# *Intermediate Lab*
## PHYS 3870

## Lecture 3

# Distribution Functions

References:  Taylor Ch.  5 (and Chs. 10 and 11 for Reference)
Taylor Ch.  6 and 7
Also refer to "Glossary of Important Terms in Error Analysis"
"Probability Cheat Sheet"

UtahState
UNIVERSITY

# *Intermediate Lab*
## PHYS 3870

# Distribution Functions

# Practical Methods to Calculate Mean and St. Deviation

**We need to develop a good way to tally, display, and think about a collection of repeated measurements of the same quantity.**
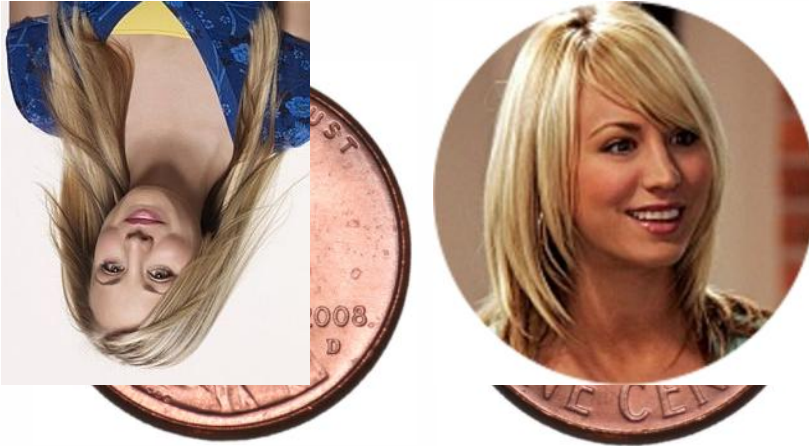
**Here is where we are headed:**
• **Develop the notion of a probability distribution function, a distribution to describe the probable outcomes of a measurement**
• **Define what a distribution function is, and its properties**
• **Look at the properties of the most common distribution function, the Gaussian distribution for purely random events**
• **Introduce other probability distribution functions**

**We will develop the mathematical basis for:**
• **Mean**
• **Standard deviation**
• **Standard deviation of the mean (SDOM)**
• **Moments and expectation values**
• **Error propagation formulas**
• **Addition of errors in quadrature (for independent and random measurements)**
• **Schwartz inequality (i.e., the uncertainty principle)  (next lecture)**
• **Numerical values for confidence limits (t-test)**
• **Principle of maximal likelihood**
• **Central limit theorem**

# Two Practical Exercises in Probabilities

Flip a penny 50 times and record the results



Roll a pair of dice 50 times and record the results



Grab a partner and a set of instructions and complete the exercise.

Intermediate 3870
Fall 2013

DISTRIBUTION FUNCTIONS

Lecture 3  Slide 4

UtahState
UNIVERSITY

# Two Practical Exercises in Probabilities

**Flip a penny 50 times and record the results**

**Group Two Instructions**
1. Flip penny 50 times
2. Record each results as "H" or "T" in list below

— — — — — — — — — —
— — — — — — — — — —
— — — — — — — — — —
— — — — — — — — — —
— — — — — — — — — —

**Group One Instructions**
1. Flip penny 50 times
2. Tally results on list below

Heads:

Tails:

**What is the asymmetry of the results?**

UtahState
UNIVERSITY

# Two Practical Exercises in Probabilities

**Flip a penny 50 times and record the results**

**Group Two Instructions**
1. Roll two dice 50 times
2. Record each results as "H" or "T" in list below

```
_H_ _H_ _T_ _H_ _T_ _T_ _H_ _H_ _T_ _T_
_H_ _H_ _H_ _T_ _T_ _H_ _H_ _T_ _H_ _H_
_H_ _T_ _H_ _T_ _H_ _T_ _H_ _H_ _T_ _T_
_H_ _H_ _T_ _T_ _H_ _H_ _T_ _H_ _T_ _H_
_H_ _T_ _H_ _H_ _T_ _T_ _H_ _T_ _T_ _T_
```

**Group One Instructions**
1. Flip penny 50 times
2. Tally results on list below

Heads: 54%

Tails: 46%

**??% asymmetry**                    **4% asymmetry**

**What is the asymmetry of the results?**

# Two Practical Exercises in Probabilities

**Roll a pair of dice 50 times and record the results**
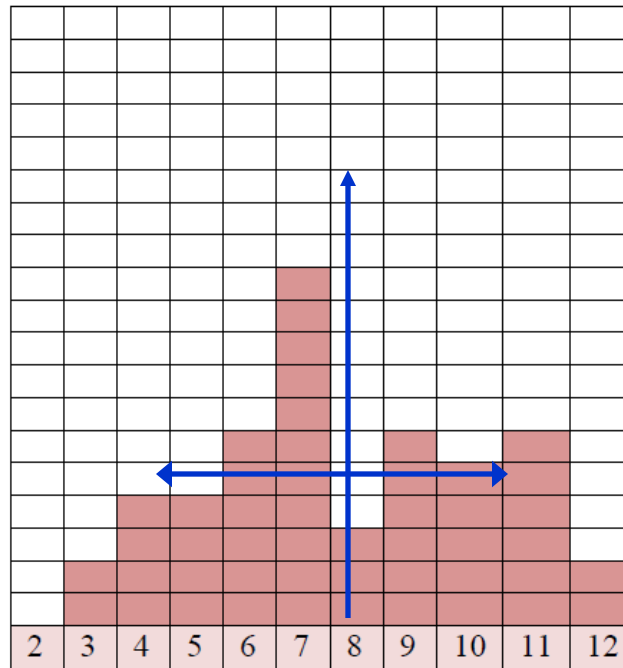
### Group One Instructions

Roll two dice 50 times

Record results on table below, checking one box for each die

| 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|----|----|----|
|   |   |   |   |   |   |   |   |    |    |    |

### Group Two Instructions

Roll two dice 50 times

Record each result in list below

— — — — — — — — — —
— — — — — — — — — —
— — — — — — — — — —
— — — — — — — — — —
— — — — — — — — — —

**What is the mean value?
The standard deviation?**

# Two Practical Exercises in Probabilities

**Roll a pair of dice 50 times and record the results**

## Group Two Instructions

Roll two dice 50 times

Record each result in list below

_3_  _4_  _10_  _6_  _9_  _7_  _9_  _3_  _11_  _10_
_6_  _6_  _9_  _7_  _11_  _12_  _6_  _11_  _6_  _7_
_7_  _7_  _8_  _4_  _11_  _10_  _12_  _11_  _7_  _4_
_12_  _8_  _9_  _9_  _7_  _10_  _9_  _11_  _4_  _7_
_7_  _6_  _10_  _5_  _8_  _7_  _5_  _5_  _5_  _7_

**What is the mean value?**

**The standard deviation?**

**What is the asymmetry (kurtosis)?**

**What is the probability of rolling a 4?**

## Group One Instructions

Roll two dice 50 times

Record results on table below, checking one box for each die



Mean = 7.3
St. Dev. = 2.8

UtahState
UNIVERSITY

# Discrete Distribution Functions

## A data set to play with

26, 24, 26, 28, 23, 24, 25, 24, 26, 25.     (5.1)

23, 24, 24, 24, 25, 25, 26, 26, 26, 28.     (5.2)

## Written in terms of "occurrence" F



**Table 5.1.** Measured lengths $x$ and their numbers of occurrences.

| Different values, $x_k$ | 23 | 24 | 25 | 26 | 27 | 28 |
|---|---|---|---|---|---|---|
| Number of times found, $n_k$ | 1 | 3 | 2 | 3 | 0 | 1 |

## The mean value

$$\bar{x} = \frac{\sum_i x_i}{N} = \frac{23 + 24 + 24 + 24 + 25 + \ldots + 28}{10}.$$

This equation is the same as

$$\bar{x} = \frac{23 + (24 \times 3) + (25 \times 2) + \ldots + 28}{10}$$

or in general

$$\sum_k (n_k) = N \quad \rightarrow \quad X = \frac{\sum_k (n_k \cdot x_k)}{N} = \frac{\sum_k (n_k \cdot x_k)}{\sum_k (n_k)}$$

## In terms of fractional expectations

**Fractional expectations**     $F_k = \dfrac{n_k}{N}$
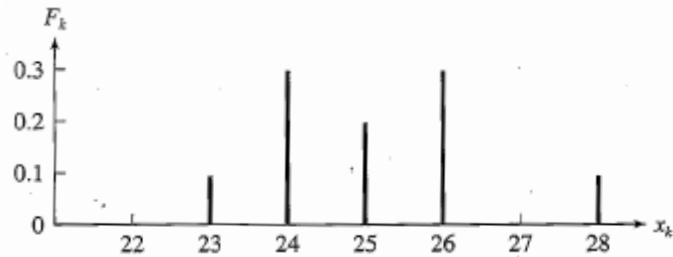
**Normalization condition**     $1 = \sum_k (F_k)$

**Mean value**     $X = \sum_k (F_k \cdot x_k)$

**(This is just a weighted sum.)**

# Limit of Discrete Distribution Functions



**Binned data sets**

**Table 5.2.** The 10 measurements (5.9) grouped in bins.

| Bin | 22 to 23 | 23 to 24 | 24 to 25 | 25 to 26 | 26 to 27 | 27 to 28 |
|---|---|---|---|---|---|---|
| Observations in bin | 1 | 3 | 1 | 4 | 1 | 0 |

26.4, 23.9, 25.1, 24.6, 22.7, 23.8, 25.1, 23.9, 25.3, 25.4.     (5.9)

**"Normalizing" data sets**   $f_k \Delta_k$ = fraction of measurements in $k$th bin.

$f_k$ ≡ fractional occurrence

$\Delta_k$ ≡ bin width

**Mean value:**   $X = \sum_k F_k \, x_k = \sum_k (f_k \, \Delta_k) x_k$
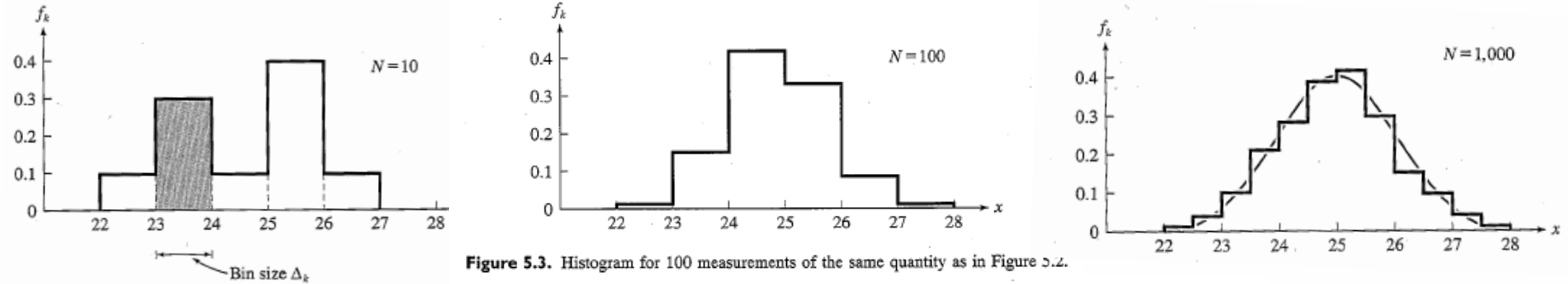
**Normalization:**   $1 = \dfrac{\sum_k (f_k \, \Delta_k) x_k}{\sum_k \Delta_k}$

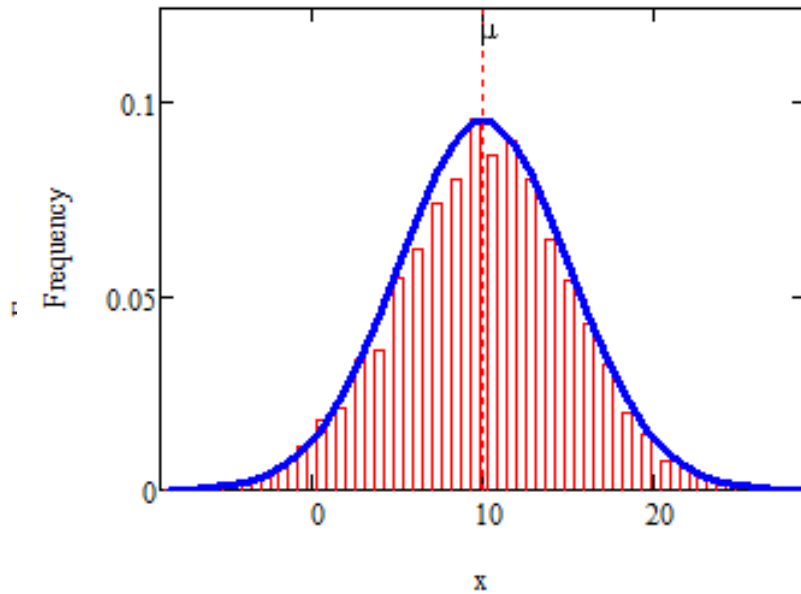**Expected value:**   $Prob(4) = \dfrac{F_4}{N} = f_4$

# Limits of Distribution Functions

**Consider the limiting distribution function as N $\rightarrow \infty$ and $\Delta_k \rightarrow 0$**

**Larger data sets**



Figure 5.3. Histogram for 100 measurements of the same quantity as in Figure 5.2.

