




Dimensionality Reduction PCA continued...

Machine Learning – CSE446
Carlos Guestrin
University of Washington

May 22, 2013
©Carlos Guestrin 2005-2013

1

Dimensionality reduction

- 
- Input data may have thousands or millions of dimensions!
 - e.g., text data has *x with 10 000 — 10 000 000 dims*
 - **Dimensionality reduction:** represent data with fewer dimensions
 - easier learning – fewer parameters
 - visualization – hard to visualize more than 3D or 4D
 - discover “intrinsic dimensionality” of data
 - high dimensional data that is truly lower dimensional

©Carlos Guestrin 2005-2013

Lower dimensional projections

- Rather than picking a subset of the features, we can new features that are combinations of existing features

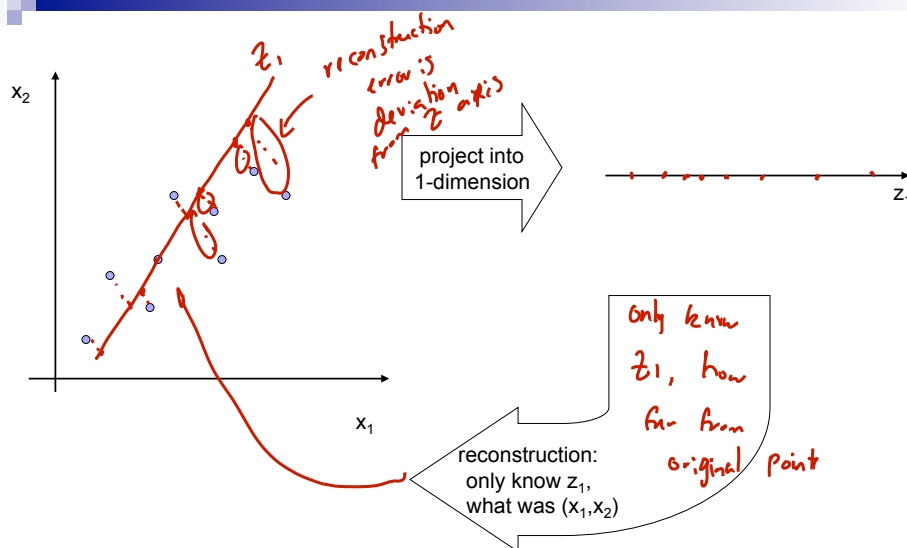
$$z_7 = 2.5x_1 + 2.9x_2 - 3.2x_3 \dots$$

x (1000) $\xrightarrow{\text{learn}}$ z (100) $\xrightarrow{\text{model}}$ $z = Ax$ $\xrightarrow{\text{loss/accuracy}}$ reconstruction error

- Let's see this in the unsupervised setting
 - just X , but no Y

©Carlos Guestrin 2005-2013

Linear projection and reconstruction



©Carlos Guestrin 2005-2013

Principal component analysis – basic idea

- Project d-dimensional data into k-dimensional space while preserving information:
 - e.g., project space of 10000 words into 3-dimensions
 - e.g., project 3-d into 2-d
- Choose projection with minimum reconstruction error

©Carlos Guestrin 2005-2013

Linear projections, a review

- Project a point into a (lower dimensional) space:
 - **point:** $\mathbf{x} = (x_1, \dots, x_d)$
 - **select a basis** – set of basis vectors – $(\mathbf{u}_1, \dots, \mathbf{u}_k)$
 - we consider orthonormal basis:
 - $\mathbf{u}_i \bullet \mathbf{u}_i = 1$, and $\mathbf{u}_i \bullet \mathbf{u}_j = 0$ for $i \neq j$
 - **select a center** – $\bar{\mathbf{x}}$, defines offset of space
 - **best coordinates** in lower dimensional space defined by dot-products: (z_1, \dots, z_k) , $z_i = (\mathbf{x} - \bar{\mathbf{x}}) \bullet \mathbf{u}_i$
 - minimum squared error

