

CSC-591: Foundations of Data Science

T/Th. 12:50-2:05pm. EBI-1005.

Ranga Raju Vatsavai

Chancellors Faculty Excellence Associate Professor in Geospatial Analytics
Department of Computer Science, North Carolina State University (NCSU)
Associate Director, Center for Geospatial Analytics, NCSU
&
Joint Faculty, Oak Ridge National Laboratory (ORNL)

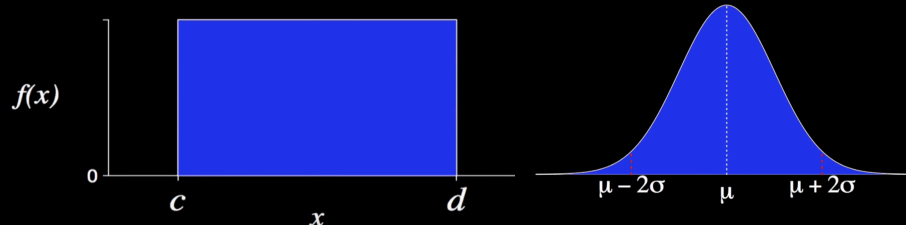
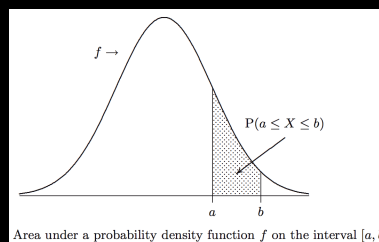
W3: 9/1-3/15

Administrative

- Books and Lecture Notes
- Introductory Statistics
 - <https://www.openintro.org/stat/>
- Linear and Matrix Algebra
 - <http://www.math.uwaterloo.ca/~hwolkowi/matrixcookbook.pdf>
 - <http://www.mathworks.com/moler/eigs.pdf>

Key points from 9/1

- Continuous distributions
- Uniform: $U(c,d)$
- Normal: $N(\mu, \sigma^2)$



9/8/15

© Raju Vatsavai

CSC-591. 3

Today

- Basic Linear Algebra and Calculus

9/8/15

© Raju Vatsavai

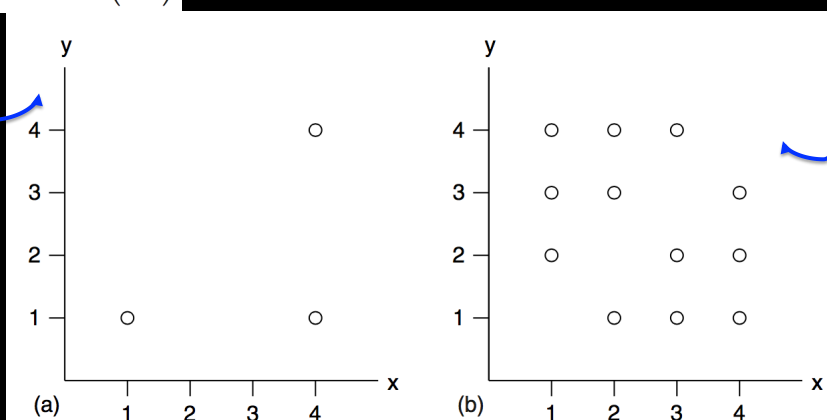
CSC-591. 4

Why?

- Let us consider the following datasets

$$\mathbf{A} = (\mathbf{x} : \mathbf{y}) = \begin{pmatrix} 1 & 1 \\ 4 & 1 \\ 4 & 4 \end{pmatrix}$$

$$\mathbf{B}' = (\mathbf{x} : \mathbf{y})' = \begin{pmatrix} \mathbf{x}' \\ \mathbf{y}' \end{pmatrix} = \begin{pmatrix} 1 & 1 & 1 & 2 & 2 & 2 & 3 & 3 & 3 & 4 & 4 & 4 \\ 2 & 3 & 4 & 1 & 3 & 4 & 1 & 2 & 4 & 1 & 2 & 3 \end{pmatrix}$$