**Frequency Distribution**



**# Data Pts:** $N \equiv 6000$

**Mean:** $\mu \equiv 10$

**Std. Dev.:** $\sigma \equiv 5$

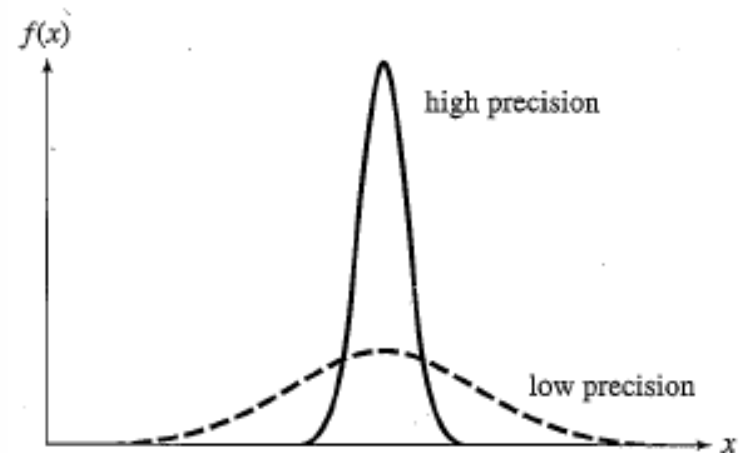**SDOM:** $\dfrac{\sigma}{\sqrt{N}} = 0.065$

**Mathcad Games:**

**Fractional Error:**

$$\left(\frac{\sigma}{\sqrt{N}}\right) \cdot \frac{1}{\mu} = 0.645\,\%$$

Intermediate 3870
Fall 2013

DISTRIBUTION FUNCTIONS

Lecture 3 Slide 11

UtahState
UNIVERSITY

# Continuous Distribution Functions

**Meaning of Distribution Interval**

$f(x)\ dx$ = fraction of measurements that fall between $x$ and $x + dx$.  (5.10)

(a)

(b)



$f(x)$

high precision

low precision

$x$

**Thus**

$$\int_a^b f(x)\ dx$$ = fraction of measurements that fall within a<x<b

**Normalization of Distribution**

$$\sum_k \left( F_k \right) = 1 \quad \rightarrow \quad \int_{-\infty}^{+\infty} f(x)dx = 1$$

**and by extension** $$\int_{-\infty}^{\infty} f(x)\ dx = 1$$

**Central (mean) Value**

$$\sum_k \left( F_k \cdot x_k \right) = X \quad \rightarrow$$

$$\int_{-\infty}^{+\infty} x\ f(x)\ dx = \bar{x} = \langle x \rangle \qquad (5.15)$$

**Width of Distribution**

$$\sum_k \left[ \left( x_k - x \right)^2 \cdot F_k \right] = \sigma_X^2 \quad \rightarrow$$

$$\int_{-\infty}^{+\infty} (x - \bar{x})^2\ f(x)\ dx = \sigma_x^2 = \langle (x - \bar{x})^2 \rangle \quad (5.16)$$

# Moments of Distribution Functions

The first moment for a **probability distribution function** is

$$\bar{x} \equiv \langle x \rangle = first\ moment = \int_{-\infty}^{+\infty} x\, f(x)\, dx$$

For a **general distribution function**,

$$\bar{x} \equiv \langle x \rangle = first\ moment = \frac{\int_{-\infty}^{+\infty} x\, g(x) dx}{\int_{-\infty}^{+\infty} g(x) dx}$$

Generalizing, the n$^{th}$ moment is

$$x_n \equiv \langle x^n \rangle = nth\ moment = \frac{\int_{-\infty}^{+\infty} x^n\, g(x) dx}{\int_{-\infty}^{+\infty} g(x) dx} = \int_{-\infty}^{+\infty} x^n\, f(x)\, dx$$

**0 (for a centered distribution)**

O$^{th}$ moment ≡ N          2$^{nd}$ moment ≡ $\langle (x - \bar{x})^2 \rangle \rightarrow \langle x^2 \rangle$

1$^{st}$ moment ≡ $\bar{x}$          3$^{rd}$ moment ≡ kurtosis

UtahState
UNIVERSITY

# Moments of Distribution Functions

Generalizing, the n[th] moment is

$$x_n \equiv \langle x^n \rangle = nth\ moment = \frac{\int_{-\infty}^{+\infty} x^n\ g(x)dx}{\int_{-\infty}^{+\infty} g(x)dx} = \int_{-\infty}^{+\infty} x^n\ f(x)\ dx$$

O[th] moment $\equiv$ N  $\qquad$ 2[nd] moment $\equiv \langle (x - \bar{x})^2 \rangle \rightarrow \langle x^2 \rangle$  **0**

1[st] moment $\equiv \bar{x}$  $\qquad$ 3[rd] moment $\equiv$ kurtosis

The n[th] moment about the mean is

$$\mu_n \equiv \langle (x - \bar{x})^n \rangle = nth\ moment\ about\ the\ mean$$

$$= \frac{\int_{-\infty}^{+\infty} (x - \bar{x})^n\ g(x)dx}{\int_{-\infty}^{+\infty} g(x)dx} = \int_{-\infty}^{+\infty} (x - \bar{x})^n\ f(x)\ dx$$

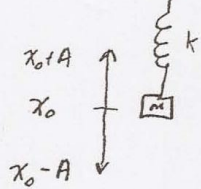The standard deviation (or second moment about the mean) is

$$\sigma_x^2 \equiv \mu_2 \equiv \langle (x - \bar{x})^2 \rangle = 2nd\ moment\ about\ the\ mean$$

$$= \frac{\int_{-\infty}^{+\infty} (x - \bar{x})^2\ g(x)dx}{\int_{-\infty}^{+\infty} g(x)dx} = \int_{-\infty}^{+\infty} (x - \bar{x})^2\ f(x)\ dx$$

## Harmonic Oscillator: Example from Mechanics

Consider a mass on a spring with frequency $\omega$ and equilibrium position $x_0$

$$\omega = \sqrt{k/m} = \frac{2\pi}{T}$$

The equations of motion are:

$$x(t) = A \sin \omega t + x_0$$
$$\dot{x}(t) = -A\omega \cos \omega t$$
$$\ddot{x}(t) = -A\omega^2 \sin \omega t$$

### Expected Values →

The expectation value of a function R(x) is

$$\langle R(x) \rangle \equiv \frac{\int_{-\infty}^{+\infty} R(x) g(x) dx}{\int_{-\infty}^{+\infty} g(x) dx}$$

$$= \int_{-\infty}^{+\infty} R(x) f(x) \, dx$$

The average time over one period is:

$$\langle t \rangle = \frac{\int_0^T t \, dt}{\int_0^T dt} = \frac{\frac{1}{2}t^2 |_0^T}{t |_0^T} = \frac{\frac{1}{2}T^2}{T} = \frac{1}{2}T = \frac{\pi}{\omega}$$

The average position over one period (and over all time) is:

$$\langle x \rangle = \frac{\int_0^T x(t) \, dt}{T} = \frac{\int_0^T [A \sin \omega t + x_0] \, dt}{T}$$
$$= \frac{[\frac{A}{\omega} \cos \omega t + x_0 t]_0^T}{T} = x_0$$

The average of the position squared is:

$$\langle x^2 \rangle = \frac{1}{T}\int_0^T [x(t)]^2 \, dt = \frac{1}{T}\int_0^T [A \sin \omega t + x_0]^2 \, dt$$
$$= \frac{1}{T}\left[ T\left(x_0^2 + \frac{A^2}{2}\right)\right] = x_0^2 - \frac{A^2}{2}$$

The average of the force is:

$$\langle F \rangle = \langle -k(x - x_0) \rangle = 0$$
$$\langle F \rangle = \langle m\ddot{x} \rangle = -m\omega^2 \langle x - x_0 \rangle = 0$$

The average velocity is:

$$\langle v \rangle = \frac{\int_0^T (-A\omega \cos \omega t) \, dt}{T} = \frac{-A\omega}{T}\left[ -\sin \omega t \right]_0^T = 0$$

The average kinetic energy is:

$$\langle KE \rangle = \frac{1}{2}m\langle v^2 \rangle = \frac{m}{2}A^2\omega^2 \int_0^T \cos^2 \omega t \, dt = \frac{mA^2\omega^2}{2}\left[ \frac{t}{2} + \frac{\sin 2\omega t}{4\omega} \right]_0^T$$
$$= \frac{mA^2\omega^2 T}{4} = \frac{\pi}{2}\sqrt{km}$$

Note: $\langle KE \rangle \neq \frac{1}{2}m\langle v \rangle^2 = 0$.

## Boltzmann Distribution: Example from Kinetic Theory

### Expected Values ➔

The expectation value of a function R(x) is

$$\langle R(x) \rangle \equiv \frac{\int_{-\infty}^{+\infty} R(x) g(x) dx}{\int_{-\infty}^{+\infty} g(x) dx}$$
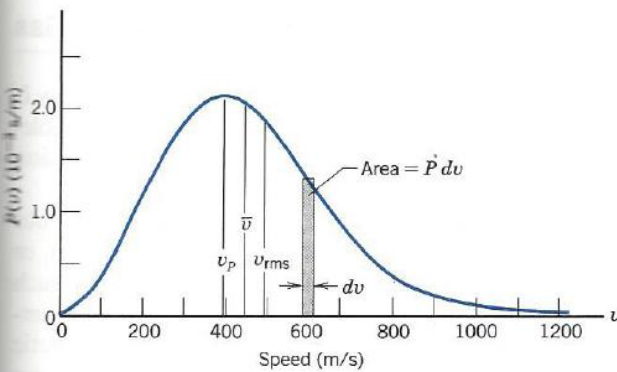
$$= \int_{-\infty}^{+\infty} R(x) f(x) dx$$

The Boltzmann distribution function for velocities of particles as a function of temperature, T is:

$$P(v; T) = 4\pi \left( \frac{M}{2\pi\, k_B T} \right)^{3/2} v^2 exp\left[ \frac{1}{2}Mv^2 \Big/ \frac{1}{2} k_B T \right]$$

Then

$$\langle v \rangle = \int_{-\infty}^{+\infty} v\, P(v)\, dv = \left[ 8\, k_B T \Big/ \pi\, M \right]^{1/2}$$

$$\langle v^2 \rangle = \int_{-\infty}^{+\infty} v^2\, P(v)\, dv = \left[ 3\, k_B T \Big/ M \right]^{1/2}$$

$$\text{implies } \langle KE \rangle = \frac{1}{2}M\langle v^2 \rangle = \frac{3}{2} k_B T$$

$$v_{peak=} \sqrt{\left[ 2\, k_B T \Big/ M \right]^{1/2}} = \left[ 2 \Big/ 3 \right]^{1/2} \langle v^2 \rangle$$



Intermediate 3870
Fall 2013

DISTRIBUTION FUNCTIONS

Lecture 3   Slide 16

UtahState
UNIVERSITY

## Fermi-Dirac Distribution: Example from Kinetic Theory

For a system of identical fermions, the average number of fermions in a single-particle state $i$, is given by the Fermi–Dirac (F–D) distribution,

$$\bar{n}_i = \frac{1}{e^{(\epsilon_i - \mu)/kT} + 1}$$

where $k_B$ is Boltzmann's constant, $T$ is the absolute temperature, $\epsilon_i$ is the energy of the single-particle state $i$, and $\mu$ is the total chemical potential.

