output type matches the input type
    Ypred = np.zeros(num test, dtype = self.ytr.dtype)
    # loop over all test rows
    for i in xrange(num test):
      # find the nearest training image to the i'th test image
      # using the L1 distance (sum of absolute value differences)
      distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
```

min\_index = np.argmin(distances) # get the index with smallest distance
Ypred[i] = self.ytr[min\_index] # predict the label of the nearest example

return Ypred

### Nearest Neighbor Classifier

#### Justin Johnson

#### Lecture 2 - 46

#### Fall 2019

```
import numpy as np
```

class NearestNeighbor: def \_\_init\_\_(self): pass

# def train(self, X, y): """ X is N x D where each row is an example. Y is 1-dimension of size N """ # the nearest neighbor classifier simply remembers all the training data self.Xtr = X self.ytr = y

def predict(self, X):
 """ X is N x D where each row is an example we wish to predict label for """
 num\_test = X.shape[0]
 # lets make sure that the output type matches the input type
 Ypred = np.zeros(num\_test, dtype = self.ytr.dtype)

```
# loop over all test rows
for i in xrange(num_test):
    # find the nearest training image to the i'th test image
    # using the L1 distance (sum of absolute value differences)
    distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
    min_index = np.argmin(distances) # get the index with smallest distance
    Ypred[i] = self.ytr[min_index] # predict the label of the nearest example
```

#### return Ypred

### Nearest Neighbor Classifier

### Memorize training data

#### Justin Johnson



impo	rt	numpy	as	np
------	----	-------	----	----

```
class NearestNeighbor:
    def __init__(self):
        pass
```

def train(self, X, y):
 """ X is N x D where each row is an example. Y is 1-dimension of size N """
 # the nearest neighbor classifier simply remembers all the training data
 self.Xtr = X
 self.ytr = y

def predict(self, X):
 """ X is N x D where each row is an example we wish to predict label for """
 num\_test = X.shape[0]
 # lets make sure that the output type matches the input type
 Ypred = np.zeros(num\_test, dtype = self.ytr.dtype)

# loop over all test rows
for i in xrange(num\_test):
 # find the nearest training image to the i'th test image
 # using the L1 distance (sum of absolute value differences)
 distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
 min\_index = np.argmin(distances) # get the index with smallest distance
 Ypred[i] = self.ytr[min\_index] # predict the label of the nearest example

For each test image: Find nearest training image Return label of nearest image

return Ypred

### Nearest Neighbor Classifier

#### Justin Johnson



```
import numpy as np
```

```
class NearestNeighbor:
    def __init__(self):
        pass
```

def train(self, X, y):
 """ X is N x D where each row is an example. Y is 1-dimension of size N """
 # the nearest neighbor classifier simply remembers all the training data
 self.Xtr = X
 self.ytr = y

def predict(self, X):
 """ X is N x D where each row is an example we wish to predict label for """
 num\_test = X.shape[0]
 # lets make sure that the output type matches the input type
 Ypred = np.zeros(num\_test, dtype = self.ytr.dtype)

```
# loop over all test rows
for i in xrange(num_test):
    # find the nearest training image to the i'th test image
    # using the L1 distance (sum of absolute value differences)
    distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
    min_index = np.argmin(distances) # get the index with smallest distance
    Ypred[i] = self.ytr[min_index] # predict the label of the nearest example
```

return Ypred

### Nearest Neighbor Classifier

**Q**: With N examples, how fast is training?

#### Justin Johnson



```
import numpy as np
```

```
class NearestNeighbor:
    def __init__(self):
        pass
```

def train(self, X, y):
 """ X is N x D where each row is an example. Y is 1-dimension of size N """
 # the nearest neighbor classifier simply remembers all the training data
 self.Xtr = X
 self.ytr = y

def predict(self, X):
 """ X is N x D where each row is an example we wish to predict label for """
 num\_test = X.shape[0]
 # lets make sure that the output type matches the input type
 Ypred = np.zeros(num\_test, dtype = self.ytr.dtype)

```
# loop over all test rows
for i in xrange(num_test):
    # find the nearest training image to the i'th test image
    # using the L1 distance (sum of absolute value differences)
    distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
    min_index = np.argmin(distances) # get the index with smallest distance
    Ypred[i] = self.ytr[min_index] # predict the label of the nearest example
```

return Ypred

### Nearest Neighbor Classifier

Q: With N examples,how fast is training?A: O(1)

#### Justin Johnson



class NearestNeighbor: def \_\_init\_\_(self): pass

def train(self, X, y):
 """ X is N x D where each row is an example. Y is 1-dimension of size N """
 # the nearest neighbor classifier simply remembers all the training data
 self.Xtr = X
 self.ytr = y

def predict(self, X):
 """ X is N x D where each row is an example we wish to predict label for """
 num\_test = X.shape[0]
 # lets make sure that the output type matches the input type
 Ypred = np.zeros(num\_test, dtype = self.ytr.dtype)

```
# loop over all test rows
for i in xrange(num_test):
    # find the nearest training image to the i'th test image
    # using the L1 distance (sum of absolute value differences)
    distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
    min_index = np.argmin(distances) # get the index with smallest distance
    Ypred[i] = self.ytr[min_index] # predict the label of the nearest example
```

Nearest Neighbor Classifier

Q: With N examples,how fast is training?A: O(1)

