

## The concept of Dispersion

1

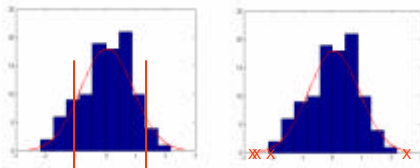
## The concept

- Law of large numbers tells us that the average result will  $\rightarrow E(X)$
- What can dispersion can we expect around  $E(X)$
- Does it have a sense to use the concept of **dispersion**?

2

## The concept

- Given a set of data or a functional distribution we would like to have a number for comparing dispersions.



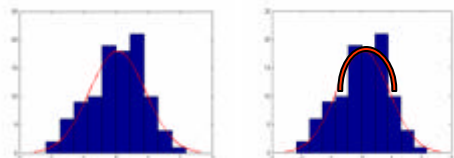
Margin that embraces 50% of all points

Extreme points

3

## The concept

- Given a set of data or a functional distribution we would like to have a number for comparing dispersions.



Mean distance to the center of mass of the distribution

Measure of curvature: Second degree polynomial

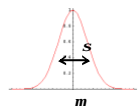
4

## The measures

- Standard deviation

$$m(X) = E(X) = \sum_{x \in X} xP(X = x)$$

$$s^2(X) = E((X - m)^2) = \sum_{x \in X} (x - m)^2 P(X = x)$$



- Mean absolute distance

$$m(X) = E(X) = \sum_{x \in X} xP(X = x)$$

$$MAD(X) = E(|X - m|) = \sum_{x \in X} |x - m| P(X = x)$$

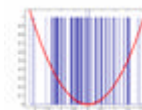
5

## The measures

- Standard deviation

- Justification:

- Samples **far** from the mean **increase quadratically** the value of the standard deviation.



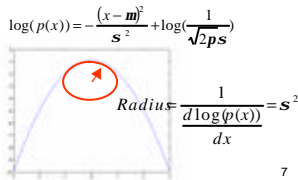
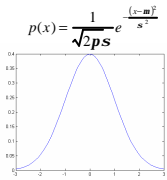
$$s^2(X) = E((X - m)^2) = \sum_{x \in X} (x - m)^2 P(X = x)$$

- Is the value of the **curvature** of the log gaussian.

6

## The measures

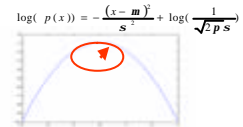
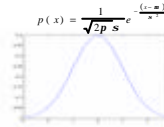
- Standard deviation
- Meaning of the **curvature of a plane curve**:
  - Radius of a circle that approximates locally the curve.



7

## The measures

- Standard deviation
- Meaning of the **curvature of a plane curve**:



- Certainty=low dispersion->Small radius->small variance
- Uncertainty=high dispersion->big radius->big variance

$$\text{Radius} = \frac{1}{\frac{d^2 \log(p(x))}{dx^2}} = s^2$$

8

## The measures

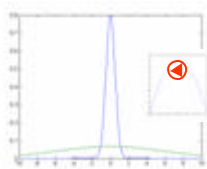
- Standard deviation
  - Certainty=low dispersion->Small radius->small variance
  - Uncertainty=high dispersion->big radius->big variance

$$\text{Radius} = \frac{1}{\frac{d^2 \log(p(x))}{dx^2}} = s^2$$

$$s = 0.5$$

$$s = 6$$

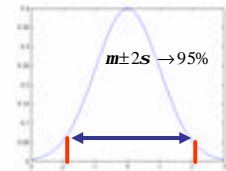
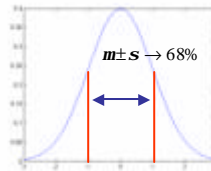
Cramer-Rao and statistics



9

## Information about Gaussian distributions.

- % of the cases for a given dispersion



$$m \pm 3s \rightarrow 99\%$$

Mensa

10

## The measures

- Mean absolute distance

$$m(X) = E(X) = \sum_{x \in X} xP(X=x)$$

$$MAD(X) = E(|X - m|) = \sum_{x \in X} |X - m|P(X=x)$$



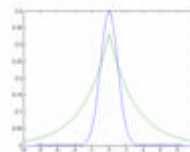
Samples far from the mean increase linearly the value of the standard deviation.  
Note: outliers do not count as much as in the standard deviation.