Since the F–D distribution was derived using the Pauli exclusion principle, which allows at most one electron to occupy each possible state, a result is that $0 < \bar{n}_i < 1$

When a quasi-continuum of energies $\epsilon$ has an associated density of states $g(\epsilon)$ (i.e. the number of states per unit energy range per unit volume) the average number of fermions per unit energy range per unit volume is,
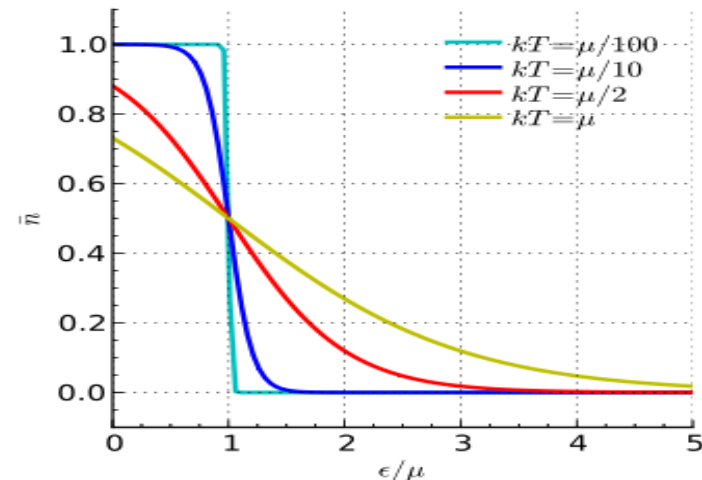
$$\bar{\mathcal{N}}(\epsilon) = g(\epsilon)\, F(\epsilon)$$

where $F(\epsilon)$ is called the Fermi function

$$F(\epsilon) = \frac{1}{e^{(\epsilon - \mu)/kT} + 1}$$

so that,

$$\bar{\mathcal{N}}(\epsilon) = \frac{g(\epsilon)}{e^{(\epsilon - \mu)/kT} + 1}$$



Intermediate 3870
Fall 2013

DISTRIBUTION FUNCTIONS

Lecture 3 Slide 17

UtahState
UNIVERSITY

# Example of Continuous Distribution Functions and Expectation Values
## Finite Square Well: Example from Quantum Mechanics

**Expectation Values**    The expectation value of a **QM operator** O(x) is   $\langle O(x) \rangle \equiv \dfrac{\int_{-\infty}^{+\infty} \Psi^*(x) O(x) \Psi(x) dx}{\int_{-\infty}^{+\infty} \Psi^*(x) \Psi(x) dx}$

For a finite square well of width L,   $\Psi_n(x) = \sqrt{2/L} \, \sin\left[\dfrac{n \pi x}{L}\right]$

$$\langle \Psi_n^*(x) | \Psi_n(x) \rangle \equiv \frac{\int_{-\infty}^{+\infty} \Psi_n^*(x) O(x) \Psi_n(x) dx}{\int_{-\infty}^{+\infty} \Psi_n^*(x) \Psi_n(x) \, dx} = 1$$

$$\langle x \rangle = \langle \Psi_n^*(x) | x | \Psi_n(x) \rangle \equiv \frac{\int_{-\infty}^{+\infty} \Psi_n^*(x) x \, \Psi_n(x) dx}{\int_{-\infty}^{+\infty} \Psi_n^*(x) \Psi_n(x) \, dx} = L/2$$

$$\langle p \rangle = \langle \Psi_n^*(x) | \frac{\hbar}{i} \frac{\partial}{\partial x} | \Psi_n(x) \rangle \equiv \frac{\int_{-\infty}^{+\infty} \Psi_n^*(x) \frac{\hbar}{i} \frac{\partial}{\partial x} \Psi_n(x) dx}{\int_{-\infty}^{+\infty} \Psi_n^*(x) \Psi_n(x) \, dx} = 0$$

$$\langle E_n \rangle = \langle \Psi_n^*(x) | i\hbar \frac{\partial}{\partial t} | \Psi_n(x) \rangle \equiv \frac{\int_{-\infty}^{+\infty} \Psi_n^*(x) i\hbar \frac{\partial}{\partial t} \Psi_n(x) dx}{\int_{-\infty}^{+\infty} \Psi_n^*(x) \Psi_n(x) \, dx} = \frac{n^2 \pi^2 \hbar^2}{2 \, m \, L^2}$$

**Probability Function (Discrete Case)**

The random variable X will be called a discrete random variable if there exists a function $f$ such that $f(x_i) \geq 0$ and $\sum_i f(x_i) = 1$ for $i = 1, 2, 3, \ldots$ and such that for any event $E$,

$$P(E) = P[X \text{ is in } E] = \sum_E f(x)$$

where $\sum_E$ means sum $f(x)$ over those values $x_i$ that are in $E$ and where $f(x) = P[X = x]$. The probability that the value of X is some real number $x$, is given by $f(x) = P[X = x]$, where $f$ is called the probability function of the random variable X.

**Cumulative Distribution Function (Discrete Case)**

The probability that the value of a random variable X is less than or equal to some real number $x$ is defined as

$$F(x) = P(X \leq x)$$
$$= \Sigma f(x_i), \qquad -\infty < x < \infty,$$

where the summation extends over those values of $i$ such that $x_i \leq x$.

**Probability Density (Continuous Case)**

The random variable X will be called a continuous random variable if there exists a function $f$ such that $f(x) \geq 0$ and $\int_{-\infty}^{\infty} f(x)\, dx = 1$ for all $x$ in interval $-\infty < x < \infty$ and such that for any event $E$

$$P(E) = P(X \text{ is in } E) = \int_E f(x)\, dx.$$

$f(x)$ is called the probability density of the random variable X. The probability that X assumes any given value of $x$ is equal to zero and the probability that it assumes a value on the interval from $a$ to $b$, including or excluding either end point, is equal to

$$\int_a^b f(x)\, dx.$$

**Cumulative Distribution Function (Continuous Case)**

The probability that the value of a random variable X is less than or equal to some real number $x$ is defined as

$$F(x) = P(X \leq x), \qquad -\infty < x < \infty$$
$$= \int_{-\infty}^{x} f(x)\, dx.$$

From the cumulative distribution, the density, if it exists, can be found from

$$f(x) = \frac{dF(x)}{dx}.$$

From the cumulative distribution

$$P(a \leq X \leq b) = P(X \leq b) - P(X \leq a)$$
$$= F(b) - F(a)$$

**Mathematical Expectation**

A. EXPECTED VALUE.

Let X be a random variable with density $f(x)$. Then the expected value of X, $E(X)$, is defined to be

$$E(X) = \sum_x x f(x)$$

if X is discrete and

$$E(X) = \int_{-\infty}^{\infty} x f(x)\, dx$$

if X is continuous. The expected value of a function $g$ of a random variable X is defined as

$$E[g(X)] = \sum_x g(x) \cdot f(x)$$

if X is discrete and

$$E[g(X)] = \int_{-\infty}^{\infty} g(x) \cdot f(x)\, dx$$

if X is continuous.

*Theorems*

1. $E[aX + bY] = aE(X) + bE(Y)$
2. $E[X \cdot Y] = E(X) \cdot E(Y)$ if X and Y are statistically independent.

B. MOMENTS

a. *Moments About the Origin.* The moments about the origin of a probability distribution are the expected values of the random variable which has the given distribution. The $r$th moment of X, usually denoted by $\mu'_r$, is defined as

$$\mu'_r = E[X^r] = \sum_x x^r f(x)$$

if X is discrete and

$$\mu'_r = E[X^r] = \int_{-\infty}^{\infty} x^r f(x)\, dx$$

if X is continuous.

The first moment, $\mu'_1$, is called the mean of the random variable X and is usually denoted by $\mu$.

b. *Moments About the Mean.* The $r$th moment about the mean, usually denoted by $\mu_r$, is defined as

$$\mu_r = E[(X - \mu)^r] = \sum_x (x - \mu)^r f(x)$$

if X is discrete and

$$\mu_r = E[(X - \mu)^r] = \int_{-\infty}^{\infty} (x - \mu)^r f(x)\, dx$$

if X is continuous.

The second moment about the mean, $\mu_2$, is given by

$$\mu_2 = E[(X - \mu)^2] = \mu'_2 - \mu^2$$

and is called the variance of the random variable X, and is denoted by $\sigma^2$. The square root of the variance, $\sigma$, is called the standard deviation.

*Theorems*

1. $\sigma^2_{cX} = c^2\sigma^2_X$
2. $\sigma^2_{c+X} = \sigma^2_X$
3. $\sigma^2_{aX+b} = a^2\sigma^2_X$

**Available on web site**

# *Intermediate Lab*
## PHYS 3870

# The Gaussian Distribution Function

References:  Taylor Ch.  5

Intermediate  3870
Fall 2013

DISTRIBUTION FUNCTIONS

Lecture  3   Slide  20

UtahState
UNIVERSITY

# Gaussian Integrals

## Factorial Approximations

$$n! \approx (2\pi n)^{1/2}\, n^n\, exp\left[-n + \frac{1}{12\,n} + O\left(\frac{1}{n^2}\right)\right]$$

$$log(n!) \approx \frac{1}{2} log(2\pi) + \left(n + \frac{1}{2}\right) log(n) - n + \frac{1}{12\,n} + O\left(\frac{1}{n^2}\right)$$

$$log(n!) \approx n\, log(n) - n \quad \text{(for all terms decreasing faster than linearly with n)}$$

## Gaussian Integrals

$$I_m = 2\int_0^\infty x^m\, exp(-x^2)\, dx \qquad ; \text{m>-1}$$

$$I_m = 2\int_0^\infty y^n\, exp(-y)\, dy \equiv \Gamma(n+1) \qquad ; x^2 \equiv y,\; 2\,dx = y^{1/2}\,dy,\quad n \equiv \tfrac{1}{2}(m-1)$$

$$I_0 = \Gamma\left(n = \tfrac{1}{2}\right) = \sqrt{\pi} \qquad ; \text{m=0},\; n = -\tfrac{1}{2}$$

$$I_{2k} = \Gamma\left(k + \tfrac{1}{2}\right) = (k - {}^1\!/_2)(k - {}^3\!/_2)\ldots({}^3\!/_2)({}^1\!/_2)\sqrt{\pi} \qquad ; \text{even m} \quad \text{m=2 k>0},\; n = k - \tfrac{1}{2}$$

$$I_{2k+1} = \Gamma(k+1) = k! \qquad ; \text{odd m} \quad \text{m=2 k+1>0},\; n = k \geq 0$$

