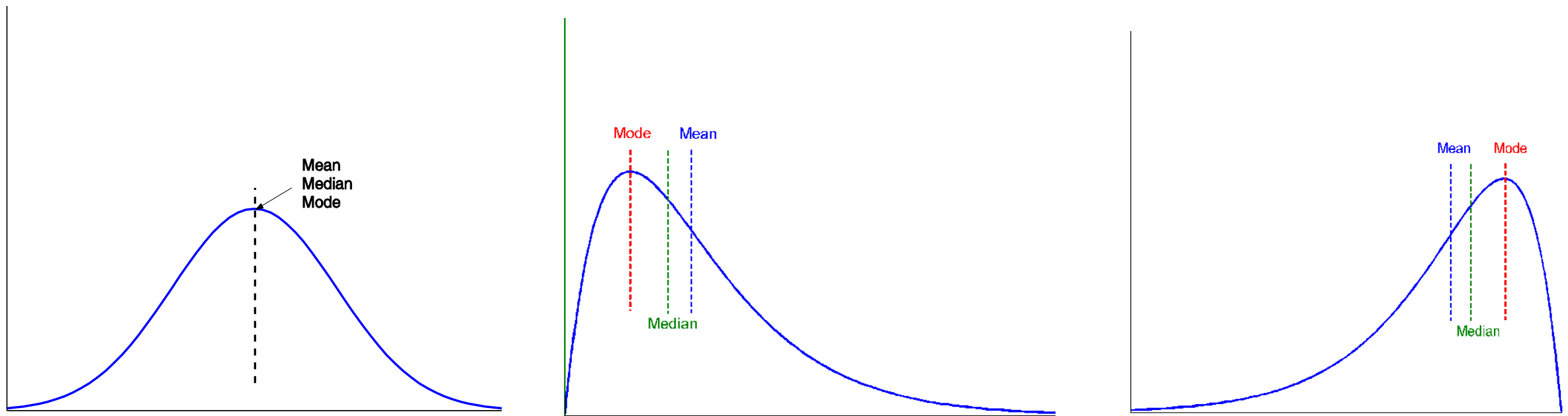


Computational Intelligence & Machine Learning

Symmetric vs Skewed Data

Median, mean and mode of symmetric, positively and negatively skewed data



Dispersion

Quartiles, outliers and boxplots

Quartiles: Q_1 (25th percentile), Q_3 (75th percentile)

Inter-quartile range: $IQR = Q_3 - Q_1$

Five number summary: min, Q_1 , median, Q_3 , max

Boxplot: ends of the box are the quartiles; median is marked; add whiskers, and plot outliers individually

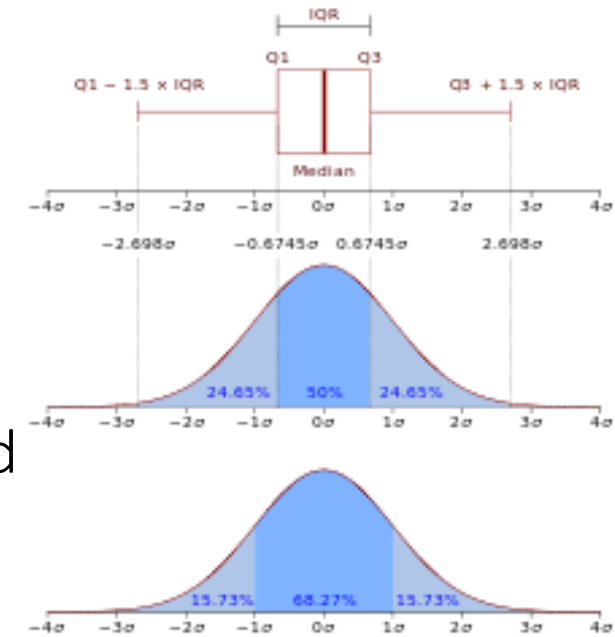
Outlier: usually, a value higher/lower than $1.5 \times IQR$

Variance and standard deviation (*sample: s , population: σ*)

Variance: (algebraic, scalable computation)

Standard deviation s (or σ) is the square root of variance s^2 (or σ^2)

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n-1} \left[\sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 \right] \quad \sigma^2 = \frac{1}{N} \sum_{i=1}^n (x_i - \mu)^2 = \frac{1}{N} \sum_{i=1}^n x_i^2 - \mu^2$$



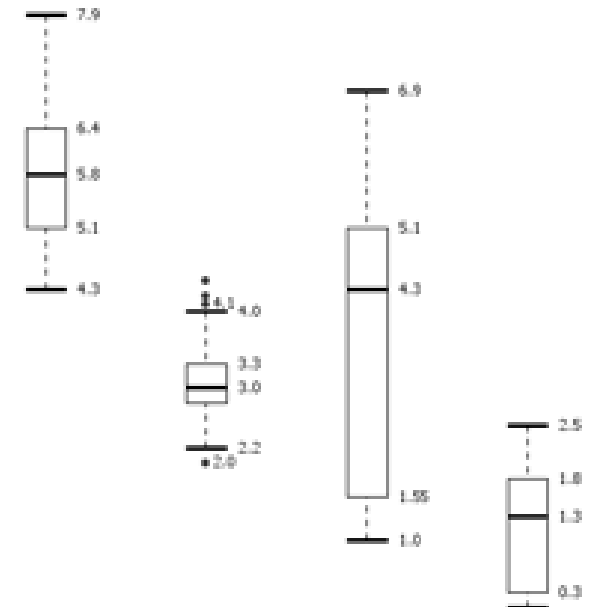
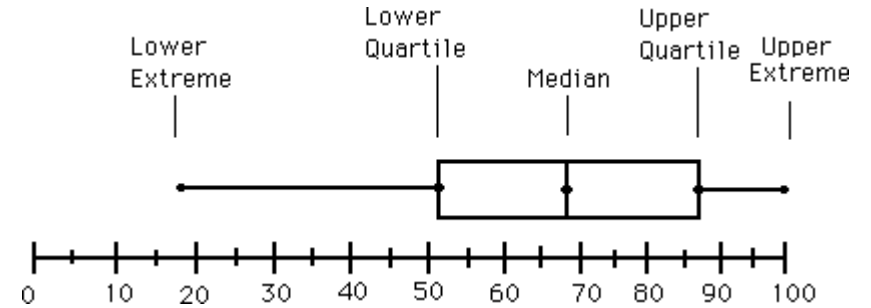
Boxplot

