

CS 2750: Machine Learning
K-NN (cont'd) + Review

Prof. Adriana Kovashka
University of Pittsburgh
February 3, 2016

Plan for today

- A few announcements
- Wrap-up K-NN
- Quizzes back + linear algebra and PCA review
- Next time: Other review + start on linear regression

Announcements

- HW2 released – due 2/24
 - Should take a lot less time
 - Overview
- Reminder: project proposal due 2/17
 - See course website for a project can be + ideas
 - “In the proposal, describe what techniques and data you plan to use, and what existing work there is on the subject.”
- Notes from board from Monday uploaded

Announcements

- Just for this week, the TA's office hours will be on Wednesday, 2-6pm

Machine Learning vs Computer Vision

- I spent 20 minutes on computer vision features
 - You will learn soon enough that how you compute your feature representation is important in machine learning
- Everything else I've shown you or that you've had to do in homework has been machine learning *applied to computer vision*
 - The goal of Part III in HW1 was to get you comfortable with the idea that each dimension in a data representation captures the result of some separate "test" applied to the data point

Machine Learning vs Computer Vision

- Half of computer vision today is *applied machine learning*
- The other half has to do with image processing, and you've seen none of that
 - See the slides for my undergraduate computer vision class

Machine Learning vs Computer Vision

- Many machine learning researchers today do computer vision as well, so these topics are not entirely disjoint
 - Look up the recent publications of Kevin Murphy (Google), the author of “Machine Learning: A Probabilistic Perspective”

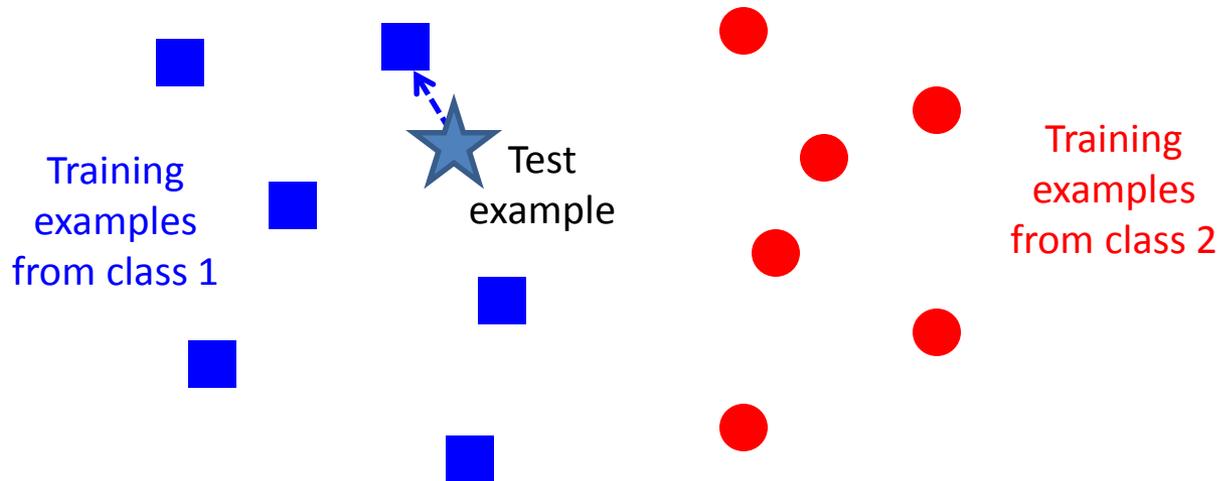
Machine Learning vs Computer Vision

- In conclusion: I will try to also use data other than images in the examples that I give, but I trust that you will have the maturity to not think that I'm teaching you computer vision just because I include examples with images 😊

Last time: Supervised Learning Part I

- Basic formulation of the simplest classifier:
K Nearest Neighbors
- Example uses
- Generalizing the distance metric and weighting neighbors differently
- Problems:
 - The curse of dimensionality
 - Picking K
 - Approximation strategies

1-Nearest Neighbor Classifier

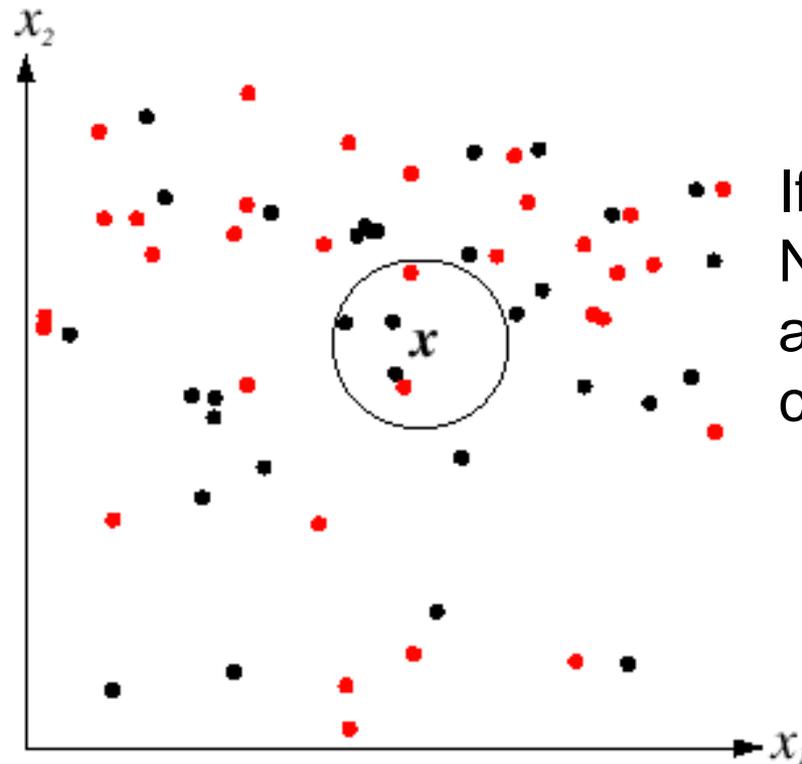


$f(\mathbf{x}) = \text{label of the training example nearest to } \mathbf{x}$

- All we need is a distance function for our inputs
- No training required!

K-Nearest Neighbors Classifier

- For a new point, find the k closest points from training data (e.g. $k=5$)
- Labels of the k points “vote” to classify



If query lands here, the 5 NN consist of 3 negatives and 2 positives, so we classify it as negative.

Black = negative
Red = positive

k-Nearest Neighbor

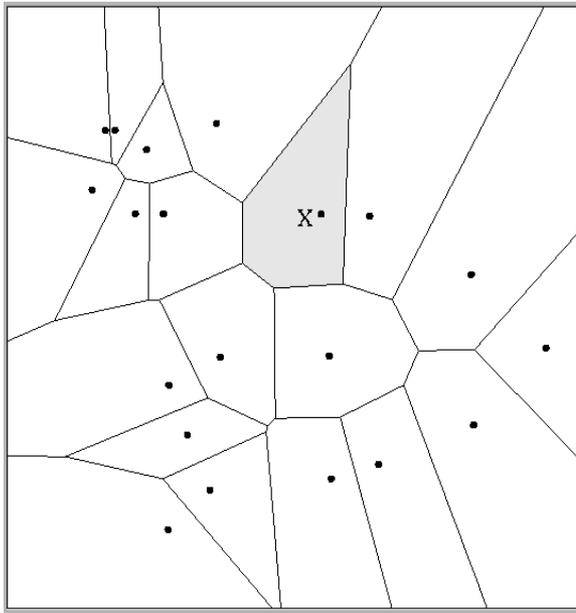
Four things make a memory based learner:

- *A distance metric*
 - **Euclidean (and others)**
- *How many nearby neighbors to look at?*
 - **k**
- *A weighting function (optional)*
 - **Not used**
- *How to fit with the local points?*
 - **Just predict the average output among the nearest neighbors**

