



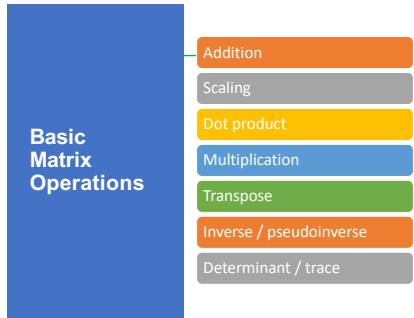
2- Linear Algebra and Matrix

Thomas W. Gyera, Assistant Professor
Computer and Information Science
University of Massachusetts Dartmouth

Courtesy of Fei-Fei Li: http://vision.stanford.edu/teaching/cs131_fall1617/

Vectors have two main uses

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



Matrix Operations

- Multiplication
- The product AB is:
- Each entry in the result is (that row of A) dot product with (that column of B)
- Many uses, which will be covered later

Outline



Vectors and Matrices
Basic operations
Special matrices



More Matrix Operations
Matrix inverse
Matrix rank
Singular Value Decomposition (SVD)
Use for image compression



Transformation Matrices
Homogeneous coordinates
Translation

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

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

T denotes the transpose operation

Vector

- We'll default to column vectors in this class

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

- You'll want to keep track of the orientation of your vectors when programming in MATLAB
- In MATLAB v' means v^T , indicating the transpose operation

Matrix

- A matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ is an array of numbers with size by, 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.

Images



$$\begin{bmatrix} 193 & 180 & 210 & 117 & 125 \\ 189 & 8 & 177 & 97 & 114 \\ 100 & 71 & 81 & 195 & 165 \\ 167 & 12 & 242 & 203 & 181 \\ 14 & 28 & 9 & 49 & 150 \end{bmatrix}$$

MATLAB represents an image as a matrix of pixel brightnesses

$$\begin{bmatrix} 193 & 180 & 210 & 117 & 125 \\ 189 & 8 & 177 & 97 & 114 \\ 100 & 71 & 81 & 195 & 165 \\ 167 & 12 & 242 & 203 & 181 \\ 14 & 28 & 9 & 49 & 150 \end{bmatrix}$$

Note that matrix coordinates are NOT Cartesian coordinates. The upper-left corner is $[y, x] = (1, 1)$

Color Images

- Grayscale images have one number per pixel, and are stored as an $m \times n$ matrix.
- Color images have 3 numbers per pixel – red, green, and blue brightnesses (RGB)
- Stored as an $m \times n \times 3$ matrix



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

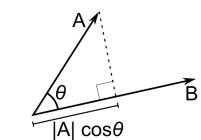
- Multiply corresponding entries of two vectors and add up the result
- $x \cdot y$ is also $|x| |y| \cos(\text{the angle between } x \text{ and } y)$

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

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



Matrix Operations

- Multiplication example:

$$A \times B$$

$$\begin{bmatrix} 0 & 2 \\ 4 & 6 \end{bmatrix} \begin{bmatrix} 10 & 14 \\ 34 & 54 \end{bmatrix}$$

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.

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

Matrix Operations

- Powers
 - By convention, we can refer to the matrix product AA as A^2 , and AAA as A^3 , etc.
 - Obviously only square matrices can be multiplied that way



Matrix Operations

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

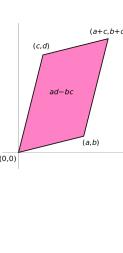
$$\begin{bmatrix} 0 & 1 & \dots \\ 2 & 3 & \\ 4 & 5 & \end{bmatrix}^T = \begin{bmatrix} 0 & 2 & 4 \\ 1 & 3 & 5 \end{bmatrix}$$

- A useful identity:

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

Matrix Operations

- Determinant
 - $\det(\mathbf{A})$ returns a scalar
 - Represents area (or volume) of the parallelogram described by the vectors in the rows of the matrix
 - For $\mathbf{A} = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$, $\det(\mathbf{A}) = ad - bc$
 - Properties: $\det(\mathbf{AB}) = \det(\mathbf{A})\det(\mathbf{B})$
 - $\det(\mathbf{A}^{-1}) = \frac{1}{\det(\mathbf{A})}$
 - $\det(\mathbf{A}^T) = \det(\mathbf{A})$
 - $\det(\mathbf{A}) = 0 \Leftrightarrow \mathbf{A}$ is singular



Outline

- Vectors and Matrices
 - Basic operations
 - Special matrices
- More Matrix Operations
 - Matrix inverse
 - Matrix rank
 - Singular Value Decomposition (SVD)
 - Use for image compression
- Transformation Matrices
 - Homogeneous coordinates
 - Translation

Matrix Operations

- Trace
 - $\text{tr}(\mathbf{A})$ = sum of diagonal elements
 - $\text{tr}\left[\begin{bmatrix} 1 & 3 \\ 5 & 7 \end{bmatrix}\right] = 1 + 7 = 8$
 - Invariant to a lot of transformations, so it's used sometimes in proofs. (Rarely in this class though.)
 - Properties:
 - $\text{tr}(\mathbf{AB}) = \text{tr}(\mathbf{BA})$
 - $\text{tr}(\mathbf{A} + \mathbf{B}) = \text{tr}(\mathbf{A}) + \text{tr}(\mathbf{B})$

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**.

Useful identities, for matrices that are invertible:

$$\begin{aligned} (\mathbf{A}^{-1})^{-1} &= \mathbf{A} \\ (\mathbf{AB})^{-1} &= \mathbf{B}^{-1}\mathbf{A}^{-1} \\ \mathbf{A}^{-T} &\triangleq (\mathbf{A}^T)^{-1} = (\mathbf{A}^{-1})^T \end{aligned}$$