Five-number summary of a distribution

Minimum, Q1, Median, Q3, Maximum

Boxplot

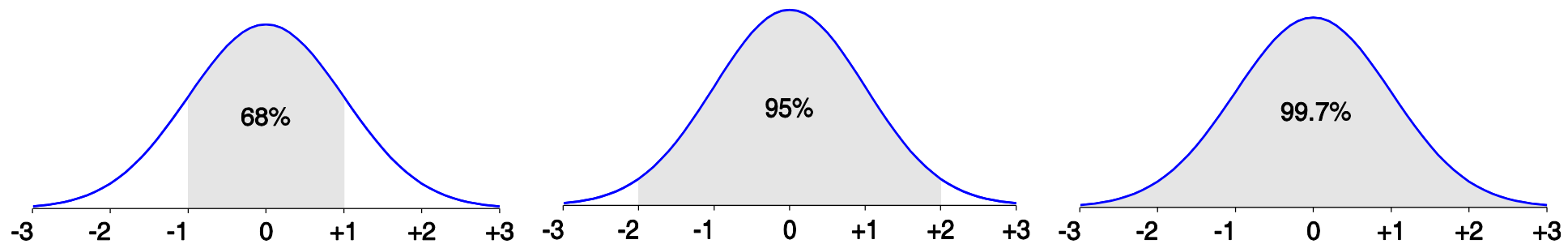
- Data are represented with a box
- The ends of the box are at the first and third quartiles, i.e., the height of the box is IQR
- The median is marked by a line within the box
- Whiskers: two lines outside the box extended to Minimum and Maximum
- Outliers: points beyond a specified outlier threshold, plotted individually



Example: Normal Distribution

The normal (distribution) curve

- From $\mu - \sigma$ to $\mu + \sigma$: contains about 68% of the measurements (μ : mean, σ : standard deviation)
- From $\mu - 2\sigma$ to $\mu + 2\sigma$: contains about 95% of it
- From $\mu - 3\sigma$ to $\mu + 3\sigma$: contains about 99.7% of it

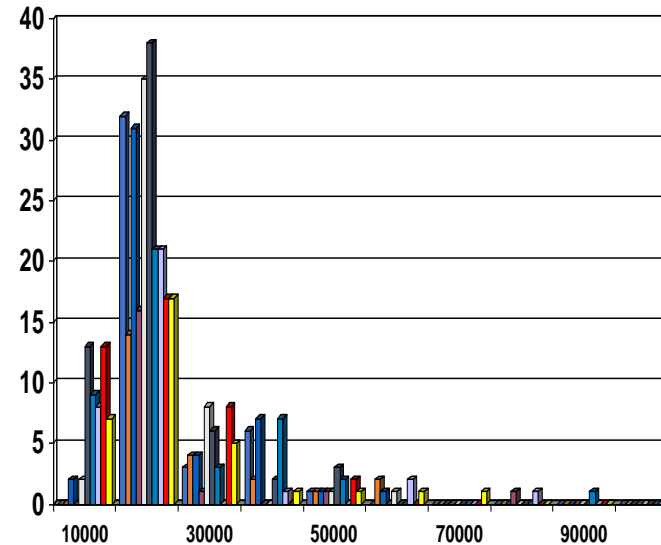


Visualization

- **Boxplot:** graphic display of five-number summary
- **Histogram:** x-axis are values, y-axis repres. frequencies
- **Quantile plot:** each value x_i is paired with f_i indicating that approximately $100 f_i \%$ of data are $\leq x_i$
- **Quantile-quantile (q-q) plot:** graphs the quantiles of one univariant distribution against the corresponding quantiles of another
- **Scatter plot:** each pair of values is a pair of coordinates and plotted as points in the plane

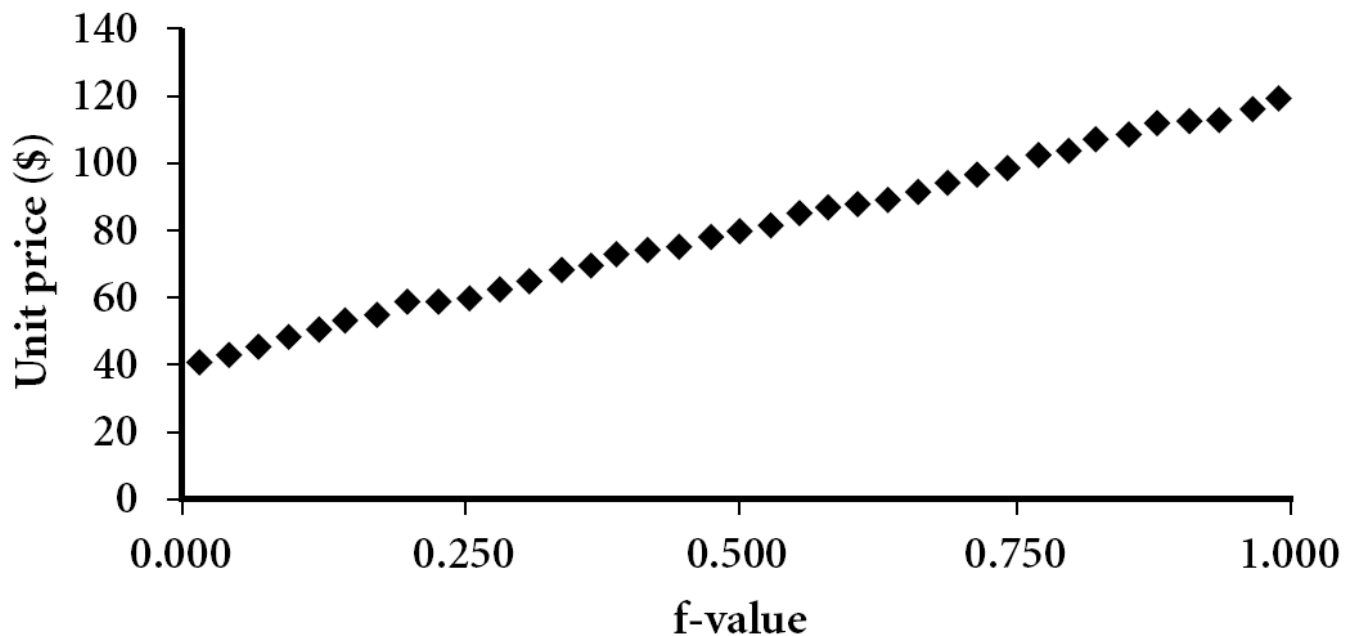
Histograms

- **Histogram:** Graph display of tabulated frequencies, shown as bars
- It shows what proportion of cases fall into each of several categories
- Differs from a bar chart in that it is the *area* of the bar that denotes the value, not the height as in bar charts, a crucial distinction when the categories are not of uniform width
- The categories are usually specified as non-overlapping intervals of some variable. The categories (bars) must be adjacent