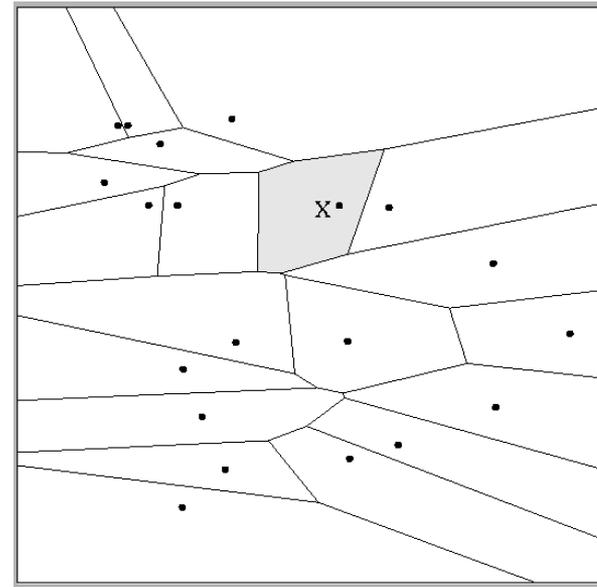
Multivariate distance metrics

Suppose the input vectors $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$ are two dimensional:

$$\mathbf{x}_1 = (x_{11}, x_{12}), \mathbf{x}_2 = (x_{21}, x_{22}), \dots, \mathbf{x}_N = (x_{N1}, x_{N2}).$$



$$Dist(\mathbf{x}_i, \mathbf{x}_j) = (x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2$$



$$Dist(\mathbf{x}_i, \mathbf{x}_j) = (x_{i1} - x_{j1})^2 + (3x_{i2} - 3x_{j2})^2$$

The relative scalings in the distance metric affect region shapes

Another generalization: Weighted K-NNs

- Neighbors weighted differently:
- Extremes
 - Bandwidth = infinity: prediction is dataset average
 - Bandwidth = zero: prediction becomes 1-NN

Kernel Regression/Classification

Four things make a memory based learner:

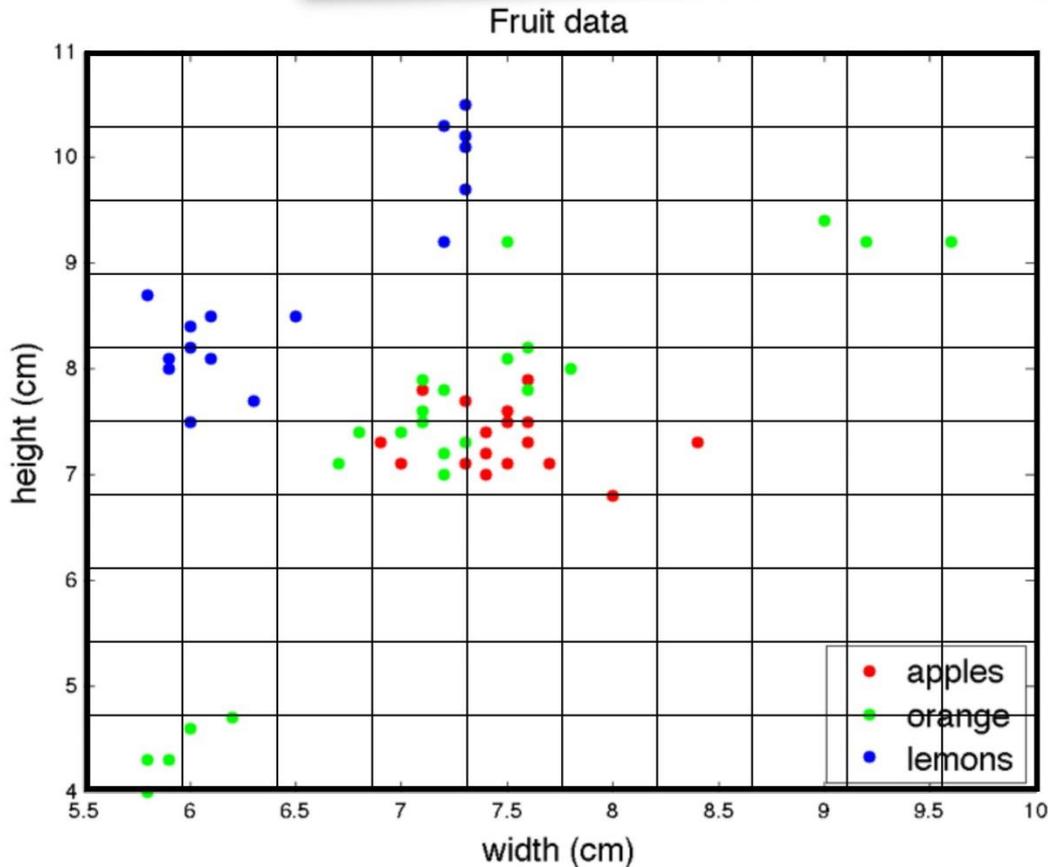
- *A distance metric*
 - **Euclidean (and others)**
- *How many nearby neighbors to look at?*
 - **All of them**
- *A weighting function (optional)*
 - $w_i = \exp(-d(x_i, \text{query})^2 / \sigma^2)$
 - Nearby points to the query are weighted strongly, far points weakly. The σ parameter is the **Kernel Width**.
- *How to fit with the local points?*
 - **Predict the weighted average of the outputs**

Problems with Instance-Based Learning

- Too many features?
 - Doesn't work well if large number of irrelevant features, distances overwhelmed by noisy features
 - Distances become meaningless in high dimensions (the **curse of dimensionality**)
- What is the impact of the **value of K**?
- Expensive
 - No learning: most real work done during testing
 - For every test sample, must search through all dataset – very slow!
 - **Must use tricks like approximate nearest neighbor search**
 - Need to store all training data

Curse of Dimensionality

How many neighborhoods are there?



$$\begin{aligned}\text{\#bins} &= 10 \times 10 \\ d &= 2\end{aligned}$$

$$\begin{aligned}\text{\#bins} &= 10^d \\ d &= 1000\end{aligned}$$

Atoms in the universe
 $\sim 10^{80}$

Curse of Dimensionality

- Consider: Sphere of radius 1 in d-dims
- Consider: an outer ε -shell in this sphere
- What is $\frac{\textit{shell volume}}{\textit{sphere volume}}$?

Curse of Dimensionality

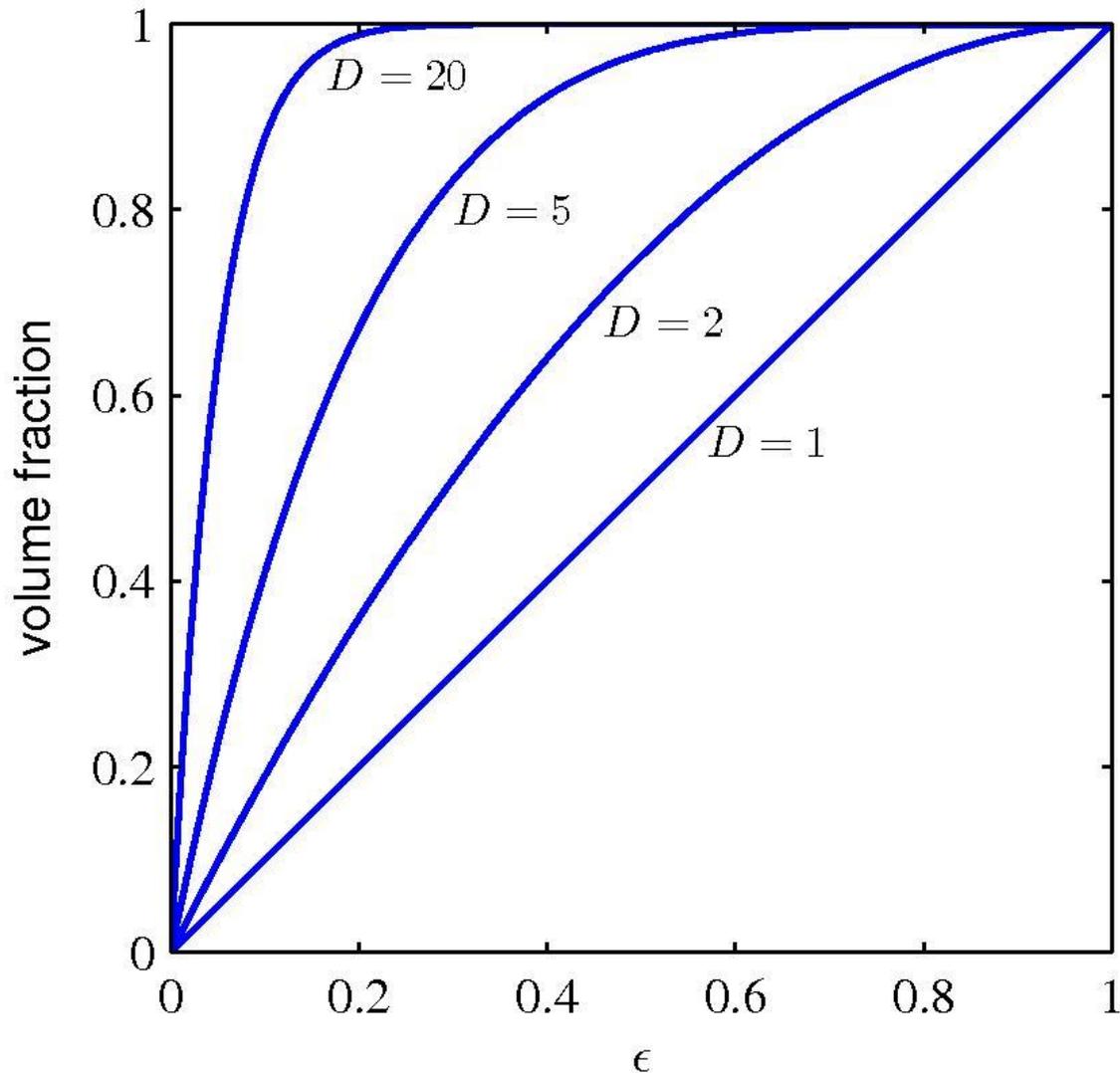


Figure 1.22 from Bishop

Curse of Dimensionality

- Problem: In very high dimensions, distances become less meaningful
- This problem applies to all types of classifiers, not just K-NN

Problems with Instance-Based Learning

- Too many features?
 - Doesn't work well if large number of irrelevant features, distances overwhelmed by noisy features
 - Distances become meaningless in high dimensions (the **curse of dimensionality**)
- What is the impact of the **value of K**?
- Expensive
 - No learning: most real work done during testing
 - For every test sample, must search through all dataset – very slow!
 - **Must use tricks like approximate nearest neighbor search**
 - Need to store all training data

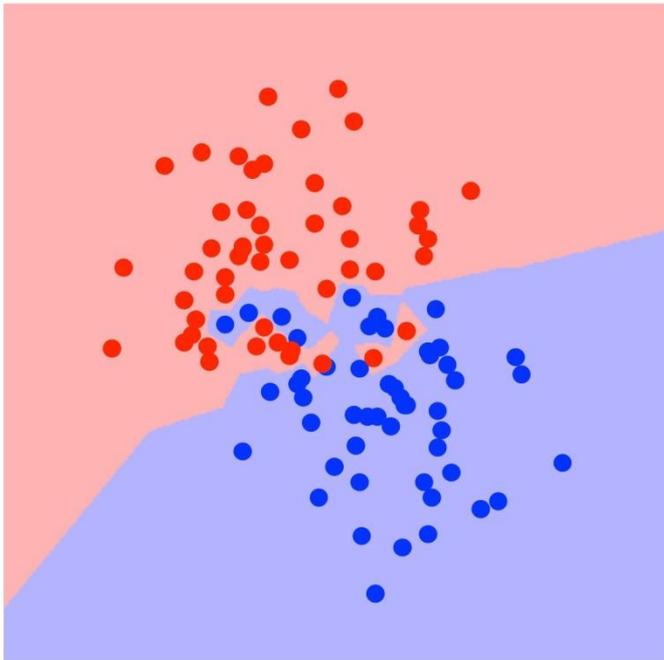
kNN Decision Boundary

- Increasing k simplifies / complicates decision boundary

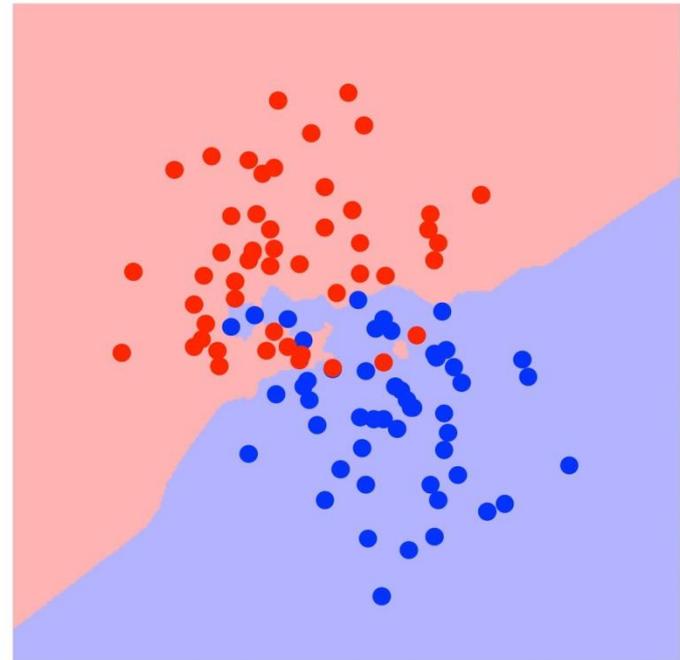
kNN Decision Boundary

- Increasing k “simplifies” decision boundary
 - Majority voting means less emphasis on individual points

$K = 1$



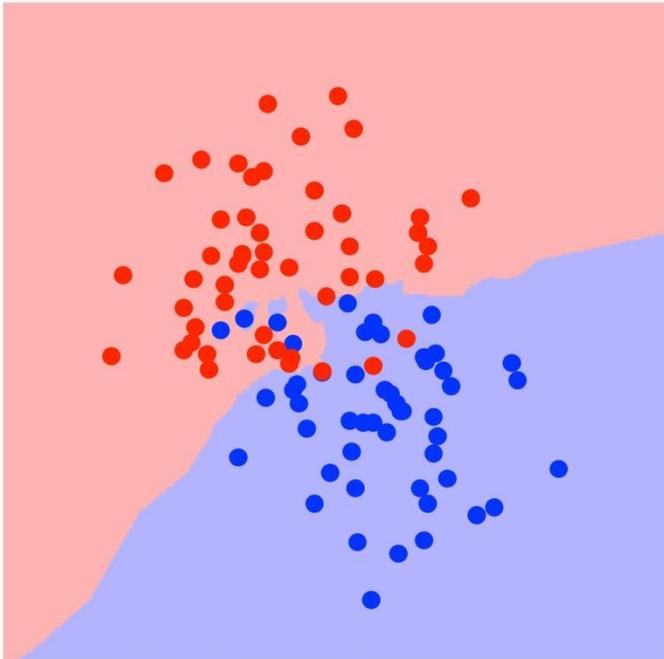
$K = 3$



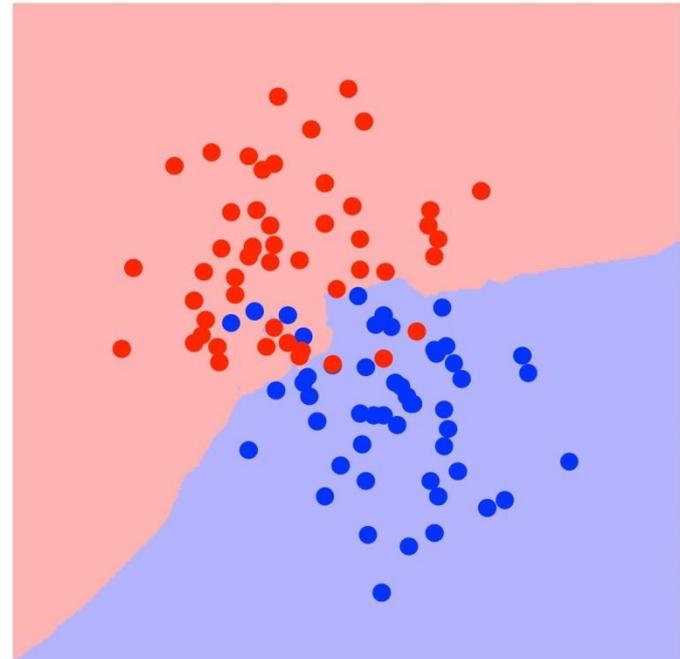
kNN Decision Boundary

- Increasing k “simplifies” decision boundary
 - Majority voting means less emphasis on individual points