**Q**: With N examples, how fast is testing?

return Ypred

#### Justin Johnson



class NearestNeighbor: def \_\_init\_\_(self): pass

def train(self, X, y):
 """ X is N x D where each row is an example. Y is 1-dimension of size N """
 # the nearest neighbor classifier simply remembers all the training data
 self.Xtr = X
 self.ytr = y

def predict(self, X):
 """ X is N x D where each row is an example we wish to predict label for """
 num\_test = X.shape[0]
 # lets make sure that the output type matches the input type
 Ypred = np.zeros(num\_test, dtype = self.ytr.dtype)

```
# loop over all test rows
for i in xrange(num_test):
    # find the nearest training image to the i'th test image
    # using the L1 distance (sum of absolute value differences)
    distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
    min_index = np.argmin(distances) # get the index with smallest distance
    Ypred[i] = self.ytr[min_index] # predict the label of the nearest example
```

Nearest Neighbor Classifier

Q: With N examples,how fast is training?A: O(1)

Q: With N examples,how fast is testing?A: O(N)

return Ypred

#### Justin Johnson



class NearestNeighbor: def \_\_init\_\_(self): pass

def train(self, X, y):
 """ X is N x D where each row is an example. Y is 1-dimension of size N """
 # the nearest neighbor classifier simply remembers all the training data
 self.Xtr = X
 self.ytr = y

def predict(self, X):
 """ X is N x D where each row is an example we wish to predict label for """
 num\_test = X.shape[0]
 # lets make sure that the output type matches the input type
 Ypred = np.zeros(num\_test, dtype = self.ytr.dtype)

# loop over all test rows
for i in xrange(num\_test):
 # find the nearest training image to the i'th test image
 # using the L1 distance (sum of absolute value differences)
 distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
 min\_index = np.argmin(distances) # get the index with smallest distance
 Ypred[i] = self.ytr[min\_index] # predict the label of the nearest example

return Ypred

Nearest Neighbor Classifier

Q: With N examples,how fast is training?A: O(1)

Q: With N examples,how fast is testing?A: O(N)

This is **bad**: We can afford slow training, but we need fast testing!

#### Justin Johnson

#### Lecture 2 - 53

#### Fall 2019

```
class NearestNeighbor:
    def __init__(self):
        pass
```

def train(self, X, y):
 """ X is N x D where each row is an example. Y is 1-dimension of size N """
 # the nearest neighbor classifier simply remembers all the training data
 self.Xtr = X
 self.ytr = y

def predict(self, X):
 """ X is N x D where each row is an example we wish to predict label for """
 num\_test = X.shape[0]
 # lets make sure that the output type matches the input type
 Ypred = np.zeros(num\_test, dtype = self.ytr.dtype)

```
# loop over all test rows
for i in xrange(num_test):
    # find the nearest training image to the i'th test image
    # using the L1 distance (sum of absolute value differences)
    distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
    min_index = np.argmin(distances) # get the index with smallest distance
    Ypred[i] = self.ytr[min_index] # predict the label of the nearest example
```

Nearest Neighbor Classifier

There are many methods for fast / approximate nearest neighbors; e.g. see

https://github.com/facebookresearch/faiss

#### return Ypred

#### Justin Johnson



## What does this look like?



Justin Johnson



## What does this look like?



Justin Johnson





	lus	ti	'n	0	h	n	S	0	n	
-			•••	 $\mathbf{}$	•••	•	<u> </u>	$\sim$		



**X**<sub>1</sub>

Nearest neighbors in two dimensions





Lecture 2 - 58



**X**<sub>0</sub>

Nearest neighbors in two dimensions

Points are training examples; colors give training labels



#### Justin Johnson

Lecture 2 - 59



**X**<sub>0</sub>

Nearest neighbors in two dimensions

Points are training examples; colors give training labels

Background colors give the category a test point would be assigned



Lecture 2 - 60

**X**<sub>0</sub>

Fall 2019

Nearest neighbors in two dimensions

Points are training examples; colors give training labels

Background colors give the category a test point would be assigned



**Decision boundary** is the boundary between two classification regions

**X**<sub>0</sub>

Nearest neighbors in two dimensions Points are training examples; colors give training labels

Background colors give the category a test point would be assigned



**Decision boundary** is the boundary between two classification regions

Decision boundaries can be noisy; affected by outliers

Fall 2019

**X**<sub>0</sub>

Nearest neighbors in two dimensions Points are training examples; colors give training labels

Background colors give the category a test point would be assigned



Decision boundary is the boundary between two classification regions

Decision boundaries can be noisy; affected by outliers

How to smooth out decision boundaries? Use more neighbors!