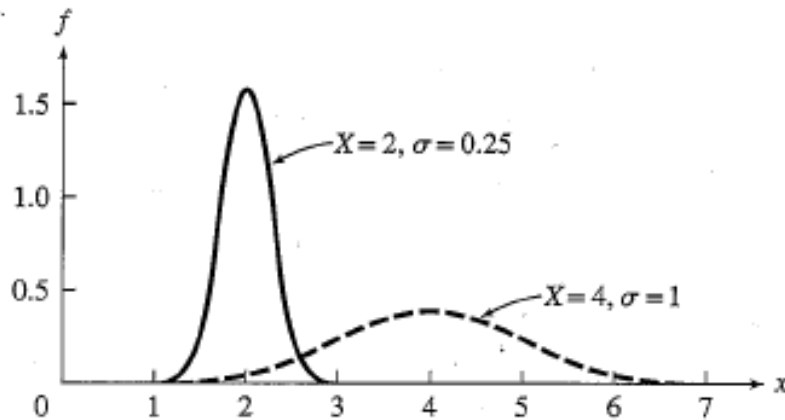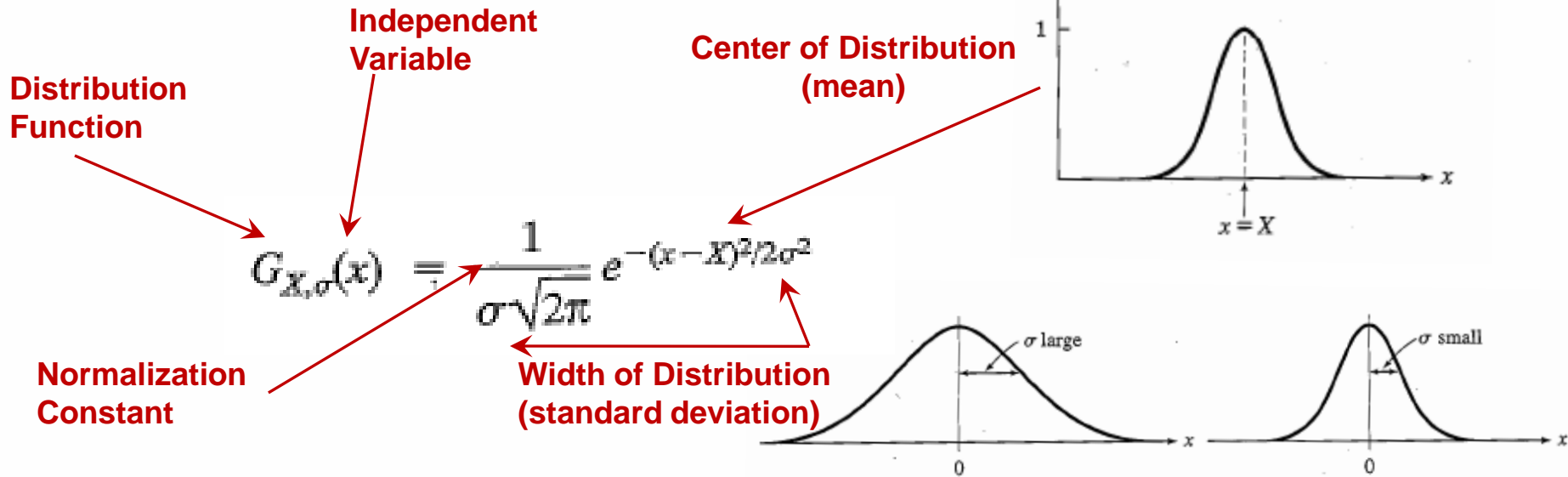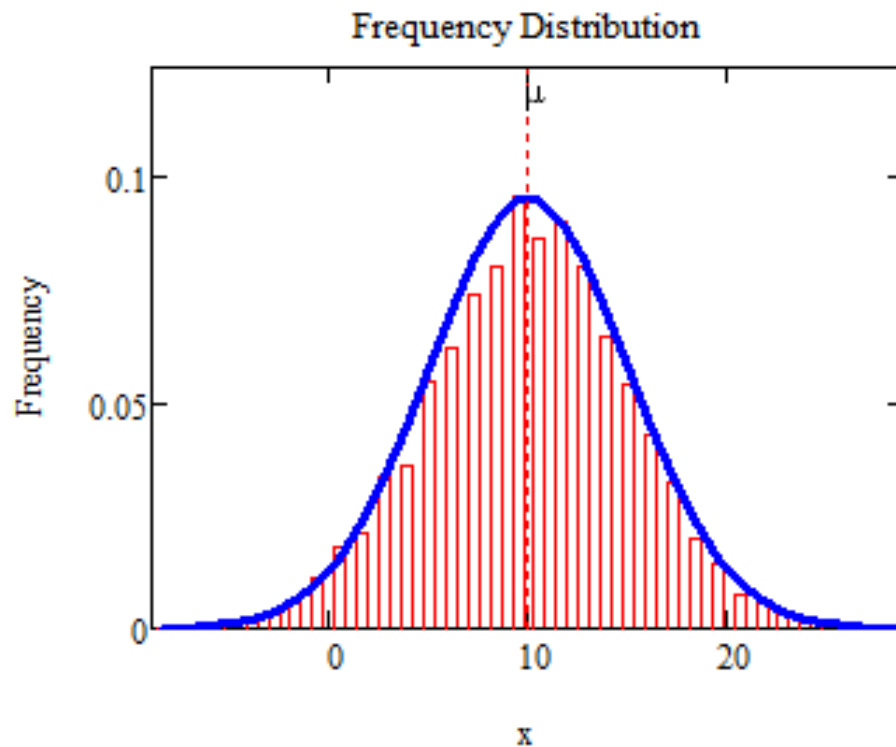# Gaussian Distribution Function

**Independent Variable**

**Center of Distribution (mean)**

**Distribution Function**



$$G_{X,\sigma}(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-X)^2/2\sigma^2}$$

**Normalization Constant**

**Width of Distribution (standard deviation)**



$\sigma$ large   $\sigma$ small

## Gaussian Distribution Function



X = 2, σ = 0.25

X = 4, σ = 1

**Figure 5.10.** Two normal, or Gauss, distributions.

# Effects of Increasing N on Gaussian Distribution Functions

**Consider the limiting distribution function as N → ∞ and dx → 0**



Frequency Distribution

# Data Pts:  $N \equiv 6000$

Mean:  $\mu \equiv 10$

Std. Dev.:  $\sigma \equiv 5$

SDOM:  $\dfrac{\sigma}{\sqrt{N}} = 0.065$

**Fractional Error:**

$$\left(\frac{\sigma}{\sqrt{N}}\right) \cdot \frac{1}{\mu} = 0.645\,\%$$

UtahState UNIVERSITY

**Defining the Gaussian distribution function in Mathcad**

## Measure Data:

Normal Distribution Parameters:

Mean: $\mu = 10.000$

Standard Deviation: $\sigma = 5.000$

Number of "Measurements": $N = 6$ $n := 0..N$

"Measurements": $x := \text{morm}(N, \mu, \sigma)$

## Histogram Calculations:

**I suggest you "investigate" these with the Mathcad sheet on the web site**

Number of intervals: $M := \text{floor}\left(N^{\frac{1}{2.5}}\right) + 2$ $M = 4.000$ $m := 0..(M)$

Interval spacing: $\Delta x := \left(\dfrac{\text{ceil}(\max(x)) - \text{floor}(\min(x))}{M}\right)$ $\Delta x = 2.500$

Calculate Intervals: $\text{Int}_m := \text{floor}(\min(x)) + m \cdot \Delta x$

Calculate Frequencies: $F := \text{hist}(\text{Int}, x)$

## Gaussian Distribution Function:

Define distribution function: $\text{Norm}(A, \mu, \sigma, x) := \dfrac{A}{\sigma \cdot \sqrt{2 \cdot \pi}} \cdot \exp\left[\dfrac{-(x - \mu)^2}{2 \cdot \sigma^2}\right]$

Maximum of distribution function: $N_{\max} := \text{Norm}(1, \mu, \sigma, \mu)$ $N_{\max} = 0.080$

UtahState
UNIVERSITY

## Alternative Distribution Function:

| Distribution function: | Normal Distribution Parameters: | | "Measurements": | Frequencies : |
|---|---|---|---|---|
| Binomial | $p := .68$ | $\nu := 15$ | $x_b := rbinom(N, \nu, p)$ | $F_b := hist(Int, x_b)$ |
| Poisson | $\lambda := 10$ | | $x_p := rpois(N, \lambda)$ | $F_p := hist(Int, x_p)$ |
| Cauchy | $1 := \mu$ | $s := \sigma$ | $x_c := rcauchy(N, 1, s)$ | $F_c := hist(Int, x_c)$ |
| Chi-squared | $d := N - 4$ | | $x_\chi := rpois(N, \lambda)$ | $F_\chi := hist(Int, x_\chi)$ |

**Using Mathcad to define other common distribution functions.**

**Consider the limiting distribution function as N →∞ and $\Delta_k$→0**



Cauchy Distribution



Chi-squared Distribution



Binomial Distribution



Poisson Distribution

UtahState
UNIVERSITY

# Gaussian Distribution Moments

**Consider the Gaussian distribution function**

$$G_{\bar{X}\sigma}(x) = \frac{1}{\sigma\sqrt{2\pi}} exp[(x-\bar{X})^2/2\sigma^2]$$

**Use the normalization condition to evaluate the normalization constant (see Taylor, p. 132)**

$$1 = \int_{-\infty}^{\infty} G_{\bar{X}\sigma}(x)\,dx = \int_{-\infty}^{\infty} N exp[-(x-\bar{X})^2/2\sigma^2]\,dx$$

$$1 \xrightarrow{y\equiv x-\bar{X},dy\equiv dx} \int_{-\infty}^{\infty} N exp[-y^2/2\sigma^2]\,dx$$

$$1 \xrightarrow{z\equiv y/\sigma,dz\equiv dy/\sigma} \int_{-\infty}^{\infty} N exp[-z^2/2]\,dz = N\sigma\sqrt{2\pi}$$

$$N = 1/(\sigma\sqrt{2\pi})$$

**The mean, $\dot{X}$, is the first moment of the Gaussian distribution function (see Taylor, p. 134)**
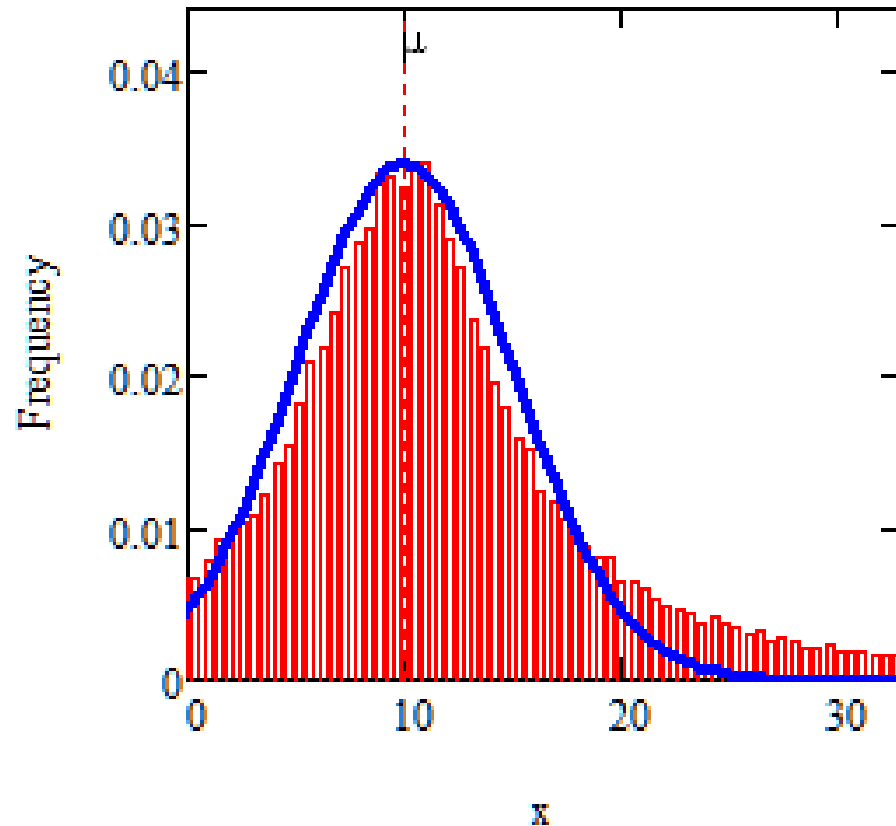
$$\langle x \rangle = \int_{-\infty}^{\infty} x\,G_{\bar{X}\sigma}(x)\,dx = \bar{X}$$

**The standard deviation, $\sigma_x$, is the standard deviation of the mean of the Gaussian distribution function (see Taylor, p. 143)**

$$\sigma_x{}^2 = \int_{-\infty}^{\infty} (x-\bar{X})^2 G_{\bar{X}\sigma}(x)\,dx = \sigma^2$$

# When is mean x not $X_{best}$?
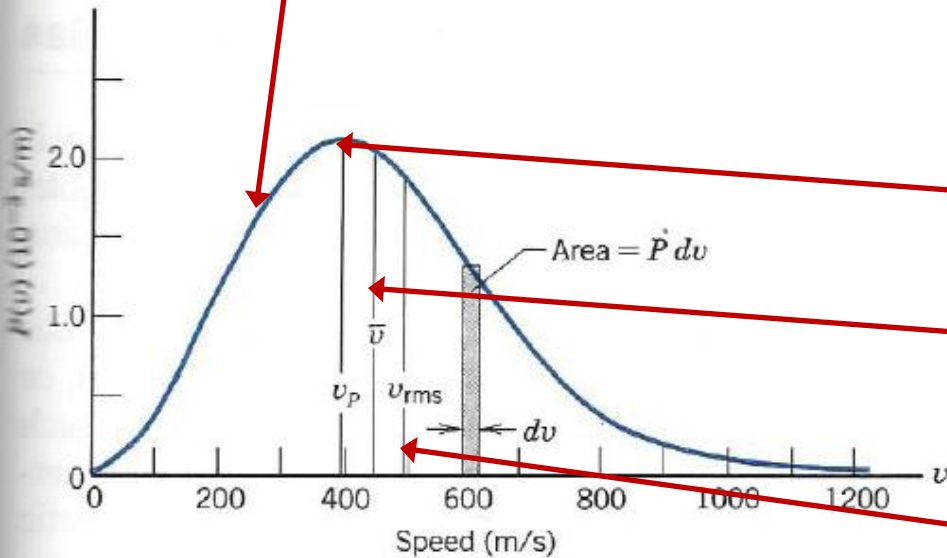


Cauchy Distribution

**Answer:  When the distribution is not symmetric about X.**

**Example:  Cauchy Distribution**

# When is mean x not X$_{best}$?

Maxwell speed distribution law is

$$P(v) = 4\pi \left( \frac{M}{2\pi RT} \right)^{3/2} v^2 e^{-Mv^2/2RT}.$$



There are three candidates for what is called the "average" value of the speed of the Maxwell speed distribution.

Firstly, by finding the maximum of the MSD (by differentiating, setting the derivative equal to zero and solving for the speed), we can determine the most probable speed. Calling this $v_{mp}$, we find that:

$$v_{mp} = \left( \frac{2kT}{m} \right)^{1/2}.$$

Second, we can find the mean value of $v$ from the MSD. Calling this $\bar{v}$:
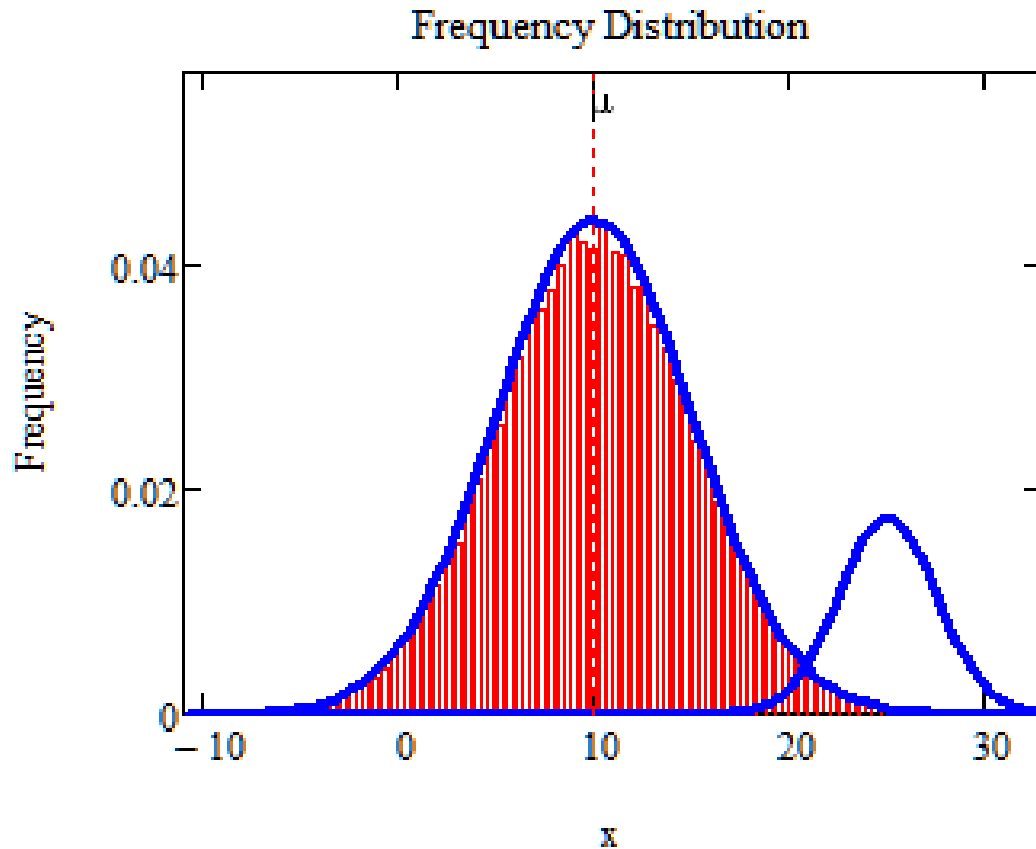
$$\bar{v} = \left( \frac{8kT}{\pi m} \right)^{1/2}.$$

Third and finally, we can find the root mean square of the speed by finding the expected value of $v^2$. (Alternatively, and much simpler, we can solve it by using the equipartition theorem.) Calling this $v_{rms}$:

$$v_{rms} = \left( \frac{3kT}{m} \right)^{1/2}.$$

Notice that $v_{mp} < \bar{v} < v_{rms}$.

These are three different ways of defining the average velocity, and they are not numerically the same. It is important to be clear about which quantity is being used.

# When is mean x not $X_{best}$?



**Answer: When the distribution is has more than one peak.**

UtahState
UNIVERSITY

# *Intermediate Lab*

## PHYS 3870

## The Gaussian Distribution Function and Its Relation to Errors

UtahState
UNIVERSITY

# A Review of Probabilities in Combination

**1 head     AND     1 Four**

**$P(H,4) = P(H) * P(4)$**

**1 Six     OR     1 Four**

**$P(6,4) = P(6) + P(4)$**

**(true for a "mutually exclusive" single role)**

**1 head     OR     1 Four**

**$P(H,4) = P(H) + P(4) - P(H \text{ and } 4)$**

**(true for a "non-mutually exclusive" events)**

**NOT   1   Six**

**$P(NOT\ 6) = 1 - P(6)$**

**Probability of a data set of N like measurements, $(x_1, x_2, \ldots x_N)$**

**$$P(x_1, x_2, \ldots x_N) = P(x)_1 * P(x_2) * \ldots P(x_N)$$**

# The Gaussian Distribution Function and Its Relation to Errors

**We will use the Gaussian distribution as applied to random variables to develop the mathematical basis for:**

- **Mean**
- **Standard deviation**
- **Standard deviation of the mean (SDOM)**
- **Moments and expectation values**
- **Error propagation formulas**
- **Addition of errors in quadrature (for independent and random measurements)**
- **Numerical values for confidence limits (t-test)**
- **Principle of maximal likelihood**
- **Central limit theorem**
- **Weighted distributions and Chi squared**
- **Schwartz inequality (i.e., the uncertainty principle) (next lecture)**

# Gaussian Distribution Moments

**Consider the Gaussian distribution function**

$$G_{\bar{X}\sigma}(x) = \frac{1}{\sigma\sqrt{2\pi}} exp[(x-\bar{X})^2/2\sigma^2]$$

**Use the normalization condition to evaluate the normalization constant (see Taylor, p. 132)**

$$1 = \int_{-\infty}^{\infty} G_{\bar{X}\sigma}(x)\,dx = \int_{-\infty}^{\infty} N exp[-(x-\bar{X})^2/2\sigma^2]\,dx$$

$$1 \xrightarrow{y\equiv x-\bar{X},dy\equiv dx} \int_{-\infty}^{\infty} N exp[-y^2/2\sigma^2]\,dx$$

$$1 \xrightarrow{z\equiv y/\sigma,dz\equiv dy/\sigma} \int_{-\infty}^{\infty} N exp[-z^2/2]\,dz = N\sigma\sqrt{2\pi}$$

$$N = 1/(\sigma\sqrt{2\pi})$$

**The mean, Ẋ, is the first moment of the Gaussian distribution function (see Taylor, p. 134)**

$$\langle x \rangle = \int_{-\infty}^{\infty} x\,G_{\bar{X}\sigma}(x)\,dx = \bar{X}$$

**The standard deviation, $\sigma_x$, is the standard deviation of the mean of the Gaussian distribution function (see Taylor, p. 143)**
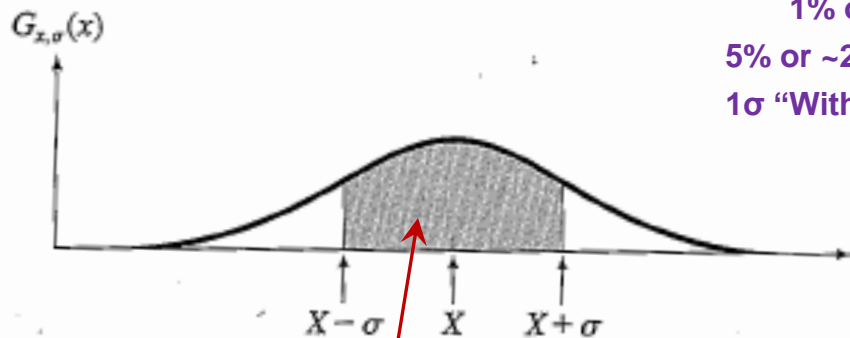
$$\sigma_x{}^2 = \int_{-\infty}^{\infty} (x-\bar{X})^2 G_{\bar{X}\sigma}(x)\,dx = \sigma^2$$

# Standard Deviation of Gaussian Distribution

$$Prob(\text{within } \sigma) = \int_{X-\sigma}^{X+\sigma} G_{X,\sigma}(x)\,dx \qquad (5.32)$$

$$= \frac{1}{\sigma\sqrt{2\pi}} \int_{X-\sigma}^{X+\sigma} e^{-(x-X)^2/2\sigma^2}\,dx. \qquad (5.33)$$
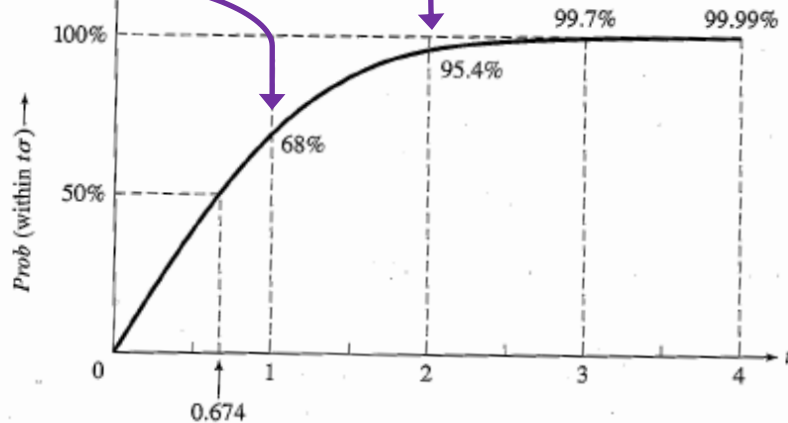
**See Sec. 10.6: Testing of Hypotheses**

**5 ppm or ~5σ "Valid for HEP"**

**1% or ~3σ "Highly Significant"**

**5% or ~2σ "Significant"**

**1σ "Within errors"**

**Area under curve (probability that −σ<x<+σ) is 68%**

**Ah, that's highly significant!**

| $t$ | 0 | 0.25 | 0.5 | 0.75 | 1.0 | 1.25 | 1.5 | 1.75 | 2.0 | 2.5 | 3.0 | 3.5 | 4.0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $Prob$ (%) | 0 | 20 | 38 | 55 | 68 | 79 | 87 | 92 | 95.4 | 98.8 | 99.7 | 99.95 | 99.99 |

**Figure 5.13.** The probability $Prob$(within $t\sigma$) that a measurement of $x$ will fall within $t$ standard deviations of the true value $x = X$. Two common names for this function are the *normal error integral* and the *error function*, erf($t$).

**More complete Table in App. A and B**

Intermediate 3870
Fall 2013

DISTRIBUTION FUNCTIONS

Lecture 3   Slide 34

UtahState
UNIVERSITY

# Error Function of Gaussian Distribution

**Error Function: (probability that –tσ<x<+tσ ).**

$$Prob(\text{within } t\sigma) = \frac{1}{\sqrt{2\pi}} \int_{-t}^{t} e^{-z^2/2} \, dz. \qquad (5.35)$$

$G_{x,\sigma}(x)$



Area under curve (probability that –tσ<x<+tσ) is given in Table at right.

**Complementary Error Function: (probability that –x<-tσ AND x>+tσ ).**
**Prob(x outside tσ) = 1 - Prob(x within tσ)**

**Probable Error: (probability that –0.67σ<x<+0.67σ) is 50%.**

**Error Function: Tabulated values-see App. A.**



| t | 0 | 0.25 | 0.5 | 0.75 | 1.0 | 1.25 | 1.5 | 1.75 | 2.0 | 2.5 | 3.0 | 3.5 | 4.0 |
|---|---|------|-----|------|-----|------|-----|------|-----|-----|-----|-----|-----|
| Prob (%) | 0 | 20 | 38 | 55 | 68 | 79 | 87 | 92 | 95.4 | 98.8 | 99.7 | 99.95 | 99.99 |

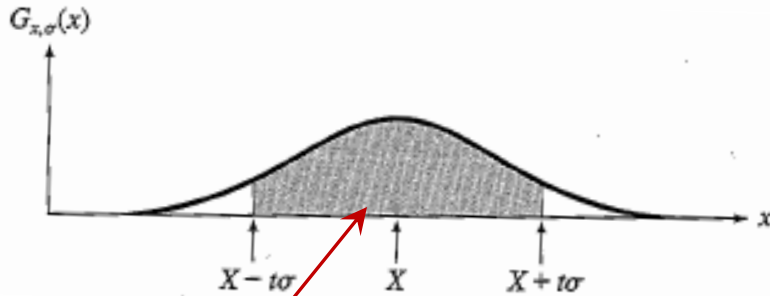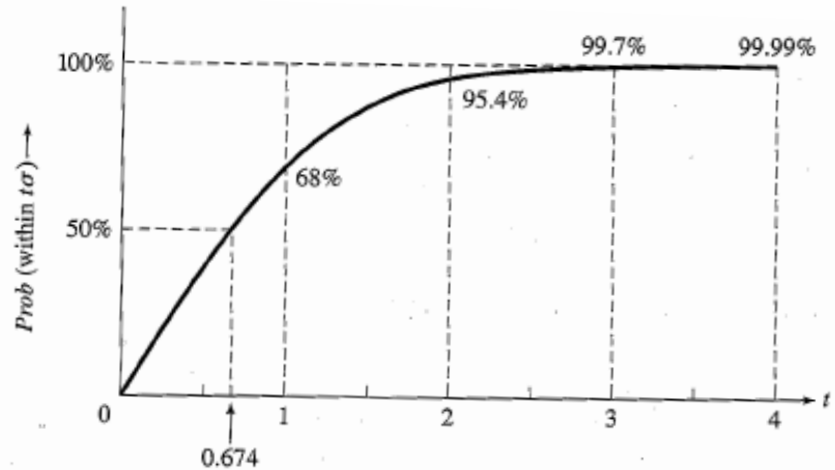**Figure 5.13.** The probability *Prob*(within *tσ*) that a measurement of *x* will fall within *t* standard deviations of the true value *x = X*. Two common names for this function are the *normal error integral* and the *error function*, erf(*t*).

**More complete Table in App. A and B**

# Useful Points on Gaussian Distribution

**Full Width at Half Maximum**
**FWHM**
**(See Prob. 5.12)**

**Points of Inflection**
**Occur at ±σ**
**(See Prob. 5.13)**

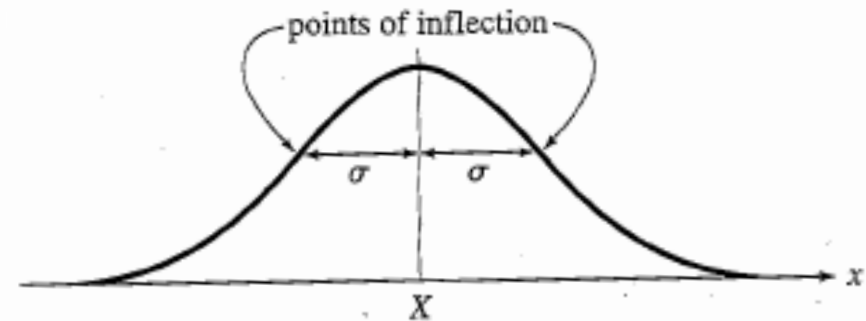$$\text{FWHM} = 2\sigma\sqrt{2 \ln 2} = 2.35\sigma.$$



**Figure 5.20.** The points $X \pm \sigma$ are the points of inflection of the Gauss curve; for Problem 5.13.

# Error Analysis and Gaussian Distribution

## Adding a Constant

$$q = x + A, \qquad (5.47)$$



$$G_{X,\sigma}(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-X)^2/2\sigma^2}$$

(probability of obtaining value $q$) $\propto e^{-[(q-A)-X]^2/2\sigma_x^2}$

$$= e^{-[q-(X+A)]^2/2\sigma_x^2}. \qquad (5.49)$$

**X→X+A**

## Multiplying by a Constant

$$q = Bx,$$



$$G_{X,\sigma}(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-X)^2/2\sigma^2}$$

(probability of obtaining value $q$) $\propto$ (probability of obtaining $x = q/B$)

$$\propto \exp\left[-\left(\frac{q}{B} - X\right)^2/2\sigma_x^2\right]$$

$$= \exp[-(q - BX)^2/2B^2\sigma_x^2]. \qquad (5.50)$$

**X→BX  and σ→B σ**

# Error Propagation: Addition

## Sum of Two Variables

**Consider the derived quantity**
**Z=X + Y**
**(with X=0 and Y=0)**

$$Prob(x) \propto \exp\left(\frac{-x^2}{2\sigma_x^2}\right) \quad (5.51)$$

$$Prob(y) \propto \exp\left(\frac{-y^2}{2\sigma_y^2}\right). \quad (5.52)$$
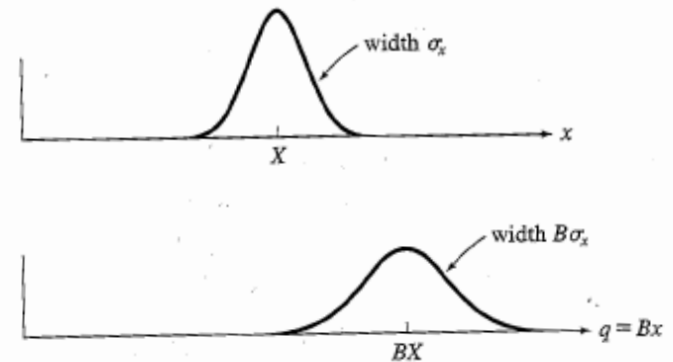
**Error in Z:**
**Multiple two probabilities**

$$Prob(x,y) = Prob(x) \cdot Prob(x) \propto exp\left[-\frac{1}{2}\left(\frac{x^2}{\sigma_x^2}+\frac{y^2}{\sigma_y^2}\right)\right]$$

$$\propto exp\left[-\frac{1}{2}\left(\frac{(x+y)^2}{(\sigma_x^2+\sigma_y^2)}-z^2\right)\right] \quad \textbf{ICBST (Eq. 5.53)}$$

$$\propto exp\left[-\frac{1}{2}\left(\frac{(x+y)^2}{(\sigma_x^2+\sigma_y^2)}\right)\right] exp\left[-\frac{z^2}{2}\right]$$

**Integrates to $\sqrt{2\pi}$**

**X+Y→Z and $\sigma_x^2 + \sigma_y^2 \rightarrow \sigma_z^2$**
**(addition in quadrature for random, independent variables!)**

Intermediate 3870
Fall 2013

DISTRIBUTION FUNCTIONS

Lecture 3   Slide 38

UtahState
UNIVERSITY

# General Formula for Error Propagation

**How do we determine the error of a derived quantity Z(X,Y,…) from errors in X,Y,…?**

General formula for error propagation  see  [Taylor, Secs. 3.5 and 3.9]

Uncertainty as a function of one variable  [Taylor, Sec. 3.5]

1. Consider a graphical method of estimating error
   a) Consider an arbitaray function q(x)
   b) Plot q(x) vs. x.
   c) On the graph, label:
      (1)  $q_{best} = q(x_{best})$
      (2)  $q_{hi} = q(x_{best} + \delta x)$
      (3)  $q_{low} = q(x_{best} - \delta x)$
   d) Making a linear approximation:

$$q_{hi} = q_{best} + slope \cdot \delta x = q_{best} + \left(\frac{\partial q}{\partial x}\right)$$

$$q_{low} = q_{best} - slope \cdot \delta x = q_{best} - \left(\frac{\partial q}{\partial x}\right)$$

   e) Therefore:

$$\delta q = \left|\frac{\partial q}{\partial x}\right| \cdot \delta x$$

   Note the absolute value.



**Figure 3.3.** Graph of $q(x)$ vs $x$. If $x$ is measured as $x_{best} \pm \delta x$, then the best estimate for $q(x)$ is $q_{best} = q(x_{best})$. The largest and smallest probable values of $q(x)$ correspond to the values $x_{best} \pm \delta x$ of $x$.

Intermediate 3870
Fall 2013
DISTRIBUTION FUNCTIONS
Lecture 3   Slide 39

UtahState
UNIVERSITY

# General Formula for Error Propagation

General formula for uncertainty of a function of one variable

$$\delta q = \left| \frac{\partial q}{\partial x} \right| \cdot \delta x \qquad \text{[Taylor, Eq. 3.23]}$$

Can you now derive for specific rules of error propagation:

1. Addition and Subtraction   **[Taylor, p. 49]**
2. Multiplication and Division  **[Taylor, p. 51]**
3. Multiplication by a constant (exact number)   **[Taylor, p. 54]**
4. Exponentiation  (powers)  **[Taylor, p. 56]**

# General Formula for Multiple Variables

Uncertainty of a function of multiple variables  [Taylor, Sec. 3.11]

1. It can easily (no, really) be shown that (see Taylor Sec. 3.11) for a function of several variables

$$\delta q(x, y, z, \ldots) = \left|\frac{\partial q}{\partial x}\right| \cdot \delta x + \left|\frac{\partial q}{\partial y}\right| \cdot \delta y + \left|\frac{\partial q}{\partial z}\right| \cdot \delta z + \ldots$$   [Taylor, Eq. 3.47]

2. More correctly, it can be shown that (see Taylor Sec. 3.11) for a function of several variables

$$\delta q(x, y, z, \ldots) \leq \left|\frac{\partial q}{\partial x}\right| \cdot \delta x + \left|\frac{\partial q}{\partial y}\right| \cdot \delta y + \left|\frac{\partial q}{\partial z}\right| \cdot \delta z + \ldots$$   [Taylor, Eq. 3.47]

where the equals sign represents an upper bound, as discussed above.

3. For a function of several *independent and random* variables

$$\delta q(x, y, z, \ldots) = \sqrt{\left(\frac{\partial q}{\partial x} \cdot \delta x\right)^2 + \left(\frac{\partial q}{\partial y} \cdot \delta y\right)^2 + \left(\frac{\partial q}{\partial z} \cdot \delta z\right)^2 + \ldots}$$  [Taylor, Eq. 3.48]

Again, the proof is left for Ch. 5.

# Error Propagation:  General Case

**How do we determine the error of a derived quantity Z(X,Y,…) from errors in X,Y,…?**

**Consider the arbitrary derived quantity
q(x,y) of two independent random variables x and y.**

**Expand q(x,y) in a Taylor series about the expected values of x and y
(i.e., at points near X and Y).**

**Fixed, shifts peak of distribution**

$$q(x, y) = q(X, Y) + \left(\frac{\partial q}{\partial x}\right)\Big|_X (x - X) + \left(\frac{\partial q}{\partial y}\right)\Big|_Y (y - Y)$$

**Fixed     Distribution centered at X with width σ$_X$**

**Product of Two Variables**

$$\delta q(x, y) = \sigma_q = \sqrt{q(X, Y) + \left[\left(\frac{\partial q}{\partial x}\right)\Big|_X \sigma_x\right]^2 + \left[\left(\frac{\partial q}{\partial y}\right)\Big|_Y \sigma_y\right]^2}$$

**0**

# SDOM of Gaussian Distribution

## Standard Deviation of the Mean

**Each measurement has similar $\sigma_{x_i} = \sigma_{\bar{x}}$**

$$\sigma_{x_1} = \cdots = \sigma_{x_N} = \sigma_x.$$

$$G_{X,\sigma}(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-X)^2/2\sigma^2}$$

**and similar partial derivatives**

$$\frac{\partial \bar{x}}{\partial x_1} = \cdots = \frac{\partial \bar{x}}{\partial x_N} = \frac{1}{N}.$$

**Thus…**

$$\sigma_{\bar{x}} = \sqrt{\left(\frac{1}{N}\sigma_x\right)^2 + \cdots + \left(\frac{1}{N}\sigma_x\right)^2}$$

$$= \sqrt{N\frac{\sigma_x^2}{N^2}} = \frac{\sigma_x}{\sqrt{N}},$$

(5.66)

width $\sigma_{\bar{x}} = \frac{\sigma_x}{\sqrt{10}}$

width $\sigma_x$

**The SDOM decreases as the square root of the number of measurements.**

**That is, the relative width, $\sigma/\dot{X}$, of the distribution gets narrower as more measurements are made.**

# Two Key Theorems from Probability

## Central Limit Theorem

For random, independent measurements (each with a well-define expectation value and well-defined variance), the arithmetic mean (average) will be approximately normally distributed.

## Principle of Maximum Likelihood

Given the N observed measurements, $x_1$, $x_2$,…$x_N$, the best estimates for $\dot{X}$ and $\sigma$ are those values for which the observed $x_1$, $x_2$,…$x_N$, are most likely.

# Mean of Gaussian Distribution as "Best Estimate"

## Principle of Maximum Likelihood

To find the most likely value of the mean (the best estimate of ẋ), find X that yields the highest probability for the data set.

Consider a data set $\{x_1, x_2, x_3 \ldots x_N\}$

Each randomly distributed with

$$Prob_{X,\sigma}(x_i) = G_{X,\sigma}(x_i) \equiv \frac{1}{\sigma\sqrt{2\pi}} e^{-(x_i-X)^2/2\sigma} \propto \frac{1}{\sigma} e^{-(x_i-X)^2/2\sigma}$$

The combined probability for the full data set is the product

$$Prob_{X,\sigma}(x_1, x_2 \ldots x_N) = Prob_{X,\sigma}(x_1) \times Prob_{X,\sigma}(x_2) \times \ldots \times Prob_{X,\sigma}(x_N)$$

$$\propto \frac{1}{\sigma} e^{-(x_1-X)^2/2\sigma} \times \frac{1}{\sigma} e^{-(x_2-X)^2/2\sigma} \times \ldots \times \frac{1}{\sigma} e^{-(x_N-X)^2/2\sigma} = \frac{1}{\sigma^N} e^{\sum -(x_i-X)^2/2\sigma}$$

## Best Estimate of X is from maximum probability or minimum summation

| | | | | | |
|---|---|---|---|---|---|
| **Minimize Sum** | $\displaystyle\sum_{i=1}^{N}(x_i - X)^2/\sigma$ | **Solve for derivative wrst X set to 0** | $\displaystyle\sum_{i=1}^{N}(x_i - X) = 0$ | **Best estimate of X** | $X_{best} = \dfrac{\sum x_i}{N}$ |

# Uncertainty of "Best Estimates" of Gaussian Distribution

## Principle of Maximum Likelihood

To find the most likely value of the standard deviation (the best estimate of the width of the x distribution), find $\sigma_x$ that yields the highest probability for the data set.

Consider a data set $\{x_1, x_2, x_3 \dots x_N\}$

The combined probability for the full data set is the product

$$Prob_{X,\sigma}(x_1, x_2 \dots x_N) = Prob_{X,\sigma}(x_1) \times Prob_{X,\sigma}(x_2) \times \dots \times Prob_{X,\sigma}(x_N)$$

$$\propto \frac{1}{\sigma} e^{-(x_1-X)^2/2\sigma} \times \frac{1}{\sigma} e^{-(x_2-X)^2/2\sigma} \times \dots \times \frac{1}{\sigma} e^{-(x_N-X)^2/2\sigma} = \frac{1}{\sigma^N} e^{\sum -(x_i-X)^2/2\sigma}$$

## Best Estimate of X is from maximum probability or minimum summation

**Minimize Sum** $\displaystyle\sum_{i=1}^{N} (x_i - X)^2/\sigma$   **Solve for derivative wrst X set to 0** $\displaystyle\sum_{i=1}^{N}(x_i - X) = 0$   **Best estimate of X** $X_{best} = \frac{\sum x_i}{N}$

## Best Estimate of σ is from maximum probability or minimum summation

**Minimize Sum** $\displaystyle\sum_{i=1}^{N} (x_i - X)^2/\sigma$   **Solve for derivative wrst σ set to 0** See Prob. 5.26   **Best estimate of σ** $\sigma_{best} = \sqrt{\frac{1}{N}\sum_{i=1}^{N}(x_i - X)^2/\sigma}$

UtahState UNIVERSITY

# *Intermediate Lab*
## PHYS 3870

## Combining Data Sets
## Weighted Averages

References:  Taylor Ch.  7

UtahState
UNIVERSITY

# Weighted Averages

**Question:  How can we properly combine two or more separate independent measurements of the same randomly distributed quantity to determine a best combined value with uncertainty?**

UtahState
UNIVERSITY

Intermediate  3870
Fall 2013

DISTRIBUTION FUNCTIONS

Lecture  3   Slide  48

# Weighted Averages

Consider two measurements of the same quantity, described by a random Gaussian distribution

$<x_1> \pm \sigma_{x1}$   and   $<x_2> \pm \sigma_{x2}$   **Assume negligible systematic errors.**

The probability of measuring two such measurements is

$$Prob_x(x_1, x_2) = Prob_x(x_1)\, Prob_x(x_2)$$

$$= \frac{1}{\sigma_1 \sigma_2} e^{-\chi^2/2}\ \ where\ \ \chi^2 \equiv \left[\frac{(x_1 - X)}{\sigma_1}\right]^2 + \left[\frac{(x_2 - X)}{\sigma_2}\right]^2$$

To find the best value for χ, find the maximum Prob or minimum $\chi^2$

**Note: $\chi^2$ , or Chi squared, is the sum of the squares of the deviations from the mean, divided by the corresponding uncertainty.**

**Such methods are called "Methods of Least Squares".  They follow directly from the Principle of Maximum Likelihood.**

# Weighted Averages

The probability of measuring two such measurements is

$$Prob_x(x_1, x_2) = Prob_x(x_1)\, Prob_x(x_2)$$

$$= \frac{1}{\sigma_1 \sigma_2} e^{-\chi^2/2} \quad where \ \chi^2 \equiv \left[\frac{(x_1 - X)}{\sigma_1}\right]^2 + \left[\frac{(x_2 - X)}{\sigma_2}\right]^2$$

To find the best value for $\chi$, find the maximum Prob or minimum $\chi^2$

**Best Estimate of $\chi$ is from maximum probility or minimum summation**

**Minimize Sum**            **Solve for derivative wrst $\chi$ set to 0**    **Solve for best estimate of $\chi$**

$$\chi^2 \equiv \left[\frac{(x_1 - X)}{\sigma_1}\right]^2 + \left[\frac{(x_2 - X)}{\sigma_2}\right]^2 \qquad 2\left|\frac{(x_1 - X)}{\sigma_1}\right| + 2\left|\frac{(x_2 - X)}{\sigma_2}\right| = 0 \qquad X_{best} = \left(\frac{x_1}{\sigma_1^2} + \frac{x_2}{\sigma_2^2}\right)\Big/\left(\frac{1}{\sigma_1^2} + \frac{1}{\sigma_2^2}\right)$$

This leads to

$$x_{W\_avg} = \frac{w_1 x_1 + w_2 x_2}{w_1 + w_2} = \frac{\sum w_i\, x_i}{\sum w_i} \quad where \ w_i = {1}/{(\sigma_i)^2}$$

**Note: If $w_1 = w_2$, we recover the standard result $X_{wavg} = (1/2)(x_1 + x_2)$**

Finally, the width of a weighted average distribution is

$$\sigma_{wieghted\ avg} = \frac{1}{\sum_i w_i}$$

# Weighted Averages on Steroids

**A very powerful method for combining data from different sources with different methods and uncertainties (or, indeed, data of related measured and calculated quantities) is Kalman filtering.**

**The Kalman filter, also known as linear quadratic estimation (LQE), is an algorithm that uses a series of measurements observed over time, containing noise (random variations) and other inaccuracies, and produces estimates of unknown variables that tend to be more precise than those based on a single measurement alone.**



The Kalman filter keeps track of the estimated state of the system and the variance or uncertainty of the estimate. The estimate is updated using a state transition model and measurements. $x_{k|k-1}$ denotes the estimate of the system's state at time step k before the $k^{th}$ measurement $y_k$ has been taken into account; $P_{k|k-1}$ is the corresponding uncertainty.  --Wikipedia, 2013.

**Ludger Scherliess, of USU Physics, is a world expert at using Kalman filtering for the assimilation of satellite and ground-based data and the USU GAMES model to predict  space weather .**

UtahState
UNIVERSITY

# *Intermediate Lab*
## PHYS 3870

## Rejecting Data
## Chauvenet's Criterion

References: Taylor Ch. 6

UtahState
UNIVERSITY

# Rejecting Data

## What is a good criteria for rejecting data?

**Question: When is it "reasonable" to discard a seemingly "unreasonable" data point from a set of randomly distributed measurements?**

- Never
- Whenever it makes things look better
- Chauvenet's criterion provides a (quantitative) compromise

UtahState
UNIVERSITY

# Rejecting Data

## Zallen's Criterion

**Question: When is it "reasonable" to discard a seemingly "unreasonable" data point from a set of randomly distributed measurements?**

Often in physics, experimental observations are termed "anomalous" before they are understood. Once theory succeeds in explaining and illuminating the observations, they are no longer "anomalous" and instead come to be regarded as "obvious." A crucial paper can trigger such an "anomalous → obvious" transition, and in the present case that key role was played by a 1975 paper by Scher and Montroll. That landmark paper has become basic to our understanding of a striking characteristic of carrier motion (now called *dispersive transport*) which is a common occurrence in amorphous semiconductors, though foreign to our experience with crystals.

UtahState
UNIVERSITY

# Rejecting Data

## Disney's Criterion

**Question:  When is it "reasonable" to discard a seemingly "unreasonable" data point from a set of randomly distributed measurements?**

- **Whenever it makes things look better**

**Disney's First Law**

**Wishing will make it so.**

**Disney's Second Law**

**Dreams are more colorful than reality.**

UtahState
UNIVERSITY

# Rejecting Data

## Chauvenet's Criterion

**Data may be rejected if the expected number of measurements at least as deviant as the suspect measurement is less than 50%.**

Consider a set of N measurements of a single quantity

$\{ x_1, x_2, \ldots\ldots x_N \}$

Calculate $<x>$ and $\sigma_x$ and then determine the fractional deviations from the mean of all the points:

$$x_{frac\_dev} = \frac{|x_i - \bar{x}|}{\sigma_x}$$

For the suspect point(s), $x_{suspect}$, find the probability of such a point occurring in N measurements

n = (expected number as deviant as $x_{suspect}$)

= N Prob(outside $x_{suspect} \cdot \sigma_x$)

# Error Function of Gaussian Distribution

**Error Function: (probability that –tσ<x<+tσ ).**

$$Prob(\text{within } t\sigma) = \frac{1}{\sqrt{2\pi}} \int_{-t}^{t} e^{-z^2/2} \, dz. \qquad (5.35)$$

$G_{x,\sigma}(x)$

**Error Function: Tabulated values-see App. A.**

**Area under curve (probability that –tσ<x<+tσ) is given in Table at right.**

| $t$ | 0 | 0.25 | 0.5 | 0.75 | 1.0 | 1.25 | 1.5 | 1.75 | 2.0 | 2.5 | 3.0 | 3.5 | 4.0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Prob (%) | 0 | 20 | 38 | 55 | 68 | 79 | 87 | 92 | 95.4 | 98.8 | 99.7 | 99.95 | 99.99 |

**Figure 5.13.** The probability *Prob*(within *t*σ) that a measurement of *x* will fall within *t* standard deviations of the true value *x* = *X*. Two common names for this function are the *normal error integral* and the *error function, erf(t).*

**Probable Error: (probability that –0.67σ<x<+0.67σ) is 50%.**

# Chauvenet's Criterion

The probability that a data point is likely to fall outside a given deviation is:

$$\text{Prob}(x_{test}, X, \sigma) := 1 - \int_{-|x_{test}|}^{|x_{test}|} \frac{1}{\sigma \cdot \sqrt{2 \cdot \pi}} \cdot e^{-\left[\frac{(x-X)^2}{2 \cdot \sigma^2}\right]} dx$$

x =

|   | 0 | m |
|---|------|---|
| 0 | 45.7 |   |
| 1 | 46.2 |   |
| 2 | 46.9 |   |
| 3 | 54.8 |   |
| 4 | 46.1 |   |
| 5 | 45.2 |   |
| 6 | 45.4 |   |
| 7 | 47   |   |
| 8 | 45.9 |   |
| 9 | 46.3 |   |

Frac_Dev =

|   | 0 |
|---|-------|
| 0 | 0.468 |
| 1 | 0.281 |
| 2 | 0.019 |
| 3 | 2.936 |
| 4 | 0.318 |
| 5 | 0.655 |
| 6 | 0.58  |
| 7 | 0.019 |
| 8 | 0.393 |
| 9 | 0.243 |

$\text{Prob}(x_i, x_{mean}, \sigma_x) =$

| |
|---|
| 0.68 |
| 0.61 |
| 0.507 |
| $1.66 \cdot 10^{-3}$ |
| 0.625 |
| 0.744 |
| 0.719 |
| 0.493 |
| 0.653 |
| 0.596 |

Including all data points

$x_{mean} := \text{mean}(x) = 46.95\,\text{m}$

$\sigma_x := \text{stdev}(x) = 2.673\,\text{m}$

Excluding the rejected data point

$X_{mean} := \text{mean}(X_{CH}) = 46.078\,\text{m}$

$\sigma_x := \text{stdev}(X_{CH}) = 0.577\,\text{m}$

UtahState UNIVERSITY

# Chauvenet's Criterion—Example 1

Example: Ten Measurements of a Length

A student makes 10 measurements of one length $x$ and gets the results (all in mm)

$$46, 48, 44, 38, 45, 47, 58, 44, 45, 43.$$

Noticing that the value 58 seems anomalously large, he checks his records but can find no evidence that the result was caused by a mistake. He therefore applies Chauvenet's criterion. What does he conclude?

Accepting provisionally all 10 measurements, he computes

$$\bar{x} = 45.8 \quad \text{and} \quad \sigma_x = 5.1.$$

# Chauvenet's Details (1)

The difference between the suspect value $x_{sus} = 58$ and the mean $\bar{x} = 45.8$ is 12.2, or 2.4 standard deviations; that is,

$$t_{sus} = \frac{x_{sus} - \bar{x}}{\sigma_x} = \frac{58 - 45.8}{5.1} = 2.4.$$

Referring to the table in Appendix A, he sees that the probability that a measurement will differ from $\bar{x}$ by $2.4\sigma_x$ or more is

$$Prob(\text{outside } 2.4\sigma) = 1 - Prob(\text{within } 2.4\sigma)$$
$$= 1 - 0.984$$
$$= 0.016.$$

In 10 measurements, he would therefore expect to find only 0.16 of one measurement as deviant as his suspect result. Because 0.16 is less than the number 0.5 set by Chauvenet's criterion, he should at least consider rejecting the result.

If he decides to reject the suspect 58, then he must recalculate $\bar{x}$ and $\sigma_x$ as

$$\bar{x} = 44.4 \quad \text{and} \quad \sigma_x = 2.9.$$

As you would expect, his mean changes a bit, and his standard deviation drops appreciably.

Intermediate 3870

Fall 2013

DISTRIBUTION FUNCTIONS

Lecture 3   Slide 60

UtahState
UNIVERSITY

# Chauvenet's Criterion—Details (2)

Consider the following example of the application of Chauvenet's Criterion to determine if a certain datum should be rejected.

A set of N=10 measurements of a length are made. The data are assumed to be described by a randon Gaussian distribution.

## Enter Data

Number of data points: $N := 10$

Data indices: $i := 0..(N-1)$

Enter data set: $x_i :=$

| |
|---|
| 45.7 · m |
| 46.2 · m |
| 46.9 · m |
| 54.8 · m |
| 46.1 · m |
| 45.2 · m |
| 45.4 · m |
| 47.0 · m |
| 45.9 · m |
| 46.3 · m |

Calculate mean: $x_{mean} := mean(x) = 46.95m$

Calculate standard deviation: $\sigma_x := stdev(x) = 2.673m$

Calculate fractional deviation from the mean: $Frac\_Dev_i := \left| \dfrac{x_i - x_{mean}}{\sigma_x} \right|$

To apply Chauvenet's criterion, we first sort the data x in order of ascending values of the fractional deviation from the mean. The probability that a data point is likely to fall outside a given deviation is then calculated. We then determine how many data that should be eliminated based on Chauvenet's

# Chauvenet's Criterion— Details (3)

Sort data in ascending order:

$$x_{order} := csort\left(augment\left(\frac{x}{m}, Frac\_Dev\right), 1\right)$$

The probability that a data point is likely to fall outside a given deviation is:

$$Prob(x_{test}, X, \sigma) := 1 - \int_{-|x_{test}|}^{|x_{test}|} \frac{1}{\sigma \cdot \sqrt{2 \cdot \pi}} \cdot e^{-\left[\frac{(x-X)^2}{2 \cdot \sigma^2}\right]} dx$$

Apply Chauvenet's criterion:

$$Reject(x, X, \sigma, N) := if[(N \cdot Prob(x, X, \sigma)) > 50 \cdot \%, "Keep", "Reject"]$$

Determine how many data points should be rejected:

$$N_{reject} := \sum_{i=0}^{N-1} if\left[(N \cdot Prob(x_i, x_{mean}, \sigma_x)) > 50 \cdot \%, 0, 1\right] = 1$$

| x = | | | $\frac{Frac\_Dev}{\%}$ = | | $Prob(x_i, x_{mean}, \sigma_x)$ = | $N \cdot Prob(x_i, x_{mean}, \sigma_x)$ = | $Reject(x_i, x_{mean}, \sigma_x, N)$ = | |
|---|---|---|---|---|---|---|---|---|
| | 0 | · m | | 0 | 0.68 | 6.8 | | 0 |
| 0 | 45.7 | | 0 | 46.759 | 0.61 | 6.105 | 0 | "Keep" |
| 1 | 46.2 | | 1 | 28.055 | 0.507 | 5.075 | 1 | "Keep" |
| 2 | 46.9 | | 2 | 1.87 | $1.66 \cdot 10^{-3}$ | 0.017 | 2 | "Keep" |
| 3 | 54.8 | | 3 | 293.645 | 0.625 | 6.247 | 3 | "Reject" |
| 4 | 46.1 | | 4 | 31.796 | 0.744 | 7.436 | 4 | "Keep" |
| 5 | 45.2 | | 5 | 65.462 | 0.719 | 7.19 | 5 | "Keep" |
| 6 | 45.4 | | 6 | 57.981 | 0.493 | 4.925 | 6 | "Keep" |
| 7 | 47 | | 7 | 1.87 | 0.653 | 6.528 | 7 | "Keep" |
| 8 | 45.9 | | 8 | 39.277 | 0.596 | 5.961 | 8 | "Keep" |
| 9 | 46.3 | | 9 | 24.315 | | | 9 | "Keep" |

UtahState UNIVERSITY

# Chauvenet's Criterion— Example 2

$$x =$$

| | | $\cdot$ m |
|---|---|---|
| | 0 | |
| 0 | 45.7 | |
| 1 | 46.2 | |
| 2 | 46.9 | |
| 3 | 54.8 | |
| 4 | 46.1 | |
| 5 | 45.2 | |
| 6 | 45.4 | |
| 7 | 47 | |
| 8 | 45.9 | |
| 9 | 46.3 | |

$$\frac{Frac\_Dev}{\%} =$$

| | |
|---|---|
| | 0 |
| 0 | 46.759 |
| 1 | 28.055 |
| 2 | 1.87 |
| 3 | 293.645 |
| 4 | 31.796 |
| 5 | 65.462 |
| 6 | 57.981 |
| 7 | 1.87 |
| 8 | 39.277 |
| 9 | 24.315 |

$$\text{Prob}(x_i, x_{mean}, \sigma_x) =$$

| |
|---|
| 0.68 |
| 0.61 |
| 0.507 |
| $1.66 \cdot 10^{-3}$ |
| 0.625 |
| 0.744 |
| 0.719 |
| 0.493 |
| 0.653 |
| 0.596 |

$$N \cdot \text{Prob}(x_i, x_{mean}, \sigma_x) =$$

| |
|---|
| 6.8 |
| 6.105 |
| 5.075 |
| 0.017 |
| 6.247 |
| 7.436 |
| 7.19 |
| 4.925 |
| 6.528 |
| 5.961 |

$$\text{Reject}(x_i, x_{mean}, \sigma_x, N) =$$

| | |
|---|---|
| | 0 |
| 0 | "Keep" |
| 1 | "Keep" |
| 2 | "Keep" |
| 3 | "Reject" |
| 4 | "Keep" |
| 5 | "Keep" |
| 6 | "Keep" |
| 7 | "Keep" |
| 8 | "Keep" |
| 9 | "Keep" |

Now recalculate the mean and standard deviation after rejecting $N_{reject}$ data points.

Truncated data set indices and data array:    $j := 0 .. N - 1 - N_{reject}$     $X_{CH_j} := x_{order_{j,0}}$

**The final analysis is:**

| | Including all data points | Excluding the rejected data point |
|---|---|---|
| Number data points: | $N = 10$ | $N - N_{reject} = 9$ |
| Mean: | $x_{mean} := \text{mean}(x) = 46.95 \, m$ | $X_{mean} := \text{mean}(X_{CH}) = 46.078$ |
| Standard Deviation: | $\sigma_x := \text{stdev}(x) = 2.673 \, m$ | $\sigma_x := \text{stdev}(X_{CH}) = 0.577$ |

UtahState UNIVERSITY

# *Intermediate Lab*
## PHYS 3870

# Summary of Probability Theory

Intermediate 3870
Fall 2013

DISTRIBUTION FUNCTIONS

Lecture 3   Slide 64

UtahState
UNIVERSITY

*Probabilit  action (Discrete Case)*

The random variable X will be called a discrete random variable if there exists a function $f$ such that $f(x_i) \geq 0$ and $\sum_i f(x_i) = 1$ for $i = 1, 2, 3, \ldots$ and such that for any event $E$,

$$P(E) = P[X \text{ is in } E] = \sum_E f(x)$$

where $\sum_E$ means sum $f(x)$ over those values $x_i$ that are in $E$ and where $f(x) = P[X = x]$.

The probability that the value of X is some real number $x$, is given by $f(x) = P[X = x]$, where $f$ is called the probability function of the random variable X.

*Cumulative Distribution Function (Discrete Case)*

The probability that the value of a random variable X is less than or equal to some real number $x$ is defined as

$$F(x) = P(X \leq x)$$
$$= \Sigma f(x_i), \qquad -\infty < x < \infty,$$

where the summation extends over those values of $i$ such that $x_i \leq x$.

*Probability Density (Continuous Case)*

The random variable X will be called a continuous random variable if there exists a function $f$ such that $f(x) \geq 0$ and $\int_{-\infty}^{\infty} f(x)\, dx = 1$ for all $x$ in interval $-\infty < x < \infty$ and such that for any event $E$

$$P(E) = P(X \text{ is in } E) = \int_E f(x)\, dx.$$

$f(x)$ is called the probability density of the random variable X. The probability that X assumes any given value of $x$ is equal to zero and the probability that it assumes a value on the interval from $a$ to $b$, including or excluding either end point, is equal to

$$\int_a^b f(x)\, dx.$$

# Summary of Probability Theory-I

Intermediate 3870
Fall 2013

DISTRIBUTION FUNCTIONS

Lecture 3   Slide 65

UtahState
UNIVERSITY

# Summary of Probability Theory-II

*Probability Density (Continuous Case)*

The random variable X will be called a continuous random variable if there exists a function $f$ such that $f(x) \geq 0$ and $\int_{-\infty}^{\infty} f(x)\, dx = 1$ for all $x$ in interval $-\infty < x < \infty$ and such that for any event $E$

$$P(E) = P(X \text{ is in } E) = \int_E f(x)\, dx.$$

$f(x)$ is called the probability density of the random variable X. The probability that X assumes any given value of $x$ is equal to zero and the probability that it assumes a value on the interval from $a$ to $b$, including or excluding either end point, is equal to

$$\int_a^b f(x)\, dx.$$

*Cumulative Distribution Function (Continuous Case)*

The probability that the value of a random variable X is less than or equal to some real number $x$ is defined as

$$F(x) = P(X \leq x), \qquad -\infty < x < \infty$$
$$= \int_{-\infty}^{x} f(x)\, dx.$$

From the cumulative distribution, the density, if it exists, can be found from

$$f(x) = \frac{dF(x)}{dx}.$$

From the cumulative distribution

$$P(a \leq X \leq b) = P(X \leq b) - P(X \leq a)$$
$$= F(b) - F(a)$$

# Summary of Probability Theory-III

*Mathematical Expectation*

A. EXPECTED VALUE.

Let X be a random variable with density $f(x)$. Then the expected value of X, $E(X)$, is defined to be

$$E(X) = \sum_x xf(x)$$

if X is discrete and

$$E(X) = \int_{-\infty}^{\infty} xf(x)\,dx$$

if X is continuous. The expected value of a function $g$ of a random variable X is defined as

$$E[g(X)] = \sum_x g(x) \cdot f(x)$$

if X is discrete and

$$E[g(X)] = \int_{-\infty}^{\infty} g(x) \cdot f(x)\,dx$$

if X is continuous.

*Theorems*

1. $E[aX + bY] = aE(X) + bE(Y)$
2. $E[X \cdot Y] = E(X) \cdot E(Y)$ if X and Y are statistically independent.

Intermediate 3870

Fall 2013

DISTRIBUTION FUNCTIONS

Lecture 3   Slide 67

UtahState
UNIVERSITY

# Summary of Probability Theory-IV

B. MOMENTS

  a. *Moments About the Origin.* The moments about the origin of a probability distribution are the expected values of the random variable which has the given distribution. The $r$th moment of X, usually denoted by $\mu'_r$, is defined as

$$\mu'_r = E[X^r] = \sum_x x^r f(x)$$

if X is discrete and

$$\mu'_r = E[X^r] = \int_{-\infty}^{\infty} x^r f(x)\, dx$$

if X is continuous.

  The first moment, $\mu'_1$, is called the mean of the random variable X and is usually denoted by $\mu$.

  b. *Moments About the Mean.* The $r$th moment about the mean, usually denoted by $\mu_r$, is defined as

$$\mu_r = E[(X - \mu)^r] = \sum_x (x - \mu)^r f(x)$$

if X is discrete and

$$\mu_r = E[(X - \mu)^r] = \int_{-\infty}^{\infty} (x - \mu)^r f(x)\, dx$$

if X is continuous.

  The second moment about the mean, $\mu_2$, is given by

$$\mu_2 = E[(X - \mu)^2] = \mu'_2 - \mu^2$$

and is called the variance of the random variable X, and is denoted by $\sigma^2$. The square root of the variance, $\sigma$, is called the standard deviation.

*Theorems*

1. $\quad \sigma^2_{cX} = c^2 \sigma^2_X$
2. $\quad \sigma^2_{c+X} = \sigma^2_X$
3. $\quad \sigma^2_{aX+b} = a^2 \sigma^2_X$

UtahState
UNIVERSITY