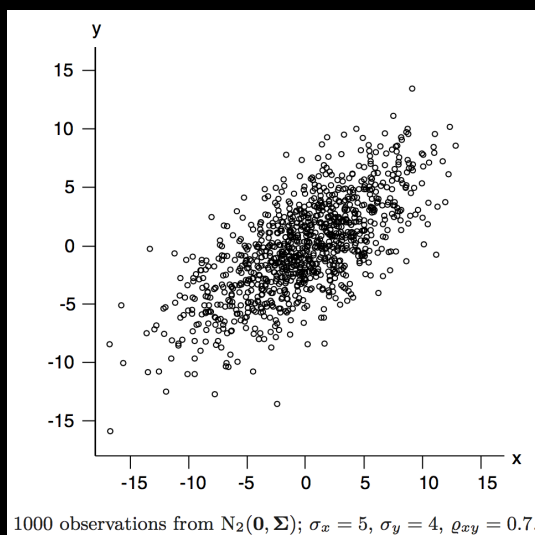
9/8/15

© Raju Vatsavai

CSC-591. 5

Why

- Let us consider the following datasets



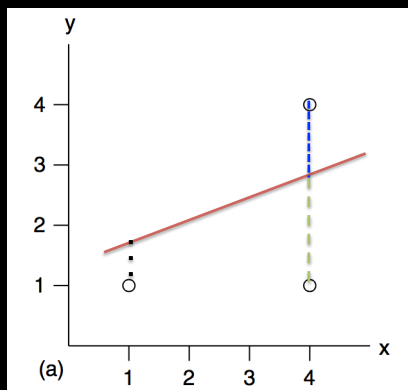
9/8/15

© Raju Vatsavai

CSC-591. 6

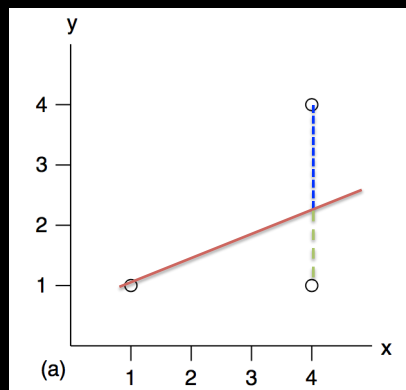
Question

- Draw a line $\hat{y} = \beta x + \varepsilon$ into each scatter plot such that the sum of squared vertical distances $(y_i - \hat{y}_i)^2$ would be as small as possible



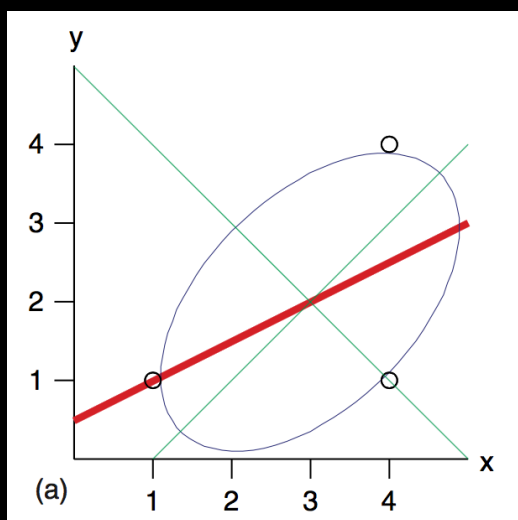
9/8/15

© Raju Vatsavai



CSC-591. 7

Solution



9/8/15

© Raju Vatsavai

CSC-591. 8

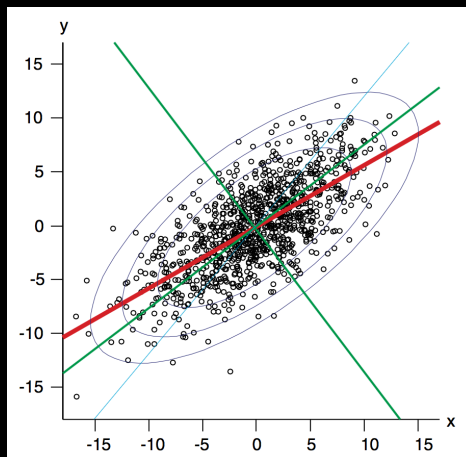
Solution

Assume Normal
Distribution:

$N_2(\mathbf{0}, \Sigma)$; $\sigma_x = 5$, $\sigma_y = 4$, $\rho_{xy} = 0.7$.