#### Justin Johnson

#### Lecture 2 - 63

Fall 2019

 $X_0$ 

Instead of copying label from nearest neighbor, take **majority vote** from K closest points

K = 3

K = 1



•

#### Justin Johnson



Using more neighbors helps smooth out rough decision boundaries

K = 1





#### Justin Johnson



Using more neighbors helps reduce the effect of outliers

K = 1



K = 3



#### Justin Johnson



When K > 1 there can be ties between classes. Need to break somehow!

K = 1





#### Justin Johnson



K-Nearest Neighbors: Distance Metric

L1 (Manhattan) distance

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



L2 (Euclidean) distance

$$d_2(I_1,I_2) = \sqrt{\sum_p \left(I_1^p - I_2^p
ight)^2}$$



Justin Johnson



## K-Nearest Neighbors: Distance Metric

L1 (Manhattan) distance

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



K = 1

L2 (Euclidean) distance

$$d_2(I_1,I_2) = \sqrt{\sum_p \left(I_1^p - I_2^p
ight)^2}$$



K = 1

#### Justin Johnson



### K-Nearest Neighbors: Distance Metric

With the right choice of distance metric, we can apply K-Nearest Neighbor to any type of data!





# K-Nearest Neighbors: Web Demo

Interactively move points around and see decision boundaries change

Play with L1 vs L2 metrics

Play with changing number of training points, value of K

http://vision.stanford.edu/teaching/cs231n-demos/knn/



Fall 2019

#### Justin Johnson

What is the best value of **K** to use? What is the best **distance metric** to use?

These are examples of **hyperparameters**: choices about our learning algorithm that we don't learn from the training data; instead we set them at the start of the learning process


What is the best value of **K** to use? What is the best **distance metric** to use?

These are examples of **hyperparameters**: choices about our learning algorithm that we don't learn from the training data; instead we set them at the start of the learning process

Very problem-dependent. In general need to try them all and see what works best for our data / task.

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

Your Dataset





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

**BAD**: K = 1 always works perfectly on training data

Your Dataset



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

**BAD**: K = 1 always works perfectly on training data

Your Dataset

#### Idea #2: Split data into train and test, choose

hyperparameters that work best on test data

train

test



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

**BAD**: K = 1 always works perfectly on training data

Your Dataset						
Idea #2: Split data into <b>train</b> and <b>test</b> , choose hyperparameters that work best on test data	<b>BAD</b> : No i will perfo	dea how algorit rm on new data	hm			
train		test				



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

**BAD**: K = 1 always works perfectly on training data

Your Dataset			
Idea #2: Split data into <b>train</b> and <b>test</b> , choose hyperparameters that work best on test data	<b>BAD</b> : No i will perfo	idea how algorit rm on new data	
train		test	
Idea #3: Split data into train, val, and test; choose Better! hyperparameters on val and evaluate on test			
train	validation	test	

Justin Johnson	Lecture 2 - 78	Fall 2019

#### Your Dataset

# Idea #4: Cross-Validation: Split data into folds, try each fold as validation and average the results

fold 1	fold 2	fold 3	fold 3 fold 4		test
fold 1	fold 2	fold 3	fold 4	fold 5	test
fold 1	fold 2	fold 3	fold 4	fold 5	test

Useful for small datasets, but (unfortunately) not used too frequently in deep learning

	1151	'In		hn	son
2	430		<u> </u>		3011



Example of 5-fold cross-validation for the value of **k**.

Each point: single outcome.

The line goes through the mean, bars indicated standard deviation

(Seems that k ~ 7 works best for this data)

#### Justin Johnson

### K-Nearest Neighbor on raw pixels is seldom used

- Very slow at test time
- Distance metrics on pixels are not informative



(all 3 images have same L2 distance to the one on the left)

Original image is CCO public domain

#### Justin Johnson



### Nearest Neighbor with ConvNet features works well!



Devlin et al, "Exploring Nearest Neighbor Approaches for Image Captioning", 2015

Justin Johnson



## Nearest Neighbor with ConvNet features works well!

#### Example: Image Captioning with Nearest Neighbor



A bedroom with a bed and a couch.



A cat sitting in a bathroom sink.



A train is stopped at a train station.



A wooden bench in front of a building.

Devlin et al, "Exploring Nearest Neighbor Approaches for Image Captioning", 2015

#### Justin Johnson



# Machine learning

- Why learning?
