

CS 147:
Computer Systems Performance Analysis
Review of Statistics

2015-06-15 CS147

CS 147:
Computer Systems Performance Analysis
Review of Statistics

Introduction to Statistics

- ▶ Concentration on applied statistics
- ▶ Especially those useful in measurement
- ▶ Today's lecture will cover 15 basic concepts
- ▶ You should already be familiar with them

2015-06-15 CS147
└ 15 Concepts

└ Introduction to Statistics

Introduction to Statistics

- Concentration on applied statistics
- Especially those useful in measurement
- Today's lecture will cover 15 basic concepts
- You should already be familiar with them

1. Independent Events

- ▶ Occurrence of one event doesn't affect probability of other
- ▶ Examples:
 - ▶ Coin flips
 - ▶ Inputs from separate users
 - ▶ "Unrelated" traffic accidents

2015-06-15 CS147
└ 15 Concepts
└ Independent Events
└ 1. Independent Events

1. Independent Events

- Occurrence of one event doesn't affect probability of other
- Examples:
 - Coin flips
 - Inputs from separate users
 - "Unrelated" traffic accidents

1. Independent Events

- ▶ Occurrence of one event doesn't affect probability of other
- ▶ Examples:
 - ▶ Coin flips
 - ▶ Inputs from separate users
 - ▶ "Unrelated" traffic accidents
- ▶ What about second basketball free throw after the player misses the first?

2015-06-15
CS147
└ 15 Concepts
└ Independent Events
└ 1. Independent Events

1. Independent Events

- Occurrence of one event doesn't affect probability of other
- Examples:
 - Coin flips
 - Inputs from separate users
 - "Unrelated" traffic accidents
- What about second basketball free throw after the player misses the first?

2. Random Variable

- ▶ Variable that takes values probabilistically
- ▶ Variable usually denoted by capital letters, particular values by lowercase
- ▶ Examples:
 - ▶ Number shown on dice
 - ▶ Network delay
 - ▶ CS 70 attendance
- ▶ What about disk seek time?

2015-06-15
CS147
└─ 15 Concepts
 └─ Random Variable
 └─ 2. Random Variable

2. Random Variable

- Variable that takes values probabilistically
- Variable usually denoted by capital letters, particular values by lowercase
- Examples:
 - Number shown on dice
 - Network delay
 - CS 70 attendance
- What about disk seek time?

3. Cumulative Distribution Function (CDF)

- ▶ Maps a value a to probability that the outcome is less than or equal to a :

$$F_x(a) = P(x \leq a)$$

- ▶ Valid for discrete and continuous variables
- ▶ Monotonically increasing
- ▶ Easy to specify, calculate, measure

2015-06-15

CS147

└ 15 Concepts

└ CDF

└ 3. Cumulative Distribution Function (CDF)

3. Cumulative Distribution Function (CDF)

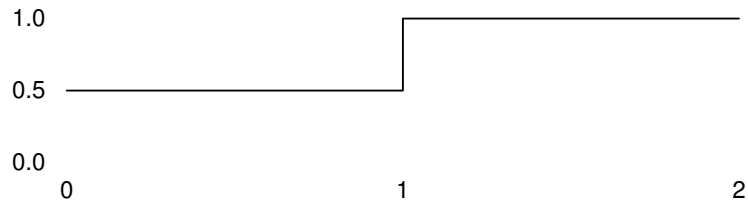
- Maps a value a to probability that the outcome is less than or equal to a :

$$F_x(a) = P(x \leq a)$$

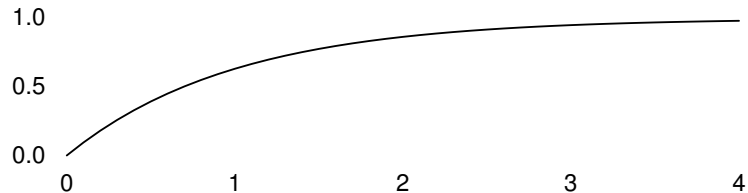
- Valid for discrete and continuous variables
- Monotonically increasing
- Easy to specify, calculate, measure

CDF Examples

- ▶ Coin flip ($T = 0, H = 1$):



- ▶ Exponential packet interarrival times:



2015-06-15

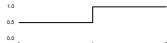
CS147

└ 15 Concepts

└ CDF

└ CDF Examples

CDF Examples

• Coin flip ($T = 0, H = 1$):