©Carlos Guestrin 2005-2013

PCA finds projection that minimizes reconstruction error

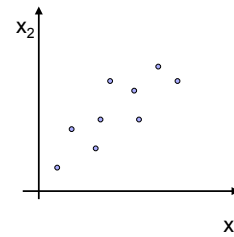
- Given m data points: $\mathbf{x}^i = (x_1^i, \dots, x_d^i)$, $i=1 \dots N$
- Will represent each point as a projection:

$$\square \hat{\mathbf{x}}^i = \bar{\mathbf{x}} + \sum_{j=1}^k z_j^i \mathbf{u}_j \quad \text{where: } \bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}^i \quad \text{and} \quad z_j^i = (\mathbf{x}^i - \bar{\mathbf{x}}) \cdot \mathbf{u}_j$$

- PCA:

- Given $k < d$, find $(\mathbf{u}_1, \dots, \mathbf{u}_k)$ minimizing reconstruction error:

$$error_k = \sum_{i=1}^N (\mathbf{x}^i - \hat{\mathbf{x}}^i)^2$$



©Carlos Guestrin 2005-2013

Understanding the reconstruction error

- Note that \mathbf{x}^i can be represented exactly by d -dimensional projection:

$$\mathbf{x}^i = \bar{\mathbf{x}} + \sum_{j=1}^d z_j^i \mathbf{u}_j$$

$$\hat{\mathbf{x}}^i = \bar{\mathbf{x}} + \sum_{j=1}^k z_j^i \mathbf{u}_j$$

$$z_j^i = (\mathbf{x}^i - \bar{\mathbf{x}}) \cdot \mathbf{u}_j$$

- Given $k < d$, find $(\mathbf{u}_1, \dots, \mathbf{u}_k)$ minimizing reconstruction error:

$$error_k = \sum_{i=1}^N (\mathbf{x}^i - \hat{\mathbf{x}}^i)^2$$

- Rewriting error:

©Carlos Guestrin 2005-2013

Reconstruction error and covariance matrix

$$error_k = \sum_{i=1}^N \sum_{j=k+1}^d [\mathbf{u}_j \cdot (\mathbf{x}^i - \bar{\mathbf{x}})]^2$$

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

©Carlos Guestrin 2005-2013

Minimizing reconstruction error and eigen vectors

- Minimizing reconstruction error equivalent to picking (ordered) orthonormal basis $(\mathbf{u}_1, \dots, \mathbf{u}_d)$ minimizing:

$$error_k = \sum_{j=k+1}^d \mathbf{u}_j^T \Sigma \mathbf{u}_j$$

- Eigen vector:
- Minimizing reconstruction error equivalent to picking $(\mathbf{u}_{k+1}, \dots, \mathbf{u}_d)$ to be eigen vectors with smallest eigen values

©Carlos Guestrin 2005-2013

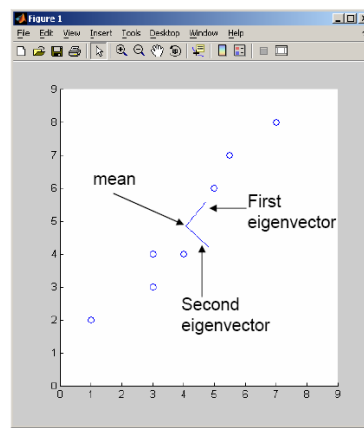
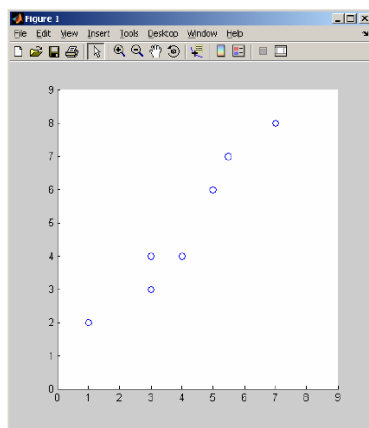
Basic PCA algorithm

- Start from m by n data matrix \mathbf{X}
- **Recenter**: subtract mean from each row of \mathbf{X}
 - $\mathbf{X}_c \leftarrow \mathbf{X} - \bar{\mathbf{X}}$
- **Compute covariance matrix**:
 - $\Sigma \leftarrow 1/N \mathbf{X}_c^T \mathbf{X}_c$
- Find **eigen vectors and values** of Σ
- **Principal components**: k eigen vectors with highest eigen values