- Nearest neighbors
- Bayesian classification
- Logistic regression
- Support vector machines
- Clustering
- Principal component analysis

## **Bayesian classification**





### **Bayesian classification**

#### 5.1.2 Bayesian classification

For some simple machine learning problems, e.g., if we have an analytic model of feature construction and noising, or if we can gather enough samples, we can determine the probability distributions of the feature vectors for each class  $p(\mathbf{x}|C_k)$  as well as the prior class likelihoods  $p(C_k)$ .<sup>4</sup> According to Bayes' Rule (4.33), the likelihood of class  $C_k$  given a feature vector  $\mathbf{x}$ is given by

$$y_{k} = p(\mathcal{C}_{k}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_{k})p(\mathcal{C}_{k})}{\sum_{j} p(\mathbf{x}|\mathcal{C}_{j})p(\mathcal{C}_{j})}$$
(5.2)

$$=\frac{\exp l_k}{\sum_j \exp l_j},\tag{5.3}$$

where the second form (using the exp functions) is known as the *normalized exponential* or *softmax* function. The quantity

$$l_k = \log p(\mathbf{x}|\mathcal{C}_k) + \log p(\mathcal{C}_k)$$
(5.4)

is the *log-likelihood* of sample x being from class  $C_k$ .<sup>5</sup> It is sometimes convenient to denote the softmax function itself as a vector-to-vector valued function,

$$\mathbf{y} = \operatorname{softmax}(\mathbf{l}). \tag{5.5}$$

Richard Szeliski

UW CSE 576 - Machine Learning

# Logistic sigmoid function



For the binary (two class) classification task, we can re-write (5.3) as

$$p(\mathcal{C}_0|\mathbf{x}) = \frac{1}{1 + \exp(-l)} = \sigma(l), \qquad (5.7)$$

where  $l = l_0 - l_1$  is the difference between the two class log likelihood and is known as the *log odds* or *logit*.

The  $\sigma(l)$  function is called the *logistic sigmoid function* (or simply the *logistic function* or *logistic curve*), where *sigmoid* means an S-shaped curve (Figure 5.6). The sigmoid was a popular *activation function* in earlier neural networks, although it has now been replaced with other functions, as discussed in Section 5.3.2.

## **Multivariate Gaussian distributions**

• How do we determine the class posterior?



### Multivariate Gaussian distribution

For Gaussians with identical covariance matrices  $\Sigma$ , we have

$$p(\mathbf{x}|\mathcal{C}_k) = \frac{1}{(2\pi)^{D/2}} \frac{1}{\|\mathbf{\Sigma}\|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \mathbf{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_k)\right\}$$
(5.8)

In the case of two classes (binary classification), we obtain (Bishop 2006, Section 4.2.1)

$$p(\mathcal{C}_0|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + b), \tag{5.9}$$

with

$$\mathbf{w} = \mathbf{\Sigma}^{-1} (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1), \quad \text{and} \tag{5.10}$$

$$b = \frac{1}{2}\boldsymbol{\mu}_0^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_0 + \frac{1}{2}\boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 + \log \frac{p(\mathcal{C}_0)}{p(\mathcal{C}_1)}.$$
 (5.11)

# Logistic regression

$$p(\mathcal{C}_0|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + b), \tag{5.9}$$

Equation (5.9), which we will revist shortly in the context of non-generative (discriminative) classification (5.18), is called *logistic regression*, since we pass the output of a linear regression formula

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b \tag{5.12}$$

through the logistic function to obtain a class probability. Figure 5.7 illustrates this in two dimensions, there the posterior likelihood of the red class  $p(C_0|\mathbf{x})$  is shown on the right side.

In linear regression (5.12), w plays the role a the *weight* vector along which we project the feature vector  $\mathbf{x}$ , and b plays the role of the *bias*, which determines where to set the





### Multiple classes & linear discriminant analysis

For K > 2 classes, the softmax function (5.3) can be applied to the linear regression log likelihoods,

$$l_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + b_k,$$

$$y_k = p(\mathcal{C}_k | \mathbf{x}) = \frac{p(\mathbf{x} | \mathcal{C}_k) p(\mathcal{C}_k)}{\sum_j p(\mathbf{x} | \mathcal{C}_j) p(\mathcal{C}_j)}$$

with

$$\mathbf{w}_{k} = \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_{k}, \quad \text{and} \quad (5.14)$$
$$b_{k} = -\frac{1}{2} \boldsymbol{\mu}_{k}^{T} \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_{k} + \log p(\mathcal{C}_{k}). \quad (5.15)$$

Because the decision boundaries along which the classification switches from one class to another are linear,

$$\mathbf{w}_k \mathbf{x} + b_k > \mathbf{w}_l \mathbf{x} + b_l, \tag{5.16}$$

the technique of classifying examples using such crieteria is known as *linear discriminant analysis* (Bishop 2006, Section 4.1; Murphy 2012, Section 4.2.2).<sup>6</sup>

**Richard Szeliski** 

### Quadratic discriminant analysis

• What are the white lines on the right called?



# Machine learning

- Why learning?
- Nearest neighbors
- Bayesian classification
- Logistic regression (linear classifiers)
- Support vector machines
- Clustering
- Principal component analysis



EECS 498-007 / 598-005 Deep Learning for Computer Vision Fall 2019

# Lecture 3: Linear Classifiers



Justin Johnson



#### Neural Network



This image is CC0 1.0 public domain





### Recall CIFAR10

airplane automobile bird cat deer dog frog horse ship truck



**50,000** training images each image is **32x32x3** 

10,000 test images.

#### Justin Johnson



#### Parametric Approach



Justin Johnson



### Parametric Approach: Linear Classifier

$$f(x,W) = Wx$$

f(x,W)

#### Image



**10** numbers giving class scores

Array of **32x32x3** numbers (3072 numbers total)

W parameters or weights





#### Justin Johnson





#### Justin Johnson



### Example for 2x2 image, 3 classes (cat/dog/ship)



$$f(x,W) = Wx + b$$

#### Justin Johnson



### Example for 2x2 image, 3 classes (cat/dog/ship)



#### Justin Johnson

#### Linear Classifier: <u>Algebraic Viewpoint</u>



#### Justin Johnson

#### Linear Classifier: Bias Trick

Add extra one to data vector; bias is absorbed into last column of weight matrix

Stretch pixels into column



Input image (2, 2)

	l .						
0.2	-0.5	0.1	2.0	1.1			
1.5	1.3	2.1	0.0	3.2			
0	0.25	0.2	-0.3	-1.2			
<b>W</b> (3, 5)							



#### Justin Johnson

Linear Classifier: Predictions are Linear!

f(x, W) = Wx (ignore bias)

$$f(cx, W) = W(cx) = c * f(x, W)$$

Justin Johnson



Linear Classifier: Predictions are Linear!

f(x, W) = Wx (ignore bias)

$$f(cx, W) = W(cx) = c * f(x, W)$$
Image Scores 0.5 \* Image 0.5 \* Scores
$$-96.8$$

$$437.8$$

$$62.0$$

$$f(cx, W) = W(cx) = c * f(x, W)$$

$$-48.4$$

$$-48.4$$

$$31.0$$

		• •								
	C'	гι	n	$\square$	n	n	C	$\frown$	n	
J	$\mathbf{S}$	LI					2	U		



#### Interpreting a Linear Classifier

#### **Algebraic Viewpoint**

$$f(x,W) = Wx + b$$



#### Justin Johnson



#### Interpreting a Linear Classifier



Justin Johnson

Lecture 3 - 108

Fall 2019
#### Interpreting an Linear Classifier





#### Justin Johnson



## Interpreting an Linear Classifier: Visual Viewpoint



```
Justin Johnson
```

Lecture 3 - 110

Fall 2019

## Interpreting an Linear Classifier: Visual Viewpoint





Fall 2019



#### Justin Johnson

## Interpreting an Linear Classifier: Visual Viewpoint

Linear classifier has one "template" per category

A single template cannot capture multiple modes of the data

e.g. horse template has 2 heads!

cat

bird



Fall 2019

Justin Johnson

car

plane

Lecture 3 - 112

deer



f(x,W) = Wx + b



Array of **32x32x3** numbers (3072 numbers total)

Justin Johnson





	ust	in J	0	hn	SO	n
-	ase					





J	U	S	ti	'n	0	h	n	S	0	n
2	<u> </u>	-	<u> </u>					-	$\overline{}$	





Justin Johnson

Lecture 3 - 116

Fall 2019



Justin Johnson

Lecture 3 - 117

Fall 2019

## Hard Cases for a Linear Classifier

**Class 1**: First and third quadrants

#### **Class 2**: Second and fourth quadrants

Class 1: 1 <= L2 norm <= 2

**Class 2**: Everything else



**Class 1**: Three modes

#### **Class 2**: Everything else



#### Justin Johnson



## Linear Classifier: Three Viewpoints

**Algebraic Viewpoint** 

f(x,W) = Wx



#### Visual Viewpoint

One template per class



#### **Geometric Viewpoint**

#### Hyperplanes cutting up space



#### Justin Johnson



## So Far: Defined a linear <u>score function</u> f(x,W) = Wx + b







airplane	-3.45	-0.51	3.42
automobile	-8.87	6.04	4.64
bird	0.09	5.31	2.65
cat	2.9	-4.22	5.1
deer	4.48	-4.19	2.64
dog	8.02	3.58	5.55
frog	3.78	4.49	-4.34
horse	1.06	-4.37	-1.5
ship	-0.36	-2.09	-4.79
truck	-0.72	-2,93	6 1 4

Given a W, we can compute class scores for an image x.