Pseudoinverse

- 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
```

Matrix Rank

- Column/row rank

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

- Column rank always equals row rank
- Matrix rank

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

$$\begin{bmatrix} 1 & 0 & 1 \\ 1 & 1 & 0 \\ 2 & 0 & 2 \end{bmatrix}$$

What is the rank of the left matrix and why?

Matrix Rank

- For transformation matrices, the rank tells you the **dimensions** of the output
- E.g., if rank of \mathbf{A} is 1, then the transformation

$\mathbf{p}' = \mathbf{Ap}$
maps points onto a line.

- Here's a matrix with rank 1:

$$\begin{bmatrix} 1 & 1 \\ 2 & 2 \end{bmatrix} \times \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} x+y \\ 2x+2y \end{bmatrix}$$

All points get mapped to the line $y=2x$

Special Matrices

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

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

- Diagonal matrix
 - Square matrix with numbers along diagonal, 0's elsewhere
 - A diagonal \times [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
 $\mathbf{A}^T = \mathbf{A}$

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

- Skew-symmetric matrix
 $\mathbf{A}^T = -\mathbf{A}$

$$\begin{bmatrix} 0 & -2 & -5 \\ 2 & 0 & -7 \\ 5 & 7 & 0 \end{bmatrix}$$

Pseudoinverse



Fortunately, there are workarounds to solve $AX=B$ in these situations. And MATLAB can do them!



Instead of taking an inverse, directly ask MATLAB to solve for X in $AX=B$, by typing $\mathbf{A}\backslash\mathbf{B}$, called "mldivide".



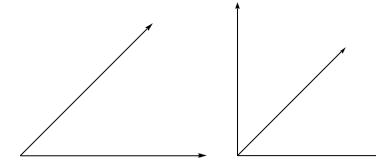
For complete instruction please check:
<http://www.mathworks.com/help/matlab/math/mldivide.html>

Pseudoinverse

- 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 pre-multiply by it: $A^{-1}Ax = A^{-1}B \rightarrow x = A^{-1}B$
- MATLAB command would be $\text{inv}(A)\mathbf{B}$
- But calculating the inverse for large matrices often brings problems with computer floating-point resolution (because it involves working with very small and very large numbers together).
- Or your matrix might not even have an inverse.

Linear Independence

Linearly independent set Not linearly independent



Linear Independence

- Suppose we have a set of vectors v_1, v_2, \dots, v_n
- If we can express v_i as a linear combination of the other vectors v_2, \dots, v_n , then v_i is linearly dependent on the other vectors
 - $v_1 = [0,1], v_2 = [0,2]$
- If no vector is linearly dependent on the rest of the set, the set is linearly independent
 - $v_1 = [1,0,0], v_2 = [0,1,0], v_3 = [0,0,1]$

There are several computer algorithms that can "factorize" a matrix, representing it as the product of some other matrices

The most useful of these is the Singular Value Decomposition (SVD).

Represent any matrix \mathbf{A} as a product of three matrices: $\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$

MATLAB command: $[\mathbf{U}, \mathbf{\Sigma}, \mathbf{V}] = \text{svd}(\mathbf{A})$

Matrix Rank

If an $m \times m$ matrix is rank m , we say it's "full rank"

- Maps an $m \times 1$ vector uniquely to another $m \times 1$ vector
- An inverse matrix can be found

If rank $< m$, we say it's "singular"

At least one dimension is getting collapsed. No way to look at the result and tell what the input was

- Inverse does not exist

Inverse also doesn't exist for non-square matrices

Singular Value Decomposition

Singular Value Decomposition

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

- Where U and V are rotation matrices, and Σ is a scaling matrix. For example:

$$\begin{bmatrix} U & \Sigma & V^T \\ -.40 & .916 & \begin{bmatrix} 5.39 & 0 \\ 0 & 3.154 \end{bmatrix} & \begin{bmatrix} .05 & .999 \\ .999 & .05 \end{bmatrix} \end{bmatrix} = \begin{bmatrix} A \\ 3 & -2 \\ 1 & 5 \\ 4 & 6 \end{bmatrix}$$

SVD Applications

$$\begin{aligned} & \begin{bmatrix} U\Sigma & V^T \\ -.367 & -.71 & 0 \\ -.88 & .30 & 0 \end{bmatrix} \times \begin{bmatrix} .42 & .57 & .70 \\ .81 & .11 & -.58 \\ .41 & .82 & .41 \end{bmatrix} = \begin{bmatrix} A_{partial} \\ 1.6 & 2.1 & 2.6 \\ 3.8 & 5.0 & 6.2 \\ 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \\ & + \begin{bmatrix} U\Sigma & V^T \\ -.367 & -.71 & 0 \\ -.88 & .30 & 0 \end{bmatrix} \times \begin{bmatrix} .42 & .57 & .70 \\ .81 & .11 & -.58 \\ .41 & .82 & .41 \end{bmatrix} = \begin{bmatrix} A_{partial} \\ 1.6 & 2.1 & 2.6 \\ 3.8 & 5.0 & 6.2 \\ 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \end{aligned}$$

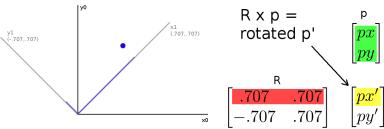
- Each product of (column i of U)-(value i from Σ)-(row i of V^T) produces a component of the final A .

Outline

- Vectors and Matrices**
 - Basic operations
 - Special matrices
- More Matrix Operations**
 - Matrix inverse
 - Matrix rank
 - Singular Value Decomposition (SVD)
 - Use for image compression
- Transformation Matrices**
 - Homogeneous coordinates
 - Translation

Rotation

- Insight: this is what happens in a matrix*vector multiplication
 - Result x coordinate is: [original vector] dot [matrix row 1]
 - So matrix multiplication can rotate a vector p:



Singular Value Decomposition

U and 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.

Σ is a diagonal matrix

- The number of nonzero entries = rank of A
- The algorithm always sorts the entries high to low

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

SVD Applications

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

SVD Applications

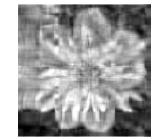
$$\begin{bmatrix} U & \Sigma & V^T \\ -.39 & -.92 & \begin{bmatrix} 9.51 & 0 & 0 \\ 0 & .77 & 0 \\ .41 & -.82 & .41 \end{bmatrix} & \begin{bmatrix} .42 & .57 & .70 \\ .81 & .11 & -.58 \\ .41 & .82 & .41 \end{bmatrix} \end{bmatrix} = \begin{bmatrix} A \\ 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}$$

- Look at how the multiplication works out, left to right:
- Column 1 of U gets scaled by the first value from Σ .

$$\begin{bmatrix} U\Sigma & V^T \\ -.367 & -.71 & 0 & \begin{bmatrix} .42 & .57 & .70 \\ .81 & .11 & -.58 \\ .41 & .82 & .41 \end{bmatrix} & \begin{bmatrix} A_{partial} \\ 1.6 & 2.1 & 2.6 \\ 3.8 & 5.0 & 6.2 \end{bmatrix} \end{bmatrix}$$

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

SVD Applications



For this image, using only the first 10 of 300 principal components produces a recognizable reconstruction

So, SVD can be used for image compression

SVD Applications

$$\begin{bmatrix} U\Sigma & V^T \\ -.367 & -.71 & 0 & \begin{bmatrix} .42 & .57 & .70 \\ .81 & .11 & -.58 \\ .41 & .82 & .41 \end{bmatrix} & \begin{bmatrix} A_{partial} \\ 1.6 & 2.1 & 2.6 \\ 3.8 & 5.0 & 6.2 \end{bmatrix} \end{bmatrix}$$

- We can call those first few columns of the data 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 V^T show how the principal components are mixed to produce the columns of the matrix

SVD Applications

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

- 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

Transformation

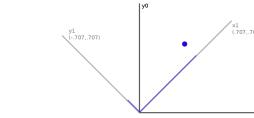
- Matrices can be used to transform vectors in useful ways, through multiplication: $x' = Ax$
- Simplest is scaling:

$$\begin{bmatrix} s_x & 0 \\ 0 & s_y \end{bmatrix} \times \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} s_x x \\ s_y y \end{bmatrix}$$

- (Verify to yourself that the matrix multiplication works out this way)

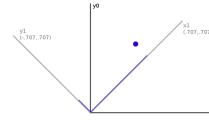
Transformation

- How can you convert a vector represented in frame "0" to a rotated coordinate frame "1"?
- Remember what a vector is: [component in direction of the frame's x axis, and component in direction of y axis]



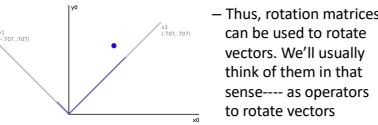
Rotation

- So to rotate it we must produce this vector: [component in direction of new x axis, component in direction of new y axis]
- We can do this easily with dot products!
- New x coordinate is [original vector] dot [the new x axis]
- New y coordinate is [original vector] dot [the new y axis]



Rotation

- Suppose we express a point in the new coordinate system which is rotated left
- If we plot the result in the original coordinate system, we have rotated the point right



2D Rotation Matrix Formula

Counter-clockwise rotation by an angle θ

$$\begin{aligned} x' &= \cos\theta x - \sin\theta y \\ y' &= \cos\theta y + \sin\theta x \\ \begin{bmatrix} x' \\ y' \end{bmatrix} &= \begin{bmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} \\ P' &= R P \end{aligned}$$

Transformation Matrices

- Multiple transformation matrices can be used to transform a point: $p' = R_2 R_1 S p$
- The effect of this is to apply their transformations one after the other, from right to left.
- In the example above, the result is $(R_1 (R_2 (S p)))$
- The result is exactly the same if we multiply the matrices first, to form a single transformation matrix.
- $p' = (R_2 R_1 S) p$

Homogeneous System

- In general, a matrix multiplication lets us linearly combine components of a vector

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} \times \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} ax + by \\ cx + dy \end{bmatrix}$$

- This is sufficient for scale, rotate, skew transformations.
- But notice, we can't add a constant! \oplus

Homogeneous System

- The (somewhat hacky) solution? Stick a "1" at the end of every vector:

$$\begin{bmatrix} a & b & c \\ d & e & f \\ 0 & 0 & 1 \end{bmatrix} \times \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \begin{bmatrix} ax + by + c \\ dx + ey + f \\ 1 \end{bmatrix}$$

- Now we can rotate, scale, and skew like before, AND translate (note how the multiplication works out, above)
- This is called "homogeneous coordinates"

Homogeneous System

- In homogeneous coordinates, the multiplication works out so the rightmost column of the matrix is a vector that gets added.

$$\begin{bmatrix} a & b & c \\ d & e & f \\ 0 & 0 & 1 \end{bmatrix} \times \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \begin{bmatrix} ax + by + c \\ dx + ey + f \\ 1 \end{bmatrix}$$

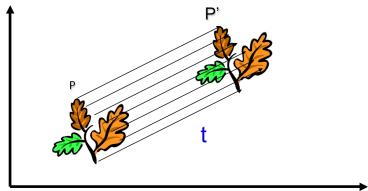
- Generally, a homogeneous transformation matrix will have a bottom row of [0 0 1], so that the result has a "1" at the bottom too.

Homogeneous System

- One more thing we might want: to divide the result by something
 - For example, we may want to divide by a coordinate, to make things scale down as they get farther away in a camera image
 - Matrix multiplication can't actually divide
 - So, by convention, in homogeneous coordinates, we'll divide the result by its last coordinate after doing a matrix multiplication

$$\begin{bmatrix} x \\ y \\ 7 \end{bmatrix} \Rightarrow \begin{bmatrix} x/7 \\ y/7 \\ 1 \end{bmatrix}$$

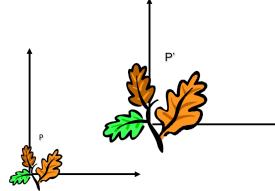
2D Translation



Using Homogeneous Coordinates

$$\begin{aligned} P &= (x, y) \rightarrow (x, y, 1) \\ t &= (t_x, t_y) \rightarrow (t_x, t_y, 1) \\ P' &\rightarrow \begin{bmatrix} x+t_x \\ y+t_y \\ 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & t_x \\ 0 & 1 & t_y \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} \\ &= \begin{bmatrix} 1 & 0 & t_x \\ 0 & 1 & t_y \\ 0 & 0 & 1 \end{bmatrix} \cdot P = T \cdot P \end{aligned}$$

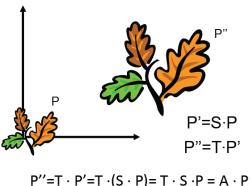
Scaling



Scaling Equation

$$\begin{aligned} P &= (x, y) \rightarrow P' = (s_x x, s_y y) \\ P &= (x, y) \rightarrow (x, y, 1) \\ P' &= (s_x x, s_y y) \rightarrow (s_x x, s_y y, 1) \\ P' &\rightarrow \begin{bmatrix} s_x x \\ s_y y \\ 1 \end{bmatrix} = \begin{bmatrix} s_x & 0 & 0 \\ 0 & s_y & 0 \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \begin{bmatrix} S & 0 \\ 0 & 1 \end{bmatrix} \cdot P = S \cdot P \end{aligned}$$

Scaling & Translation



Scaling & Translation

$$\begin{aligned} P'' &= T \cdot S \cdot P = \begin{bmatrix} 1 & 0 & t_x \\ 0 & 1 & t_y \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} s_x & 0 & 0 \\ 0 & s_y & 0 \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \\ &= \underbrace{\begin{bmatrix} s_x & 0 & t_x \\ 0 & s_y & t_y \\ 0 & 0 & 1 \end{bmatrix}}_A \cdot \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \begin{bmatrix} s_x x + t_x \\ s_y y + t_y \\ 1 \end{bmatrix} = \begin{bmatrix} S & t \\ 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = A \cdot P \end{aligned}$$

Order Matters

$$\begin{aligned} P''' &= T \cdot S \cdot P = \begin{bmatrix} 1 & 0 & t_x \\ 0 & 1 & t_y \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} s_x & 0 & 0 \\ 0 & s_y & 0 \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \begin{bmatrix} s_x & 0 & t_x \\ 0 & s_y & t_y \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \begin{bmatrix} s_x x + t_x \\ s_y y + t_y \\ 1 \end{bmatrix} \\ P''' &= S \cdot T \cdot P = \begin{bmatrix} s_x & 0 & 0 \\ 0 & s_y & 0 \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 1 & 0 & t_x \\ 0 & 1 & t_y \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \\ &= \begin{bmatrix} s_x & 0 & s_x t_x \\ 0 & s_y & s_y t_y \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \begin{bmatrix} s_x x + s_x t_x \\ s_y y + s_y t_y \\ 1 \end{bmatrix} \end{aligned}$$

Rotation Equation

Counter-clockwise rotation by an angle θ

$$\begin{aligned} x' &= \cos \theta x - \sin \theta y \\ y' &= \sin \theta x + \cos \theta y \\ \begin{bmatrix} x' \\ y' \end{bmatrix} &= \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} \\ P' &= R P \end{aligned}$$

Rotation Matrix Properties

- Transpose of a rotation matrix produces a rotation in the opposite direction

$$\begin{aligned} R \cdot R^T &= R^T \cdot R = I \\ \det(R) &= 1 \end{aligned}$$

- The rows of a rotation matrix are always mutually perpendicular (a.k.a. orthogonal) unit vectors
-(and so are its columns)

Properties

$$\begin{bmatrix} x' \\ y' \end{bmatrix} = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$$

A 2D rotation matrix is 2x2

Note: R belongs to the category of *normal* matrices and satisfies many interesting properties:

$$\begin{aligned} R \cdot R^T &= R^T \cdot R = I \\ \det(R) &= 1 \end{aligned}$$

Scaling + Translation + Rotation

$$\begin{aligned} P' &= (T R S) P \\ P' &= T \cdot R \cdot S \cdot P = \begin{bmatrix} 1 & 0 & t_x \\ 0 & 1 & t_y \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} \cos \theta & 0 & 0 \\ 0 & \sin \theta & 0 \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} s_x & 0 & 0 \\ 0 & s_y & 0 \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \\ &= \begin{bmatrix} \cos \theta & -\sin \theta & t_x \\ \sin \theta & \cos \theta & t_y \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \begin{bmatrix} s_x x + s_y t_x \\ s_y y + s_x t_y \\ 1 \end{bmatrix} \\ &= \begin{bmatrix} R & t \\ 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} S & 0 \\ 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \boxed{\begin{bmatrix} R S & t \\ 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}} \end{aligned}$$

This is the form of the general-purpose transformation matrix



3- Introduction to Data Mining

Thomas W Gyeera, Assistant Professor
Computer and Information Science
University of Massachusetts Dartmouth

Courtesy to Prof. Panayiotis Tsaparas

The data is also very complex

- Multiple types of data: tables, time series, images, graphs, etc.
- Spatial and temporal aspects
- Interconnected data of different types:
 - From the mobile phone we can collect, location of the user, friendship information, check-ins to venues, opinions through twitter, images through cameras, queries to search engines.

Example: transaction data

- Billions of real-life customers:**
- WALMART: 20M transactions per day
 - AT&T 300M calls per day
 - Credit card companies: billions of transactions per day.

The point cards allow companies to collect information about specific users

Why do we need data mining?

Huge amounts of raw data!

- In the digital age, TB of data is generated by the second
 - Mobile devices, digital photographs, web documents.
 - Facebook updates, Tweets, Blogs, User-generated content
 - Transactions, sensor data, surveillance data
 - Queries, clicks, browsing
 - Cheap storage has made possible to maintain this data

Need to analyze the raw data to extract knowledge

Why do we need data mining?

"The data is the computer"

- Large amounts of data can be more powerful than complex algorithms and models
- Google has solved many Natural Language Processing problems, simply by
 - Example: misspellings, synonyms
- Data is power
 - Today, the collected data is one of the biggest assets of an online company
 - Data of users
 - The friendship and updates of Facebook
 - Tweets and follows of Twitter
 - Amazon transactions

We need a way to harness the collective intelligence

Example: genomic sequences

<http://www.1000genomes.org/page.php>

Full sequence of 1000 individuals

3×10^9 nucleotides per person $\rightarrow 3 \times 10^{12}$ nucleotides

Lots more data in fact: medical history of the persons, gene expression data

Climate data (just an example)

- A database of temperature, precipitation and pressure records managed by the National Climatic Data Center, Arizona State University and the Carbon Dioxide Information Analysis Center
- 6000 temperature stations, 7500 precipitation stations, 2000 pressure stations*
- Spatiotemporal data

Example: environmental data

Mobile phones today record a large amount of information about the user behavior

- GPS records position
- Camera products images
- Communication via phone and SMS
- Search engines and browser
- Association with entities via check-ins

Amazon collects all the items that you browsed, placed into your basket, read reviews about, purchased.

Google and Bing record all your browsing activity via toolbar plugins. They also record the queries you asked, the pages you saw and the clicks you did.

Data collected for millions of users on a daily basis

Mobile phones today record a large amount of information about the user behavior

- GPS records position
- Camera products images
- Communication via phone and SMS
- Search engines and browser
- Association with entities via check-ins

Amazon collects all the items that you browsed, placed into your basket, read reviews about, purchased.

Google and Bing record all your browsing activity via toolbar plugins. They also record the queries you asked, the pages you saw and the clicks you did.

Data collected for millions of users on a daily basis

Behavioral data

Mobile phones today record a large amount of information about the user behavior

- GPS records position
- Camera products images
- Communication via phone and SMS
- Search engines and browser
- Association with entities via check-ins

Amazon collects all the items that you browsed, placed into your basket, read reviews about, purchased.

Google and Bing record all your browsing activity via toolbar plugins. They also record the queries you asked, the pages you saw and the clicks you did.

Data collected for millions of users on a daily basis

Mobile phones today record a large amount of information about the user behavior

- GPS records position
- Camera products images
- Communication via phone and SMS
- Search engines and browser
- Association with entities via check-ins

Amazon collects all the items that you browsed, placed into your basket, read reviews about, purchased.

Google and Bing record all your browsing activity via toolbar plugins. They also record the queries you asked, the pages you saw and the clicks you did.

Data collected for millions of users on a daily basis

Types of Attributes

- There are different types of attributes
 - Categorical
 - Examples: eye color, zip codes, words, rankings (e.g., good, fair, bad), height in (tall, medium, short)
 - Nominal (no order or comparison) vs Ordinal (order but not comparable)
 - Numeric
 - Examples: dates, temperature, time, length, value, count.
 - Discrete (counts) vs Continuous (temperature)
 - Special case: Binary attributes (yes/no, exists/not exists)

Projection of x Load

| Projection of x Load | Projection of y Load | Distance | Load | Thickness |
|------------------------|------------------------|----------|------|-----------|
| 10.23 | 5.27 | 15.22 | 2.7 | 1.2 |
| 12.65 | 6.25 | 16.22 | 2.2 | 1.1 |

Numeric Record Data

- If data objects have the same **fixed set** of **numeric attributes**, then the data objects can be thought of as **points** in a multi-dimensional space, where each dimension represents a distinct attribute
- Such data set can be represented by an **n-by-d data matrix**, where there are **n** rows, one for each object, and **d** columns, one for each attribute

| Projection of x Load | Projection of y Load | Distance | Load | Thickness |
|------------------------|------------------------|----------|------|-----------|
| 10.23 | 5.27 | 15.22 | 2.7 | 1.2 |
| 12.65 | 6.25 | 16.22 | 2.2 | 1.1 |

Categorical Data

- Data that consists of a collection of records, each of which consists of a **fixed set of categorical attributes**

| Id | Refund | Marital Status | Taxable Income | Cheat |
|----|--------|----------------|----------------|-------|
| 1 | Yes | Single | High | No |
| 2 | No | Married | Medium | No |
| 3 | No | Single | Low | No |
| 4 | Yes | Married | High | No |
| 5 | No | Divorced | Medium | Yes |
| 6 | No | Married | Low | No |
| 7 | Yes | Divorced | Medium | No |
| 8 | No | Single | Medium | Yes |
| 9 | No | Married | Medium | No |
| 10 | No | Single | Medium | Yes |

Document Data

- Each document becomes a 'term' vector, each term is a component (attribute) of the vector, the value of each component is the number of times the corresponding term occurs in the document.
- Bag-of-words representation – no ordering

| | Brain | Computer | Art | Big | Machine | Learning | AI | ML | Deep | Neural | TensorFlow | PyTorch |
|------------|-------|----------|-----|-----|---------|----------|----|----|------|--------|------------|---------|
| Document 1 | 3 | 0 | 5 | 1 | 2 | 6 | 1 | 0 | 2 | 0 | 2 | 1 |
| Document 2 | 0 | 7 | 0 | 2 | 1 | 0 | 0 | 3 | 0 | 0 | 0 | 0 |
| Document 3 | 0 | 1 | 0 | 0 | 1 | 2 | 2 | 0 | 3 | 0 | 0 | 0 |



Transaction Data

- Each record (transaction) is a set of items.

| TID | Items |
|-----|---------------------------|
| 1 | Bread, Coke, Milk |
| 2 | Beer, Bread |
| 3 | Beer, Coke, Diaper, Milk |
| 4 | Beer, Bread, Diaper, Milk |
| 5 | Coke, Diaper, Milk |

- A set of items can also be represented as a binary vector, where each attribute is an item.
- A document can also be represented as a set of words (no counts)

Sparisty: average number of products bought by a customer

What can you do with the data?

- Suppose that you are the owner of a supermarket and you have collected billions of market basket data. What information would you extract from it and how would you use it?

| TID | Items |
|-----|---------------------------|
| 1 | Bread, Coke, Milk |
| 2 | Beer, Bread |
| 3 | Beer, Coke, Diaper, Milk |
| 4 | Beer, Bread, Diaper, Milk |
| 5 | Coke, Diaper, Milk |

- What if this was an online store?

Ordered Data

- Genomic sequence data