• Exponential packet interarrival times:



4. Probability Density Function (pdf)

- ▶ Derivative of (continuous) CDF:

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

- ▶ Usable to find probability of a range:

$$\begin{aligned} P(x_1 < x \leq x_2) &= F(x_2) - F(x_1) \\ &= \int_{x_1}^{x_2} f(x) dx \end{aligned}$$

2015-06-15

CS147

└ 15 Concepts

└ pdf

└ 4. Probability Density Function (pdf)

4. Probability Density Function (pdf)

- Derivative of (continuous) CDF:

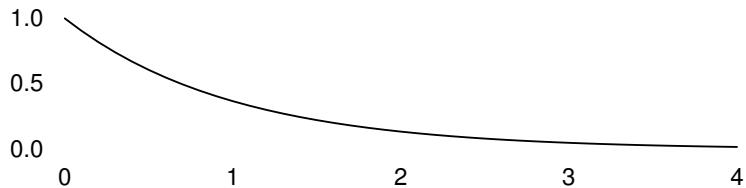
$$f(x) = \frac{dF(x)}{dx}$$

- Usable to find probability of a range:

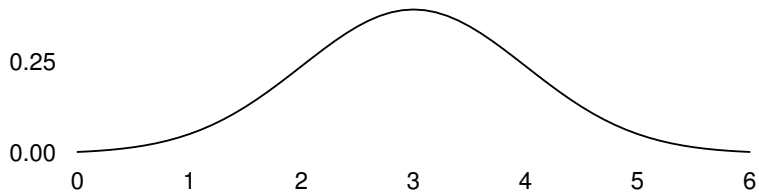
$$\begin{aligned} P(x_1 < x \leq x_2) &= F(x_2) - F(x_1) \\ &= \int_{x_1}^{x_2} f(x) dx \end{aligned}$$

Examples of pdf

- ▶ Exponential interarrival times:



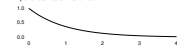
- ▶ Gaussian (normal) distribution:



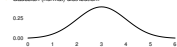
2015-06-15
CS147
└ 15 Concepts
└ pdf
└ Examples of pdf

Examples of pdf

• Exponential interarrival times:



• Gaussian (normal) distribution:



5. Probability Mass Function (pmf)

- ▶ CDF not differentiable for discrete random variables
- ▶ pmf serves as replacement: $f(x_i) = p_i$ where p_i is the probability that x will take on the value x_i :

$$\begin{aligned}P(x_1 < x \leq x_2) &= F(x_2) - F(x_1) \\ &= \sum_{x_1 < x \leq x_2} p_i\end{aligned}$$

2015-06-15

CS147

└ 15 Concepts

└ pmf

└ 5. Probability Mass Function (pmf)

5. Probability Mass Function (pmf)

- CDF not differentiable for discrete random variables
- pmf serves as replacement: $f(x_i) = p_i$ where p_i is the probability that x will take on the value x_i :

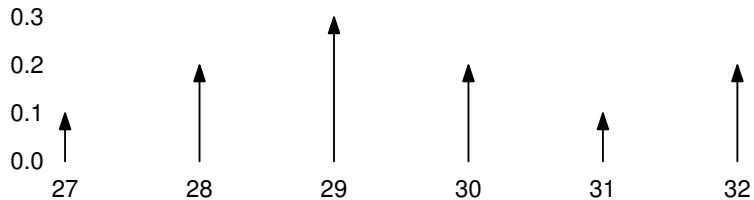
$$\begin{aligned}P(x_1 < x \leq x_2) &= F(x_2) - F(x_1) \\ &= \sum_{x_1 < x \leq x_2} p_i\end{aligned}$$

Examples of pmf

► Coin flip:



► Typical CS grad class size:



2015-06-15

CS147

└ 15 Concepts

└ pmf

└ Examples of pmf

Examples of pmf

► Coin flip:

1.0

0.5

0.0

0

1

► Typical CS grad class size:

0.3

0.2

0.1

0.0

27

28

29

30

31

32

6. Expected Value (Mean)

- ▶ Mean:

$$\mu = E(x) = \sum_{i=1}^n p_i x_i = \int_{-\infty}^{\infty} x f(x) dx$$

- ▶ Summation if discrete
- ▶ Integration if continuous

2015-06-15
CS147
└ 15 Concepts
└ Mean
└ 6. Expected Value (Mean)