Regression line has the
slope $\beta \approx \rho_{xy}\sigma_y/\sigma_x$

Direction of first major
axis of contour ellipse is
determined by the first
eigenvector of Σ



9/8/15

© Raju Vatsavai

CSC-591. 9

Linear Algebra

- Linear algebra is the study of linear maps on finite-dimensional vector spaces.

9/8/15

© Raju Vatsavai

CSC-591. 10

Complex Numbers

Definition *complex numbers*

- A **complex number** is an ordered pair (a, b) , where $a, b \in \mathbf{R}$, but we will write this as $a + bi$.
- The set of all complex numbers is denoted by \mathbf{C} :

$$\mathbf{C} = \{a + bi : a, b \in \mathbf{R}\}.$$

- **Addition and multiplication** on \mathbf{C} are defined by

$$(a + bi) + (c + di) = (a + c) + (b + d)i,$$

$$(a + bi)(c + di) = (ac - bd) + (ad + bc)i;$$

here $a, b, c, d \in \mathbf{R}$.

Example

- $(2+3i)(4+5i)$
- $-7 + 22i$

Properties of Complex Arithmetic

commutativity

$\alpha + \beta = \beta + \alpha$ and $\alpha\beta = \beta\alpha$ for all $\alpha, \beta \in \mathbb{C}$;

associativity

$(\alpha + \beta) + \lambda = \alpha + (\beta + \lambda)$ and $(\alpha\beta)\lambda = \alpha(\beta\lambda)$ for all $\alpha, \beta, \lambda \in \mathbb{C}$;

identities

$\lambda + 0 = \lambda$ and $\lambda 1 = \lambda$ for all $\lambda \in \mathbb{C}$;

additive inverse

for every $\alpha \in \mathbb{C}$, there exists a unique $\beta \in \mathbb{C}$ such that $\alpha + \beta = 0$;

multiplicative inverse

for every $\alpha \in \mathbb{C}$ with $\alpha \neq 0$, there exists a unique $\beta \in \mathbb{C}$ such that $\alpha\beta = 1$;

distributive property

$\lambda(\alpha + \beta) = \lambda\alpha + \lambda\beta$ for all $\lambda, \alpha, \beta \in \mathbb{C}$.

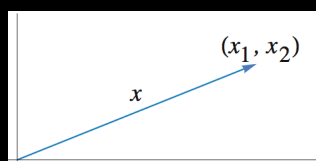
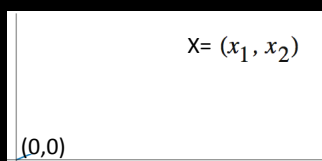
9/8/15

© Raju Vatsavai

CSC-591. 13

Vector

- $[1.2, 2.5]$
 - Is a 2-vector over \mathbb{R} , written as \mathbb{R}^2
 - General notation, \mathbb{R}^d
- A vector can be thought as a function
- $F^d : \{0, 1, 2, \dots, d-1\} \rightarrow F$
- Sparse vector
 - If most elements are 0's



Elements of \mathbb{R}^2 can be thought of as points or vectors

9/8/15

© Raju Vatsavai

CSC-591. 14

What kind of data

- Document as bag of words
 - $f : \text{WORDS} \rightarrow \mathbb{R}$
- Collections of attributes
 - Physical characteristics of persons
 - Demographic records of customers
- Probability distribution, e.g., $\{1:1/6, 2:1/6, \dots, 6:1/6\}$
- Images
 - $\{r, g, b\}$

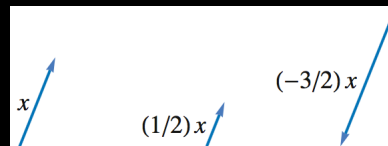
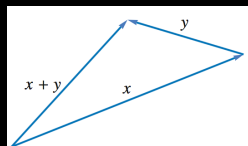
9/8/15