11

## The measures

- Mean absolute distance
  - Uses: measure of dispersion for distributions with longer tail than gaussians.

$$MAD(X) = E(|X - m|) = \sum_{x \in X} |X - m|P(X=x)$$

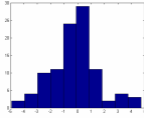


**Laplace vs. Gaussian**  
Note: That high values are much more probable in the case of an exponential random variable

12

## The measures

- Mean absolute distance
  - Outliers do not change so much the value of the estimation.

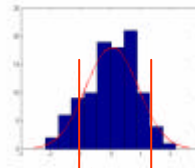


$$MAD(X) = E(|X - m|) = \sum_{x \in X} |x - m| P(X = x)$$

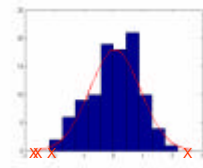
13

## The concept

- Given a set of data or a functional distribution we would like to have a number for comparing dispersions.



Margin that embraces 50% of all points



Extreme points

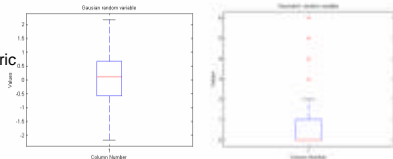
14

## The Box Plot

- Gives information about the dispersion summarizing the information about:
  - The interquartile size.  $x_{25}$ ,  $x_{75}$
  - Median
  - Sample nearest to the 1.5 times the interquartile margin
  - The outliers (points  $>1.5$  IQR)

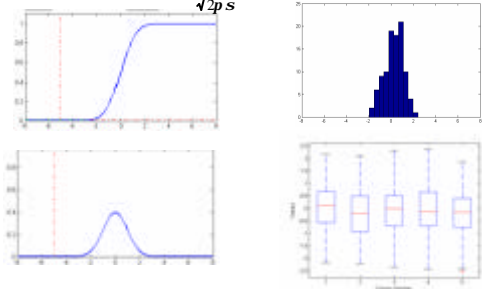
- Examples:  
100 points

Gaussian vs. geometric.



## Description of a random variable

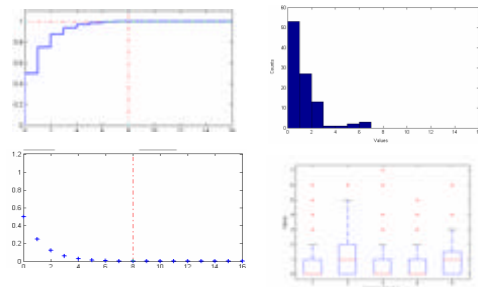
- Gaussian  $p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-m)^2}{2\sigma^2}}$



16

## Description of a random variable

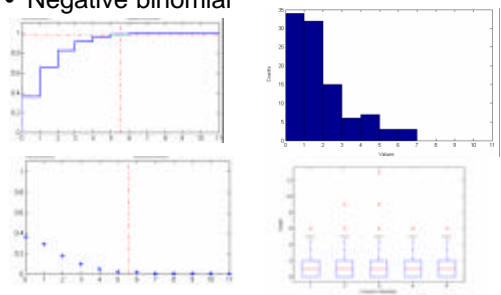
- Geometric  $P(X = i) = p^{i-1}(1-p)$  for  $i = 1, 2, 3, \dots$