©Carlos Guestrin 2005-2013

PCA example

$$\hat{\mathbf{x}}^i = \bar{\mathbf{x}} + \sum_{j=1}^k z_j^i \mathbf{u}_j$$

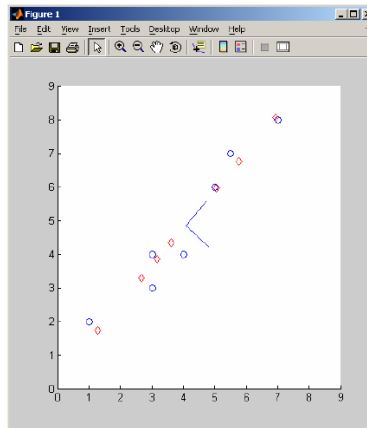
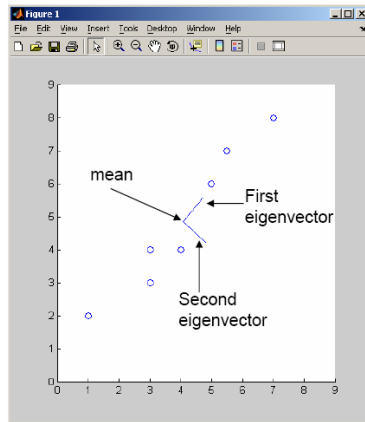


©Carlos Guestrin 2005-2013

PCA example – reconstruction

$$\hat{x}^i = \bar{x} + \sum_{j=1}^k z_j^i u_j$$

only used first principal component



©Carlos Guestrin 2005-2013

Eigenfaces [Turk, Pentland '91]

■ Input images:



■ Principal components:



©Carlos Guestrin 2005-2013

Eigenfaces reconstruction

- Each image corresponds to adding 8 principal components:



©Carlos Guestrin 2005-2013

Scaling up

- Covariance matrix can be really big!
 - Σ is d by d
 - Say, only 10000 features
 - finding eigenvectors is very slow...
- Use singular value decomposition (SVD)
 - finds to k eigenvectors
 - great implementations available, e.g., R or Matlab svd

©Carlos Guestrin 2005-2013

SVD

- Write $\mathbf{X} = \mathbf{W} \mathbf{S} \mathbf{V}^T$
 - $\mathbf{X} \leftarrow$ data matrix, one row per datapoint
 - $\mathbf{W} \leftarrow$ weight matrix, one row per datapoint – coordinate of \mathbf{x}^i in eigenspace
 - $\mathbf{S} \leftarrow$ singular value matrix, diagonal matrix
 - in our setting each entry is eigenvalue λ_j
 - $\mathbf{V}^T \leftarrow$ singular vector matrix
 - in our setting each row is eigenvector \mathbf{v}_j

©Carlos Guestrin 2005-2013

PCA using SVD algorithm

- Start from m by n data matrix \mathbf{X}
- **Recenter**: subtract mean from each row of \mathbf{X}
 - $\mathbf{X}_c \leftarrow \mathbf{X} - \bar{\mathbf{X}}$
- Call SVD algorithm on \mathbf{X}_c – ask for k singular vectors
- **Principal components**: k singular vectors with highest singular values (rows of \mathbf{V}^T)
 - **Coefficients** become:

©Carlos Guestrin 2005-2013

What you need to know

- Dimensionality reduction
 - why and when it's important
- Simple feature selection
- Principal component analysis
 - minimizing reconstruction error
 - relationship to covariance matrix and eigenvectors
 - using SVD

©Carlos Guestrin 2005-2013

Bayes optimal classifier Naïve Bayes

Machine Learning – CSE446
Carlos Guestrin
University of Washington

May 22, 2013

©Carlos Guestrin 2005-2013

20

Classification

- **Learn:** $h: \mathbf{X} \mapsto Y$
 - \mathbf{X} – features
 - Y – target classes
- Suppose you know $P(Y|\mathbf{X})$ exactly, how should you classify?
 - Bayes optimal classifier:

Bayes Rule

$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$$

Which is shorthand for:

$$(\forall i, j) P(Y = y_i | X = x_j) = \frac{P(X = x_j | Y = y_i) P(Y = y_i)}{P(X = x_j)}$$

How hard is it to learn the optimal classifier?

- Data =

Sky	Temp	Humid	Wind	Water	Forecast	EnjoySpt
Sunny	Warm	Normal	Strong	Warm	Same	Yes
Sunny	Warm	High	Strong	Warm	Same	Yes
Rainy	Cold	High	Strong	Warm	Change	No
Sunny	Warm	High	Strong	Cool	Change	Yes

- How do we represent these? How many parameters?

- Prior, $P(Y)$:

- Suppose Y is composed of k classes

- Likelihood, $P(\mathbf{X}|Y)$:

- Suppose \mathbf{X} is composed of d binary features

- **Complex model ! High variance with limited data!!!**

©Carlos Guestrin 2005-2013

23

Conditional Independence

- X is **conditionally independent** of Y given Z , if the probability distribution governing X is independent of the value of Y , given the value of Z
 $(\forall i, j, k) P(X = i | Y = j, Z = k) = P(X = i | Z = k)$

- e.g., $P(\text{Thunder} | \text{Rain}, \text{Lightning}) = P(\text{Thunder} | \text{Lightning})$

- Equivalent to:

$$P(X, Y | Z) = P(X | Z)P(Y | Z)$$

©Carlos Guestrin 2005-2013

24

What if features are independent?

- Predict Thunder
- From two **conditionally Independent** features
 - Lightning
 - Rain

The Naïve Bayes assumption

- Naïve Bayes assumption:
 - Features are independent given class:

$$\begin{aligned}P(X_1, X_2|Y) &= P(X_1|X_2, Y)P(X_2|Y) \\ &= P(X_1|Y)P(X_2|Y)\end{aligned}$$

- More generally:

$$P(X_1 \dots X_d | Y) = \prod_i P(X_i | Y)$$

- How many parameters now?
 - Suppose \mathbf{X} is composed of d binary features

The Naïve Bayes Classifier

- Given:

- Prior $P(Y)$
- d conditionally independent features \mathbf{X} given the class Y
- For each X_i , we have likelihood $P(X_i|Y)$

- Decision rule:

$$\begin{aligned} y^* = h_{NB}(\mathbf{x}) &= \arg \max_y P(y) P(x_1, \dots, x_d | y) \\ &= \arg \max_y P(y) \prod_i P(x_i | y) \end{aligned}$$

- If assumption holds, NB is optimal classifier!

©Carlos Guestrin 2005-2013

27

MLE for the parameters of NB

- Given dataset

- $\text{Count}(A=a, B=b)$ == number of examples where $A=a$ and $B=b$

- MLE for NB, simply:

- Prior: $P(Y=y) =$

- Likelihood: $P(X_i=x_i|Y=y) =$

©Carlos Guestrin 2005-2013

28

Subtleties of NB classifier 1 – Violating the NB assumption

- Usually, features are not conditionally independent:

$$P(X_1 \dots X_d | Y) \neq \prod_i P(X_i | Y)$$

- Actual probabilities $P(Y|\mathbf{X})$ often biased towards 0 or 1
- Nonetheless, NB is the single most used classifier out there
 - NB often performs well, even when assumption is violated
 - [Domingos & Pazzani '96] discuss some conditions for good performance

©Carlos Guestrin 2005-2013

29

Subtleties of NB classifier 2 – Insufficient training data

- What if you never see a training instance where $X_1=a$ when $Y=b$?
 - e.g., $Y=\{\text{SpamEmail}\}$, $X_1=\{\text{'Enlargement'}\}$
 - $P(X_1=a | Y=b) = 0$
- Thus, no matter what the values X_2, \dots, X_d take:
 - $P(Y=b | X_1=a, X_2, \dots, X_d) = 0$
- “Solution”: smoothing
 - Add “fake” counts, usually uniformly distributed
 - Equivalent to Bayesian Learning

©Carlos Guestrin 2005-2013

30

Text classification

- Classify e-mails
 - $Y = \{\text{Spam}, \text{NotSpam}\}$
- Classify news articles
 - $Y = \{\text{what is the topic of the article?}\}$
- Classify webpages
 - $Y = \{\text{Student, professor, project, ...}\}$
- What about the features **X**?
 - The text!

©Carlos Guestrin 2005-2013

31

Features **X** are entire document – X_i for i^{th} word in article

Article from rec.sport.hockey

Path: cantaloupe.srv.cs.cmu.edu!das-news.harvard.e
From: xxx@yyy.zzz.edu (John Doe)
Subject: Re: This year's biggest and worst (opinion)
Date: 5 Apr 93 09:53:39 GMT

I can only comment on the Kings, but the most obvious candidate for pleasant surprise is Alex Zhitnik. He came highly touted as a defensive defenseman, but he's clearly much more than that. Great skater and hard shot (though wish he were more accurate). In fact, he pretty much allowed the Kings to trade away that huge defensive liability Paul Coffey. Kelly Hrdey is only the biggest disappointment if you thought he was any good to begin with. But, at best, he's only a mediocre goaltender. A better choice would be Tomas Sandstrom, though not through any fault of his own, but because some things in Toronto decided

32

NB for Text classification

- $P(\mathbf{X}|Y)$ is huge!!!