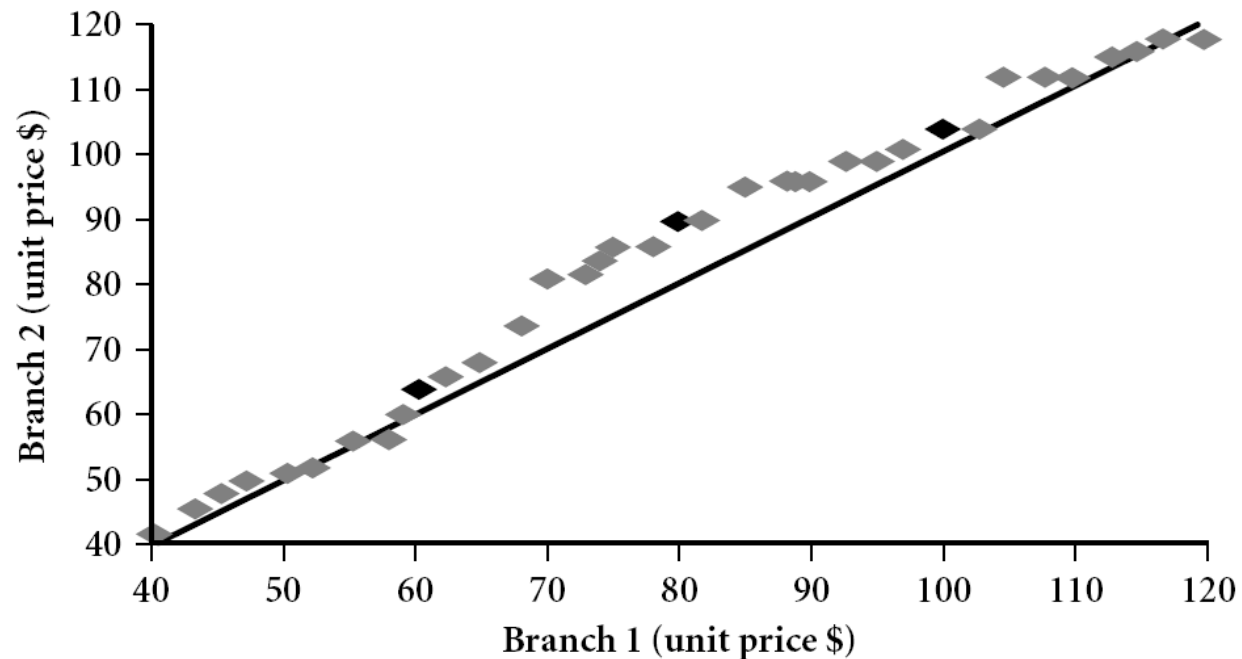
Quantile

- Displays all the data (allowing the user to assess both the overall behavior and unusual occurrences)
- Plots **quantile** information
- For a data x_i data sorted in increasing order, f_i indicates that approximately $100 \times f_i\%$ of the data are below or equal to the value x_i



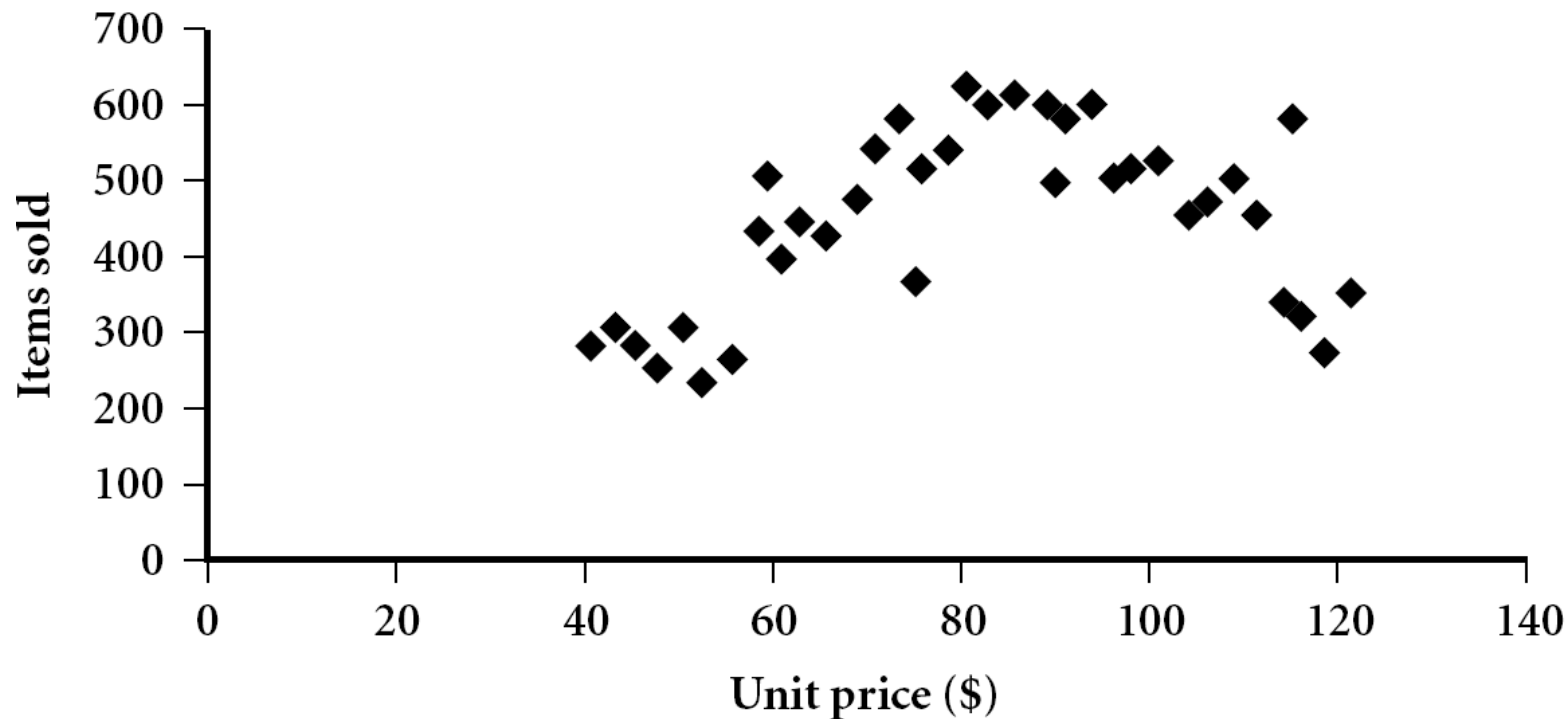
Quantile-Quantile (Q-Q) Plot

- Graphs the quantiles of one univariate distribution against the corresponding quantiles of another
- View: Is there a shift in going from one distribution to another?
- Example shows unit price of items sold at Branch 1 vs. Branch 2 for each quantile. Unit prices of items sold at Branch 1 tend to be lower than those at Branch 2.