```
GGTTCCGCCCTTCAGCCCCGGCGC  
CGCAGGGGCCGCCGGCGCGCGCG  
GAGAAAGGCCGCCGGCGCTGGCGGGG  
GGGGGAGGCCGCCGGCGCGCGCG  
CCAAACCGAGTCGCACCCAGTGGCC  
CCCTCTGCTCGGCCCTAGACCTGA  
GCTCATTAAGCGGCAGCGACAG  
GCCAAGTAGAACACCGCGAACGCG  
TGGCGTGGCTGCGACCGGG
```

- Data is a long ordered string

Ordered Data

- Time series
 - Sequence of ordered (over "time") numeric values.



What can you do with the data?

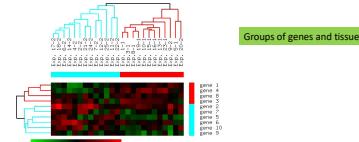
- Suppose you are a search engine and you have a toolbar log consisting of
 - pages browsed,
 - queries,
 - pages clicked,
 - ads clicked

Ad click prediction
Query reformulations

- each with a user id and a timestamp. What information would you like to get out of the data?

What can you do with the data?

- Suppose you are biologist who has microarray expression data: thousands of genes, and their expression values over thousands of different settings (e.g. tissues). What information would you like to get out of your data?



What can you do with the data?

- You are the owner of a social network, and you have full access to the social graph, what kind of information do you want to get out of your graph?

- Who is the most important node in the graph?
- What is the shortest path between two nodes?
- How many friends two nodes have in common?
- How does information spread on the network?

Why data mining?

Commercial point of view
Data has become the key competitive advantage of companies.
Examples: Facebook, Google, Amazon
Being able to extract useful information out of the data is key for exploiting them commercially.

Scientific point of view
Scientists are at an unprecedented position where they can collect TB of information.
Cloud of sensor data, astronomy data, social network data, gene data
We need the tools to analyze such data to get a better understanding of the world and advance science

Scale (in data size and feature dimension)
Why not use traditional analytic methods?
Enormity of data, curse of dimensionality
The amount and the complexity of data does not allow for manual processing of the data. We need automated techniques.

What is Data Mining again?

Data mining is the analysis of (often large) observational data sets to find unsuspected relationships and to summarize the data in novel ways that are both understandable and useful to the data analyst. (Hand, Mannila, Smyth)

Data mining is the discovery of models for data (Rajaraman, Ullman)

We can have the following types of models:
Models that explain the data (e.g., a single instance).
Models that predict the data.
Models that summarize the data.
Models that extract the most prominent features of the data.

Frequent Itemsets and Association Rules

- Given a set of records each of which contain some number of items from a given collection;
- Identify sets of items (itemsets) occurring frequently together
- Produce dependency rules which will predict occurrence of an item based on occurrences of other items.

| TID | Items |
|-----|---------------------------|
| 1 | Bread, Coke, Milk |
| 2 | Beer, Bread |
| 3 | Beer, Coke, Diaper, Milk |
| 4 | Beer, Bread, Diaper, Milk |
| 5 | Coke, Diaper, Milk |

Itemsets Discovered:
(Milk, Coke)
(Diaper, Milk)

Rules Discovered:
(Milk → (Coke))
(Diaper, Milk) → (Beer)

Frequent Item sets: Applications

Text mining: finding associated phrases in text

- There are lots of documents that contain the phrases "association rules", "data mining" and "efficient algorithm"
- Users who buy this item often buy this item as well
- Users who watched James Bond movies, also watched Jason Bourne movies.
- Recommendations make use of item and user similarity

Association Rule Discovery: Application

- Supermarket shelf management.
 - Goal: To identify items that are bought together by sufficiently many customers.
 - Approach: Process the point-of-sale data collected with barcode scanners to find dependencies among items.
 - A classic rule --
 - If a customer buys diaper and milk, then he is very likely to buy beer.
 - So, don't be surprised if you find six-packs stacked next to diapers!

Types of data

- Numeric data: Each object is a point in a multidimensional space
- Categorical data: Each object is a vector of categorical values
- Set data: Each object is a set of values (with or without counts)
 - Sets can also be represented as binary vectors, or vectors of counts
- Ordered sequences: Each object is an ordered sequence of values.
- Graph data

What can you do with the data?

- Suppose you are a stock broker and you observe the fluctuations of multiple stocks over time. What information would you like to get out of your data?



What can we do with data mining?

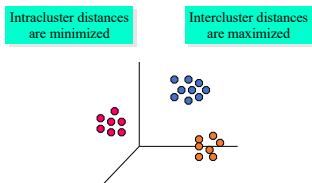
- Some examples:
 - Frequent itemsets and Association Rules
 - Clustering
 - Classification
 - Ranking
 - Exploratory analysis

Clustering Definition

- Given a set of data points, each having a set of attributes, and a similarity measure among them, find clusters such that
 - Data points in one cluster are more similar to one another.
 - Data points in separate clusters are less similar to one another.
- Similarity Measures?
 - Euclidean Distance if attributes are continuous.
 - Other Problem-specific Measures.

Illustrating Clustering

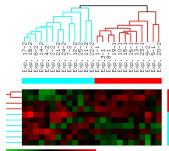
Euclidean Distance Based Clustering in 3-D space.



Tan, M. Steinbach and V. Kumar, Introduction to Data Mining

Clustering: Application 1

- Bioinformatics applications:
- Goal: Group genes and tissues together such that genes are coexpressed on the same tissues



Clustering: Application 2

- Document Clustering:
 - Goal: To find groups of documents that are similar to each other based on the important terms appearing in them.
 - Approach: To identify frequently occurring terms in each document. Form a similarity measure based on the frequencies of different terms. Use it to cluster.
 - Gain: Information Retrieval can utilize the clusters to relate a new document or search term to clustered documents.

Clustering of S&P 500 Stock Data

- Observe Stock Movements every day.
- Cluster stocks if they change similarly over time.

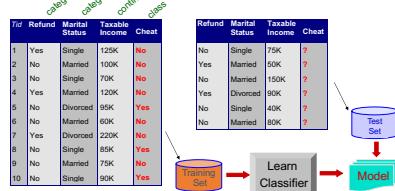
| Discussed Clusters | Industry Group |
|---|------------------|
| Apple-Microsoft-Samsung-Nvidia-AMD-AT&T-Verizon- Cisco-Sys-DOWN-INTEL-DOWN-MSFT-Logi-DOYEN- Motorola-Comcast-NTT- Natl-Semiconductors-DOWN-Google-DOYEN-SG-DOYEN- | Technology1-DOWN |
| Apple-Cisco-DOYEN-Avaya-DOYEN-DEC-DOYEN- CSCO-DOYEN-DOYEN-DOYEN-DOYEN-DOYEN- Compaq-DOYEN-EMC-DOYEN-Gate-DOYEN- Motorola-DOYEN-NTT- Natl-Semiconductors-DOYEN-DOYEN-DOYEN | Technology2-DOWN |
| Fannie-Mae-DOYEN-Pet-Homes-Louis-DOYEN- MBSNA-Cap-DOYEN-Morgan-Stanley-DOYEN | Financial-DOWN |
| Risk-Master-UP-Dowm-Index-UP-Halliburton-HEI-DOYEN- Lambert-Lockheed-Martin-DOYEN-DOYEN-DOYEN- Schlumberger-DOYEN-DOYEN-DOYEN | Oil-UP |

Tan, M. Steinbach and V. Kumar, Introduction to Data Mining

Classification: Definition



Classification Example



Tan, M. Steinbach and V. Kumar, Introduction to Data Mining

Classification: Application 1

- Ad Click Prediction
 - Goal: Predict if a user that visits a web page will click on a displayed ad. Use it to target users with high click probability.
- Approach:
 - Collected data for users over a period of time and record who clicks and who does not. The (click, no click) information forms the class attribute.
 - Use the history of the user (web pages browsed, queries issued) as the features.
 - Learn a classifier model and test on new users.

Classification: Application 2

- Fraud Detection
 - Goal: Predict fraudulent cases in credit card transactions.
- Approach:
 - Use credit card transactions and the information on its account-holder as attributes.
 - When does a customer buy, what does he buy, how often he pays on time, etc.
 - Label past transactions as fraud or fair transactions. This forms the class attribute.
 - Learn a model for the class of the transactions.
 - Use this model to detect fraud by observing credit card transactions on an account.

Tan, M. Steinbach and V. Kumar, Introduction to Data Mining

Link Analysis Ranking

Given a collection of web pages that are linked to each other, rank the pages according to importance (authoritativeness) in the graph

- Intuition: A page gains authority if it is linked to by another page.

Application: When retrieving pages, the authoritativeness is factored in the ranking.

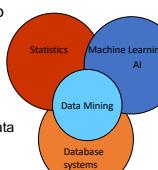
Exploratory Analysis

Trying to understand the data as a physical phenomenon, and describe them with simple metrics

- What does the web graph look like?
- How often do people repeat the same query?
- Are friends in facebook also friends in twitter?

The important thing is to find the right metrics and ask the right questions

It helps our understanding of the world and can lead to models of the phenomena we observe.



Connections of Data Mining with others

- Draws ideas from machine learning/AI, pattern recognition, statistics, and database systems
- Traditional Techniques may be unsuitable due to
 - Enormity of data
 - High dimensionality of data
 - Heterogeneous, distributed nature of data
 - Emphasis on the use of data

Cultures

- Databases:** concentrate on large-scale (non-main-memory) data.
- AI (machine-learning):** concentrate on complex methods, small data.
 - In today's world data is more important than algorithms
- Statistics:** concentrate on models.

CS345A: Data Mining on the Web: Anand Rajaraman, Jeff Ullman

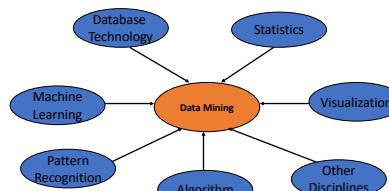
Models vs. Analytic Processing

- To a database person, data-mining is an extreme form of **analytic processing** – queries that examine large amounts of data.
 - Result is the query answer.
- To a statistician, data-mining is the inference of models.
 - Result is the parameters of the model.

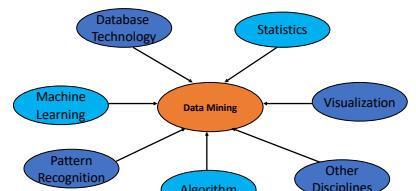
(Way too Simple) Example

- Given a billion numbers, a DB person would compute their average and standard deviation.
- A statistician might fit the billion points to the best Gaussian distribution and report the mean and standard deviation of *that distribution*.

Data Mining: Confluence of Multiple Disciplines



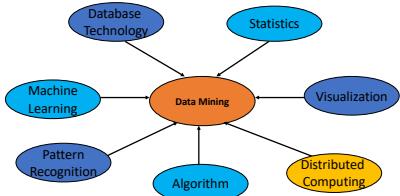
Data Mining: Confluence of Multiple Disciplines



CS345A: Data Mining on the Web: Anand Rajaraman, Jeff Ullman

CS345A: Data Mining on the Web: Anand Rajaraman, Jeff Ullman

Data Mining: Confluence of Multiple Disciplines



Sampling ...

- The key principle for effective sampling is the following:
 - using a sample will work almost as well as using the entire data sets, if the sample is representative
- A sample is representative if it has approximately the same property (of interest) as the original set of data

The data analysis pipeline

- Mining is not the only step in the analysis process
 - Data Preprocessing:** real data is noisy, incomplete and inconsistent. **Data cleaning** is required to make sense of the data.
 - Techniques: Sampling, Dimensionality Reduction, Feature selection.
 - A dirty work, but it is often the most important step for the analysis.
 - Post-Processing:** Make the data actionable and useful to the user
 - Statistical analysis of importance
 - Visualization.
 - Pre- and Post-processing are often data mining tasks as well



Data Quality

- Examples of data quality problems:
 - Noise and outliers
 - missing values
 - duplicate data

Sampling

- Sampling is the main technique employed for data selection.
 - It is often used for both the preliminary investigation of the data and the final data analysis.
- Statisticians sample because obtaining the entire set of data of interest is too expensive or time consuming.
- Sampling is used in data mining because processing the entire set of data of interest is too expensive or time consuming.

Types of Sampling

- Simple Random Sampling**
 - There is an equal probability of selecting any particular item
- Sampling without replacement**
 - As each item is selected, it is removed from the population
- Sampling with replacement**
 - Objects are not removed from the population as they are selected for the sample
 - In sampling with replacement, the same object can be picked up more than once
- Stratified sampling**
 - Split the data into several partitions; then draw random samples from each partition

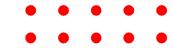
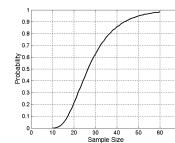
Sample Size



8000 points 2000 Points 500 Points

Sample Size

- What sample size is necessary to get at least one object from each of 10 groups.





4- Pre- and Post-Processing

Thomas W Gyeera, Assistant Professor
Computer and Information Science
University of Massachusetts Dartmouth

Courtesy to Prof. Panayiotis Tsaparas

What is Data Mining?

Data mining is the use of efficient techniques for the analysis of very large collections of data and the extraction of useful and possibly unexpected patterns in data.

"Data mining is the analysis of (often large) observational data sets to find unsuspected relationships and to summarize the data in novel ways that are both interesting and useful to the data analyst" (Hand, Mannila, Smyth)

"Data mining is the discovery of models" (Rajaraman, Ullman)

- What are the different types of models?
 - Models that result in rules (e.g., a single function)
 - Models that predict the future outcomes.
 - Models that summarize the data
 - Models that extract the most prominent features of the data.

Why do we need data mining?

- Really huge amounts of complex data generated from multiple sources are collected in different places