6. Expected Value (Mean)

- Mean: $\mu = E(x) = \sum_{i=1}^n p_i x_i = \int_{-\infty}^{\infty} x f(x) dx$
- Summation if discrete
- Integration if continuous

7. Variance

- ▶ Variance:

$$\begin{aligned}\text{Var}(x) = E[(x - \mu)^2] &= \sum_{i=1}^n p_i (x_i - \mu)^2 \\ &= \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx\end{aligned}$$

- ▶ Often easier to calculate equivalent $E(x^2) - E(x)^2$
- ▶ Usually denoted σ^2 ; square root σ is called *standard deviation*

2015-06-15
 CS147
 └ 15 Concepts
 └ └ Variance
 └ └ └ 7. Variance

7. Variance

• Variance:

$$\begin{aligned}\text{Var}(x) = E[(x - \mu)^2] &= \sum_{i=1}^n p_i (x_i - \mu)^2 \\ &= \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx\end{aligned}$$

- Often easier to calculate equivalent $E(x^2) - E(x)^2$
- Usually denoted σ^2 ; square root σ is called *standard deviation*

8. Coefficient of Variation (C.O.V. or C.V.)

- ▶ Ratio of standard deviation to mean:

$$\text{C.V.} = \frac{\sigma}{\mu}$$

- ▶ Indicates how well mean represents the variable

2015-06-15

CS147

└ 15 Concepts

└ Coefficient of Variation

└ 8. Coefficient of Variation (C.O.V. or C.V.)

8. Coefficient of Variation (C.O.V. or C.V.)

- Ratio of standard deviation to mean:

$$\text{C.V.} = \frac{\sigma}{\mu}$$

- Indicates how well mean represents the variable

9. Covariance

- ▶ Given x, y with means μ_x and μ_y , their covariance is:

$$\begin{aligned}\text{Cov}(x, y) = \sigma_{xy}^2 &= E[(x - \mu_x)(y - \mu_y)] \\ &= E(xy) - E(x)E(y)\end{aligned}$$

- ▶ Two typos on p.181 of book
- ▶ High covariance implies y departs from mean whenever x does

2015-06-15
 CS147
 └ 15 Concepts
 └ Covariance
 └ 9. Covariance

9. Covariance

- Given x, y with means μ_x and μ_y , their covariance is:

$$\begin{aligned}\text{Cov}(x, y) = \sigma_{xy}^2 &= E[(x - \mu_x)(y - \mu_y)] \\ &= E(xy) - E(x)E(y)\end{aligned}$$
- Two typos on p.181 of book
- High covariance implies y departs from mean whenever x does

Covariance (cont'd)

- ▶ For independent variables, $E(xy) = E(x)E(y)$ so $\text{Cov}(x, y) = 0$
- ▶ Reverse isn't true: $\text{Cov}(x, y) = 0$ does **NOT** imply independence
- ▶ If $y = x$, covariance reduces to variance

2015-06-15 CS147
└ 15 Concepts
└ Covariance
└ Covariance (cont'd)

Covariance (cont'd)

- For independent variables, $E(xy) = E(x)E(y)$ so $\text{Cov}(x, y) = 0$
- Reverse isn't true: $\text{Cov}(x, y) = 0$ does **NOT** imply independence
- If $y = x$, covariance reduces to variance

10. Correlation Coefficient

- ▶ Normalized covariance:

$$\text{Correlation}(x, y) = \rho_{xy} = \frac{\sigma_{xy}^2}{\sigma_x \sigma_y}$$

- ▶ Always lies between -1 and 1
- ▶ Correlation of 1 $\Rightarrow x \sim y$, -1 $\Rightarrow x \sim \frac{1}{y}$

2015-06-15

CS147

└ 15 Concepts

└ Correlation Coefficient

└ 10. Correlation Coefficient

10. Correlation Coefficient

• Normalized covariance:

$$\text{Correlation}(x, y) = \rho_{xy} = \frac{\sigma_{xy}^2}{\sigma_x \sigma_y}$$

• Always lies between -1 and 1

• Correlation of 1 $\Rightarrow x \sim y$, -1 $\Rightarrow x \sim \frac{1}{y}$

11. Mean and Variance of Sums

- ▶ For any random variables,

$$E(a_1x_1 + \cdots + a_kx_k) = a_1E(x_1) + \cdots + a_kE(x_k)$$

- ▶ For independent variables,

$$\text{Var}(a_1x_1 + \cdots + a_kx_k) = a_1^2\text{Var}(x_1) + \cdots + a_k^2\text{Var}(x_k)$$