$K = 5$



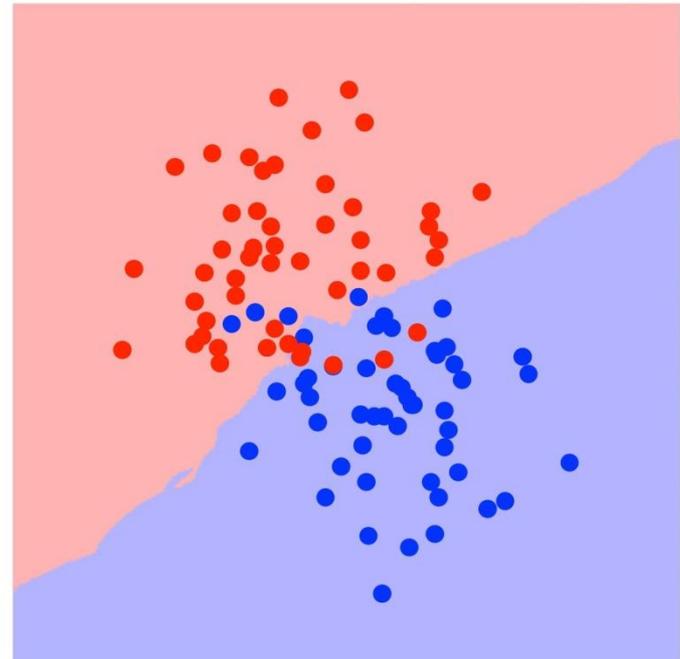
$K = 7$



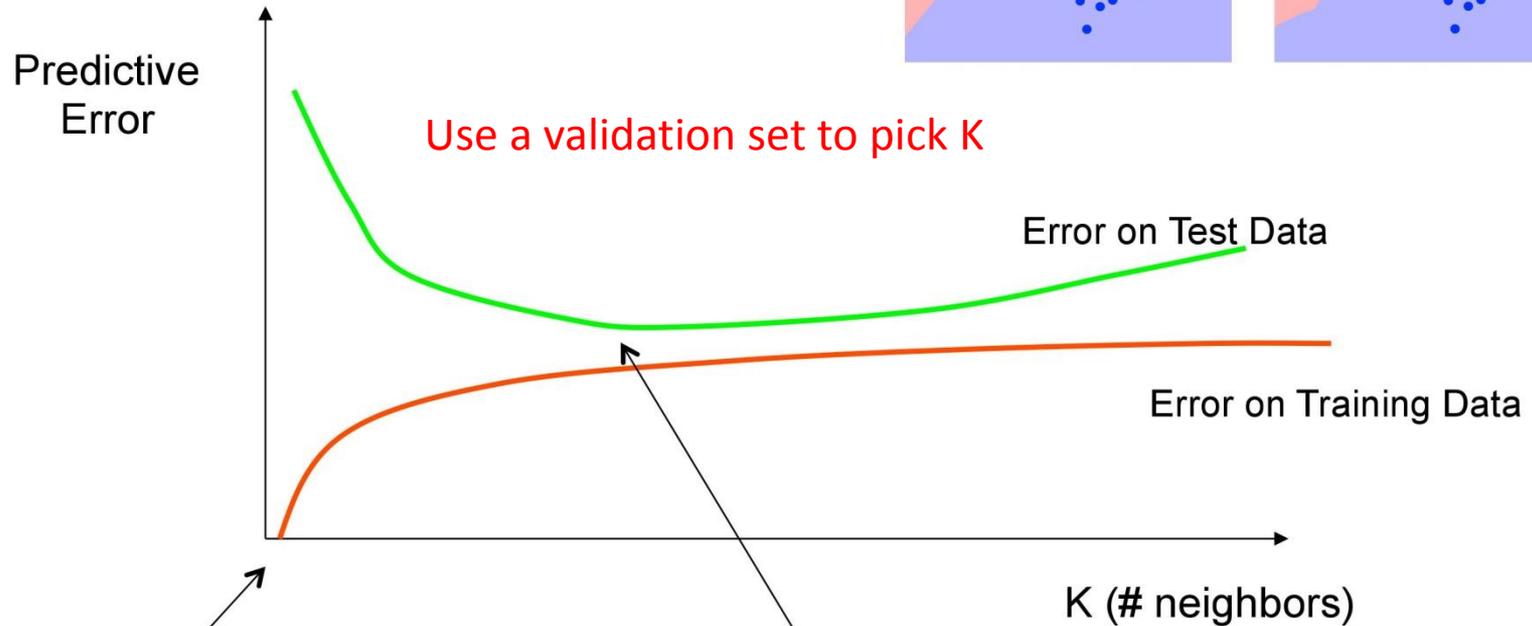
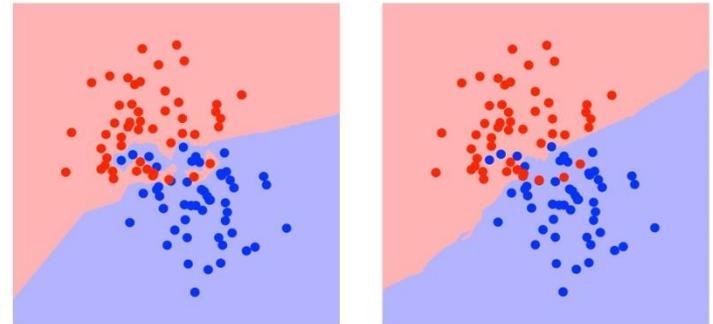
kNN Decision Boundary

- Increasing k “simplifies” decision boundary
 - Majority voting means less emphasis on individual points

$K = 25$



Error rates and K



K=1? Zero error!
Training data have been memorized...

Too complex

Best value of K

Problems with Instance-Based Learning

- Too many features?
 - Doesn't work well if large number of irrelevant features, distances overwhelmed by noisy features
 - Distances become meaningless in high dimensions (the **curse of dimensionality**)
- What is the impact of the **value of K**?
- Expensive
 - No learning: most real work done during testing
 - For every test sample, must search through all dataset – very slow!
 - **Must use tricks like approximate nearest neighbor search**
 - Need to store all training data