- Scientific data from different disciplines
 - Weather, astronomy, physics, biological microarrays, genetics
 - Human genome
 - The Web, scientific articles, news, tweets, Facebook postings.
- Business data
 - Retail store records, credit card records
 - Behavioral data
 - Mobile phone data, query logs, browsing behavior, ad clicks
 - Health care data
- All these types of data can be combined in many ways:
 - E-commerce, medical test, images, user behavior, ad transactions.

- We need to analyze this data to extract knowledge
- Knowledge can be used for commercial or scientific purposes.
- Our solution should scale to the size of the data

The data analysis pipeline

- Mining is not the only step in the analysis process



- **Pre-Processing:** Real data is noisy, incomplete, and inconsistent.
 - Data Cleaning: The first step of the process
 - Techniques: Sampling, Dimensionality Reduction, Feature selection.
 - A dirty work, but it is often the most important step for the analysis.

- **Post-Processing:** Make the data actionable and useful to the user
 - Statistical analysis of importance
 - Visualizations
- **Pre- and Post-processing** are often data mining tasks as well

Data Quality

- Examples of data quality problems:
 - Noise and outliers
 - Missing values
 - Duplicate data

| Tid | Refund | Marital Status | Taxable Income | Cheat |
|-----|--------|----------------|----------------|-------|
| 1 | Yes | Single | 125K | No |
| 2 | No | Married | 100K | No |
| 3 | No | Single | 70K | No |
| 4 | Yes | Married | 120K | No |
| 5 | No | Divorced | 1000K | Yes |
| 6 | No | NULL | 60K | No |
| 7 | Yes | Divorced | 220K | NULL |
| 8 | No | Single | 85K | Yes |
| 9 | No | Married | 90K | No |
| 9 | No | Single | 90K | No |

A mistake or a millionaire?

Missing values

Inconsistent duplicate entries

Sampling

Sampling is the main technique employed for data selection.

- It is often used for both the preliminary investigation of the data and the final data analysis.

Statisticians sample because obtaining the entire set of data of interest is too expensive or time consuming.

- Example: What is the average height of a person in UMD?
- We cannot measure the height of everybody

Sampling is used in data mining because processing the entire set of data of interest is too expensive or time consuming.

- Example: We have 1M documents. What fraction has at least 100 words in common?
- Computing number of common words for all pairs requires $O(n^2)$ operations.
- Example: What fraction of tweets in a year contain the word "apple"?
- 300M tweets per day, if 100 characters on average, $85.5\% \text{ to store all tweets}$

Sampling

The key principle for effective sampling is the following:

- using a sample will work almost as well as using the entire data sets, if the sample is representative

A sample is representative if it has approximately the same property (of interest) as the original set of data

- Otherwise, we say that the sample introduces some bias

What happens if we take a sample from the university campus to compute the average height of a person at UMD?

Types of Sampling

- Simple Random Sampling
 - There is an equal probability of selecting any particular item

- Sampling without replacement
 - As each item is selected, it is removed from the population

- Sampling with replacement
 - Objects are not removed from the population as they are selected for the sample.

- In sampling with replacement, the same object can be picked up more than once. This means that the probability of picking the same item again is $P(W|W)$.
- e.g., we have 100 people, 51 are women $P(W) = 0.51$, 49 men $P(M) = 0.49$. If I pick two persons what is the probability $P(W|W)$ that both are women?

- Sampling with replacement: $P(W|W) = 0.51^2$

- Sampling without replacement: $P(W|W) = 51/100 * 50/99$

Types of Sampling

- **Stratified sampling**
 - Split the data into several groups; then draw random samples from each group.
 - Ensures that both groups are represented.

Example 1. I want to understand the differences between legitimate and fraudulent credit card transactions. 0.1% of transactions are fraudulent. What happens if I select 1000 transactions at random?

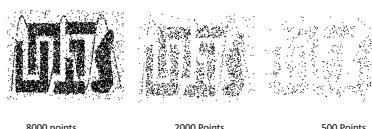
- I get 1 fraudulent transaction (in expectation). Not enough to draw any conclusions. Solution: sample 100 legitimate and 1000 fraudulent transactions

Probability Reminder: If an event has probability p of happening and I do N trials, the expected number of times the event occurs is pk .

Example 2. I want to answer the question: Do web pages that are linked have on average more words in common than those that are not? I have 1M pages, and 1M links, what happens if I select 10K pairs of pages at random?

- Most likely I will not get any links. Solution: sample 10K random pairs, and 10K links

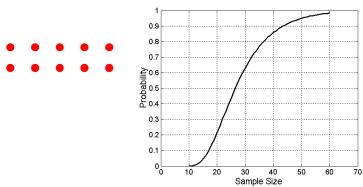
Sample Size



8000 points 2000 Points 500 Points

Sample Size

- What sample size is necessary to get at least one object from each of 10 groups.



Reservoir sampling

- Algorithm: With probability $1/n$ select the n -th item of the stream and replace the previous choice.

Claim: Every item has probability $1/N$ to be selected after N items have been read.

- Proof
 - What is the probability of the n -th item to be selected?
 - $\frac{1}{n}$
 - What is the probability of the n -th item to survive for $N-n$ rounds?
 - $(1 - \frac{1}{n+1})(1 - \frac{1}{n+2}) \cdots (1 - \frac{1}{N})$

A (detailed) data preprocessing example



Suppose we want to mine the comments/reviews of people on Yelp and Foursquare.

Data Collection



- Today there is an abundance of data online

• Facebook, Twitter, Wikipedia, Web, etc...

- We can extract interesting information from this data, but first we need to collect it

• Customized crawlers, use of public APIs

• Additional cleaning/processing to parse out the useful parts

• Respect of crawling etiquette

Mining Task

- Collect all reviews for the top-10 most reviewed restaurants in NY in Yelp

- Find five terms that best describe the restaurants.

- Algorithm?

- Standard interview question for many companies

Measures of Spread: Range and Variance

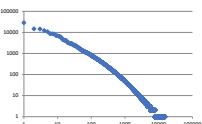
- Range is the difference between the max and min
- The variance or standard deviation is the most common measure of the spread of a set of points.

$$\text{var}(x) = \frac{1}{m} \sum_{i=1}^m (x - \bar{x})^2$$

$$\sigma(x) = \sqrt{\text{var}(x)}$$

Zipf's law

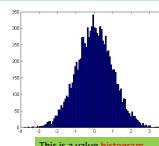
- Power laws can be detected by a linear relationship in the log-log space for the rank-frequency plot



- $f(r)$: Frequency of the r -th most frequent word
 $f(r) = r^{-\beta}$

Normal Distribution

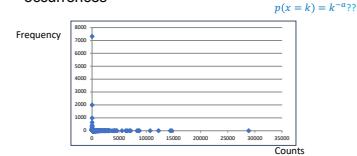
$$\cdot \phi(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



- An important distribution that characterizes many quantities and has a central role in probabilities and statistics.
- Appears also in the central limit theorem
- Fully characterized by the mean μ and standard deviation σ

Not everything is normally distributed

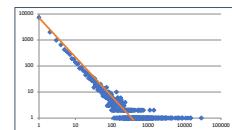
- Plot of y number of words with x number of occurrences



- If this was a normal distribution, we would not have a count as large as 28K

Power-law distribution

- We can understand the distribution of words if we take the log-log plot

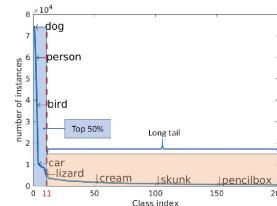


- Linear relationship in the log-log space for:
 $p(x=k) = k^{-\alpha}$

Post-processing

- Visualization
 - The human eye is a powerful analytical tool
 - If we visualize the data properly, we can discover patterns
 - Visualization is the way to present the data so that patterns can be seen
 - E.g., histograms and plots are a form of visualization
 - There are multiple techniques (a field on its own)

The Long Tail



Meaningfulness of Answers

A big data-mining risk is that you will "discover" patterns that are meaningless.

Statisticians call it Bonferroni's principle: (roughly) if you look in more places for interesting patterns than your amount of data will support, you are bound to find crap.

The Rhine Paradox: a great example of how not to conduct scientific research.

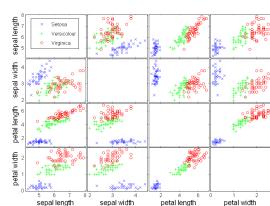
Rhine Paradox – (1)

Joseph Rhine was a parapsychologist in the 1950's who hypothesized that some people had Extra-Sensory Perception (ESP).

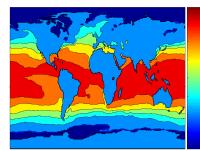
He devised (something like) an experiment where subjects were asked to guess 10 hidden cards – red or blue.

He discovered that almost 1 in 1000 had ESP – they were able to get all 10 right!

Scatter Plot Array of Iris Attributes



Contour Plot Example: SST Dec, 1998



CS345A Data Mining on the Web: Anand Rajaraman, Jeff Ullman

Rhine Paradox – (2)

He told these people they had ESP and called them in for another test of the same type.

Alas, he discovered that almost all of them had lost their ESP.

What did he conclude?

• Answer on next slide.

Rhine Paradox – (3)

- He concluded that you shouldn't tell people they have ESP; it causes them to lose it.
- The facts behind...

$$\begin{aligned} P(\text{at least one wins}) &= 1 - P(\text{no one wins}) \\ &= 1 - \left(1 - \frac{1}{2}\right)^{1000} \\ &\approx 0.6235762 \end{aligned}$$

45

CS345A Data Mining on the Web: Anand Rajaraman, Jeff Ullman

46

47

CS345A Data Mining on the Web: Anand Rajaraman, Jeff Ullman

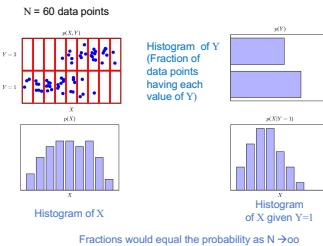


5- Probability Theory

Thomas W Gyeera, Assistant Professor
Computer and Information Science
University of Massachusetts Dartmouth

Courtesy to Prof. Sargur N. Srihari

Joint Distribution over Two Variables



Independent Variables

- If $p(X, Y) = p(X)p(Y)$ then X and Y are said to be independent
- Why?
- From product rule $p(Y|X) = \frac{p(X,Y)}{p(X)} = p(Y)$
- In fruit example if each box contained same fraction of apples and oranges then $p(F|B) = p(F)$

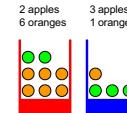


Expectation: an Example

- Assume the final grade is represented by a random variable X
- Blow is the marginal probability of X in class
 - $p(X = A) = 0.3; p(X = B) = 0.4; p(X = C) = 0.3$
- We further assume there is a mapping f between final grade and future salary level
 - $f(A) = 100K; f(B) = 90K; f(C) = 80K$
- Question: what is the expected salary level of this class?
 - $100K * 0.3 + 90K * 0.4 + 80K * 0.3 = 90K$

Probabilities of Interest

- Marginal Probability
 - What is the probability of an apple?
- Conditional Probability
 - Given that we have an orange what is the probability that we chose the blue box?
- Joint Probability
 - What is the probability of orange AND blue box?



Let $p(B = r) = 4/10$
and $p(B = b) = 6/10$

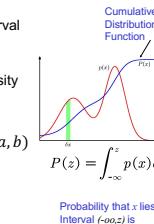
Bayes Theorem

- Think about relation between mid-term score and final grade:
 - Everyone wants to know, if my mid-term is in the range [80-90], what is the probability of getting A?
 - Assume the range of score is represented by random variable X , and final grade is by Y
 - So, we are modeling: $p(Y = A|X \in [80,90])$
- Below is something we may know from the course records of year 2018:
 - $p(Y = A) = 0.3;$
 - $p(X \in [80,90]) = 0.2$
 - $p(X \in [80,90]|Y = A) = 0.1$

Probability Densities Function (PDF)

- Continuous Variables
 - If probability that x falls in interval $(x, x + \delta x)$ is given by $p(x)dx$ for $\delta x \rightarrow 0$
 - then $p(x)$ is a Probability Density Function (PDF) of x
- Probability x lies in interval (a, b) is

$$p(x \in (a, b)) = \int_a^b p(x)dx$$



Expectation of $f(x)$

- Expectation $E[f]$ is average value of some function $f(x)$ under the probability distribution $p(x)$
- For a discrete distribution: $E[f] = \sum p(x)f(x)$
- For a continuous distribution: $E[f] = \int p(x)f(x)dx$
- If there are N points drawn from a PDF, then expectation can be approximated as: $E[f] = \frac{1}{N} \sum f(x)$

Sum Rule of Probability Theory

- Consider two random variables
 - X can take on values $x_i, i = 1, \dots, M$
 - Y can take on values $y_j, j = 1, \dots, L$
 - N trials sampling both X and Y
 - No. of trials with $X = x_i$ and $Y = y_j$ is n_{ij}



$$\text{Joint Probability } p(X = x_i, Y = y_j) = \frac{n_{ij}}{N}$$

$$p(X = x_i) = \frac{c_i}{N}$$

$$\text{Since } c_i = \sum_j n_{ij} \quad p(X = x_i) = \sum_{j=1}^L p(X = x_i, Y = y_j)$$

Bayes Theorem

- From the product rule together with the symmetry property $p(X, Y) = p(Y, X)$ we get

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

- which is called Bayes' theorem

- Sum and product rule for $p(X)$

$$p(X) = \sum_Y p(X|Y)p(Y)$$

Normalization constant to ensure sum of conditional probability equal to 1 over all values of Y



Bayes rule applied to Fruit Problem

- Probability that box is red given that fruit picked is orange

$$p(B = r|F = o) = \frac{p(F = o|B = r)p(B = r)}{p(F = o)}$$

$$= \frac{\frac{3}{5} \times \frac{1}{2}}{\frac{9}{10}} = \frac{2}{5} = 0.44$$

The *a posteriori* probability of 0.66 is different from the *a priori* probability of 0.4.

- Probability that fruit is orange
 - From sum and product rules

$$p(F = o) = p(F = o, B = r) + p(F = o, B = b)$$

$$= p(F = o|B = r)p(B = r) + p(F = o|B = b)p(B = b)$$

$$= \frac{6}{8} \times \frac{4}{10} + \frac{1}{8} \times \frac{6}{10} = \frac{9}{20} = 0.45$$

The marginal probability of 0.45 is lower than average probability of 7/12=0.58

Sum, Product, Bayes for Continuous

- Rules apply for continuous, or combinations of discrete and continuous variables