17

## Description of a random variable

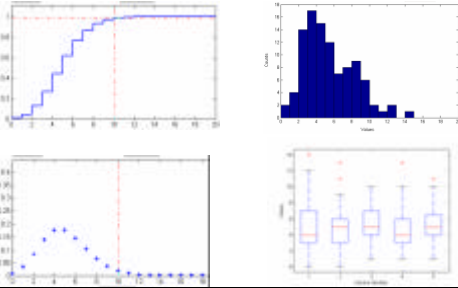
- Negative binomial



18

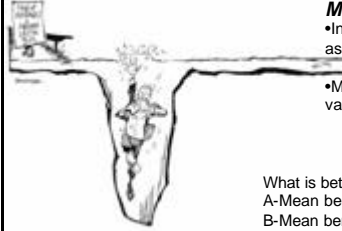
# Description of a random variable

- Poisson



# Some intuitions

- The *flaw* of averages



### Markowitz's Idea:

- Introduce variability when assessing the value of an asset.
- Maximize mean, minimizing variance.

What is better?:  
 A-Mean benefit of 500 plus minus 400  
 B-Mean benefit of 200 plus minus 50

# The flaw of averages

- An investment problem

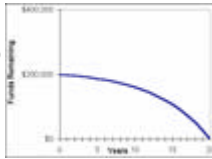
- Suppose you want your \$200,000 retirement fund invested in the Standard & Poor's 500 index to last 20 years. **How much can you withdraw per year?**
- The return of the S&P has varied over the years but has **averaged about 14** percent per year since 1952.
- If you do this you will be pleased to find that you can withdraw \$32,000 per year.

$$A(1+r)^{20} - \sum_{k=0}^{19} x(1+r)^k = 0$$

$$A = 200,000$$

$$r = 14\%$$

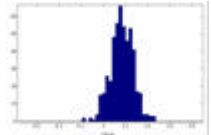
$$x = \frac{rA(1+r)^{20}}{(1+r)^{20} - 1}$$



# The flaw of averages

- Model of the return:

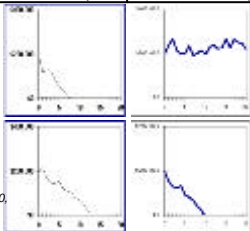
- $r$  % with probability  $p$
- *Histogram of the return*
- *Average value  $r$ , but can fluctuate.*
  - Sometimes gives benefits
  - Sometimes losses
- Note that each month a fixed quantity is subtracted independently of  $r$



# The flaw of averages

- Simulations on real data:

Start: 1973 Avg. Return 14% Tanks in 8 yrs.	Start: 1974 Avg. Return 15.4% Goes the distance.
Start: 1975 Avg. Return 15.4% Tanks in 13 yrs.	Start: 1976 Avg. Return 15.3% Tanks in 10 yrs



### The Flaw of Averages

BY SAM SAVAGE  
 Published Sunday, October 8, 2000  
 in the San Jose Mercury News

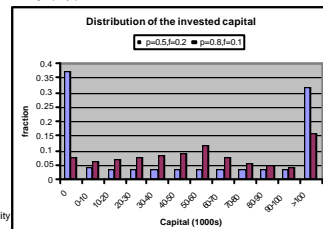
# The flaw of averages

- Model of the return:

- $r$  % with probability  $p$
- $(1+f)r$  % with probability  $(1-p)/2$
- $(1-f)r$  % with probability  $(1-p)/2$

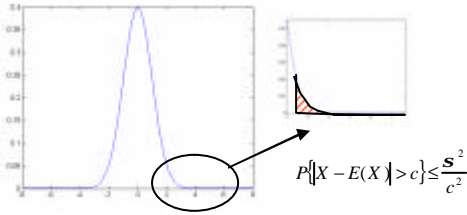
- Simulation.

- 4,000,000 runs.



## Tchebychev Inequality

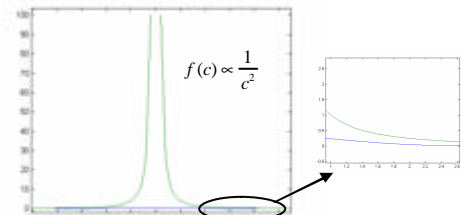
- A bound on the upper probability.



25

## Tchebychev Inequality

- Geometrical meaning of  $P\{|X - E(X)| > c\} \leq \frac{s^2}{c^2}$



26

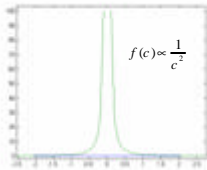
## Tchebychev Inequality

- Geometrical meaning of  $P\{|X - E(X)| > c\} \leq \frac{s^2}{c^2}$

- What happens with?

$$P(X=x) = \frac{1}{p(1+x^2)}$$

$$P(X=x) = \frac{1}{x^a}$$



Only valid on distributions that have finite variance!

27

## Tchebychev Inequality

- For a given probability distribution of a random variable X, with finite variance we have:

$$P\{|X - E(X)| > ks\} \leq \frac{1}{k^2}$$

- for any  $k > 0$  or equivalently

$$P\{|X - E(X)| > c\} \leq \frac{s^2}{c^2}$$

28

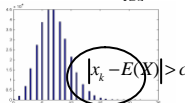
## Tchebychev Inequality

- Proof

- Given an ordered set  $\{x_1, x_2, x_3, x_4, x_5, \dots\}$
- We define the subset  $A = \{k \mid |x_k - E(X)| > c\}$
- then

$$s^2 = \sum_k (x_k - E(X))^2 p_k \geq \sum_{k \in A} (x_k - E(X))^2 p_k \geq c^2 \sum_{k \in A} p_k$$

$$s^2 \geq c^2 P\{|X - E(X)| > c\}$$



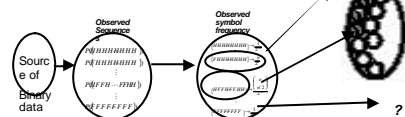
29

## Tchebychev Inequality

- Note that the inequality can be rough, and highly inexact for high values of c

- Uses:

- Information theory, bounds on probabilities and events highly improbable.



30

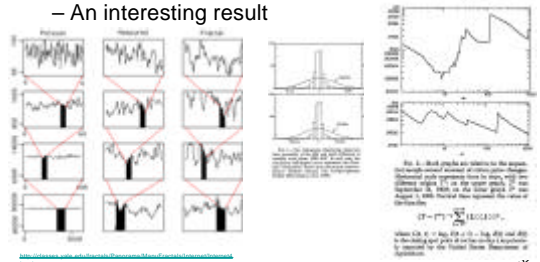
## Random variables without variance

- Family known as
  - Pareto Stable or Mandelbrot Levy
- Models:
  - Internet traffic
  - Processes in unix systems
  - Speculative prices/Pluviometric Data

31

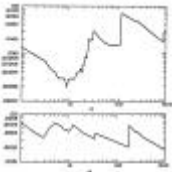
## Speculative Prices

- Mandelbrot's paper on long tail densities
  - An interesting result



## Syndrom of infinite variance

- Mandelbrot's paper on long tail densities



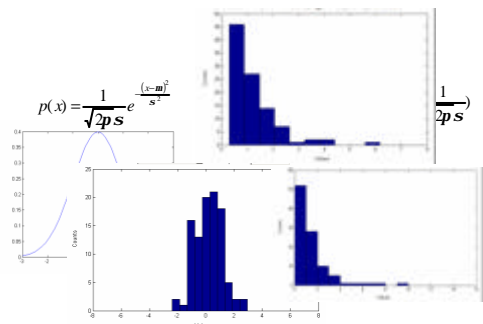
$$m(X) = E(X) = \sum_{x \in X} xP(X = x)$$

$$s^2(X) = E((X - m)^2) = \sum_{x \in X} (x - m)^2 P(X = x)$$

$$P(X = x) = \frac{1}{x^\alpha}$$

Note that for  $\alpha < 2$  the sum diverges

33



34