But how can we actually choose a good W?

Cat image by Nikita is licensed under CC-BY 2.0; Car image is CCO 1.0 public domain; Frog image is in the public domain

#### Justin Johnson



## Choosing a good W







airplane	-3.45	-0.51	3.42
automobile	-8.87	6.04	4.64
bird	0.09	5.31	2.65
cat	2.9	-4.22	5.1
deer	4.48	-4.19	2.64
dog	8.02	3.58	5.55
frog	3.78	4.49	-4.34
horse	1.06	-4.37	-1.5
ship	-0.36	-2.09	-4.79
truck	-0.72	-2.93	6.14

TODO:

- 1. Use a **loss function** to quantify how good a value of W is
- 2. Find a W that minimizes the loss function (optimization)

#### Justin Johnson

Lecture 3 - 121

9



f(x,W) = Wx + b

A **loss function** tells how good our current classifier is

Low loss = good classifier High loss = bad classifier

(Also called: **objective function**; **cost function**)





A **loss function** tells how good our current classifier is

Low loss = good classifier High loss = bad classifier

(Also called: **objective function**; **cost function**)

Negative loss function sometimes called **reward function**, **profit function**, **utility function**, **fitness function**, etc

#### Justin Johnson



A **loss function** tells how good our current classifier is

Low loss = good classifier High loss = bad classifier

(Also called: **objective function**; **cost function**)

Negative loss function sometimes called **reward function**, **profit function**, **utility function**, **fitness function**, etc Given a dataset of examples

$$\{(x_i, y_i)\}_{i=1}^N$$

Where  $oldsymbol{x_i}$  is image and  $oldsymbol{y_i}$  is (integer) label

#### Justin Johnson



A **loss function** tells how good our current classifier is

Low loss = good classifier High loss = bad classifier

(Also called: **objective function**; **cost function**)

Negative loss function sometimes called **reward function**, **profit function**, **utility function**, **fitness function**, etc Given a dataset of examples

$$\{(x_i, y_i)\}_{i=1}^N$$

Where  $oldsymbol{x_i}$  is image and  $oldsymbol{y_i}$  is (integer) label

Loss for a single example is  $L_i(f(x_i, W), y_i)$ 

#### Justin Johnson



A **loss function** tells how good our current classifier is

Low loss = good classifier High loss = bad classifier

(Also called: **objective function**; **cost function**)

Negative loss function sometimes called **reward function**, **profit function**, **utility function**, **fitness function**, etc Given a dataset of examples

$$\{(x_i, y_i)\}_{i=1}^N$$

Where  $oldsymbol{x_i}$  is image and  $oldsymbol{y_i}$  is (integer) label

Loss for a single example is  $L_i(f(x_i, W), y_i)$ Loss for the dataset is average of

per-example losses:

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

"The score of the correct class should be higher than all the other scores"





"The score of the correct class should be higher than all the other scores"



among other classes

Justin Johnson



"The score of the correct class should be higher than all the other scores"





"The score of the correct class should be higher than all the other scores"



Given an example  $(x_i, y_i)$ ( $x_i$  is image,  $y_i$  is label)

Let  $s = f(x_i, W)$  be scores

Then the SVM loss has the form:

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

#### Justin Johnson



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

**Data loss**: Model predictions should match training data





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

**Data loss**: Model predictions should match training data

**Regularization**: Prevent the model from doing *too* well on training data





$$L(W) = \frac{1}{N} \sum_{i=1}^{N} L_i(f(x_i, W), y_i) + \lambda R(W) \qquad \begin{array}{l} \lambda_i = \text{regularization strength} \\ \text{(hyperparameter)} \end{array}$$

**Data loss**: Model predictions should match training data

**Regularization**: Prevent the model from doing *too* well on training data

#### Justin Johnson



$$L(W) = \frac{1}{N} \sum_{i=1}^{N} L_i(f(x_i, W), y_i) + \lambda R(W) \qquad \begin{array}{l} \lambda_i = \text{regularization strength} \\ \text{(hyperparameter)} \end{array}$$

**Data loss**: Model predictions should match training data

**Regularization**: Prevent the model from doing *too* well on training data

# Simple examplesMore complex:L2 regularization: $R(W) = \sum_k \sum_l W_{k,l}^2$ DropoutL1 regularization: $R(W) = \sum_k \sum_l |W_{k,l}|$ Batch normalizationElastic net (L1 + L2): $R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$ Cutout, Mixup, Stochastic depth, etc...