$$\text{Marginalization } p(x) = \int p(x, y) dy$$

$$\text{Product } p(x, y) = p(y|x)p(x)$$

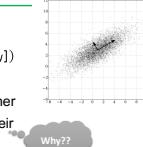
$$\text{Bayes } p(y|x) = \frac{p(x|y)p(y)}{p(x)}$$

- Formal justification of sum, product rules for continuous variables requires measure theory

Covariance

- For two random variables x and y covariance is defined as

$$\begin{aligned} cov[x, y] &= E_{xy}(x - E[x])(y - E[y]) \\ &= E_{xy}[xy] - E[x]E[y] \end{aligned}$$



- Expresses how x and y vary together

- If x and y are independent, then their covariance vanishes

- If we consider covariance of components of vector x with each other then we denote it as

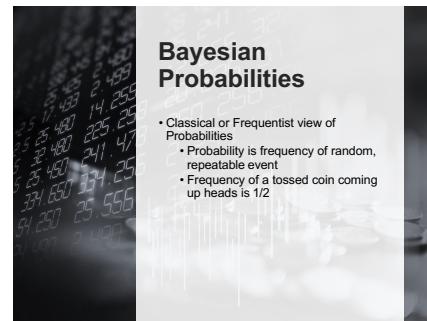
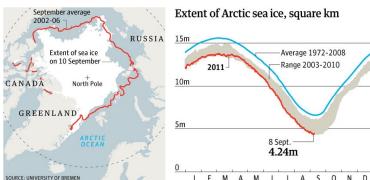
$$var[x] = cov[x, x]$$

Covariance Matrix

- If x and y are **two vectors** of random variables, then the covariance is a matrix
- Example:
 - assume random variable vectors $x = [x_1, x_2, x_3, x_4]; y = [y_1, y_2, y_3, y_4]$

$$\text{cov}(x, y) = \begin{pmatrix} \text{cov}(x_1, y_1) & \text{cov}(x_1, y_2) & \text{cov}(x_1, y_3) & \text{cov}(x_1, y_4) \\ \text{cov}(x_2, y_1) & \text{cov}(x_2, y_2) & \text{cov}(x_2, y_3) & \text{cov}(x_2, y_4) \\ \text{cov}(x_3, y_1) & \text{cov}(x_3, y_2) & \text{cov}(x_3, y_3) & \text{cov}(x_3, y_4) \\ \text{cov}(x_4, y_1) & \text{cov}(x_4, y_2) & \text{cov}(x_4, y_3) & \text{cov}(x_4, y_4) \end{pmatrix}$$

Arctic Ice (2011)



Bayesian Probabilities

- Classical or Frequentist view of Probabilities
 - Probability is frequency of random, repeatable event
 - Frequency of a tossed coin coming up heads is 1/2



Bayesian Probabilities

Bayesian View

- Probability is a quantification of uncertainty
- Degree of belief in propositions that do not involve random variables

Examples of uncertain events as probabilities:

- Whether Shakespeare's plays were written by Francis Bacon
- Whether moon was once in its own orbit around the sun
- Whether Thomas Jefferson had a child by one of his slaves
- Whether a signature on a check is genuine

Examples of Uncertain Events

Probability that Mr. M was the murderer of Mrs. M given the evidence

Whether Arctic ice cap will disappear by end of century

- We have some idea of how quickly polar ice is melting
- Revise it on the basis of fresh evidence (satellite observations)
- Assessment will affect actions we take (to reduce greenhouse gases)

All can be achieved by general Bayesian interpretation

Role of Likelihood Function

- Uncertainty in w expressed as

$$p(w|D) = \frac{p(D|w)p(w)}{p(D)} \quad \text{where } p(D) = \int p(D|w)p(w)dw$$

- Denominator is normalization factor
- Involves marginalization over w
- Bayes theorem in words

$$\text{posterior} \propto \text{likelihood} \times \text{prior}$$



- Likelihood Function plays central role in both Bayesian and frequentist paradigms
 - Frequentist: w is a **fixed parameter** determined by an estimator; error bars are estimates obtained from possible data sets $\mathcal{D} = \{D_1, D_2\}$
 - Bayesian: there is a single data set D ; uncertainty is expressed as **probability distribution** over w

Maximum Likelihood Approach

- In frequentist setting, w is considered to be a fixed parameter
 - w is set to value that maximizes likelihood function $p(D|w)$
- In ML, negative log of likelihood function is called error function, since maximizing likelihood is equivalent to minimizing error



- Bootstrap approach to creating L data sets
 - From N data points new data sets are created by drawing n points at random with replacement
 - Repeat L times to generate L data sets
 - Accuracy of parameter estimate can be evaluated by variability of predictions between different bootstrap sets

Bayesian vs. Frequentist Approach

- Inclusion of prior knowledge arises naturally
- Coin Toss Example
 - Fair looking coin is tossed three times and lands Head each time
 - Classical MLE of the probability of landing heads is 1 implying all future tosses will land Heads
 - Bayesian approach with reasonable prior will lead to less extreme conclusion



- Marginalization over whole parameter space is required to make predictions or compare models
- Factors making it practical:
 - Sampling Methods such as Markov Chain Monte Carlo methods
 - Increased speed and memory of computers
- Deterministic approximation schemes such as Variational Bayes and Expectation propagation are alternatives to sampling

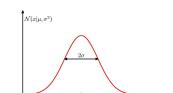
The Gaussian Distribution

- For single real-valued variable x

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{1}{2\sigma^2}(x - \mu)^2\right\}$$

- Parameters:

- Mean μ , variance σ^2
- Standard deviation σ
- Precision $= \frac{1}{\sigma^2}$
- $E[x] = \mu$
- $\text{var}[x] = \sigma^2$



Maximum of a distribution is its mode
For a Gaussian, mode coincides with its mean

Likelihood Function for Gaussian

- Given N observations $x_i, i = 1, \dots, n$

- Independent and identically distributed (i.i.d.)

- Probability of data set is given by likelihood function

$$p(\mathbf{x}|\mu, \sigma^2) = \prod_{n=1}^N \mathcal{N}(x_n|\mu, \sigma^2).$$

Data: black points
Likelihood= product of blue values. Pick mean and variance to maximize this product

- Log-likelihood function is

$$\ln p(\mathbf{x}|\mu, \sigma^2) = -\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 - \frac{N}{2} \ln \sigma^2 - \frac{N}{2} \ln(2\pi).$$

Likelihood Function for Gaussian

- Log-likelihood function is

$$\ln p(\mathbf{x}|\mu, \sigma^2) = -\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 - \frac{N}{2} \ln \sigma^2 - \frac{N}{2} \ln(2\pi).$$

- Maximum likelihood solutions are given by:

$$\mu_{ML} = \frac{1}{N} \sum_{n=1}^N x_n$$

$$\sigma_{ML}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{ML})^2$$

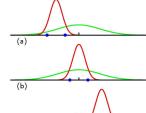
Data: black points
Likelihood= product of blue values. Pick mean and variance to maximize this product

Bias in Maximum Likelihood

- Maximum likelihood systematically underestimates variance

$$E[\mu_{ML}] = \mu$$

$$E[\sigma_{ML}^2] = \left(\frac{N-1}{N}\right)\sigma^2$$



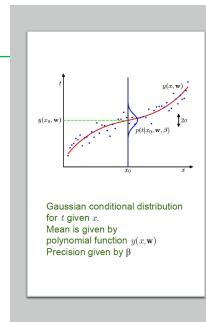
Averaged across three data sets
mean is correct
Variance is underestimated
because it is estimated relative
to sample mean and not true mean

- Not an issue as N increases
- Problem is related to overfitting problem

Curve Fitting (Probabilistic)

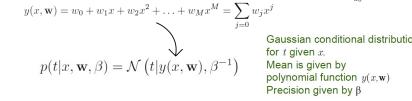
- Goal is to predict for target variable t given a new value of the input variable x

- Given N input values $\mathbf{x} = (x_1, \dots, x_N)^T$ and corresponding target values $\mathbf{t} = (t_1, \dots, t_N)^T$



Curve Fitting (Probabilistic)

- Assume given value of x , value of t has a Gaussian distribution with mean equal to $y(x, w)$ of polynomial curve



$$y(t|x, w, \beta) = N(t|y(x, w), \beta^{-1})$$

Precision Parameter with MLE

- Maximum likelihood can also be used to determine β of Gaussian conditional distribution

Maximizing likelihood w.r.t. β gives

$$\frac{1}{\beta_{ML}} = \frac{1}{N} \sum_{n=1}^N \{y(x_n, w_{ML}) - t_n\}^2.$$

- First determine parameter vector w_{ML} governing the mean and subsequently use this to find precision β_{ML}

Curve Fitting with Maximum Likelihood

- Can omit last two terms -- don't depend on w
- Can replace $\beta/2$ with β
 - since it is constant w.r.t. w
- Minimize negative log-likelihood
- Identical to sum-of-squares error function

$$\ln p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = -\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi).$$

Curve Fitting with Maximum Likelihood

- Likelihood Function is

$$p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | y(x_n, \mathbf{w}), \beta^{-1}).$$

- Logarithm of the Likelihood function is

$$\ln p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = -\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi).$$

- To find maximum likelihood solution for polynomial coefficients \mathbf{w}_{ML}

- Maximize w.r.t. \mathbf{w}

Predictive Distribution

- Knowing parameters \mathbf{w} and β
- Predictions for new values of x can be made using

- Instead of a point estimate we are now giving a probability distribution over t

$$p(t|x, w_{ML}, \beta_{ML}) = \mathcal{N}(t|y(x, w_{ML}), \beta_{ML}^{-1}).$$



6- Decision Trees

Thomas W Gyeera, Assistant Professor
Computer and Information Science
University of Massachusetts Dartmouth

Courtesy to Prof. Panayiotis Tsaparas

Catching tax-evasion

| Tid | Refund | Marital Status | Taxable Income | Cheat |
|-----|--------|----------------|----------------|-------|
| 1 | Yes | Single | 125K | No |
| 2 | No | Married | 100K | No |
| 3 | No | Single | 70K | No |
| 4 | Yes | Married | 120K | No |
| 5 | No | Divorced | 95K | Yes |
| 6 | No | Married | 60K | No |
| 7 | Yes | Divorced | 220K | No |
| 8 | No | Single | 85K | Yes |
| 9 | No | Married | 75K | No |
| 10 | No | Single | 90K | Yes |

Tax-return data for year 2011

A new tax return for 2012

Is this a cheating tax return?

| Refund | Marital Status | Taxable Income | Cheat |
|--------|----------------|----------------|-------|
| No | Married | 80K | ? |

An instance of the classification problem: learn a method for discriminating between records of different classes (cheaters vs non-cheaters)

Examples of Classification Tasks

| | |
|---------------|---|
| Predicting | Predicting tumor cells as benign or malignant |
| Classifying | Classifying credit card transactions as legitimate or fraudulent |
| Categorizing | Categorizing news stories as finance, weather, entertainment, sports, etc |
| Identifying | Identifying spam email, spam web pages, adult content |
| Understanding | Understanding if a web query has commercial intent or not |

General approach to classification

- Training set consists of records with known class labels
- Training set is used to build a classification model
- A labeled test set of previously unseen data records is used to evaluate the quality of the model.
- The classification model is applied to new records with unknown class labels

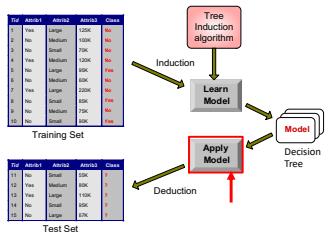
Classification Techniques



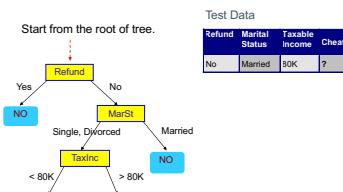
Decision Tree

- Decision tree
 - A flow-chart-like tree structure
 - Internal node denotes a test on an attribute
 - Branch represents an outcome of the test
 - Leaf nodes represent class labels or class distribution

Decision Tree Classification Task



Apply Model to Test Data



| Tid | Refund | Marital Status | Taxable Income | Cheat |
|-----|--------|----------------|----------------|-------|
| 1 | Yes | Single | 125K | No |
| 2 | No | Married | 100K | No |
| 3 | No | Single | 70K | No |
| 4 | Yes | Married | 120K | No |
| 5 | No | Divorced | 95K | Yes |
| 6 | No | Married | 60K | No |
| 7 | Yes | Divorced | 220K | No |
| 8 | No | Single | 85K | Yes |
| 9 | No | Married | 75K | No |
| 10 | No | Single | 90K | Yes |

What is classification?

- Classification is the task of learning a target function f that maps attribute set x to one of the predefined class labels y

| Tid | Refund | Marital Status | Taxable Income | Cheat |
|-----|--------|----------------|----------------|-------|
| 1 | Yes | Single | 125K | No |
| 2 | No | Married | 100K | No |
| 3 | No | Single | 70K | No |
| 4 | Yes | Married | 120K | No |
| 5 | No | Divorced | 95K | Yes |
| 6 | No | Married | 60K | No |
| 7 | Yes | Divorced | 220K | No |
| 8 | No | Single | 85K | Yes |
| 9 | No | Married | 75K | No |
| 10 | No | Single | 90K | Yes |

One of the attributes is the class attribute
In this case: Cheat

Two class labels (or classes): Yes (1), No (0)

Figure 4.2. Classification as the task of mapping an input attribute set x into its class label y .

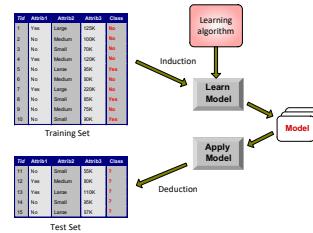
Why classification?

The target function f is known as a classification model

Descriptive modeling: Explanatory tool to distinguish between objects of different classes (e.g., understand why people cheat on their taxes)

Predictive modeling: Predict a class of a previously unseen record

Illustrating Classification Task



Evaluation of classification models

- Counts of test records that are correctly (or incorrectly) predicted by the classification model

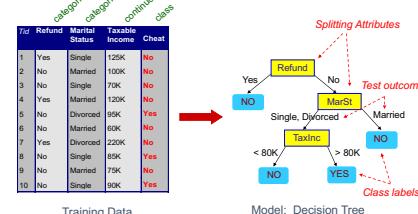
Confusion matrix

| Actual Class | Predicted Class | |
|--------------|-----------------|-----------|
| | Class = 1 | Class = 0 |
| Class = 1 | f_{11} | f_{10} |
| Class = 0 | f_{01} | f_{00} |

$$\text{Accuracy} = \frac{\# \text{ correct predictions}}{\text{total \# of predictions}} = \frac{f_{11} + f_{00}}{f_{11} + f_{10} + f_{01} + f_{00}}$$

$$\text{Error rate} = \frac{\# \text{ wrong predictions}}{\text{total \# of predictions}} = \frac{f_{10} + f_{01}}{f_{11} + f_{10} + f_{01} + f_{00}}$$

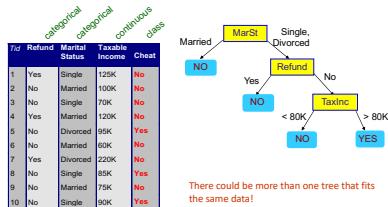
Example of a Decision Tree



Training Data

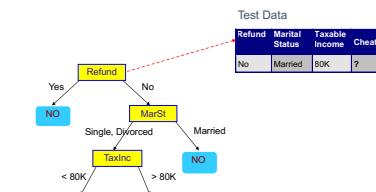
Model: Decision Tree

Another Example of Decision Tree

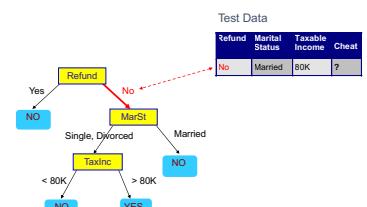


There could be more than one tree that fits the same data!

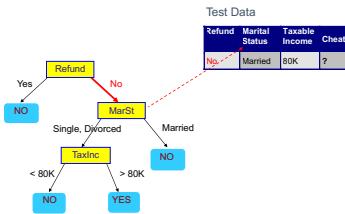
Apply Model to Test Data



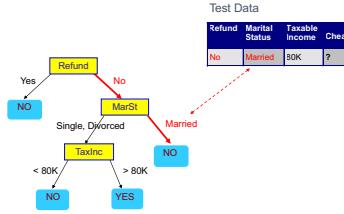
Apply Model to Test Data



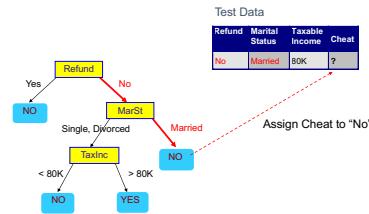
Apply Model to Test Data



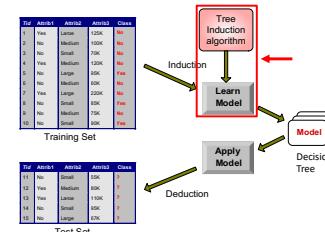
Apply Model to Test Data



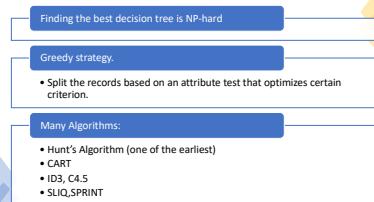
Apply Model to Test Data



Decision Tree Classification Task



Tree Induction



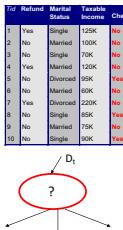
General Structure of Hunt's Algorithm

Let D_t be the set of training records that reach a node t