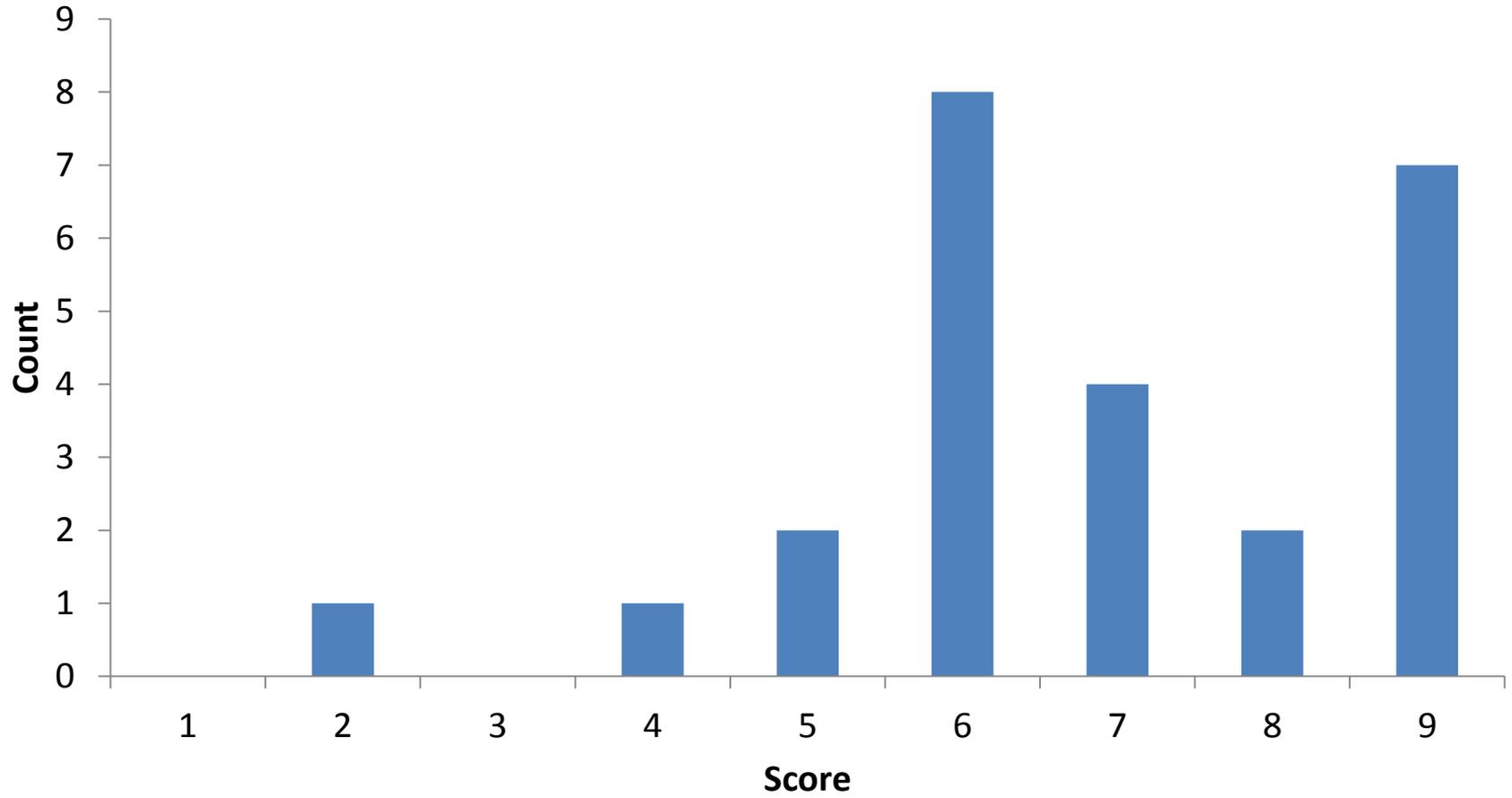
Approximate distance methods

- Build a balanced tree of data points (*kd*-trees), splitting data points along different dimensions
- Using tree, find “current best” guess of nearest neighbor
- Intelligently eliminate parts of the search space if they cannot contain a better “current best”
- Only search for neighbors up until some budget exhausted

Summary

- K-Nearest Neighbor is the most basic and simplest to implement classifier
- Cheap at training time, expensive at test time
- Unlike other methods we'll see later, naturally works for any number of classes
- Pick K through a validation set, use approximate methods for finding neighbors
- Success of classification depends on the amount of data and the meaningfulness of the distance function

Quiz 1



Review

- Today:
 - Linear algebra
 - PCA / dimensionality reduction
- Next time:
 - Mean shift vs k-means
 - Regularization and bias/variance

Other topics

- Hierarchical agglomerative clustering
 - Read my slides, can explain in office hours
- Graph cuts
 - Won't be on any test, I can explain in office hours
- Lagrange multipliers
 - Will cover in more depth later if needed
- HW1 Part I
 - We'll release solutions

Terminology

- Representation: The vector x_i for your data point i
- Dimensionality: the value of d in your $1 \times d$ vector representation x_i
- Let x_i^j be the j -th dimension of x_i , for $j = 1, \dots, d$
- Usually the different x_i^j are independent, can think of them as responses to different “tests” applied to x_i

Terminology

- For example: In HW1 Part V, I was asking you to try different combinations of RGB and gradients
- This means you can have a d -dimensional representation where d can be anything from 1 to 5
- There are 5 “tests” total (R, G, B, G_x, G_y)

The background features a large, faint watermark of the Stanford University seal. The seal is circular and contains a redwood tree in the center, with the text 'STANFORD UNIVERSITY' around the top and '1891' at the bottom. The words 'FREIHEIT WEHT' are also visible on the left side of the seal.

Linear Algebra Primer

Professor Fei-Fei Li
Stanford

Another, very in-depth linear algebra review from CS229 is available here:

<http://cs229.stanford.edu/section/cs229-linalg.pdf>

And a video discussion of linear algebra from EE263 is here (lectures 3 and 4):

<http://see.stanford.edu/see/lecturelist.aspx?coll=17005383-19c6-49ed-9497-2ba8bfcfe5f6>

Vectors and Matrices

- Vectors and matrices are just collections of ordered numbers that represent something: movements in space, scaling factors, word counts, movie ratings, pixel brightnesses, etc. We'll define some common uses and standard operations on them.

Vector

- A column vector $\mathbf{v} \in \mathbb{R}^{n \times 1}$ where

$$\mathbf{v} = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix}$$

- A row vector $\mathbf{v}^T \in \mathbb{R}^{1 \times n}$ where

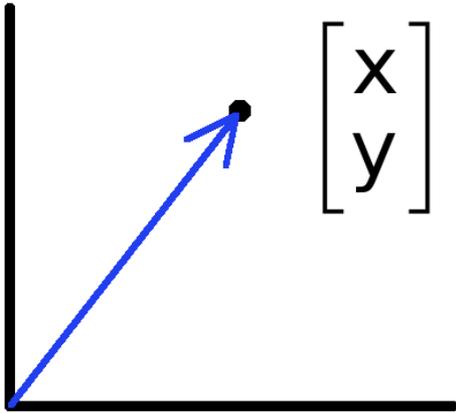
$$\mathbf{v}^T = [v_1 \quad v_2 \quad \dots \quad v_n]$$

T denotes the transpose operation

Vector

- You'll want to keep track of the orientation of your vectors when programming in MATLAB.
- You can transpose a vector V in MATLAB by writing V' .

Vectors have two main uses



- Vectors can represent an offset in 2D or 3D space
- Points are just vectors from the origin
- Data can also be treated as a vector
- Such vectors don't have a geometric interpretation, but calculations like "distance" still have value

Matrix

- A matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ is an array of numbers with size $m \downarrow$ by $n \rightarrow$, i.e. m rows and n columns.

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & a_{13} & \dots & a_{1n} \\ a_{21} & a_{22} & a_{23} & \dots & a_{2n} \\ \vdots & & & & \vdots \\ a_{m1} & a_{m2} & a_{m3} & \dots & a_{mn} \end{bmatrix}$$

- If $m = n$, we say that \mathbf{A} is square.

Matrix Operations

- Addition

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} + \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} = \begin{bmatrix} a + 1 & b + 2 \\ c + 3 & d + 4 \end{bmatrix}$$

- Can only add a matrix with matching dimensions, or a scalar.

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} + 7 = \begin{bmatrix} a + 7 & b + 7 \\ c + 7 & d + 7 \end{bmatrix}$$

- Scaling

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} \times 3 = \begin{bmatrix} 3a & 3b \\ 3c & 3d \end{bmatrix}$$

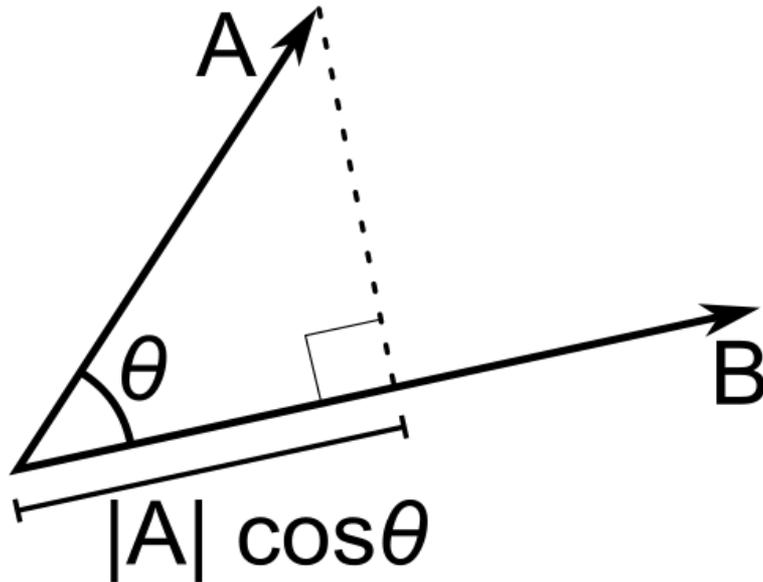
Matrix Operations

- Inner product (dot product) of vectors
 - Multiply corresponding entries of two vectors and add up the result
 - $\mathbf{x} \cdot \mathbf{y}$ is also $|\mathbf{x}| |\mathbf{y}| \cos(\text{the angle between } \mathbf{x} \text{ and } \mathbf{y})$

$$\mathbf{x}^T \mathbf{y} = \begin{bmatrix} x_1 & \dots & x_n \end{bmatrix} \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = \sum_{i=1}^n x_i y_i \quad (\text{scalar})$$