#### Justin Johnson



$$L(W) = \frac{1}{N} \sum_{i=1}^{N} L_i(f(x_i, W), y_i) + \lambda R(W) \qquad \begin{array}{l} \lambda_i = \text{regularization strength} \\ \text{(hyperparameter)} \end{array}$$

**Data loss**: Model predictions should match training data

**Regularization**: Prevent the model from doing *too* well on training data

#### **Purpose of Regularization:**

- Express preferences in among models beyond "minimize training error"
- Avoid **overfitting**: Prefer simple models that generalize better
- Improve optimization by adding curvature









Justin Johnson





Justin Johnson



Justin Johnson

Want to interpret raw classifier scores as probabilities



- cat **3.2**
- car 5.1

frog -1.7



Want to interpret raw classifier scores as probabilities



$$s = f(x_i; W)$$
 F

$$P(Y=k|X=x_i)=rac{e^{s_k}}{\sum_j e^{s_j}}$$
 Softmax

frog -1.7

Justin Johnson



Want to interpret raw classifier scores as probabilities



$$s = f(x_i; W)$$

$$P(Y=k|X=x_i)=rac{e^{s_k}}{\sum_j e^{s_j}}$$
 Softmax function



Unnormalized logprobabilities / logits

Justin Johnson



Want to interpret raw classifier scores as probabilities



$$P(Y=k|X=x_i)=rac{e^{s_k}}{\sum_j e^{s_j}}$$
 Softmax

#### Justin Johnson



Want to interpret raw classifier scores as probabilities



Justin Johnson


Want to interpret raw classifier scores as probabilities



Justin Johnson



Want to interpret raw classifier scores as probabilities



Justin Johnson



Want to interpret raw classifier scores as probabilities

Fall 2019

			s = f(x)	$e_i; W)$	P(Y =	$=k X=x_i)=rac{e^{s_k}}{\sum_j e^{s_j}}$	Softmax function
		P n	<pre>Probabilities nust be &gt;= C</pre>	Pr mu	obabilities st sum to 1	$L_i = -\log P(Y=y)$	$_i X=x_i)$
cat	3.2		24.5		0.13	→ Compare ←	1.00
car	5.1	exp	164.0	normalize	0.87		0.00
frog	-1.7		0.18		0.00		0.00
Unnormalized log- probabilities / logitsunnormalized probabilitiesCorrect probabilitiesprobabilities / logitsprobabilitiesprobabilities							

Lecture 3 - 164

Justin Johnson

Want to interpret raw classifier scores as probabilities

Fall 2019

			s = f(x)	$e_i; W)$	P(Y =	$=k X=x_i)=rac{e^{s_k}}{\sum_j e^{s_j}}$	Softmax function	
	Probabilities must be >= 0Probabilities must sum to 1 $L_i = -\log P(Y = y_i)$							
cat	3.2		24.5		0.13	🔶 Compare ←	1.00	
car	5.1	exp	164.0	normalize	0.87	Kullback–Leibler divergence	0.00	
frog	-1.7		0.18		0.00	$D_{KL}(P \  Q) =$	0.00	
Unnormalized log- probabilities / logits probabilities probabilities $\sum_{y} P(y) \log \frac{P(y)}{Q(y)}$								

Justin Johnson

Want to interpret raw classifier scores as probabilities

			s = f(x)	(i;W)	P(Y =	$=k X=x_i)=rac{e^{s_k}}{\sum_j e^{s_j}}$	Softmax function	
		n P	nust be >= 0	mu	st sum to 1	$L_i = -\log P(Y = y_i)$	$_{i} X=x_{i})$	
cat	3.2		24.5		0.13	🔶 Compare ←	1.00	
car	5.1	exp	164.0	normalize	0.87	Cross Entropy	0.00	
frog	-1.7		0.18		0.00	H(P,Q) =	0.00	
Unnormalized log- probabilities / logits probabilities probabilities $H(p) + D_{KL}(P  Q)$								

Want to interpret raw classifier scores as probabilities



3.2

$$s=f(x_i;W)$$
  $P(Y=k|X=x_i)=rac{e^{s_k}}{\sum_j e^{s_j}}$ 

Maximize probability of correct class

 $L_i = -\log P(Y = y_i | X = x_i)$ 

Putting it all together:

$$L_i = -\log(rac{e^{sy_i}}{\sum_j e^{s_j}})$$

car 5.1

cat

frog -1.7



Softmax

function

Want to interpret raw classifier scores as probabilities



3.2

cat

$$s=f(x_i;W)$$
  $P(Y=k|X=x_i)=rac{e^{s_k}}{\sum_j e^{s_j}}$  Softmax function

Maximize probability of correct class

Maximize probability of correct class Putting it all together:  

$$L_i = -\log P(Y = y_i | X = x_i)$$
  $L_i = -\log(\frac{e^{sy_i}}{e^{sy_i}})$ 

$$L_i = -\log(rac{e^{sy_i}}{\sum_j e^{s_j}})$$



Want to interpret raw classifier scores as probabilities



3.2

cat

$$s=f(x_i;W)$$
  $P(Y=k|X=x_i)=rac{e^{s_k}}{\sum_j e^{s_j}}$  Softmax function

Maximize probability of correct class

Putting it all together:

$$L_i = -\log P(Y=y_i|X=x_i) \qquad L_i = -\log(rac{e^{sy_i}}{\sum_j e^{s_j}})$$

**A**: Min 0, max +infinity



Want to interpret raw classifier scores as probabilities



3.2

cat

$$s=f(x_i;W)$$
  $P(Y=k|X=x_i)=rac{e^{s_k}}{\sum_j e^{s_j}}$  Softmax function

Maximize probability of correct class

$$L_i = -\log P(Y=y_i|X=x_i)$$

$$L_i = -\log(rac{e^{sy_i}}{\sum_j e^{s_j}})$$

car5.1Q: If all scores arefrog-1.7small random values,what is the loss?