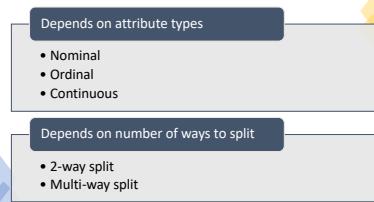
General Procedure:

- If D_t contains records that belong to the same class y_t , then t is a leaf node labeled as y_t .
- If D_t contains records with the same attribute values, then t is a leaf node labeled with the majority class y_t .
- If D_t is an empty set, then t is a leaf node labeled by the default class y_d .
- If D_t contains records that belong to more than one class, use an attribute test to split the data into smaller subsets.

Recursively apply the procedure to each subset.



How to Specify Test Condition?

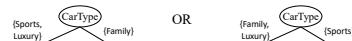


Splitting Based on Nominal Attributes

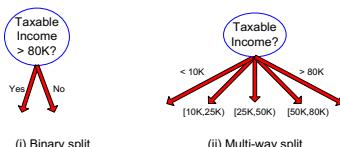
- Multi-way split:** Use as many partitions as distinct values.



- Binary split:** Divides values into two subsets. Need to find optimal partitioning.

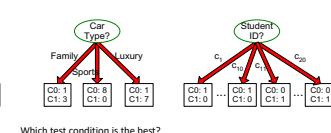


Splitting Based on Continuous Attributes



How to determine the Best Split

Before Splitting: 10 records of class 0, 10 records of class 1



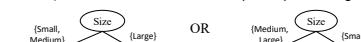
Which test condition is the best?

Splitting Based on Ordinal Attributes

- Multi-way split:** Use as many partitions as distinct values.



- Binary split:** Divides values into two subsets – respects the order. Need to find optimal partitioning.



- What about this split?



Splitting Based on Continuous Attributes

- Different ways of handling:
 - Discretization** to form an **ordinal** categorical attribute
 - Static** – discretize once at the beginning
 - Dynamic** – ranges can be found by equal interval bucketing, equal frequency bucketing (percentiles), or clustering
 - Binary Decision:** $(A < v)$ or $(A \geq v)$
 - consider all possible splits and finds the best cut
 - can be more compute intensive

How to determine the Best Split

- Greedy approach:**

- Nodes with **homogeneous** class distribution are preferred

- Need a measure of node **impurity**:

C0: 5
C1: 5

Non-homogeneous,
High degree of impurity

C0: 9
C1: 1

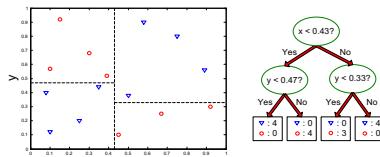
Homogeneous,
Low degree of impurity

- Ideas?

Measuring Node Impurity

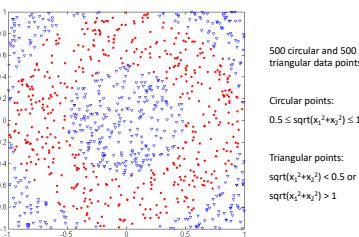
- $p(i|t)$: fraction of records associated with node t belonging to class i
- Entropy(t) = $-\sum_{i=1}^c p(i|t) \log p(i|t)$
- Used in ID3 and C4.5
- Gini(t) = $1 - \sum_{i=1}^c [p(i|t)]^2$
- Used in CART, SLIQ, SPRINT.
- Classification error(t) = $1 - \max_i [p(i|t)]$

Decision Boundary



- Border line between two neighboring regions of different classes is known as decision boundary
- Decision boundary is parallel to axes because test condition involves a single attribute at-a-time

Underfitting and Overfitting (Example)



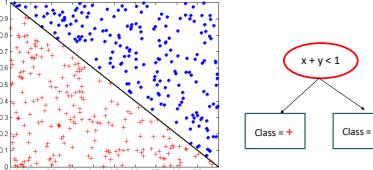
Notes on Overfitting

- Overfitting results in decision trees that are more complex than necessary
- Training error no longer provides a good estimate of how well the tree will perform on previously unseen records
 - The model does not generalize well
- Need new ways for estimating errors

Expressiveness

- Decision tree provides expressive representation for learning discrete-valued function
- But they do not generalize well to certain types of Boolean functions
- Example: parity function:
 - Class = 1 if there is an even number of Boolean attributes with truth value = True
 - Class = 0 if there is an odd number of Boolean attributes with truth value = True
 - For accurate modeling, must have a complete tree
- Less expressive for modeling continuous variables
 - Particularly when test condition involves only a single attribute at-a-time

Oblique Decision Trees

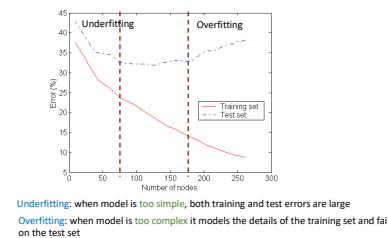


- Test condition may involve multiple attributes
- More expressive representation
- Finding optimal test condition is computationally expensive

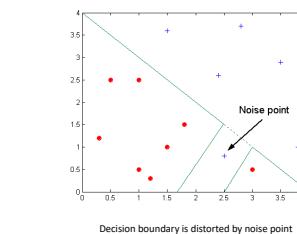
Practical Issues of Classification

- Underfitting and Overfitting
- Evaluation

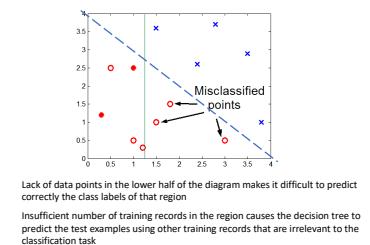
Underfitting and Overfitting



Underfitting and Overfitting



Overfitting due to Insufficient Examples



Estimating Generalization Errors

- Re-substitution errors: error on training ($\Sigma e(t)$)
- Generalization errors: error on testing ($\Sigma e'(t)$)
- Methods for estimating generalization errors:
 - Optimistic approach: $e'(t) = e(t)$
 - Pessimistic approach:
 - For each leaf node: $e'(t) = (e(t) + 0.5)$
 - Total errors: $e'(T) = e(T) + N \times 0.5$ (N : number of leaf nodes)
 - Penalize large trees
 - For a tree with 30 leaf nodes and 10 errors on training (out of 1000 instances)
 - Training error = 10/1000 = 1%
 - Generalization error = $(10 + 30 \times 0.5)/1000 = 2.5\%$
- Using validation set:
 - Split data into training, validation, test
 - Use validation dataset to estimate generalization error
 - Drawback: less data for training

Occam's Razor

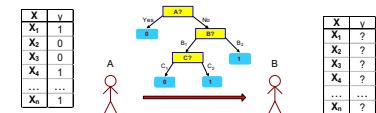


GIVEN TWO MODELS OF SIMILAR GENERALIZATION ERRORS, ONE SHOULD PREFER THE SIMPLER MODEL OVER THE MORE COMPLEX MODEL

FOR COMPLEX MODELS, THERE IS A GREATER CHANCE THAT IT WAS FITTED ACCIDENTALLY BY ERRORS IN DATA

THEFORE, ONE SHOULD INCLUDE MODEL COMPLEXITY WHEN EVALUATING A MODEL

Minimum Description Length (MDL)



- $\text{Cost}(\text{Model}, \text{Data}) = \text{Cost}(\text{Data}| \text{Model}) + \text{Cost}(\text{Model})$
- Search for the least costly model.

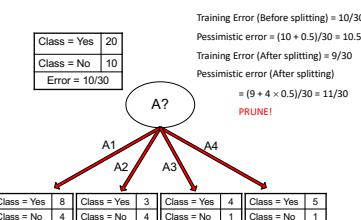
How to Address Overfitting

- Pre-Pruning (Early Stopping Rule)**
 - Stop the algorithm before it becomes a fully-grown tree
 - Typical stopping conditions for a node:
 - Stop if all instances belong to the same class
 - Stop if all the attribute values are the same
- More **restrictive** conditions:
 - Stop if **number of instances** is less than some user-specified threshold
 - Stop if class distribution of instances are **independent** of the available features (e.g., using χ^2 test)
 - Stop if expanding the current node **does not improve impurity** measures (e.g., Gini or information gain).

How to Address Overfitting...

- Post-pruning**
 - Grow decision tree to its entirety
 - Trim the nodes of the decision tree in a **bottom-up** fashion
 - If generalization error improves after trimming, replace sub-tree by a leaf node.
 - Class label of leaf node is determined from majority class of instances in the sub-tree
 - Can use **MDL** for post-pruning

Example of Post-Pruning



Model Evaluation

- Metrics for Performance Evaluation
- How to evaluate the performance of a model?
- Methods for Performance Evaluation
- How to obtain reliable estimates?
- Methods for Model Comparison
- How to compare the relative performance among competing models?

Metrics for Performance Evaluation

- Focus on the **predictive capability** of a model
 - Rather than how fast it takes to classify or build models, scalability, etc.
- Confusion Matrix:

| | | PREDICTED CLASS | |
|--------------|-----------|-----------------|----------|
| | | Class=Yes | Class=No |
| ACTUAL CLASS | Class=Yes | a | b |
| | Class=No | c | d |
| | | | |

a: TP (true positive)
 b: FN (false negative)
 c: FP (false positive)
 d: TN (true negative)

Metrics for Performance Evaluation...

| ACTUAL CLASS | PREDICTED CLASS | | |
|--------------|-----------------|-----------|-----------|
| | Class=Yes | | Class=No |
| | Class=Yes | a (TP) | b (FN) |
| Class=No | c (FP) | d (TN) | |

- Most widely-used metric:

$$\text{Accuracy} = \frac{a + d}{a + b + c + d} = \frac{TP + TN}{TP + TN + FP + FN}$$

Computing Cost of Classification

| Cost Matrix | | PREDICTED CLASS | |
|--------------|-------|-----------------|----------|
| | | Class=Yes | Class=No |
| ACTUAL CLASS | C(ij) | + | - |
| | + | 100 | |
| - | -1 | 0 | |

| Model M ₁ | | PREDICTED CLASS | |
|----------------------|---|-----------------|----------|
| | | Class=Yes | Class=No |
| ACTUAL CLASS | + | 150 | 40 |
| | - | 60 | 250 |

Accuracy = 80%
Cost = 3910

Model M₂

| Model M ₂ | | PREDICTED CLASS | |
|----------------------|---|-----------------|----------|
| | | Class=Yes | Class=No |
| ACTUAL CLASS | + | 250 | 45 |
| | - | 5 | 200 |

Accuracy = 90%
Cost = 4255

Cost vs Accuracy

| Count | | PREDICTED CLASS | |
|--------------|-------|-----------------|----------|
| | | Class=Yes | Class=No |
| ACTUAL CLASS | C(ij) | + | - |
| | + | a | b |
| - | - | c | d |

Accuracy is proportional to cost if
1. C(Yes|No)=C(No|Yes) = q
2. C(Yes|Yes)=C(No|No) = p

$$N = a + b + c + d$$

$$\text{Accuracy} = (a + d)/N$$

| Cost | | PREDICTED CLASS | |
|--------------|-----------|-----------------|----------|
| | | Class=Yes | Class=No |
| ACTUAL CLASS | Cost | p | q |
| | Class=Yes | p | q |
| Class=No | q | p | |

$$\begin{aligned} \text{Cost} &= p(a+d) + q(b+c) \\ &= p(a+d) + q(N - p(a+d)) \\ &= q(N - p(a+d)) + p(a+d) \\ &= N(q-p) + \text{Accuracy} \end{aligned}$$

Limitation of Accuracy

- Consider a 2-class problem
 - Number of Class 0 examples = 9990
 - Number of Class 1 examples = 10

- If model predicts everything to be class 0, accuracy is $9990/10000 = 99.9\%$
- Accuracy is misleading because model does not detect any class 1 example

Cost Matrix

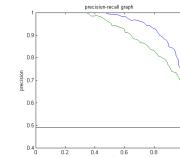
| | | PREDICTED CLASS | | |
|--------------|-----------|-----------------|-----------|----------|
| | | C(ij) | Class=Yes | Class=No |
| ACTUAL CLASS | C(ij) | C(Yes Yes) | C(No Yes) | |
| | C(Yes No) | C(Yes No) | C(No No) | |

$C(ij)$: Cost of classifying class i example as class j

$$\text{Weighted Accuracy} = \frac{w_a + w_d}{w_a + w_b + w_c + w_d}$$

Precision-Recall plot

Usually for parameterized models, it controls the precision/recall tradeoff



Dealing with class imbalance



We can modify the optimization criterion by using a cost-sensitive criterion:
We can balance the class distribution:
• Sample from the larger class so that the size of the two classes is the same
• Replicate the data of the class of interest so that the classes are balanced -> over-fitting issues

Model Evaluation

| Metrics for Performance Evaluation | • How to evaluate the performance of a model? |
|------------------------------------|---|
| Metrics for Performance Evaluation | • How to obtain reliable estimates? |
| Methods for Model Comparison | • How to compare the relative performance among competing models? |

Methods for Performance Evaluation

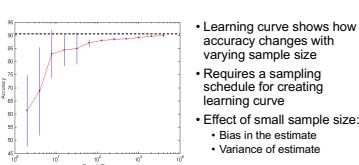


Performance of a model may depend on other factors besides the learning algorithm:
Class distribution
Cost of misclassification
Size of training and test sets

Methods of Estimation

| Holdout | • Reserve 2/3 for training and 1/3 for testing |
|--------------------|---|
| Random subsampling | • One sample may be biased – Repeated holdout |
| Cross validation | • Partition data into k disjoint subsets • k-fold: train on k-1 partitions, test on the remaining one • Leave-one-out: k=n • Guarantees that each record is used the same number of times for training and testing |
| Bootstrap | • Sampling with replacement • ~63% of records used for training, ~27% for testing |

Learning Curve



- Learning curve shows how accuracy changes with varying sample size
- Requires a sampling schedule for creating learning curve
- Effect of small sample size:
 - Bias in the estimate
 - Variance of estimate

Model Evaluation

| Metrics for Performance Evaluation | • How to evaluate the performance of a model? |
|------------------------------------|---|
| Metrics for Performance Evaluation | • How to obtain reliable estimates? |
| Methods for Model Comparison | • How to compare the relative performance among competing models? |

ROC (Receiver Operating Characteristic)

- Developed in 1950s for signal detection theory to analyze noisy signals
- Characterize the trade-off between positive hits and false alarms

ROC curve plots TPR (on the y-axis) against FPR (on the x-axis)

| | | PREDICTED CLASS | |
|--------|-----|-----------------|-----------|
| | | Yes | No |
| Actual | Yes | a (TP) | b (FN) |
| | No | c (FP) | d (TN) |

Fraction of positive instances predicted correctly

$$TPR = \frac{TP}{TP + FN}$$

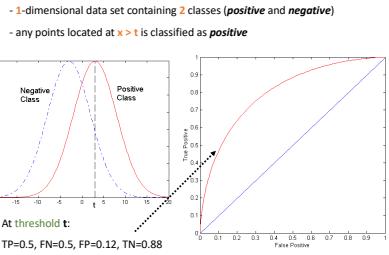
 Fraction of negative instances predicted incorrectly

$$FPR = \frac{FP}{FP + TN}$$

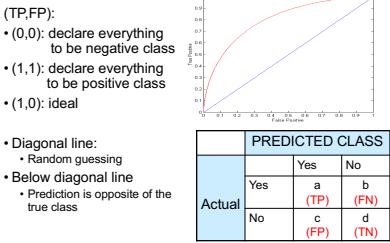
ROC (Receiver Operating Characteristic)

| | | Performance of a classifier represented as a point on the ROC curve | |
|--|--|--|--|
| | | Changing some parameter of the algorithm, sample distribution or cost matrix changes the location of the point | |

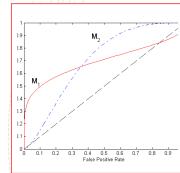
ROC Curve



ROC Curve

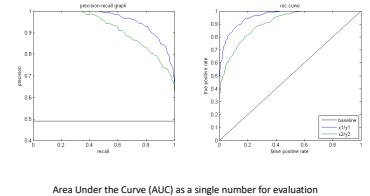


Using ROC for Model Comparison



- No model consistently outperform the others
 - M1 is better for small FPR
 - M2 is better for large FPR
- Area Under the ROC curve (AUC)
 - Ideal: Area = 1
 - Random guess: Area = 0.5

ROC curve vs Precision-Recall curve





7- Similarity and Distance

Thomas W Gyeera, Assistant Professor
Computer and Information Science
University of Massachusetts Dartmouth

Courtesy to Prof. Panayiotis Tsaparas

Similarity: Intersection

- Number of words in common



- $\text{Sim}(D,D) = 3$, $\text{Sim}(D,D) = \text{Sim}(D,D) = 2$
- What about this document?

Vera releases new book with
apple pie recipes

$$\text{Sim}(D,D) = \text{Sim}(D,D) = 3$$

Example

| document | Apple | Microsoft | Obama | Election |
|----------|-------|-----------|-------|----------|
| D1 | 1/3 | 2/3 | 0 | 0 |
| D2 | 1/3 | 2/3 | 0 | 0 |
| D3 | 0 | 0 | 1/3 | 2/3 |

Documents D1, D2 are in the "same direction"

Document D3 is orthogonal to these two

Cosine Similarity



Figure 2.16. Geometric illustration of the cosine measure.

- $\text{Sim}(X,Y) = \cos(X,Y)$: The cosine of the angle between X and Y
- If the vectors are aligned (correlated) angle is zero degrees and $\cos(X,Y) = 1$
- If the vectors are orthogonal (no common coordinates) angle is 90 degrees and $\cos(X,Y) = 0$
- Cosine is commonly used for comparing documents, where we assume that the vectors are normalized by the document length.

Distance

- Numerical measure of how different two data objects are
 - A function that maps pairs of objects to real values
 - Lower when objects are more alike
- Minimum distance is 0, when comparing an object with itself.
- Upper limit varies

Distance Metric

- A distance function d is a **distance metric** if it is a function from pairs of objects to real numbers such that:
 - $d(x,y) \geq 0$. (non-negativity)
 - $d(x,y) = 0$ iff $x = y$. (identity)
 - $d(x,y) = d(y,x)$. (symmetry)
 - $d(x,y) \leq d(x,z) + d(z,y)$ (triangle inequality).

Similarity and Distance

- For many different problems we need to quantify how close two objects are.
- Examples:
 - For an item bought by a customer, find other similar items
 - Group together the customers of site so that similar customers are shown the same ad.
 - Group together web documents so that you can separate the ones that talk about politics and the ones that talk about sports.
 - Find all the near-duplicate mirrored web documents.
 - Find credit card transactions that are very different from previous transactions.
- To solve these problems, we need a definition of **similarity**, or **distance**.
 - The definition depends on the type of data that we have

Similarity

- Numerical measure of how alike two data objects are.
 - A function that maps pairs of objects to real values
 - Higher when objects are more alike.
- Often falls in the range $[0,1]$, sometimes in $[-1,1]$
- Desirable properties for similarity
 - $\text{JSim}(p,q) = 1$ (or maximum similarity) only if $p = q$. (Identity)
 - $\text{JSim}(p,q) = \text{JSim}(q,p)$ for all p and q. (Symmetry)

Similarity between sets

- Consider the following documents

apple
releases
new ipod

apple
releases
new ipad

new
apple pie
recipe

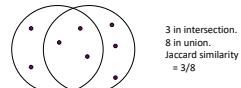
• Which ones are more similar?

• How would you quantify their similarity?

Jaccard Similarity

- The **Jaccard similarity** (Jaccard coefficient) of two sets S_1, S_2 is the size of their intersection divided by the size of their union.

$$\text{JSim}(C_1, C_2) = |C_1 \cap C_2| / |C_1 \cup C_2|$$



- Extreme behavior:
 - $\text{JSim}(X,Y) = 1$ iff $X = Y$
 - $\text{JSim}(X,Y) = 0$ iff X, Y have no elements in common
- JSim is symmetric

Similarity: Intersection

- Number of words in common



- $\text{JSim}(D,D) = 3/5$
- $\text{JSim}(D,D) = \text{JSim}(D,D) = 2/6$
- $\text{JSim}(D,D) = \text{JSim}(D,D) = 3/9$

Similarity between vectors

Documents (and sets in general) can also be represented as vectors

| document | Apple | Microsoft | Obama | Election |
|----------|-------|-----------|-------|----------|
| D1 | 10 | 20 | 0 | 0 |
| D2 | 30 | 60 | 0 | 0 |
| D3 | 0 | 0 | 10 | 20 |

How do we measure the similarity of two vectors?

How well are the two vectors aligned?

Cosine Similarity - math

- If d_1 and d_2 are two vectors, then

$$\text{cost}(d_1, d_2) = (d_1 \bullet d_2) / \|d_1\| \|d_2\|$$
 where \bullet indicates vector dot product and $\|d\|$ is the length of vector d .
- Example:

$$d_1 = 3 \ 2 \ 0 \ 5 \ 0 \ 0 \ 0 \ 2 \ 0$$

$$d_2 = 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0$$

$$d_1 \bullet d_2 = 3 \cdot 1 + 2 \cdot 0 + 0 \cdot 0 + 5 \cdot 0 + 0 \cdot 0 + 0 \cdot 0 + 2 \cdot 1 + 0 \cdot 0 + 0 \cdot 2 = 5$$

$$\|d_1\| = (\sqrt{3^2 + 2^2 + 0^2 + 0^2 + 5^2 + 0^2 + 0^2 + 2^2 + 0^2})^{0.5} = (42)^{0.5} = 6.481$$

$$\|d_2\| = (\sqrt{1^2 + 0^2 + 0^2 + 0^2 + 0^2 + 0^2 + 0^2 + 1^2 + 0^2})^{0.5} = (6)^{0.5} = 2.425$$

$$\cos(d_1, d_2) = 3/150$$

Similarity between vectors

| document | Apple | Microsoft | Obama | Election |
|----------|-------|-----------|-------|----------|
| D1 | 10 | 20 | 0 | 0 |
| D2 | 30 | 60 | 0 | 0 |
| D3 | 0 | 0 | 10 | 20 |

$$\cos(D1, D2) = 1$$

$$\cos(D1, D3) = \cos(D2, D3) = 0$$

Distances for real vectors

- Vectors $x = (x_1, \dots, x_d)$ and $y = (y_1, \dots, y_d)$
- L_p norms or Minkowski distance:

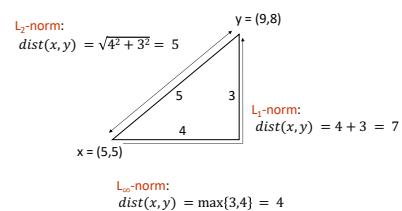
$$L_p(x,y) = [|x_1 - y_1|^p + \dots + |x_d - y_d|^p]^{1/p}$$
- L_2 norm: Euclidean distance:

$$L_2(x,y) = \sqrt{|x_1 - y_1|^2 + \dots + |x_d - y_d|^2}$$
- L_1 norm: Manhattan distance:

$$L_1(x,y) = |x_1 - y_1| + \dots + |x_d - y_d|$$
- L_∞ norm:

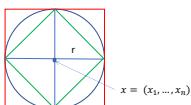
$$L_\infty(x,y) = \max\{|x_1 - y_1|, \dots, |x_d - y_d|\}$$
 - The limit of L_p as p goes to infinity.

Example of Distances



L_p norms are known to be distance metrics

Example



Green: All points y at distance $L_p(x,y) = r$ from point x

Blue: All points y at distance $L_p(x,y) = r$ from point x

Red: All points y at distance $L_\infty(x,y) = r$ from point x

Hamming Distance

- Hamming distance is the number of positions in which bit-vectors differ.

- Example

• $p_1 = 10101$, $p_2 = 10011$, $d(p_1, p_2) = 2$ because the bit-vectors differ in the 3rd and 4th positions, the L_1 norm for the binary vectors

- Hamming distance between two vectors of categorical attributes is the number of positions in which they differ.

- Example:

• $x = (\text{married, low income, cheat})$, $y = (\text{single, low income, not cheat})$, $d(x,y) = 2$

21

Variant Edit Distances

- Allow insert, delete, and mutate.
 - Change one character into another.
- Minimum number of inserts, deletes, and mutates also forms a distance measure.
- Same for any set of operations on strings.
 - Example: substring reversal or block transposition OK for DNA sequences
 - Example: character transposition is used for spelling

25

22

L_p distances for sets

- We can apply all the L_p distances to the cases of sets of attributes, with or without counts, if we represent the sets as vectors
 - E.g., a transaction is a 0/1 vector
 - E.g., a document is a vector of counts.

• E.g., a transaction is a 0/1 vector

• E.g., a document is a vector of counts.

Similarities into distances

- Jaccard distance:

$$JDist(x,y) = 1 - JSim(x,y)$$

- Jaccard Distance is a metric

- Cosine distance:

$$Dist(x,y) = 1 - \cos(x,y)$$

- Cosine distance is a metric

Why Jaccard Distance is a Distance Metric?

- $JDist(x,x) = 0$
 - since $JSim(x,x) = 1$
- $JDist(x,y) = JDist(y,x)$
 - by symmetry of intersection
- $JDist(x,y) \geq 0$
 - since intersection of x,y cannot be bigger than the union.
- **Triangle inequality:**
 - Follows from the fact that $JSim(x,y)$ is the probability of randomly selected element from the union of x and y to belong to the intersection

23

Distance between strings

- How do we define similarity between strings?

| | |
|-------------|---------------|
| weird | wierd |
| intelligent | unintelligent |
| Athena | Athina |

- Important for recognizing and correcting typing errors and analyzing DNA sequences.

Edit Distance for strings

- The edit distance of two strings is the number of inserts and deletes of characters needed to turn one into the other.
- Example: $x = abcd e$; $y = bcd u v e$
 - Turn x into y by deleting a, then inserting u and v after d.
 - Edit distance = 3.
- Minimum number of operations can be computed using dynamic programming
- Common distance measure for comparing DNA sequences

24

Why Edit Distance is a Distance Metric?

- $d(x,x) = 0$ because 0 edits suffice.
- $d(x,y) = d(y,x)$ because insert/delete are inverses of each other.
- $d(x,y) \geq 0$: no notion of negative edits.
- **Triangle inequality:** changing x to z and then to y is one way to change x to y. The minimum is no more than that

25

Distances between distributions

- We can view a document as a distribution over the words

| document | Apple | Microsoft | Obama | Election |
|----------|-------|-----------|-------|----------|
| D1 | 0.35 | 0.5 | 0.1 | 0.05 |
| D2 | 0.4 | 0.4 | 0.1 | 0.1 |

- KL-divergence (Kullback-Leibler) for distributions P,Q

$$D_{KL}(P\|Q) = \sum_x p(x) \log \frac{p(x)}{q(x)}$$

- KL-divergence is asymmetric. We can make it symmetric by taking the average of both sides

- JS-divergence (Jensen-Shannon)

$$JS(P, Q) = \frac{1}{2} D_{KL}(P\|Q) + \frac{1}{2} D_{KL}(Q\|P)$$

26

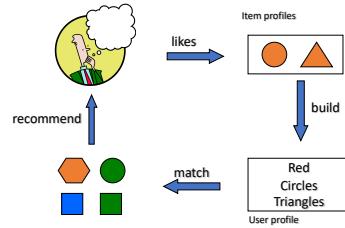


8- Min Hashing and Locality Sensitive Hashing

Thomas W Gyeera, Assistant Professor
Computer and Information Science
University of Massachusetts Dartmouth

Courtesy to Prof. Panayiotis Tsaparas

Plan of action



Recommendation Systems (III)

- Collaborative Filtering (item-item)**
 - For item s, find other similar items
 - Estimate rating for item based on ratings for similar items
 - Can use same similarity metrics and prediction functions as in user-user model
- In practice, it has been observed that item-item often works better than user-user

Finding similar items

- Both the problems we described have a common component
 - We need a quick way to find **highly similar** items to a **query** item
 - OR, we need a method for finding **all pairs** of items that are **highly similar**.
- Also known as the **Nearest Neighbor** problem, or the **All Nearest Neighbors** problem
- We will examine it for the case of near-duplicate web documents.

Why is similarity important?

- We see many definitions of similarity and distance
- How do we make use of similarity in practice?
- What issues do we have to deal with?

An important problem

- Recommendation systems**
 - When a user buys an **item** (initially books) we want to recommend other items that the user may like
 - When a user rates a **movie**, we want to recommend movies that the user may like
 - When a user likes a **song**, we want to recommend other songs that they may like
- A big success of data mining
- Exploits the long tail

Recommendation systems

- Content-based:**
 - Represent the items into a **feature space** and recommend items to customer C **similar** to previous items rated highly by C
 - Movie recommendations: recommend movies with same actor(s), director, genre, ...
 - Websites, blogs, news: recommend other sites with "similar" content

Limitations of content-based approach

- Finding the appropriate features
 - e.g., images, movies, music
- Overspecialization
 - Never recommends items outside user's content profile
 - People might have multiple interests
- Recommendations for new users
 - How to build a profile?

Recommendation Systems (II)

- Collaborative Filtering (user-user)**
 - Consider user c
 - Find set D of other users whose ratings are "similar" to c's ratings
 - Estimate user's ratings based on ratings of users in D

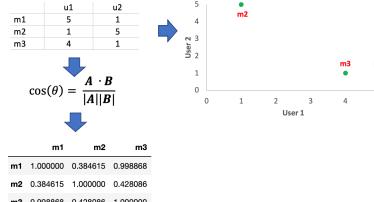
Recommendation Systems (II)

- Collaborative Filtering (user-user)**



Recommendation Systems (III)

Collaborative Filtering (item-item)



Pros and cons of collaborative filtering

- Works for any kind of item
 - No feature selection needed
- New user problem
- New item problem
- Sparsity of rating matrix
 - Cluster-based smoothing?

Another important problem

- Find **duplicate** and **near-duplicate** documents from a web crawl.
- Why is it important:
 - Identify **mirrored web pages**, and avoid indexing them, or serving them multiple times
 - Find **replicated news stories** and cluster them under a single story.
 - Identify plagiarism
- What if we wanted exact duplicates?

Main issues

- What is the **right representation** of the document when we check for similarity?
 - E.g., representing a document as a set of characters will not do (why?)
- When we have billions of documents, keeping the full text in memory is not an option.
 - We need to find a **shorter representation**
- How do we do **pairwise comparisons** of billions of documents?
 - If exact match was the issue, it would be ok, can we replicate this idea?

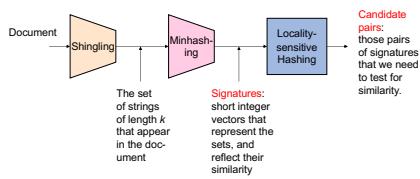
Three Essential Techniques

- Shingling:** convert documents, emails, etc., to sets.
- Minhashing:** convert large sets to short signatures, while preserving similarity.
- Locality-Sensitive Hashing (LSH):** focus on pairs of signatures likely to be similar.

Motivating problem

- Find **duplicate** and **near-duplicate** documents from a web crawl.
- If we wanted exact duplicates, we could do this by hashing
 - We will see how to adapt this technique for **near duplicate** documents

The Big Picture



Shingles: Compression Option

- To compress long shingles, we can hash them to (say) 4 bytes.
- Represent a doc by the set of **hash values** of its k -shingles.
- From now on we will assume that shingles are integers
- Collisions are possible, but very rare

Shingles

- A **k -shingle** (or **k -gram**) for a document is a sequence of k characters that appears in the document.
- Example: document = **abcab**, $k=2$
 - Set of 2-shingles = {**ab**, **bc**, **ca**}.
 - Option: regard shingles as a **bag**, and count **ab** twice.
- Represent a document by its set of k -shingles.

Shingling

- Shingle: a sequence of k contiguous characters

a rose is a rose is a rose
 a rose is
 rose is a
 rose is a
 ose is a r
 se is a ro
 e is a ros
 is a rose
 is a rose i
 a rose is
 a rose is

Fingerprinting

- Hash shingles to 64-bit integers

| Set of Shingles | Hash function (Rabin's fingerprints) | Set of 64-bit integers |
|-----------------|--------------------------------------|------------------------|
| a rose is | 1111 | |
| rose is a | 2222 | |
| rose is a | 3333 | |