- Article at least 1000 words, $\mathbf{X}=\{X_1, \dots, X_{1000}\}$
- X_i represents i^{th} word in document, i.e., the domain of X_i is entire vocabulary, e.g., Webster Dictionary (or more), 10,000 words, etc.

- NB assumption helps a lot!!!

- $P(X_i=x_i|Y=y)$ is just the probability of observing word x_i in a document on topic y

$$h_{NB}(\mathbf{x}) = \arg \max_y P(y) \prod_{i=1}^{LengthDoc} P(x_i|y)$$

©Carlos Guestrin 2005-2013

33

Bag of words model

- Typical additional assumption – **Position in document doesn't matter**: $P(X_i=x_i|Y=y) = P(X_k=x_i|Y=y)$

- “Bag of words” model – order of words on the page ignored
- Sounds really silly, but often works very well!

$$P(y) \prod_{i=1}^{LengthDoc} P(x_i|y)$$

When the lecture is over, remember to wake up the person sitting next to you in the lecture room.

©Carlos Guestrin 2005-2013

34

Bag of words model

- Typical additional assumption – **Position in document doesn't matter**: $P(X_i=x_i|Y=y) = P(X_k=x_k|Y=y)$

- “Bag of words” model – order of words on the page ignored
- Sounds really silly, but often works very well!

$$P(y) \prod_{i=1}^{LengthDoc} P(x_i|y)$$

in is lecture lecture next over person remember room
sitting the the the to to up wake when you

Bag of Words Approach



The screenshot shows the TOTAL company website. The header includes the TOTAL logo and a navigation menu with links: Global Activities, Corporate Structure, TOTAL's Story, Upstream Strategy, Downstream Strategy, Chemicals Strategy, TOTAL Foundation, and Homepage. The main content area is titled "all about the company" and contains three paragraphs of text about the company's operations and strengths.

aardvark	0
about	2
all	2
Africa	1
apple	0
anxious	0
...	
gas	1
...	
oil	1
...	
Zaire	0

NB with Bag of Words for text classification

■ Learning phase:

- Prior $P(Y)$
 - Count how many documents you have from each topic (+ prior)
- $P(X_i|Y)$
 - For each topic, count how many times you saw word in documents of this topic (+ prior)

■ Test phase:

- For each document
 - Use naïve Bayes decision rule

$$h_{NB}(x) = \arg \max_y P(y) \prod_{i=1}^{LengthDoc} P(x_i|y)$$

©Carlos Guestrin 2005-2013

37

Twenty News Groups results

Given 1000 training documents from each group
Learn to classify new documents according to
which newsgroup it came from

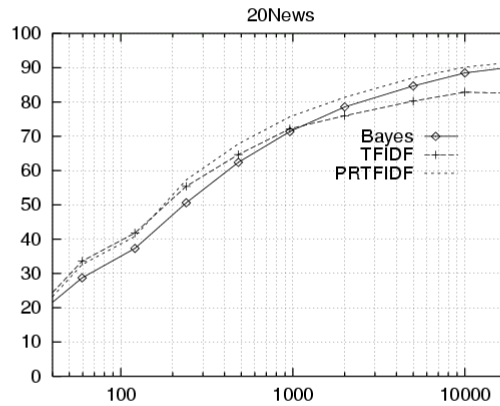
comp.graphics	misc.forsale
comp.os.ms-windows.misc	rec.autos
comp.sys.ibm.pc.hardware	rec.motorcycles
comp.sys.mac.hardware	rec.sport.baseball
comp.windows.x	rec.sport.hockey
alt.atheism	sci.space
soc.religion.christian	sci.crypt
talk.religion.misc	sci.electronics
talk.politics.mideast	sci.med
talk.politics.misc	
talk.politics.guns	

Naive Bayes: 89% classification accuracy

©Carlos Guestrin 2005-2013

38

Learning curve for Twenty News Groups



Accuracy vs. Training set size (1/3 withheld for test)