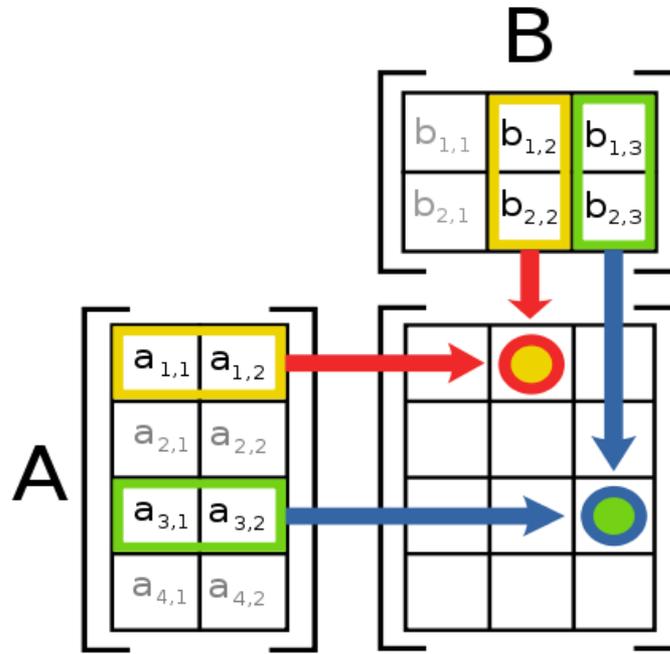
Matrix Operations

- Inner product (dot product) of vectors
 - If B is a unit vector, then $A \cdot B$ gives the length of A which lies in the direction of B (projection)



Matrix Operations

- Multiplication
- The product AB is:



- Each entry in the result is (that row of A) dot product with (that column of B)

Matrix Operations

- Multiplication example:

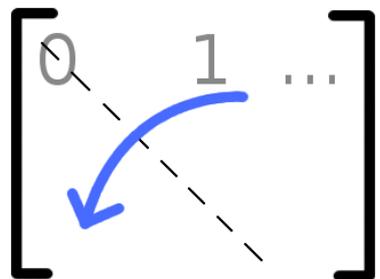
$$\begin{array}{ccc} A & \times & B \\ \downarrow & & \swarrow \\ \begin{bmatrix} 0 & 2 \\ 4 & 6 \end{bmatrix} & & \begin{bmatrix} 1 & 3 \\ 5 & 7 \end{bmatrix} \\ & & \begin{bmatrix} \square & 14 \\ \square & \square \end{bmatrix} \end{array}$$

$$0 \cdot 3 + 2 \cdot 7 = 14$$

- Each entry of the matrix product is made by taking the dot product of the corresponding row in the left matrix, with the corresponding column in the right one.

Matrix Operations

- Transpose – flip matrix, so row 1 becomes column 1


$$\begin{bmatrix} 0 & 1 & \dots \\ 2 & 3 & \\ 4 & 5 & \end{bmatrix}^T = \begin{bmatrix} 0 & 2 & 4 \\ 1 & 3 & 5 \end{bmatrix}$$

- A useful identity:

$$(ABC)^T = C^T B^T A^T$$

Special Matrices

- Identity matrix \mathbf{I}
 - Square matrix, 1's along diagonal, 0's elsewhere
 - $\mathbf{I} \cdot [\text{another matrix}] = [\text{that matrix}]$

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

- Diagonal matrix
 - Square matrix with numbers along diagonal, 0's elsewhere
 - A diagonal \cdot [another matrix] scales the rows of that matrix

$$\begin{bmatrix} 3 & 0 & 0 \\ 0 & 7 & 0 \\ 0 & 0 & 2.5 \end{bmatrix}$$

Special Matrices

- Symmetric matrix

$$\mathbf{A}^T = \mathbf{A}$$

$$\begin{bmatrix} 1 & 2 & 5 \\ 2 & 1 & 7 \\ 5 & 7 & 1 \end{bmatrix}$$

Inverse

- Given a matrix \mathbf{A} , its inverse \mathbf{A}^{-1} is a matrix such that $\mathbf{AA}^{-1} = \mathbf{A}^{-1}\mathbf{A} = \mathbf{I}$
- E.g.
$$\begin{bmatrix} 2 & 0 \\ 0 & 3 \end{bmatrix}^{-1} = \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & \frac{1}{3} \end{bmatrix}$$
- Inverse does not always exist. If \mathbf{A}^{-1} exists, \mathbf{A} is *invertible* or *non-singular*. Otherwise, it's *singular*.

Pseudo-inverse

- Say you have the matrix equation $AX=B$, where A and B are known, and you want to solve for X
- You could use MATLAB to calculate the inverse and premultiply by it: $A^{-1}AX=A^{-1}B \rightarrow X=A^{-1}B$
- MATLAB command would be **inv(A)*B**
- But calculating the inverse for large matrices often brings problems with computer floating-point resolution, or your matrix might not even have an inverse. Fortunately, there are workarounds.
- Instead of taking an inverse, directly ask MATLAB to solve for X in $AX=B$, by typing **A\B**
- MATLAB will try several appropriate numerical methods (including the pseudoinverse if the inverse doesn't exist)
- MATLAB will return the value of X which solves the equation
 - If there is no exact solution, it will return the closest one
 - If there are many solutions, it will return the smallest one

Matrix Operations

- MATLAB example:

$$AX = B$$

$$A = \begin{bmatrix} 2 & 2 \\ 3 & 4 \end{bmatrix}, B = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

```
>> x = A\B
```

```
x =
```

```
    1.0000
```

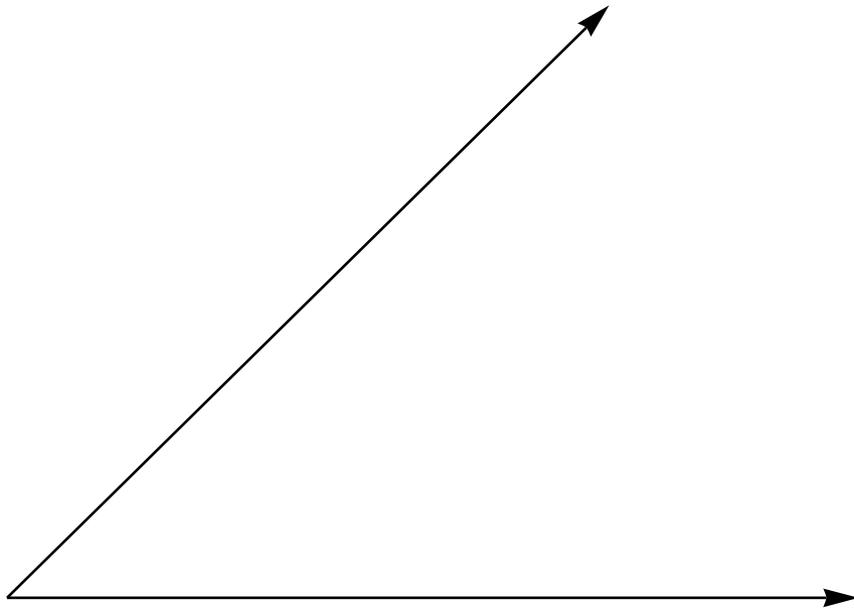
```
   -0.5000
```

Linear independence

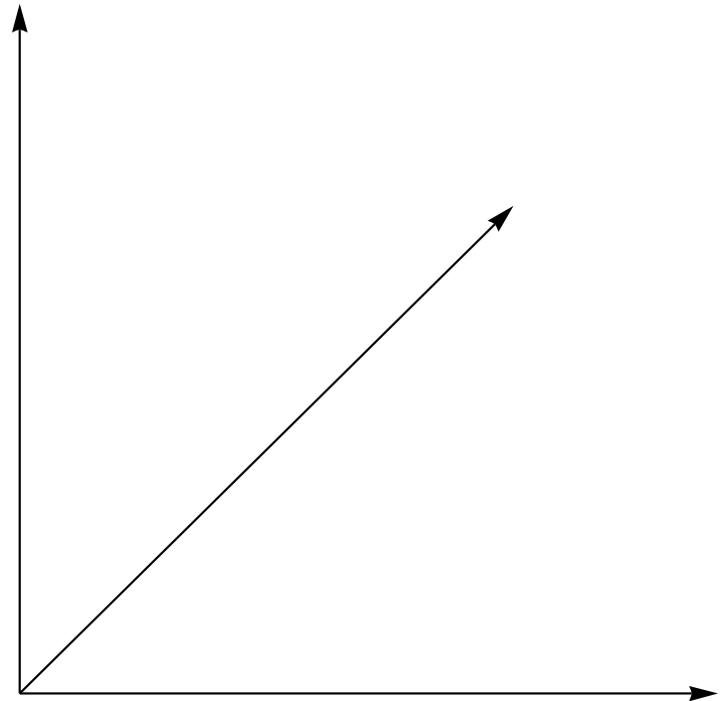
- Suppose we have a set of vectors $\mathbf{v}_1, \dots, \mathbf{v}_n$
- If we can express \mathbf{v}_1 as a linear combination of the other vectors $\mathbf{v}_2 \dots \mathbf{v}_n$, then \mathbf{v}_1 is linearly *dependent* on the other vectors.
 - The direction \mathbf{v}_1 can be expressed as a combination of the directions $\mathbf{v}_2 \dots \mathbf{v}_n$. (E.g. $\mathbf{v}_1 = .7 \mathbf{v}_2 - .7 \mathbf{v}_4$)
- If no vector is linearly dependent on the rest of the set, the set is linearly *independent*.
 - Common case: a set of vectors $\mathbf{v}_1, \dots, \mathbf{v}_n$ is always linearly independent if each vector is perpendicular to every other vector (and non-zero)