| ose is a r | 4444 | |
| se is a ro | 5555 | |
| e is a ros | 6666 | |
| is a rose | 7777 | |
| s a rose i | 8888 | |
| a rose is | 9999 | |
| a rose is | 0000 | |

Basic Data Model: Sets

- Document: A document is represented as a **set** shingles (more accurately, hashes of shingles)
- Document similarity: **Jaccard** similarity of the sets of shingles.
 - Common shingles over the union of shingles
 - $\text{Sim}(C_1, C_2) = |C_1 \cap C_2| / |C_1 \cup C_2|$
- Although we use the documents as our driving example the techniques we will describe apply to any kind of sets.
 - E.g., similar customers or items.

From Sets to Boolean Matrices

- Represent the data as a boolean matrix M
 - Rows = the universe of all possible set elements
 - In our case, shingle fingerprints take values in $[0 \dots 2^{64}-1]$
 - Columns = the sets
 - In our case, documents, sets of shingle fingerprints
 - $M(r, S) = 1$ in row r and column S , if and only if r is a member of S .
- Typical matrix is **sparse**.
 - We do not really materialize the matrix

Example

- Universe: $U = \{A, B, C, D, E, F, G\}$

$$\begin{aligned} X &= \{A, B, F, G\} \\ Y &= \{A, E, F, G\} \end{aligned}$$

$$\text{Sim}(X, Y) = \frac{3}{5}$$

| | x | y |
|---|---|---|
| A | 1 | 1 |
| B | 1 | 0 |
| C | 0 | 0 |
| D | 0 | 0 |
| E | 0 | 1 |
| F | 1 | 1 |
| G | 1 | 1 |

Example

- Universe: $U = \{A, B, C, D, E, F, G\}$

$$\begin{aligned} X &= \{A, B, F, G\} \\ Y &= \{A, E, F, G\} \end{aligned}$$

$$\text{Sim}(X, Y) = \frac{3}{5}$$

| | x | y |
|---|---|---|
| A | 1 | 1 |
| B | 1 | 0 |
| C | 0 | 0 |
| D | 0 | 0 |
| E | 0 | 1 |
| F | 1 | 1 |
| G | 1 | 1 |

At least one of the columns has value 1

Minhashing

- Pick a random permutation of the rows (the universe U).
- Define “hash” function for set S
 - $h(S)$ is the index of the first row (in the permuted order) in which column S has 1.
 - OR
 - $h(S)$ is the index of the first element of S in the permuted order.
- Use k (e.g., $k = 100$) independent random permutations to create a signature.

Example of minhash signatures

- Input matrix

| | s_1 | s_2 | s_3 | s_4 |
|---|-------|-------|-------|-------|
| A | 1 | 0 | 1 | 0 |
| B | 1 | 0 | 0 | 1 |
| C | 0 | 1 | 0 | 1 |
| D | 0 | 1 | 0 | 1 |
| E | 0 | 1 | 0 | 1 |
| F | 1 | 0 | 1 | 0 |
| G | 1 | 0 | 1 | 0 |

A → C → F → D → B → E → G →

1 2 3 4

Example of minhash signatures

- Input matrix

| | s_1 | s_2 | s_3 | s_4 | s_5 |
|---|-------|-------|-------|-------|-------|
| A | 1 | 0 | 1 | 0 | |
| B | 1 | 0 | 0 | 1 | |
| C | 0 | 1 | 0 | 1 | |
| D | 0 | 1 | 0 | 1 | |
| E | 0 | 1 | 0 | 1 | |
| F | 1 | 0 | 1 | 0 | |
| G | 1 | 0 | 1 | 0 | |

A → B → C → D → E → F → G →

2 3 4 5 1

Working Assumption

- Documents that have lots of shingles in common have similar text, even if the text appears in different order.

- Careful:** you must pick k large enough, or most documents will have most shingles.
 - Extreme case $k = 1$: all documents are the same
 - $k = 5$ is OK for short documents; $k = 10$ is better for long documents.
- Alternative ways to define shingles:
 - Use words instead of characters
 - Anchor on stop words (to avoid templates)

Signatures

- Problem:** shingle sets are too large to be kept in memory.
- Key idea:** “hash” each set S to a small **signature** $\text{Sig}(S)$, such that:
 - $\text{Sig}(S)$ is small enough that we can fit a signature in main memory for each set.
 - $\text{Sim}(S_1, S_2)$ is (almost) the same as the “similarity” of $\text{Sig}(S_1)$ and $\text{Sig}(S_2)$. (Signature preserves similarity).
- Warning:** This method can produce **false negatives**, and **false positives** (if an additional check is not made).
 - False negatives:** Similar items deemed as non-similar
 - False positives:** Non-similar items deemed as similar

Example

- Universe: $U = \{A, B, C, D, E, F, G\}$

$$\begin{aligned} X &= \{A, B, F, G\} \\ Y &= \{A, E, F, G\} \end{aligned}$$

| | x | y |
|---|---|---|
| A | 1 | 1 |
| B | 1 | 0 |
| C | 0 | 0 |
| D | 0 | 0 |
| E | 0 | 1 |
| F | 1 | 1 |
| G | 1 | 1 |

Both columns have value 1

Example of minhash signatures

- Input matrix

| | s_1 | s_2 | s_3 | s_4 | s_5 |
|---|-------|-------|-------|-------|-------|
| A | 1 | 0 | 1 | 0 | |
| B | 1 | 0 | 0 | 1 | |
| C | 0 | 1 | 0 | 1 | |
| D | 0 | 1 | 0 | 1 | |
| E | 0 | 1 | 0 | 1 | |
| F | 1 | 0 | 1 | 0 | |
| G | 1 | 0 | 1 | 0 | |

C → G → A → B → D → E → F →

3 4 1 2 5 0

Example of minhash signatures

- Input matrix

| | S ₁ | S ₂ | S ₃ | S ₄ | S ₅ |
|---|----------------|----------------|----------------|----------------|----------------|
| A | 1 | 0 | 0 | 1 | 0 |
| B | 1 | 0 | 0 | 0 | 1 |
| C | 0 | 1 | 0 | 1 | 0 |
| D | 0 | 1 | 0 | 1 | 0 |
| E | 0 | 1 | 0 | 1 | 0 |
| F | 1 | 0 | 1 | 0 | 0 |
| G | 1 | 0 | 1 | 0 | 0 |

| | S ₁ | S ₂ | S ₃ | S ₄ | S ₅ |
|----------------|----------------|----------------|----------------|----------------|----------------|
| h ₁ | 1 | 2 | 1 | 3 | 2 |
| h ₂ | 2 | 1 | 3 | 1 | 3 |
| h ₃ | 1 | 3 | 2 | 1 | 3 |

≈

Signature matrix

- $\text{Sig}(S)$ = vector of hash values
 - e.g., $\text{Sig}(S_1) = [2, 1, 1]$
- $\text{Sig}(S_i)$ = value of the i -th hash function for set S
 - E.g., $\text{Sig}(S_2, 3) = 1$

Hash function Property

$$\Pr(h(S_1) = h(S_2)) = \text{Sim}(S_1, S_2)$$

where the probability is over all choices of permutations.

Why?

- The first row where one of the two sets has value 1 belongs to the union.
 - Recall that union contains rows with at least one 1.
- We have equality if both sets have value 1, and this row belongs to the intersection

Example

- Universe: $U = \{A, B, C, D, E, F, G\}$
- $X = \{A, B, F, G\}$
- $Y = \{A, E, F, G\}$

| | X | Y | |
|---|---|---|---|
| A | 1 | 1 | |
| B | 1 | 0 | * |
| C | 0 | 0 | * |
| D | 0 | 0 | * |
| E | 0 | 1 | * |
| F | 1 | 1 | * |
| G | 1 | 1 | * |

The question is what is the value of the first * element

- Union = $\{A, B, E, F, G\}$
- Intersection = $\{A, F, G\}$

Example

- Universe: $U = \{A, B, C, D, E, F, G\}$
- $X = \{A, B, F, G\}$
- $Y = \{A, E, F, G\}$

| | X | Y | |
|---|---|---|---|
| A | 1 | 1 | |
| B | 1 | 0 | * |
| C | 0 | 0 | * |
| D | 0 | 0 | * |
| E | 0 | 1 | * |
| F | 1 | 1 | * |
| G | 1 | 1 | * |

If it belongs to the intersection, then $h(X) = h(Y)$

- Union = $\{A, B, E, F, G\}$
- Intersection = $\{A, F, G\}$

Is it now feasible?

- Assume a billion rows
- Hard to pick a random permutation of 1 billion...
- Even representing a random permutation requires 1 billion entries!!!**
- How about accessing rows in permuted order?
⊗

Being more practical

- Instead of permuting the rows we will apply a **hash function** that maps the rows to a new (possibly larger) space
 - The value of the hash function is the position of the row in the new order (permutation).
 - Each set is represented by the smallest hash value among the elements in the set
- The space of the hash functions should be such that if we select one at random each element (row) has equal probability to have the smallest value
- Min-wise independent hash functions

Algorithm – All sets, k hash functions

| x | Row | S ₁ | S ₂ | h(x) | g(x) | Sig ₁ | Sig ₂ |
|---|-----|----------------|----------------|------|------|------------------|------------------|
| 0 | A | 1 | 0 | 1 | 3 | h(0) = 1 | |
| 1 | B | 0 | 1 | 2 | 0 | g(0) = 3 | 3 |
| 2 | C | 1 | 1 | 3 | 2 | | |
| 3 | D | 1 | 0 | 4 | 4 | h(1) = 2 | 1 |
| 4 | E | 0 | 1 | 0 | 1 | g(1) = 0 | 0 |

$h(x) = x \times 1 \bmod 5$
 $g(x) = 2x \bmod 5$

| h(Row) | Row | S ₁ | S ₂ | g(Row) | Row | S ₁ | S ₂ |
|--------|-----|----------------|----------------|--------|-----|----------------|----------------|
| 0 | E | 0 | 1 | 0 | 8 | 0 | 1 |
| 1 | A | 1 | 0 | 1 | E | 0 | 1 |
| 2 | B | 0 | 1 | 2 | C | 1 | 0 |
| 3 | C | 1 | 1 | 3 | A | 1 | 1 |
| 4 | D | 1 | 0 | 4 | D | 1 | 0 |

$h(3) = 4$
 $g(3) = 4$

Implementation

- Often, data is given by column, not row.
 - E.g., columns = documents, rows = shingles.
- If so, sort matrix once so it is by row.
- And **always** compute $h_i(r)$ only once for each row.

Example

- Universe: $U = \{A, B, C, D, E, F, G\}$

- $X = \{A, B, F, G\}$

- $Y = \{A, E, F, G\}$

Rows CD could be anywhere they do not affect the probability

- Union = $\{A, B, E, F, G\}$

- Intersection = $\{A, F, G\}$

| X | Y |
|---|---|
| A | 1 |
| B | 0 |
| C | 0 |
| D | 0 |
| E | 0 |
| F | 1 |
| G | 1 |

| X | Y |
|---|---|
| D | * |
| C | * |
| B | * |
| A | * |
| F | * |
| E | * |
| G | * |

| X | Y |
|---|---|
| D | 0 |
| C | 0 |
| B | 0 |
| A | 0 |
| F | 0 |
| E | 0 |
| G | 0 |

Example

- Universe: $U = \{A, B, C, D, E, F, G\}$
- $X = \{A, B, F, G\}$
- $Y = \{A, E, F, G\}$

Every element of the union is equally likely to be the * element

$$\Pr(h(X) = h(Y)) = \frac{|[A, F, G]|}{|[A, B, E, F, G]|} = \frac{3}{5} = \text{Sim}(X, Y)$$

- Union = $\{A, B, E, F, G\}$
- Intersection = $\{A, F, G\}$

| X | Y |
|---|---|
| A | 1 |
| B | 0 |
| C | 0 |
| D | 0 |
| E | 0 |
| F | 1 |
| G | 1 |

| X | Y |
|---|---|
| D | * |
| C | * |
| B | * |
| A | * |
| F | * |
| E | * |
| G | * |

| X | Y |
|---|---|
| D | 0 |
| C | 0 |
| B | 0 |
| A | 0 |
| F | 0 |
| E | 0 |
| G | 0 |

Example

- Universe: $U = \{A, B, C, D, E, F, G\}$

- $X = \{A, B, F, G\}$

- $Y = \{A, E, F, G\}$

The * rows belong to the union

| X | Y |
|---|---|
| A | 1 |
| B | 0 |
| C | 0 |
| D | 0 |
| E | 0 |
| F | 1 |
| G | 1 |

| X | Y |
|---|---|
| D | * |
| C | * |
| B | * |
| A | * |
| F | * |
| E | * |
| G | * |

| X | Y |
|---|---|
| D | 0 |
| C | 0 |
| B | 0 |
| A | 0 |
| F | 0 |
| E | 0 |
| G | 0 |

Similarity for Signatures

The **similarity of signatures** is the fraction of the hash functions in which they agree.

| S ₁ | S ₂ | Actual | Sig |
|----------------|----------------|--------|-----|
| 1 | 2 | 0 | 0 |
| 2 | 1 | 1/2 | 0 |
| 0 | 0 | 0 | 0 |
| 1 | 1 | 3/4 | 1 |
| 0 | 1 | 0 | 0 |

Zero similarity is preserved
High similarity is well approximated

Algorithm – All sets, k hash functions

Pick $k=100$ hash functions (h_1, \dots, h_k)

| | |
|---|---|
| In practice this means selecting the hash function parameters | In practice only the rows (shingles) that appear in the data |
| for each row r | compute $h_i(r)$ |
| if column S that has 1 in row r | $h_i(r)$ = index of row r in permutation |
| if $h_i(r)$ is a smaller value than $\text{Sig}(S, i)$ then | $Sig(S, i) = h_i(r)$; Find the row r with minimum index |
| | $Sig(S, i)$ will become the smallest value of $h_i(r)$ among all rows (shingles) for which column S has value 1 (shingle belongs in S); i.e., $h_i(r)$ gives the min index for the i -th permutation |

Locality-Sensitive Hashing

What we want: a function $f(X, Y)$ that tells whether or not X and Y are a **candidate pair**: a pair of elements whose similarity must be evaluated.

A simple idea: X and Y are a candidate pair if they have the same min-hash signature.

Easy to test by hashing the signatures.

Similar sets are more likely to have the same signature.

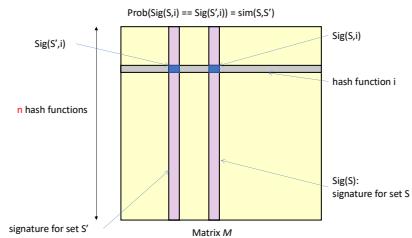
Likely to produce many false negatives.

Requiring full match of signature is strict, some similar sets will be lost.

Improvement: Compute multiple signatures; candidate pairs should have at least one common signature.

Reduce the probability for false negatives.

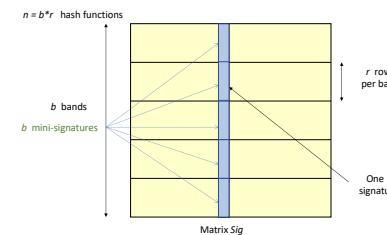
Signature matrix reminder



Partition into Bands – (1)

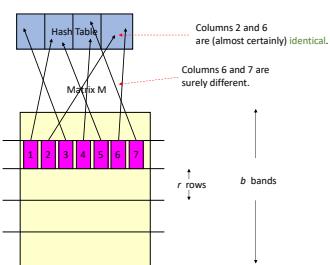
- Divide the signature matrix Sig into b bands of r rows.
 - Each band is a mini-signature with r hash functions.

Partition into Bands – (1)



Partition into Bands – (2)

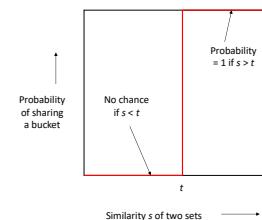
- Divide the signature matrix Sig into b bands of r rows.
 - Each band is a mini-signature with r hash functions.
- For each band, hash the mini-signature to a hash table with k buckets.
 - Make k as large as possible so that mini-signatures that hash to the same bucket are almost certainly identical.
 - Make k as large as possible so that mini-signatures that hash to the same bucket are almost certainly identical.



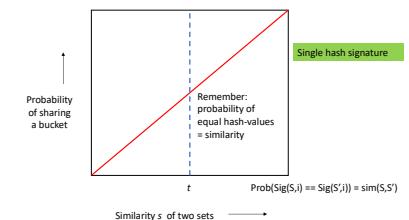
Partition into Bands – (3)

- Divide the signature matrix Sig into b bands of r rows.
 - Each band is a mini-signature with r hash functions.
- For each band, hash the mini-signature to a hash table with k buckets.
 - Make k as large as possible so that mini-signatures that hash to the same bucket are almost certainly identical.
- Candidate** column pairs are those that hash to the same bucket for at least 1 band.
- Tune b and r to catch most similar pairs, but few non-similar pairs.

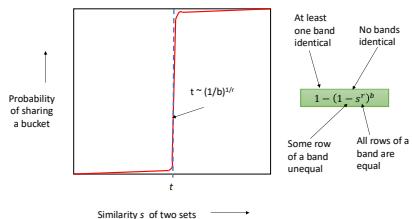
Analysis of LSH – What we want



What One Band of One Row Gives You

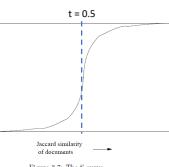


What b Bands of r Rows Gives You



Example: $b = 20; r = 5$

| s | $1 - (1 - s)^b$ |
|-----|-----------------|
| .2 | .006 |
| .3 | .047 |
| .4 | .186 |
| .5 | .470 |
| .6 | .802 |
| .7 | .975 |
| .8 | .996 |



Suppose S_1, S_2 are 80% Similar

- We want all 80%-similar pairs. Choose 20 bands of 5 integers/band.
- Probability S_1, S_2 identical in one particular band: $(0.8)^5 = 0.328$.
- Probability S_1, S_2 are not similar in any of the 20 bands: $(1 - 0.328)^{20} = 0.00035$
- i.e., about 1/3000-th of the 80%-similar column pairs are false negatives.
- Probability S_1, S_2 are similar in at least one of the 20 bands: $1 - 0.00035 = 0.999$

Suppose S_1, S_2 Only 40% Similar

- Probability S_1, S_2 identical in any one particular band: $(0.4)^5 = 0.01$.
- Probability S_1, S_2 identical in at least 1 of 20 bands: $\leq 20 * 0.01 = 0.2$.
- But false positives much lower for similarities << 40%.