Want to interpret raw classifier scores as probabilities



$$s = f(x_i; W)$$

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$
Softmax function  
Maximize probability of correct class
$$Putting it all together:$$

$$L_i = -\log P(Y = y_i | X = x_i)$$

$$L_i = -\log(\frac{e^{sy_i}}{\sum_j e^{s_j}})$$
5.1
$$Q: \text{ If all scores are small random values, what is the loss?}$$

$$A: -\log(1/C) \log(10) \approx 2.3$$

cat

car

frog



Softmax

### Recap: Three ways to think about linear classifiers

**Algebraic Viewpoint** 

f(x,W) = Wx



#### Visual Viewpoint

One template per class



#### **Geometric Viewpoint**

Hyperplanes cutting up space



#### Justin Johnson



### Recap: Loss Functions quantify preferences

- We have some dataset of (x, y)
- We have a **score function**:
- We have a **loss function**:

$$s = f(x;W) = Wx$$
Linear classifier



#### Justin Johnson



### Recap: Loss Functions quantify preferences

- We have some dataset of (x, y)
- We have a **score function**:
- We have a **loss function**:

**Q**: How do we find the best W?

$$s = f(x; W) = Wx$$
Linear classifier

### A: Next lecture



#### Justin Johnson

# Machine learning

- Why learning?
- Nearest neighbors
- Bayesian classification
- Logistic regression
- Support vector machines & random forests (lightening speed)
- Clustering
- Principal component analysis

# Support vector machines (SVMs)

- Maximize the *margin* between the decision surfaces
- The only points that matter are the circled *support vectors*



n margin classifier, we need to  $l_i = \mathbf{w} \cdot \mathbf{x}_i + b$  (5.17) have an note this more compactly, let

 $\hat{t}_i = 2t_i - 1, \quad \hat{t}_i \in \{-1, 1\}$ 

an now re-write the inequality c

 $\hat{t}_i(\mathbf{w} \cdot \mathbf{x}_i + b) \ge 1.$ 

simply find the smallest norm ization problem

 $\arg\min_{\mathbf{w},b}\|\mathbf{w}\|^2$ 

assic quadratic programming 1

### Kernelized support vector machines

 Replace linear function with a sum of *kernel functions* (e.g., Gaussian bumps)

$$l = f(\mathbf{x}) = \sum_{k} \mathbf{w}_{k} \phi(\|\mathbf{x} - \mathbf{x}_{k}\|).$$

 Circled points are support vectors, lie on f = ± 1 surfaces (f = 0 is the dark curve)



### Hinge loss vs. logistic regression

Figure 7.5 Plot of the 'hinge' error function used in support vector machines, shown in blue, along with the error function for logistic regression, rescaled by a factor of  $1/\ln(2)$  so that it passes through the point (0,1), shown in red. Also shown are the misclassification error in black and the squared error in green.

remaining points we have  $\xi_n = 1 - y_n t_n$ . Thus the objective function (7.21) can be written (up to an overall multiplicative constant) in the form

$$\sum_{n=1}^{N} E_{\rm SV}(y_n t_n) + \lambda \|\mathbf{w}\|^2$$
(7.44)

where  $\lambda = (2C)^{-1}$ , and  $E_{SV}(\cdot)$  is the *hinge* error function defined by

$$E_{\rm SV}(y_n t_n) = [1 - y_n t_n]_+ \tag{7.45}$$

where  $[\cdot]_+$  denotes the positive part. The hinge error function, so-called because of its shape, is plotted in Figure 7.5. It can be viewed as an approximation to the misclassification error, i.e., the error function that ideally we would like to minimize, which is also shown in Figure 7.5.



#### **Richard Szeliski**

# Decision trees (lighting speed)



## Random forest



T = # trees

# Random forest





*ρ* = 5

 $\rho$  = # samples at each node at construction time

## Application: Kinect body pose estimation



## Lecture recap: machine learning

- Why learning?
- Nearest neighbors
- Bayesian classification
- Logistic regression
- Support vector machines
- Unsupervised learning:
  - Clustering
  - Principal component analysis