Linear independence

Linearly independent set



Not linearly independent



Matrix rank

- Column/row rank

col-rank(\mathbf{A}) = the maximum number of linearly independent column vectors of \mathbf{A}
row-rank(\mathbf{A}) = the maximum number of linearly independent row vectors of \mathbf{A}

- Column rank always equals row rank

- Matrix rank $\text{rank}(\mathbf{A}) \triangleq \text{col-rank}(\mathbf{A}) = \text{row-rank}(\mathbf{A})$

- If a matrix is not full rank, inverse doesn't exist
 - Inverse also doesn't exist for non-square matrices

Singular Value Decomposition (SVD)

- There are several computer algorithms that can “factor” a matrix, representing it as the product of some other matrices
- The most useful of these is the Singular Value Decomposition
- Represents any matrix **A** as a product of three matrices: **$U\Sigma V^T$**
- MATLAB command: **$[U,S,V] = \text{svd}(A);$**

Singular Value Decomposition (SVD)

$$\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T = \mathbf{A}$$

- Where \mathbf{U} and \mathbf{V} are rotation matrices, and $\mathbf{\Sigma}$ is a scaling matrix. For example:

$$\begin{matrix} U & & \Sigma & & V^T & & A \\ \begin{bmatrix} -.40 & .916 \\ .916 & .40 \end{bmatrix} & \times & \begin{bmatrix} 5.39 & 0 \\ 0 & 3.154 \end{bmatrix} & \times & \begin{bmatrix} -.05 & .999 \\ .999 & .05 \end{bmatrix} & = & \begin{bmatrix} 3 & -2 \\ 1 & 5 \end{bmatrix} \end{matrix}$$

Singular Value Decomposition (SVD)

- In general, if \mathbf{A} is $m \times n$, then \mathbf{U} will be $m \times m$, $\mathbf{\Sigma}$ will be $m \times n$, and \mathbf{V}^T will be $n \times n$.

$$\begin{matrix} U \\ \begin{bmatrix} -.39 & -.92 \\ -.92 & .39 \end{bmatrix} \end{matrix} \times \begin{matrix} \Sigma \\ \begin{bmatrix} 9.51 & 0 & 0 \\ 0 & .77 & 0 \end{bmatrix} \end{matrix} \times \begin{matrix} V^T \\ \begin{bmatrix} -.42 & -.57 & -.70 \\ .81 & .11 & -.58 \\ .41 & -.82 & .41 \end{bmatrix} \end{matrix} = \begin{matrix} A \\ \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \end{matrix}$$

Singular Value Decomposition

(SVD)

- \mathbf{U} and \mathbf{V} are always rotation matrices.
 - Geometric rotation may not be an applicable concept, depending on the matrix. So we call them “unitary” matrices – each column is a unit vector.
- $\mathbf{\Sigma}$ is a diagonal matrix
 - The number of nonzero entries = rank of \mathbf{A}
 - The algorithm always sorts the entries high to low

$$\begin{matrix} U & & \Sigma & & V^T & & A \\ \begin{bmatrix} -.39 & -.92 \\ -.92 & .39 \end{bmatrix} & \times & \begin{bmatrix} 9.51 & 0 & 0 \\ 0 & .77 & 0 \end{bmatrix} & \times & \begin{bmatrix} -.42 & -.57 & -.70 \\ .81 & .11 & -.58 \\ .41 & -.82 & .41 \end{bmatrix} & = & \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \end{matrix}$$

Singular Value Decomposition (SVD)

$$M = U\Sigma V^T$$

Illustration from Wikipedia

SVD Applications

- We've discussed SVD in terms of geometric transformation matrices
- But SVD of a data matrix can also be very useful
- To understand this, we'll look at a less geometric interpretation of what SVD is doing

SVD Applications

$$\begin{matrix} U & & \Sigma & & V^T & & A \\ \left[\begin{array}{cc} -.39 & -.92 \\ -.92 & .39 \end{array} \right] & \times & \left[\begin{array}{ccc} 9.51 & 0 & 0 \\ 0 & .77 & 0 \end{array} \right] & \times & \left[\begin{array}{ccc} -.42 & -.57 & -.70 \\ .81 & .11 & -.58 \\ .41 & -.82 & .41 \end{array} \right] & = & \left[\begin{array}{ccc} 1 & 2 & 3 \\ 4 & 5 & 6 \end{array} \right]
 \end{matrix}$$

- Look at how the multiplication works out, left to right:
- Column 1 of \mathbf{U} gets scaled by the first value from

$$\begin{matrix} \swarrow U\Sigma \\ \left[\begin{array}{ccc} -3.67 & -.71 & 0 \\ -8.8 & .30 & 0 \end{array} \right]
 \end{matrix}$$

- The resulting vector gets scaled by row 1 of \mathbf{V}^T to produce a contribution to the columns of \mathbf{A}

SVD Applications

$$\begin{array}{l}
 \begin{array}{c} U\Sigma \\ \left[\begin{array}{ccc} -3.67 & -0.71 & 0 \\ -8.8 & 0.30 & 0 \end{array} \right] \times \begin{array}{c} V^T \\ \left[\begin{array}{ccc} -0.42 & -0.57 & -0.70 \\ 0.81 & 0.11 & -0.58 \\ 0.41 & -0.82 & 0.41 \end{array} \right] \end{array} \\
 \\
 + \begin{array}{c} U\Sigma \\ \left[\begin{array}{ccc} -3.67 & -0.71 & 0 \\ -8.8 & 0.30 & 0 \end{array} \right] \times \begin{array}{c} V^T \\ \left[\begin{array}{ccc} -0.42 & -0.57 & -0.70 \\ 0.81 & 0.11 & -0.58 \\ 0.41 & -0.82 & 0.41 \end{array} \right] \end{array} \\
 \\
 = \begin{array}{c} A_{\text{partial}} \\ \left[\begin{array}{ccc} 1.6 & 2.1 & 2.6 \\ 3.8 & 5.0 & 6.2 \end{array} \right] \\
 \\
 A_{\text{partial}} \\ \left[\begin{array}{ccc} -0.6 & -0.1 & 0.4 \\ 0.2 & 0 & -0.2 \end{array} \right] \\
 \\
 \left[\begin{array}{ccc} 1 & 2 & 3 \\ 4 & 5 & 6 \end{array} \right]
 \end{array}
 \end{array}$$

- Each product of (*column i of \mathbf{U}*)·(*value i from Σ*)·(*row i of \mathbf{V}^T*) produces a component of the final \mathbf{A} .

SVD Applications

$$\begin{array}{c}
 U\Sigma \\
 \begin{bmatrix} -3.67 & -0.71 & 0 \\ -8.8 & 0.30 & 0 \end{bmatrix} \times \begin{array}{c} V^T \\ \begin{bmatrix} -0.42 & -0.57 & -0.70 \\ 0.81 & 0.11 & -0.58 \\ 0.41 & -0.82 & 0.41 \end{bmatrix} \end{array} \\
 \begin{array}{c} A_{\text{partial}} \\ \begin{bmatrix} 1.6 & 2.1 & 2.6 \\ 3.8 & 5.0 & 6.2 \end{bmatrix} \end{array}
 \end{array}
 \qquad
 \begin{array}{c}
 A \\
 \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}
 \end{array}$$

$$\begin{array}{c}
 U\Sigma \\
 \begin{bmatrix} -3.67 & -0.71 & 0 \\ -8.8 & 0.30 & 0 \end{bmatrix} \times \begin{array}{c} V^T \\ \begin{bmatrix} -0.42 & -0.57 & -0.70 \\ 0.81 & 0.11 & -0.58 \\ 0.41 & -0.82 & 0.41 \end{bmatrix} \end{array} \\
 \begin{array}{c} A_{\text{partial}} \\ \begin{bmatrix} -0.6 & -0.1 & 0.4 \\ 0.2 & 0 & -0.2 \end{bmatrix} \end{array}
 \end{array}$$

- We're building \mathbf{A} as a linear combination of the columns of \mathbf{U}
- Using all columns of \mathbf{U} , we'll rebuild the original matrix perfectly
- But, in real-world data, often we can just use the first few columns of \mathbf{U} and we'll get something close (e.g. the first $\mathbf{A}_{\text{partial}}$, above)

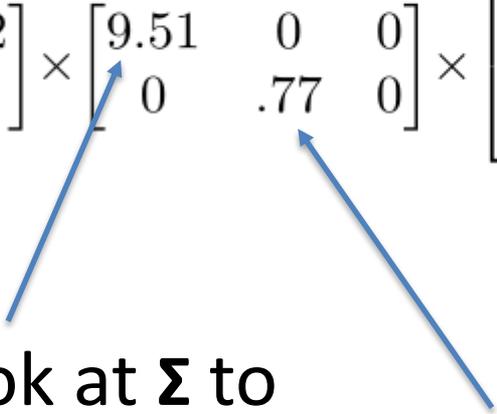
SVD Applications

$$\begin{array}{c}
 U\Sigma \\
 \begin{bmatrix} -3.67 & -.71 & 0 \\ -8.8 & .30 & 0 \end{bmatrix} \times \begin{array}{c} V^T \\ \begin{bmatrix} -.42 & -.57 & -.70 \\ .81 & .11 & -.58 \\ .41 & -.82 & .41 \end{bmatrix} \end{array} \\
 \begin{array}{c} A_{\text{partial}} \\ \begin{bmatrix} 1.6 & 2.1 & 2.6 \\ 3.8 & 5.0 & 6.2 \end{bmatrix} \end{array}
 \end{array}
 \qquad
 \begin{array}{c}
 A \\
 \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}
 \end{array}$$

$$\begin{array}{c}
 U\Sigma \\
 \begin{bmatrix} -3.67 & -.71 & 0 \\ -8.8 & .30 & 0 \end{bmatrix} \times \begin{array}{c} V^T \\ \begin{bmatrix} -.42 & -.57 & -.70 \\ .81 & .11 & -.58 \\ .41 & -.82 & .41 \end{bmatrix} \end{array} \\
 \begin{array}{c} A_{\text{partial}} \\ \begin{bmatrix} -.6 & -.1 & .4 \\ .2 & 0 & -.2 \end{bmatrix} \end{array}
 \end{array}$$

- We can call those first few columns of \mathbf{U} the *Principal Components* of the data
- They show the major patterns that can be added to produce the columns of the original matrix
- The rows of \mathbf{V}^T show how the *principal components* are mixed to produce the columns of the matrix

SVD Applications

$$\begin{matrix} U & & \Sigma & & V^T & & A \\ \begin{bmatrix} -.39 & -.92 \\ -.92 & .39 \end{bmatrix} & \times & \begin{bmatrix} 9.51 & 0 & 0 \\ 0 & .77 & 0 \end{bmatrix} & \times & \begin{bmatrix} -.42 & -.57 & -.70 \\ .81 & .11 & -.58 \\ .41 & -.82 & .41 \end{bmatrix} & = & \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \end{matrix}$$


We can look at Σ to see that the first column has a large effect

while the second column has a much smaller effect in this example

Principal Component Analysis

$$\begin{matrix} U\Sigma \\ \begin{bmatrix} -3.67 & -.71 & 0 \\ -8.8 & .30 & 0 \end{bmatrix} \end{matrix} \times \begin{matrix} V^T \\ \begin{bmatrix} -.42 & -.57 & -.70 \\ .81 & .11 & -.58 \\ .41 & -.82 & .41 \end{bmatrix} \end{matrix} = \begin{matrix} A_{\text{partial}} \\ \begin{bmatrix} 1.6 & 2.1 & 2.6 \\ 3.8 & 5.0 & 6.2 \end{bmatrix} \end{matrix}$$

- Remember, columns of \mathbf{U} are the *Principal Components* of the data: the major patterns that can be added to produce the columns of the original matrix
- One use of this is to construct a matrix where each column is a separate data sample
- Run SVD on that matrix, and look at the first few columns of \mathbf{U} to see patterns that are common among the columns
- This is called *Principal Component Analysis* (or PCA) of the data samples

Principal Component Analysis

$$\begin{matrix} U\Sigma \\ \begin{bmatrix} -3.67 & -.71 & 0 \\ -8.8 & .30 & 0 \end{bmatrix} \end{matrix} \times \begin{matrix} V^T \\ \begin{bmatrix} -.42 & -.57 & -.70 \\ .81 & .11 & -.58 \\ .41 & -.82 & .41 \end{bmatrix} \end{matrix} = \begin{matrix} A_{\text{partial}} \\ \begin{bmatrix} 1.6 & 2.1 & 2.6 \\ 3.8 & 5.0 & 6.2 \end{bmatrix} \end{matrix}$$

- Often, raw data samples have a lot of redundancy and patterns
- PCA can allow you to represent data samples as weights on the principal components, rather than using the original raw form of the data
- By representing each sample as just those weights, you can represent just the “meat” of what’s different between samples
- This minimal representation makes machine learning and other algorithms much more efficient

Addendum: How is SVD computed?

- For this class: tell MATLAB to do it. Use the result.
- But, if you're interested, one computer algorithm to do it makes use of Eigenvectors
 - The following material is presented to make SVD less of a “magical black box.” But you will do fine in this class if you treat SVD as a magical black box, as long as you remember its properties from the previous slides.

Eigenvector definition

- Suppose we have a square matrix \mathbf{A} . We can solve for vector x and scalar λ such that $Ax = \lambda x$
- In other words, find vectors where, if we transform them with \mathbf{A} , the only effect is to scale them with no change in direction.
- These vectors are called eigenvectors (German for “self vector” of the matrix), and the scaling factors λ are called eigenvalues
- An $m \times m$ matrix will have $\leq m$ eigenvectors where λ is nonzero

Finding eigenvectors

- Computers can find an x such that $Ax = \lambda x$ using this iterative algorithm:
 - $x = \text{random unit vector}$
 - while(x hasn't converged)
 - $x = Ax$
 - normalize x
- x will quickly converge to an eigenvector
- Some simple modifications will let this algorithm find all eigenvectors

Finding SVD

- Eigenvectors are for square matrices, but SVD is for all matrices
- To do $\text{svd}(A)$, computers can do this:
 - Take eigenvectors of AA^T (matrix is always square).
 - These eigenvectors are the columns of \mathbf{U} .
 - Square root of eigenvalues are the singular values (the entries of $\mathbf{\Sigma}$).
 - Take eigenvectors of $A^T A$ (matrix is always square).
 - These eigenvectors are columns of \mathbf{V} (or rows of \mathbf{V}^T)

Finding SVD

- Moral of the story: SVD is fast, even for large matrices
- It's useful for a lot of stuff
- There are also other algorithms to compute SVD or part of the SVD
 - MATLAB's `svd()` command has options to efficiently compute only what you need, if performance becomes an issue

A detailed geometric explanation of SVD is here:

<http://www.ams.org/samplings/feature-column/fcarc-svd>

Singular Value Decomposition

- Alternative method to calculate (still subtract mean 1st)
- Decompose $X = U S V^T$
 - Orthogonal: $X^T X = V S S V^T = V D V^T$
 - $X X^T = U S S U^T = U D U^T$
- $U \cdot S$ matrix provides coefficients
 - Example $x_i = (U_{i,1} S_{11})v_1 + (U_{i,2} S_{22})v_2 + \dots$
- Gives the least-squares approximation to X of this form

$$\boxed{\begin{matrix} X \\ N \times D \end{matrix}} \approx \boxed{\begin{matrix} U \\ N \times K \end{matrix}} \boxed{\begin{matrix} S \\ K \times K \end{matrix}} \boxed{\begin{matrix} V^T \\ K \times D \end{matrix}}$$

“Eigen-faces”

- “Eigen-X” = represent X using PCA
- Ex: Viola Jones data set
 - 24x24 images of faces = 576 dimensional measurements
 - Take first K PCA components