2015-06-15

CS147

└ 15 Concepts

└ Mean and Variance of Sums

└ 11. Mean and Variance of Sums

11. Mean and Variance of Sums

- For any random variables,

$$E(a_1x_1 + \cdots + a_kx_k) = a_1E(x_1) + \cdots + a_kE(x_k)$$

- For independent variables,

$$\text{Var}(a_1x_1 + \cdots + a_kx_k) = a_1^2\text{Var}(x_1) + \cdots + a_k^2\text{Var}(x_k)$$

12. Quantile

- ▶ x value at which CDF takes a value α is called α -quantile or 100α -percentile, denoted by x_α

$$P(x \leq x_\alpha) = F(x_\alpha) = \alpha$$

- ▶ If 90th-percentile score on GRE was 1500, then 90% of population got 1500 or less

2015-06-15
CS147
└ 15 Concepts
└ Quantile
└ 12. Quantile

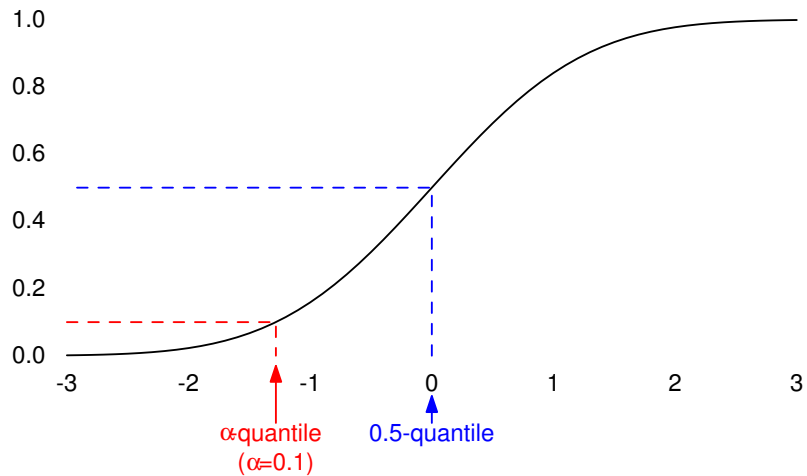
12. Quantile

• x value at which CDF takes a value α is called α -quantile or 100α -percentile, denoted by x_α .

$$P(x \leq x_\alpha) = F(x_\alpha) = \alpha$$

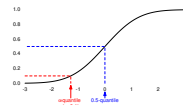
• If 90th-percentile score on GRE was 1500, then 90% of population got 1500 or less

Quantile Example



2015-06-15
CS147
└ 15 Concepts
└ Quantile
└ Quantile Example

Quantile Example



13. Median

- ▶ 50th percentile (0.5-quantile) of a random variable
- ▶ Alternative to mean
- ▶ By definition, 50% of population is below median, 50% above
 - ▶ Lots of bad (good) drivers
 - ▶ Lots of smart (stupid) people

2015-06-15
CS147
└─ 15 Concepts
 └─ Median
 └─ 13. Median

13. Median

- 50th percentile (0.5-quantile) of a random variable
- Alternative to mean
- By definition, 50% of population is below median, 50% above
 - Lots of bad (good) drivers
 - Lots of smart (stupid) people

14. Mode

- ▶ Most likely value, i.e., x_i with highest probability p_i , or x at which pdf/pmf is maximum
- ▶ Not necessarily defined (e.g., tie)
- ▶ Some distributions are bi-modal (e.g., human height has one mode for males and one for females)

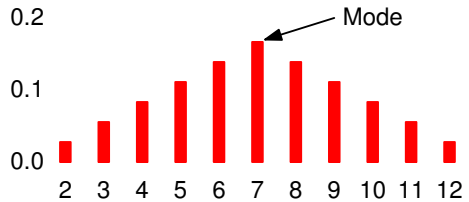
2015-06-15
CS147
└─ 15 Concepts
 └─ Mode
 └─ 14. Mode

14. Mode

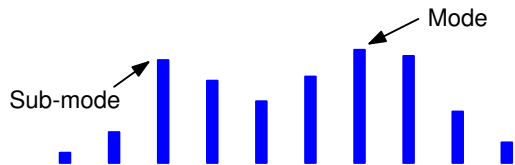
- Most likely value, i.e., x with highest probability p , or x at which pdf/pmf is maximum
- Not necessarily defined (e.g., tie)
- Some distributions are bi-modal (e.g., human height has one mode for males and one for females)