© Raju Vatsavai

CSC-591. 15

Vector Operations

- Complex numbers, translation is achieved by adding a complex number. $f(z) = z + (1+2i)$
- Vector addition: add element-wise $[u_1, u_2] + [v_1, v_2] = [u_1+v_1, u_2+v_2]$
- Scalar-vector multiplication (is a scaling operation): αv



9/8/15

© Raju Vatsavai

CSC-591. 16

Convex Combinations

- For given $0 \leq \alpha \leq 1$, and $0 \leq \beta \leq 1$, and $\alpha + \beta = 1$; an expression of form: $\alpha u + \beta v$ is called convex combination of u and v
- Application (average of two images)
- Lets say u and v are two images, and $\alpha = 0.5$ and $\beta = 0.5$, the $\alpha u + \beta v$ is average of two images

Vector Operations

- Dot product of two D -vectors is sum of the product of corresponding entries
 - $u \cdot v = \sum_{k \in D} u[k]v[k]$
- What do you think the output is?
- Examples
 - Total cost of a product
 - Measuring similarity

Vector Space

A **vector space** is a set V along with an addition on V and a scalar multiplication on V such that the following properties hold:

commutativity

$$u + v = v + u \text{ for all } u, v \in V;$$

associativity

$$(u + v) + w = u + (v + w) \text{ and } (ab)v = a(bv) \text{ for all } u, v, w \in V \text{ and all } a, b \in \mathbf{F};$$

additive identity

$$\text{there exists an element } 0 \in V \text{ such that } v + 0 = v \text{ for all } v \in V;$$

additive inverse

$$\text{for every } v \in V, \text{ there exists } w \in V \text{ such that } v + w = 0;$$

multiplicative identity

$$1v = v \text{ for all } v \in V;$$

distributive properties

$$a(u + v) = au + av \text{ and } (a + b)v = av + bv \text{ for all } a, b \in \mathbf{F} \text{ and all } u, v \in V.$$

Vector Space \mathbf{F}^S

- If S is a set, then \mathbf{F}^S denotes the set of functions from S to \mathbf{F} .

- For $f, g \in \mathbf{F}^S$, the **sum** $f + g \in \mathbf{F}^S$ is the function defined by

$$(f + g)(x) = f(x) + g(x)$$

for all $x \in S$.

- For $\lambda \in \mathbf{F}$ and $f \in \mathbf{F}^S$, the **product** $\lambda f \in \mathbf{F}^S$ is the function defined by

$$(\lambda f)(x) = \lambda f(x)$$

for all $x \in S$.

- \mathbf{F}^S is a vector space
- *In general, a vector space is an abstract entity whose elements might be lists, functions, or weird objects.*

Subspace

- A subset U of V is called a *subspace* of V if U is also a vector space (using the same addition and scalar multiplication as on V).
- A subset U of V is a subspace of V if and only if U satisfies the following three conditions:
 - Additive identity: $0 \in U$
 - Closed under addition: $u, w \in U \Rightarrow u + w \in U$
 - Closed under scalar multiplication
 $a \in F, u \in U \Rightarrow au \in U$

Important properties of subspaces

Sum of subspaces is the smallest containing subspace

Suppose U_1, \dots, U_m are subspaces of V . Then $U_1 + \dots + U_m$ is the smallest subspace of V containing U_1, \dots, U_m .

Direct Sum of Subspaces

Suppose U_1, \dots, U_m are subspaces of V .

- The sum $U_1 + \dots + U_m$ is called a **direct sum** if each element of $U_1 + \dots + U_m$ can be written in only one way as a sum $u_1 + \dots + u_m$, where each u_j is in U_j .
- If $U_1 + \dots + U_m$ is a direct sum, then $U_1 \oplus \dots \oplus U_m$ denotes $U_1 + \dots + U_m$, with the \oplus notation serving as an indication that this is a direct sum.

Linear combination of vectors

- An expression: $\alpha_1 v_1 + \dots + \alpha_n v_n$ is a linear combination of vectors $v_1 \dots v_n$ and scalars $\alpha_1 \dots \alpha_n$ are coefficients of the linear combinations
- Examples
 - Given a set of raw materials, different products can be made with different combinations of materials, then total resource utilization is linear combination of materials for each product ($u = \alpha_1 v_1 + \dots + \alpha_n v_n$)
 - In F^3 $(17, 4, 2)$ is a linear combination of $(2, 1, -3)$, $(1, -2, 4)$ because
 - $(17, -4, 2) = 6(2, 1, -3) + 5(1, -2, 4)$

9/8/15

© Raju Vatsavai

CSC-591. 23

Span

The set of all linear combinations of a list of vectors v_1, \dots, v_m in V is called the **span** of v_1, \dots, v_m , denoted $\text{span}(v_1, \dots, v_m)$. In other words,

$$\text{span}(v_1, \dots, v_m) = \{a_1 v_1 + \dots + a_m v_m : a_1, \dots, a_m \in F\}.$$

The span of the empty list $()$ is defined to be $\{0\}$.

- Previous example shows that in F^3
 - $(17, -4, 2) \in \text{span}((2, 1, -3), (1, -2, 4))$
- Spans
 - If $\text{span}(v_1, \dots, v_m)$ equals V , we say that v_1, \dots, v_m spans V

9/8/15

© Raju Vatsavai

CSC-591. 24

Linearly independent

- A list v_1, \dots, v_m of vectors in V is called **linearly independent** if the only choice of $a_1, \dots, a_m \in \mathbb{F}$ that makes $a_1 v_1 + \dots + a_m v_m$ equal 0 is $a_1 = \dots = a_m = 0$.
- The empty list $()$ is also declared to be linearly independent.

Bases

- A basis of V is a list of vectors in V that is linearly independent and spans V .
- Examples
 - The list $(1,0, \dots, 0), (0,1, \dots, 0), \dots, (0, \dots, 0, 1)$ is a basis of F^n , called the standard basis of F^n
 - The list $(1,2), (3,5)$ is a basis of F^2

Criteria for basis

- A list v_1, \dots, v_n of vectors in V is a basis of V if and only if for every $v \in V$ can be written uniquely in the form $v = a_1 v_1 + \dots + a_n v_n$, where $a_1, \dots, a_n \in F$.
- Every spanning list in a vector space can be reduced to basis of the vector space
- Every finite-dimensional vector space has a basis

9/8/15

© Raju Vatsavai

CSC-591. 27

Matrix Algebra

- Addition
- Multiplication
- Inverse

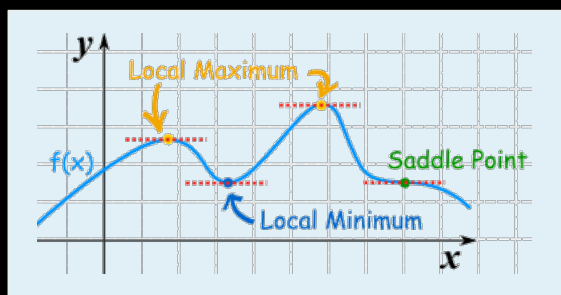
9/8/15

© Raju Vatsavai

CSC-591. 28

Basic Calculus

- How to finding maximum and minim of a function



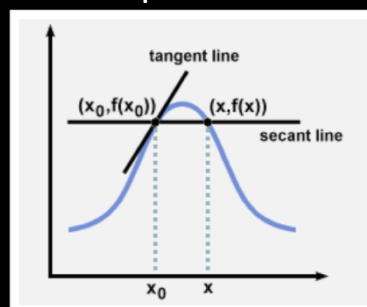
9/8/15

© Raju Vatsavai

CSC-591. 29

Derivative

- The derivative of a function $f(x)$ at $x = x_0$, denoted $f'(x_0)$ or (x_0) , can be naively defined as the slope of the graph of f at $x = x_0$.
- Derivative essentially finds the slope of a function



9/8/15

© Raju Vatsavai

CSC-591. 30

Acknowledgements

- Vipin Kumar (Minnesota)
- Jiawei Han (UIUC)
- Hanspeter Pfister (Harvard)
- Larry Wasserman (CMU)