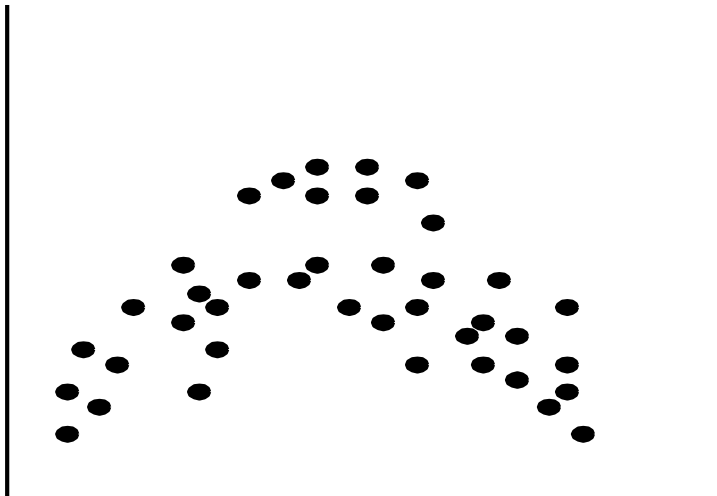
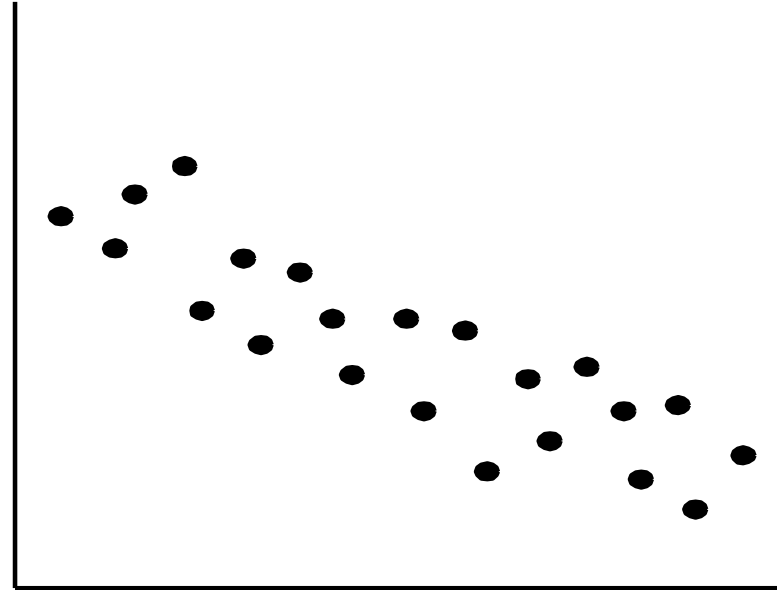
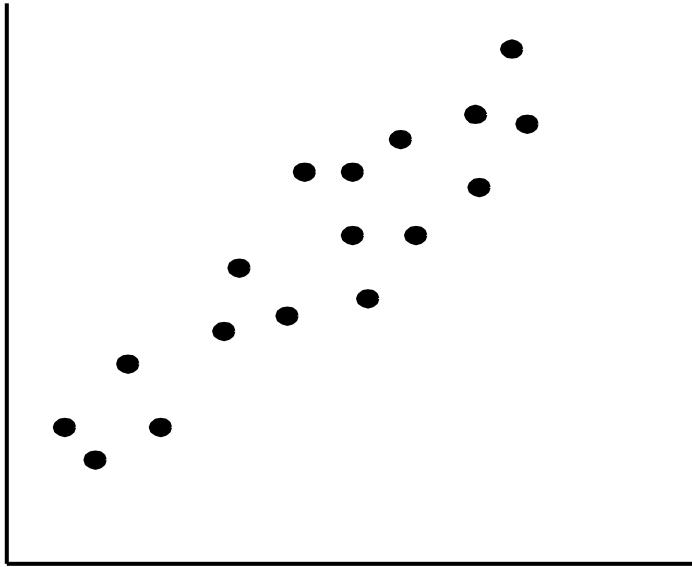


Scatter

- Provides a first look at bivariate data to see clusters of points, outliers, etc
- Each pair of values is treated as a pair of coordinates and plotted as points in the plane



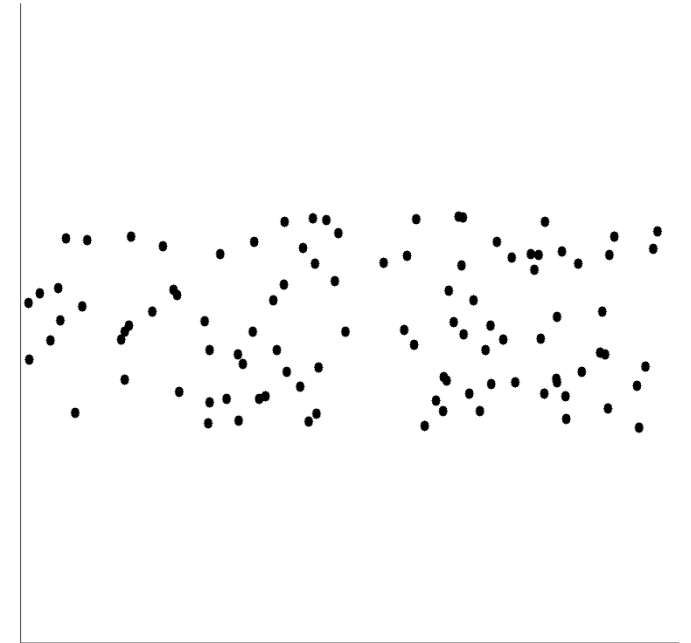
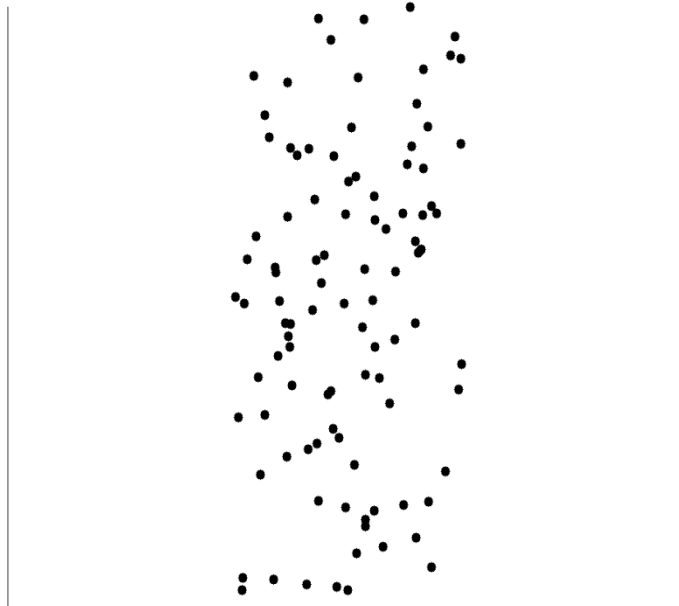
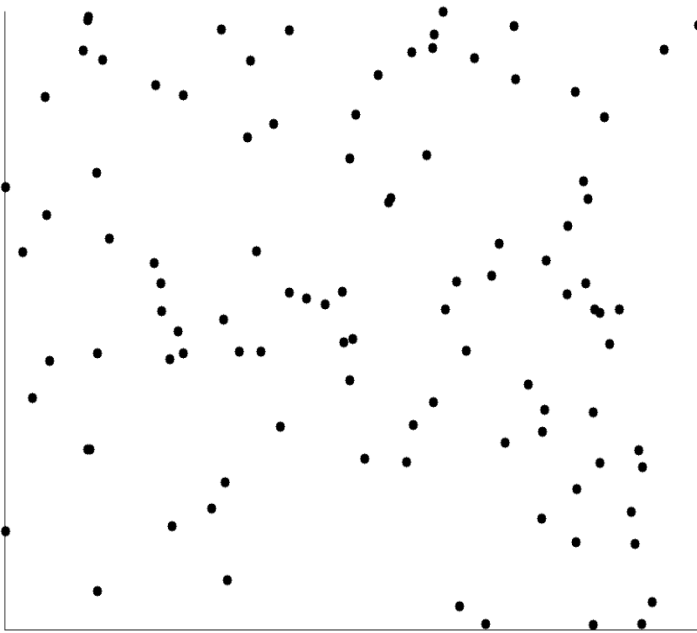
Correlation



The left half fragment is positively correlated

The right half is negative correlated

Uncorrelated Data



Similarity and Dissimilarity

Similarity

- Numerical measure of how alike two data objects are
- Value is higher when objects are more alike
- Often falls in the range $[0,1]$

Dissimilarity (e.g., distance)

- Numerical measure of how different two data objects are
- Lower when objects are more alike
- Minimum dissimilarity is often 0
- Upper limit varies

Proximity refers to a similarity or dissimilarity

Proximity for Nominal Attributes

Can take 2 or more states, e.g., red, yellow, blue, green
(generalization of a binary attribute)

Method 1: Simple matching

m : # of matches, p : total # of variables

$$d(i, j) = \frac{p - m}{p}$$

Method 2: Use a large number of binary attributes creating a new binary attribute for each of the M nominal states

Proximity for Nominal Attributes

- Example

Object Identifier	test-1 (nominal)	test-2 (ordinal)	test-3 (numeric)
1	code A	excellent	45
2	code B	fair	22
3	code C	good	64
4	code A	excellent	28

$$d(i, j) = \frac{p - m}{p}$$

- Dissimilarity Matrix (p=1)

$$\begin{bmatrix} 0 & & & \\ d(2, 1) & 0 & & \\ d(3, 1) & d(3, 2) & 0 & \\ d(4, 1) & d(4, 2) & d(4, 3) & 0 \end{bmatrix} \rightarrow \begin{bmatrix} 0 & & & \\ 1 & 0 & & \\ 1 & 1 & 0 & \\ 0 & 1 & 1 & 0 \end{bmatrix}$$

- $\text{sim}(i, j) = 1 - d(i, j) = m/p$

Proximity for Binary Attributes

- A contingency table for binary data
- Distance measure for symmetric binary variables:
- Distance measure for asymmetric binary variables:
- Jaccard coefficient (*similarity* measure for *asymmetric* binary variables):

		Object j		
		1	0	sum
Object i	1	q	r	$q+r$
	0	s	t	$s+t$
sum		$q+s$	$r+t$	p

$$d(i, j) = \frac{r + s}{q + r + s + t}$$

$$d(i, j) = \frac{r + s}{q + r + s}$$

$$sim_{Jaccard}(i, j) = \frac{q}{q + r + s}$$

Dissimilarity between Binary Variables

Example

Name	Gender	Fever	Cough	Test-1	Test-2	Test-3	Test-4
Jack	M	Y	N	P	N	N	N
Mary	F	Y	N	P	N	P	N
Jim	M	Y	P	N	N	N	N

Gender is a symmetric attribute

The remaining attributes are asymmetric binary

Let the values Y and P be 1, and the value N be 0

$$d(\text{jack}, \text{mary}) = \frac{0 + 1}{2 + 0 + 1} = \mathbf{0.33}$$

$$d(\text{jack}, \text{jim}) = \frac{1 + 1}{1 + 1 + 1} = \mathbf{0.67}$$

$$d(\text{jim}, \text{mary}) = \frac{1 + 2}{1 + 1 + 2} = \mathbf{0.75}$$

	1	0	sum
1	q	r	$q + r$
0	s	t	$s + t$
sum	$q + s$	$r + t$	p

$$d(i, j) = \frac{r + s}{q + r + s + t}$$

$$d(i, j) = \frac{r + s}{q + r + s}$$

Distance of Numeric Data

Minkowski distance: A popular distance measure

$$d(i, j) = \sqrt[h]{|x_{i1} - x_{j1}|^h + |x_{i2} - x_{j2}|^h + \dots + |x_{ip} - x_{jp}|^h}$$

where $i = (x_{i1}, x_{i2}, \dots, x_{ip})$ and $j = (x_{j1}, x_{j2}, \dots, x_{jp})$ are two p -dimensional data objects, and h is the order (the distance so defined is also called L- h norm)

Properties

$d(i, j) > 0$ if $i \neq j$, and $d(i, i) = 0$ (Positive definiteness)

$d(i, j) = d(j, i)$ (Symmetry)

$d(i, j) \leq d(i, k) + d(k, j)$ (Triangle Inequality)

A distance that satisfies these properties is a metric

Special Cases

$h = 1$: Manhattan (city block, L_1 norm) distance

E.g., the Hamming distance: the number of bits that are different between two binary vectors

$$d(i, j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{ip} - x_{jp}|$$

$h = 2$: (L_2 norm) Euclidean distance

$$d(i, j) = \sqrt{(|x_{i1} - x_{j1}|^2 + |x_{i2} - x_{j2}|^2 + \dots + |x_{ip} - x_{jp}|^2)}$$

$h \rightarrow \infty$. “supremum” (L_{\max} norm, L_∞ norm) distance.

This is the maximum difference between any component (attribute) of the vectors

$$d(i, j) = \lim_{h \rightarrow \infty} \left(\sum_{f=1}^p |x_{if} - x_{jf}|^h \right)^{\frac{1}{h}} = \max_f |x_{if} - x_{jf}|$$

$$d(i, j) = \sqrt[h]{|x_{i1} - x_{j1}|^h + |x_{i2} - x_{j2}|^h + \dots + |x_{ip} - x_{jp}|^h}$$

Examples

Manhattan (L_1)

L	x1	x2	x3	x4
x1	0			
x2	5	0		
x3	3	6	0	
x4	6	1	7	0

Euclidean (L_2)

L2	x1	x2	x3	x4
x1	0			
x2	3.61	0		
x3	2.24	5.1	0	
x4	4.24	1	5.39	0

Supremum

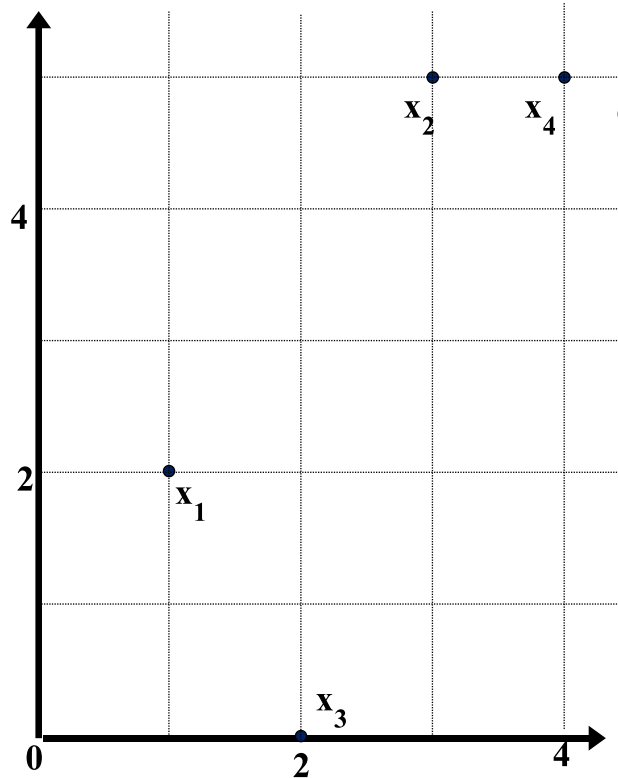
L_∞	x1	x2	x3	x4
x1	0			
x2	3	0		
x3	2	5	0	
x4	3	1	5	0

$$d(i, j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{ip} - x_{jp}|$$

$$d(i, j) = \sqrt{(|x_{i1} - x_{j1}|^2 + |x_{i2} - x_{j2}|^2 + \dots + |x_{ip} - x_{jp}|^2)}$$

$$d(i, j) = \lim_{h \rightarrow \infty} \left(\sum_{f=1}^p |x_{if} - x_{jf}|^h \right)^{\frac{1}{h}} = \max_f |x_{if} - x_{jf}|$$

point	attribute 1	attribute 2
x1	1	2
x2	3	5
x3	2	0
x4	4	5



Cosine Similarity

A **document** can be represented by thousands of attributes, each recording the *frequency* of a particular word (such as keywords) or phrase in the document.

<i>Document</i>	<i>teamcoach</i>	<i>hockey</i>	<i>baseball</i>	<i>soccer</i>	<i>penalty</i>	<i>score</i>	<i>win</i>	<i>loss</i>	<i>season</i>	
Document1	5	0	3	0	2	0	0	2	0	0
Document2	3	0	2	0	1	1	0	1	0	1
Document3	0	7	0	2	1	0	0	3	0	0
Document4	0	1	0	0	1	2	2	0	3	0

Other vector objects: gene features in micro-arrays, ...

Applications: information retrieval, biologic taxonomy, gene feature mapping, ...

Cosine measure: If d_1 and d_2 are two vectors (e.g., term-frequency vectors), then

$$\cos(d_1, d_2) = (d_1 \bullet d_2) / (||d_1|| ||d_2||),$$

where \bullet indicates vector dot product, $||d||$: the length of vector d

$$\text{sim}(d_1, d_2) = \frac{\vec{V}(d_1) \cdot \vec{V}(d_2)}{|\vec{V}(d_1)||\vec{V}(d_2)|} = \frac{\sum_{i=1}^M x_i y_i}{\sqrt{\sum_{i=1}^M \vec{V}_i^2(d)}}$$

Example

$$\cos(d_1, d_2) = (d_1 \bullet d_2) / (||d_1|| ||d_2||),$$

where \bullet indicates vector dot product, $||d||$: the length of vector d

Ex: Find the **similarity** between documents 1 and 2.

$$d_1 = (5, 0, 3, 0, 2, 0, 0, 2, 0, 0)$$

$$d_2 = (3, 0, 2, 0, 1, 1, 0, 1, 0, 1)$$

$$d_1 \bullet d_2 = 5*3 + 0*0 + 3*2 + 0*0 + 2*1 + 0*1 + 0*1 + 2*1 + 0*0 + 0*1 = 25$$

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

$$||d_2|| = (3*3 + 0*0 + 2*2 + 0*0 + 1*1 + 1*1 + 0*0 + 1*1 + 0*0 + 1*1)^{0.5} = (17)^{0.5} = 4.12$$

$$\cos(d_1, d_2) = 0.94$$

<i>Document</i>	<i>team</i>	<i>coach</i>	<i>hockey</i>	<i>baseball</i>	<i>soccer</i>	<i>penalty</i>	<i>score</i>	<i>win</i>	<i>loss</i>	<i>season</i>
Document1	5	0	3	0	2	0	0	2	0	0
Document2	3	0	2	0	1	1	0	1	0	1
Document3	0	7	0	2	1	0	0	3	0	0
Document4	0	1	0	0	1	2	2	0	3	0

Probability Distributions

Random Variable

- A random variable x takes on a defined set of values with different probabilities.
 - For example, if you roll a die, the outcome is random (not fixed) and there are 6 possible outcomes, each of which occur with probability one-sixth.
 - For example, if you poll people about their voting preferences, the percentage of the sample that responds “Yes on Proposition 100” is also a random variable (the percentage will be slightly differently every time you poll).
- Roughly, probability is how frequently we expect different outcomes to occur if we repeat the experiment over and over (“frequentist” view)

Random Variables

Discrete random variables have a countable number of outcomes

Examples: Dead/alive, treatment/placebo, dice, counts, etc.

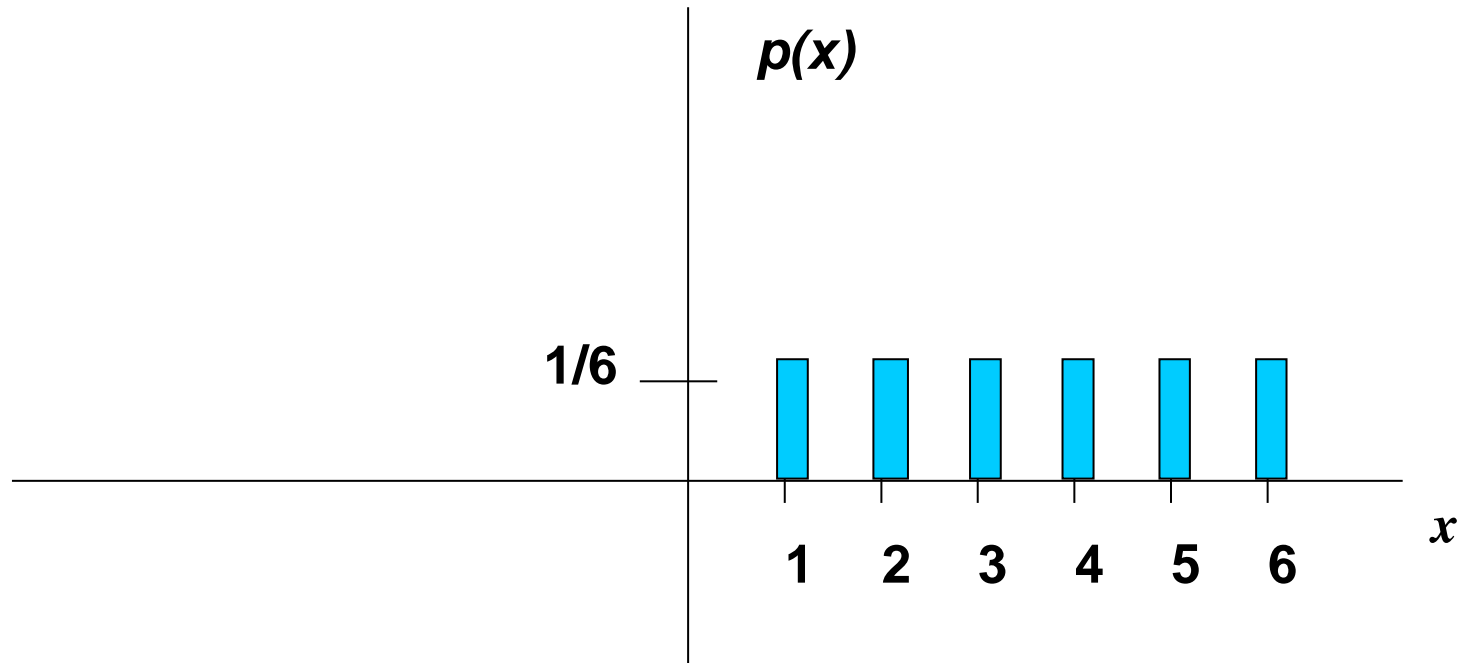
Continuous random variables have an infinite continuum of possible values.

Examples: blood pressure, weight, the speed of a car, the real numbers from 1 to 6.

Probability Functions

- A probability function maps the possible values of x against their respective probabilities of occurrence, $p(x)$
- $p(x)$ is a number from 0 to 1.0.
- The area under a probability function is always 1.

Example: roll a die



$$\sum_{\text{all } x} P(x) = 1$$

Cumulative Distribution Function

$$F_X(x) = \sum_{x_k \leq x} P_X(x_k)$$

x	$P(x \leq A)$
1	$P(x \leq 1) = 1/6$
2	$P(x \leq 2) = 2/6$
3	$P(x \leq 3) = 3/6$
4	$P(x \leq 4) = 4/6$
5	$P(x \leq 5) = 5/6$
6	$P(x \leq 6) = 6/6$

Important Discrete Distributions

Binomial

n draws of a Bernoulli distribution

Random variable X stands for the number of times that experiments are successful.

$$\Pr(X = x) = p_{\theta}(x) = \begin{cases} \binom{n}{x} p^x (1-p)^{n-x} & x = 1, 2, \dots, n \\ 0 & \text{otherwise} \end{cases}$$

$E[X] = np, \text{Var}(X) = np(1-p)$

Important Discrete Distributions

Poisson

Coming from Binomial distribution

Fix the expectation $\lambda=np$

Let the number of trials $n \rightarrow \infty$

A Binomial distribution will become a Poisson distribution

$$\Pr(X = x) = p_{\theta}(x) = \begin{cases} \frac{\lambda^x}{x!} e^{-\lambda} & x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

$$E[X] = \lambda, \text{Var}(X) = \lambda$$

Continuous Variables

- The probability function that accompanies a continuous random variable is a continuous mathematical function that integrates to 1.
- The probabilities associated with continuous functions are just areas under the curve (integrals!).
- Probabilities are given for a range of values, rather than a particular value (e.g., the probability of getting a math SAT score between 700 and 800 is 2%).

Example

- For example, recall the negative exponential function (in probability, this is called an “exponential distribution”):

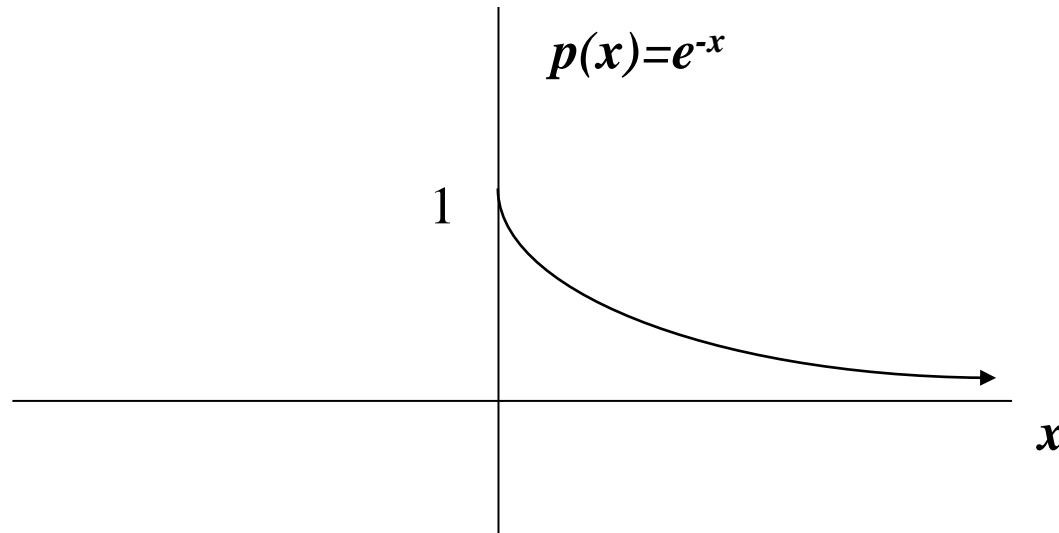
$$f(x) = e^{-x}$$

- This function integrates to 1:

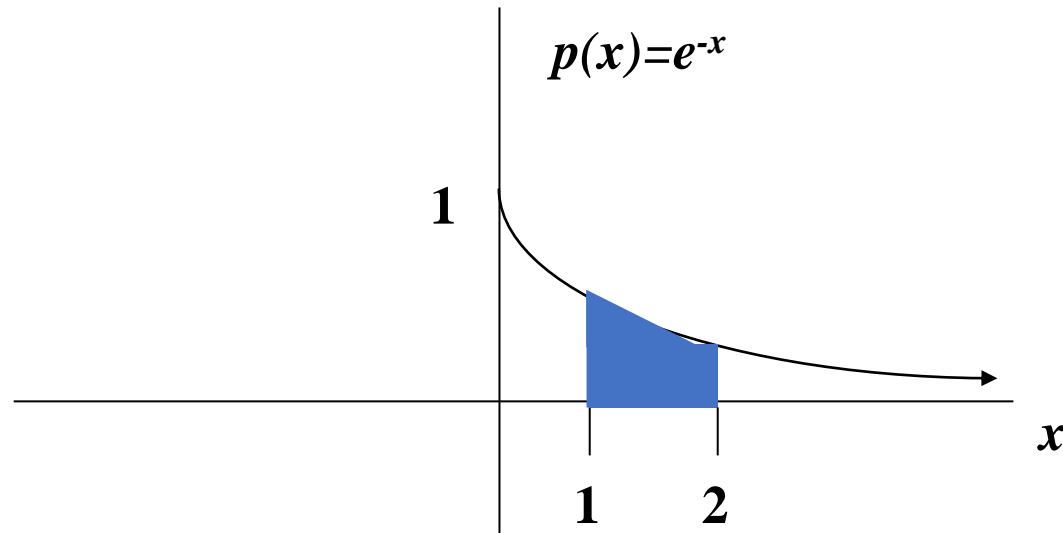
$$\int_0^{+\infty} e^{-x} = -e^{-x} \Big|_0^{+\infty} = 0 + 1 = 1$$

Probability Density Function

The probability that x is any exact particular value (such as 1.9976) is 0; we can only assign probabilities to possible ranges of x .



Example



$$P(1 \leq x \leq 2) = \int_1^2 e^{-x} = -e^{-x} \Big|_1^2 = -e^{-2} - (-e^{-1}) = -.135 + .368 = .23$$

Cumulative Distribution Function

As in the discrete case, we can specify the “cumulative distribution function” (CDF):

The CDF here = $P(x \leq A) =$

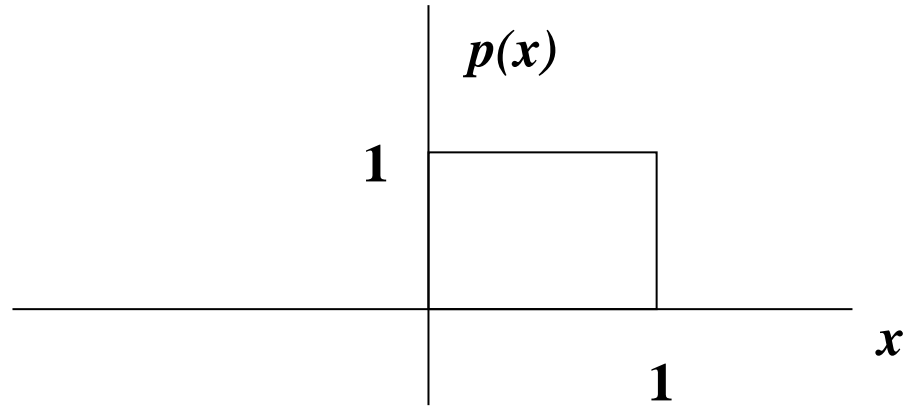
$$\int_0^A e^{-x} = -e^{-x} \Big|_0^A = -e^{-A} - (-e^0) = -e^{-A} + 1 = 1 - e^{-A}$$

Uniform Distribution

The uniform distribution: all values are equally likely

The uniform distribution:

$$f(x) = 1, \text{ for } 1 \geq x \geq 0$$



We can see it's a probability distribution because it integrates to 1 (the area under the curve is 1):

$$\int_0^1 1 = x \Big|_0^1 = 1 - 0 = 1$$

Normal Distribution

$$X \sim N(\mu, \sigma)$$

$$p_{\theta}(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$$

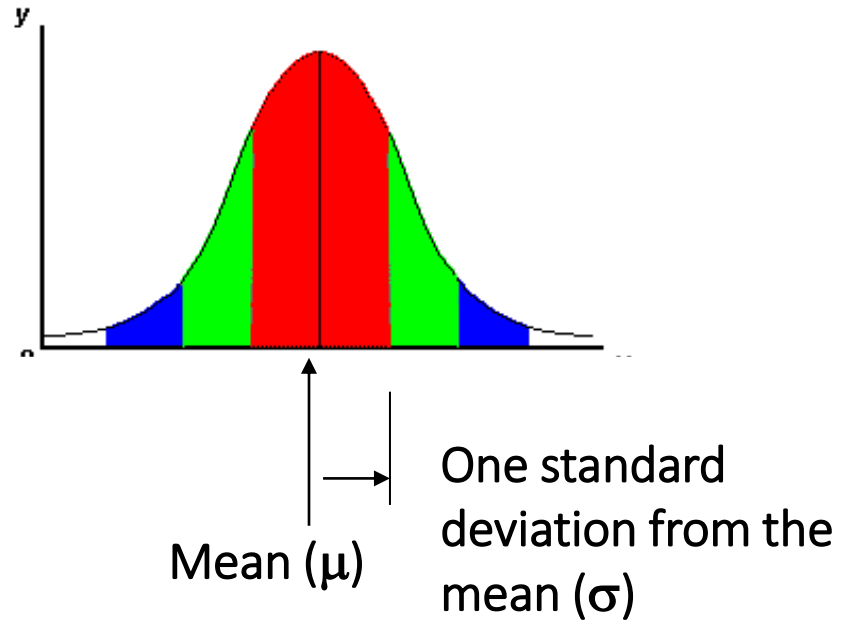
$$\Pr(a \leq X \leq b) = \int_a^b p_{\theta}(x) dx = \int_a^b \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\} dx$$

$$E[X] = \mu, \text{Var}(X) = \sigma^2$$

Expected Value and Variance

All probability distributions are characterized by an expected value and a variance (standard deviation squared).

Example: Normal Distribution



Expected Value

If we understand the underlying probability function of a certain phenomenon, then we can make informed decisions based on how **we expect x to behave on-average** over the long-run...(so called “frequentist” theory of probability).

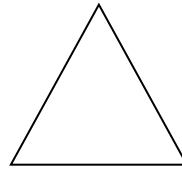
Expected value is just the weighted average or mean (μ) of random variable x. Imagine placing the masses $p(x)$ at the points X on a beam; the balance point of the beam is the expected value of x .

$$E[X] = \sum_{i=1}^k x_i p_i = x_1 p_1 + x_2 p_2 + \cdots + x_k p_k$$

$$E[X] = \int_{\mathbb{R}} x f(x) dx$$

Example

x	10	11	12	13	14
$P(x)$.4	.2	.2	.1	.1



$$\sum_{i=1}^5 x_i p(x) = 10(.4) + 11(.2) + 12(.2) + 13(.1) + 14(.1) = 11.3$$

Operators

If c = a constant number (i.e., not a variable) and X and Y are any random variables...

$$E(c) = c$$

$$E(cX) = cE(X)$$

$$E(c + X) = c + E(X)$$

$$E(X+Y) = E(X) + E(Y)$$

Variance/Deviation

“The average (expected) squared distance (or deviation) from the mean”

$$\sigma^2 = \text{Var}(x) = E[(x - \mu)^2] = \sum_{\text{all } x} (x_i - \mu)^2 p(x_i)$$

Variance

Discrete case:

$$\text{Var}(X) = \sigma^2 = \sum_{\text{all } x} (x_i - \mu)^2 p(x_i)$$

Continuous case:

$$\text{Var}(X) = \sigma^2 = \int_{-\infty}^{\infty} (x_i - \mu)^2 p(x_i) dx$$

Sample Variance

The variance of a sample: $s^2 =$

$$\frac{\sum_{i=1}^N (x_i - \bar{x})^2}{n - 1} = \sum_{i=1}^N (x_i - \bar{x})^2 \left(\frac{1}{n - 1}\right)$$

Conditional Probability

If A and B are events with $\Pr(A) > 0$, the *conditional probability of B given A* is

$$\Pr(B | A) = \frac{\Pr(AB)}{\Pr(A)}$$

Example: Drug test

	Women	Men
Success	200	1800
Failure	1800	200

	Women		Men	
	Drug I	Drug II	Drug I	Drug II
Success	200	10	19	1000
Failure	1800	190	1	1000

Bayes' Rule

Given two events A and B and suppose that $\Pr(A) > 0$. Then

$$\Pr(B | A) = \frac{\Pr(AB)}{\Pr(A)} = \frac{\Pr(A | B)\Pr(B)}{\Pr(A)}$$

Example:

$$\Pr(R) = 0.8$$

$\Pr(W R)$	R	$\neg R$
W	0.7	0.4
$\neg W$	0.3	0.6

R: It is a rainy day

W: The grass is wet

$$\Pr(R | W) = ?$$