Examples of Mode

► Dice throws:



► Adult human weight:



2015-06-15

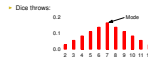
CS147

└ 15 Concepts

└ Mode

└ Examples of Mode

Examples of Mode



Adult human weight:



15. Normal (Gaussian) Distribution

- ▶ Most common distribution in data analysis
- ▶ pdf is:

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

- ▶ $-\infty \leq x \leq +\infty$
- ▶ Mean is μ , standard deviation σ

2015-06-15

CS147

└ 15 Concepts

└ Normal Distribution

└ 15. Normal (Gaussian) Distribution

15. Normal (Gaussian) Distribution

- Most common distribution in data analysis
- pdf is:

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

- $-\infty \leq x \leq +\infty$
- Mean is μ , standard deviation σ

Notation for Gaussian Distributions

- ▶ Often denoted $N(\mu, \sigma)$
- ▶ Unit normal is $N(0, 1)$
- ▶ If x has $N(\mu, \sigma)$, $\frac{x-\mu}{\sigma}$ has $N(0, 1)$
- ▶ The α -quantile of unit normal $z \sim N(0, 1)$ is denoted z_α so that

$$\left\{ P\left(\frac{X - \mu}{\sigma} \leq z_\alpha\right) \right\} = \{P(X) \leq \mu + z_\alpha \sigma\} = \alpha$$

2015-06-15

CS147

└ 15 Concepts

└ Normal Distribution

└ Notation for Gaussian Distributions

Notation for Gaussian Distributions

- Often denoted $N(\mu, \sigma)$
- Unit normal is $N(0, 1)$
- If x has $N(\mu, \sigma)$, $\frac{x-\mu}{\sigma}$ has $N(0, 1)$
- The α -quantile of unit normal $z \sim N(0, 1)$ is denoted z_α , so that

$$\left\{ P\left(\frac{X - \mu}{\sigma} \leq z_\alpha\right) \right\} = \{P(X) \leq \mu + z_\alpha \sigma\} = \alpha$$

Why Is Gaussian So Popular?

- ▶ We've seen that if $x_i \sim N(\mu_i, \alpha_i)$ and all x_i independent, then $\sum \alpha_i x_i$ is normal with mean $\sum \alpha_i \mu_i$ and variance $\sigma^2 = \sum \alpha_i^2 \sigma_i^2$
- ▶ Sum of large number of independent observations from any distribution is itself normal (Central Limit Theorem)
 \Rightarrow Experimental errors can be modeled as normal distribution.

2015-06-15

CS147

└ 15 Concepts

└ Normal Distribution

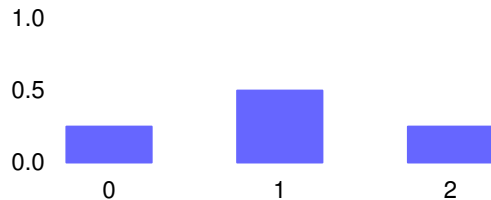
└ Why Is Gaussian So Popular?

Why Is Gaussian So Popular?

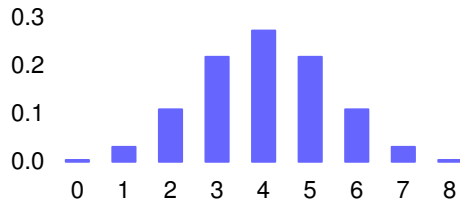
- We've seen that if $x_i \sim N(\mu_i, \alpha_i)$ and all x_i independent, then $\sum \alpha_i x_i$ is normal with mean $\sum \alpha_i \mu_i$ and variance $\sigma^2 = \sum \alpha_i^2 \sigma_i^2$
- Sum of large number of independent observations from any distribution is itself normal (Central Limit Theorem)
 \Rightarrow Experimental errors can be modeled as normal distribution.

Central Limit Theorem

- ▶ Sum of 2 coin flips ($H=1, T=0$):



- ▶ Sum of 8 coin flips:



2015-06-15

CS147

└ 15 Concepts

└ Normal Distribution

└ Central Limit Theorem

Central Limit Theorem

▶ Sum of 2 coin flips ($H=1, T=0$):

▶ Sum of 8 coin flips:

