

# CS 584 – MACHINE LEARNING

## TOPIC: PROBABILITY THEORY



**Mustafa Bilgic**



<http://www.cs.iit.edu/~mbilgic>



<https://twitter.com/bilgicm>

# MOTIVATION

- Learning
  - Statistics, expectations, etc.
  - E.g., decision trees, naïve Bayes, logistic regression, ...
- Evaluation
  - Statistics, expectations, variance, etc.
  - E.g., expected error using a sample

# SOME QUESTIONS

- Given a domain with  $n$  variables,  $X_1, X_2, \dots, X_n$ , each of which has  $v_1, v_2, \dots, v_n$  possible values, what is the size of the instance space?
- Given a sample dataset and a hypothesis  $h$ , calculate the mean, variance, and 95% confidence interval for the error rate of  $h$
- Type I error has cost  $c_1$  and type II error has cost  $c_2$ . Correct decisions have no cost. The probability of the object belonging to Positive class is  $p$ . Should it be classified as Positive or Negative?
- Given  $P(\text{Symptoms} | \text{Diseases})$ ,  $P(\text{Diseases})$  and  $P(\text{Symptoms})$ , calculate  $P(\text{Diseases} | \text{Symptoms})$
- Given a probability distribution  $p_1, p_2, \dots, p_k$ , calculate its entropy

# RANDOM VARIABLES

- Pick variables of interest
  - Medical diagnosis
    - Age, gender, weight, temperature, LT1, LT2, ...
  - Loan application
    - Income, wealth, payment history, ...
- Every variable has a domain
  - Binary (True/False)
  - Categorical
  - Real-valued
- Possible world
  - An assignment to all variables of interest

# PROBABILITY MODEL

- A **probability model** associates a numerical probability  $P(w)$  with each possible world  $w$ 
  - $P(w)$  sums to 1 over all possible worlds
- An **event** is the set of possible worlds where a given predicate is true
  - Roll two dice
    - The possible worlds are (1,1), (1,2), ..., (6,6); 36 possible worlds
    - Predicate = two dice sum to 10
    - Event = {(4,6), (5,5), (6,4)}
  - Toothache and cavity
    - Four possible worlds:  $(t, c), (t, \sim c), (\sim t, c), (\sim t, \sim c)$
    - Some worlds are more likely than others
    - Predicate can be anything about these variables:  $t \wedge c, t, t \vee \sim c,$

# AXIOMS OF PROBABILITY

1. The probability  $P(a)$  of a proposition  $a$  is a real number between 0 and 1

2.  $P(\text{true}) = 1$ ,  $P(\text{false}) = 0$

3.  $P(a \vee b) = P(a) + P(b) - P(a \wedge b)$



## $P(\neg a)$

- $P(a \vee \neg a) = P(a) + P(\neg a) - P(a \wedge \neg a)$

← Axiom #3

- $P(\text{true}) = P(a) + P(\neg a) - P(\text{false})$

- $1 = P(a) + P(\neg a) - 0$

← Axiom #2

- $P(\neg a) = 1 - P(a)$

- Intuitive explanation:

- The probability of all possible worlds is 1
- Either  $a$  or  $\neg a$  holds in one world
- The worlds that  $a$  holds and the worlds that  $\neg a$  holds are mutually exclusive and exhaustive

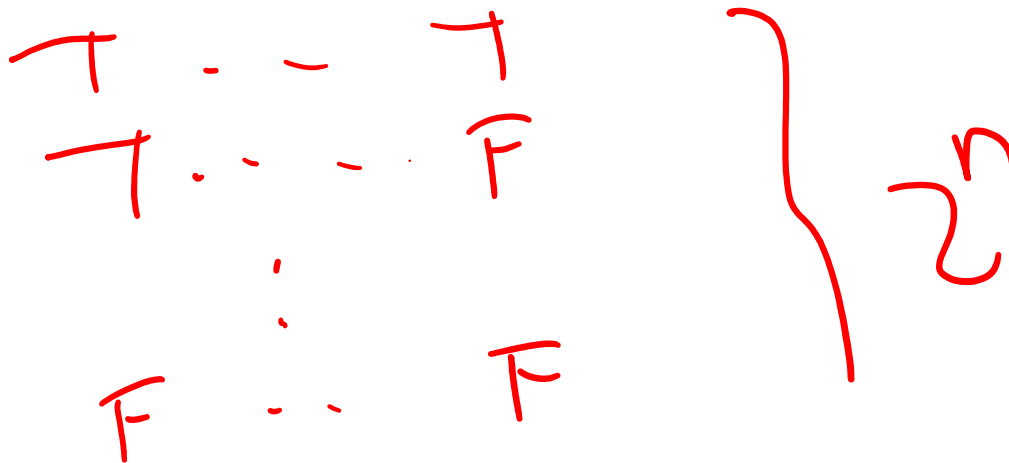
# RANDOM VARIABLES – NOTATION

- Capital:  $X$ : variable
- Lowercase:  $x$ : a particular value of  $X$
- $\text{Val}(X)$ : the set of values  $X$  can take
- Bold Capital:  $\mathbf{X}$ : a set of variables
- Bold lowercase:  $\mathbf{x}$ : an assignment to all variables in  $\mathbf{X}$
- $P(X=x)$  will be shortened as  $P(x)$
- $P(X=x \cap Y=y)$  will be shortened as  $P(x,y)$



# JOINT DISTRIBUTION

- We have  $n$  random variables,  $X_1, X_2, \dots, X_n$
- We are interested in the probability of a possible world, where
  - $X_1=\text{low}, X_2=\text{red}, \dots, X_n=\text{circle}$
- $P(X_1, X_2, \dots, X_n)$  associates a probability for each possible world  $\equiv$  the **joint distribution**
- How many entries are there, if we assume the variables are all binary?



# TOOTHACHE EXAMPLE

Feeling	X-Ray	P(F,X)
toothache	cavity	0.15
toothache	$\neg$ cavity	0.10
$\neg$ toothache	cavity	0.05
$\neg$ toothache	$\neg$ cavity	0.70

# MARGINALIZATION

- Given a distribution over  $n$  variables, you can calculate the distribution over any subset of the variables by summing out the irrelevant ones
- For example
  - Given  $P(A, B, C, D)$
  - Calculate
    - $P(A) = \sum_B \sum_C \sum_D P(A, B, C, D)$
    - $P(A, C) = \sum_B \sum_D P(A, B, C, D)$
    - ... (any subset)

## LET'S ANSWER A FEW QUERIES

Feeling	X-Ray	P(F,X)
toothache	cavity	0.15
toothache	$\neg$ cavity	0.10
$\neg$ toothache	cavity	0.05
$\neg$ toothache	$\neg$ cavity	0.70

- $P(\text{cavity}) = ?$   $0.15 + 0.05 = 0.20$
- $P(\neg \text{cavity}) = ?$   $0.80$
- $P(\text{toothache}) = ?$   $0.15 + 0.10 = 0.25$
- $P(\neg \text{toothache}) = ?$   $0.75$

# CONDITIONAL DISTRIBUTION

- $P(A, B, C \mid D, E, F, G) = \frac{P(A, B, C, D, E, F, G)}{P(D, E, F, G)}$

$$P(MP \mid I=l)$$
$$\langle \overset{\sim}{0.6}, \overset{\sim}{0.4} \rangle$$
$$P(MP \mid I=h)$$
$$\langle \overset{\sim}{0.05}, \overset{\sim}{0.95} \rangle$$

# LET'S ANSWER A FEW QUERIES

Feeling	X-Ray	P(F,X)
toothache	cavity	0.15
toothache	¬cavity	0.10
¬toothache	cavity	0.05
¬toothache	¬cavity	0.70

- $P(\text{cavity} \mid \text{toothache}) = ?$
- $P(\text{cavity} \mid \neg \text{toothache}) = ?$
- $P(\neg \text{cavity} \mid \text{toothache}) = ?$
- $P(\neg \text{cavity} \mid \neg \text{toothache}) = ?$
- $P(\text{toothache} \mid \text{cavity}) = ?$
- $P(\neg \text{toothache} \mid \text{cavity}) = ?$
- $P(\text{toothache} \mid \neg \text{cavity}) = ?$
- $P(\neg \text{toothache} \mid \neg \text{cavity}) = ?$

$$\frac{P(\neg t, c)}{P(c)} = \frac{0.05}{0.20} = \frac{1}{4}$$

$$\frac{P(\neg t, \neg c)}{P(\neg c)} = \frac{0.70}{0.80} = \frac{7}{8}$$

# BAYES' RULE

- $P(B|A) = \frac{P(A|B)*P(B)}{P(A)}$
- Example use
  - $P(\text{cause} | \text{effect}) = P(\text{effect} | \text{cause}) * P(\text{cause}) / P(\text{effect})$
- Why is this useful?
  - Because in practice it is easier to get probabilities for  $P(\text{effect} | \text{cause})$  and  $P(\text{cause})$  than for  $P(\text{cause} | \text{effect})$ 
    - E.g.,  $P(\text{disease} | \text{symptoms}) = P(\text{symptoms} | \text{disease}) * P(\text{disease}) / P(\text{symptoms})$
    - It is easier to know what symptoms diseases cause. It is harder to diagnose a disease given symptoms

# BAYES RULE

- Can we compute  $P(\alpha|\beta)$  from  $P(\beta|\alpha)$ ?



# CLASS EXAMPLE

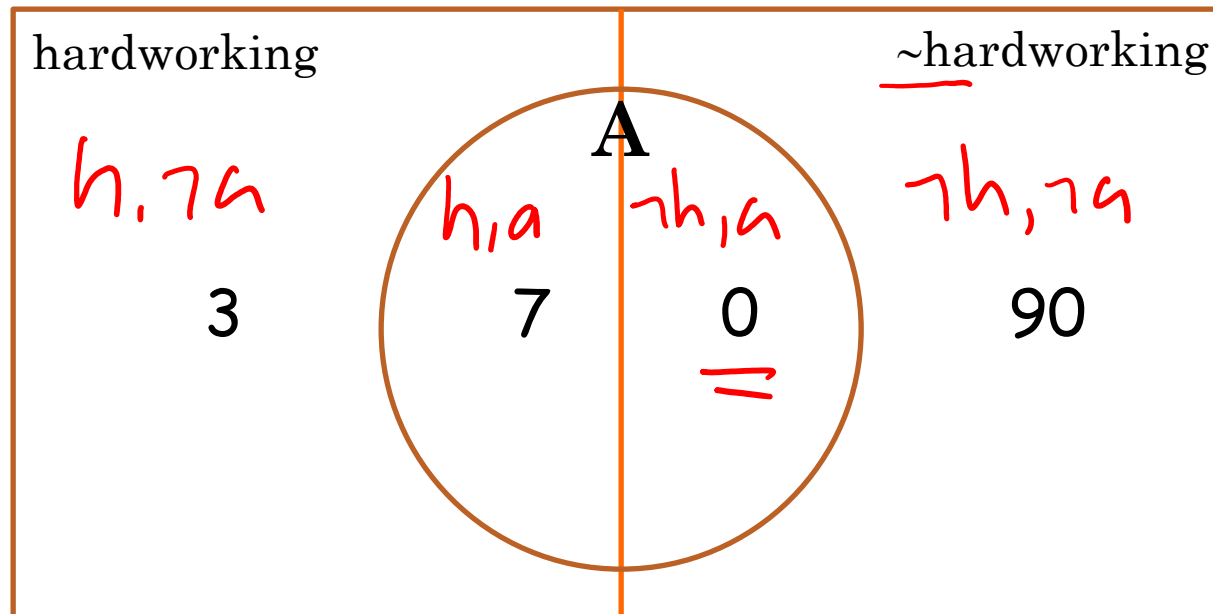
- Let's say there are 100 students in the class
- Let's say 10 of them work hard (h), 90 do not ( $\sim h$ )
- Probability of a randomly picked student being hardworking
  - $P(h) = 0.1$
- We are told that 70% of the hardworking students got an A.
  - $P(a | h) = 0.7$
  - 7 hardworking students got an A; 3 did not get an A.

○ What is  $P(h|a) = ?$

$$P(h|a) = \frac{P(a|h)P(h)}{P(a)}$$

$$p(h) = 0.10$$

## VERY HARD CLASS

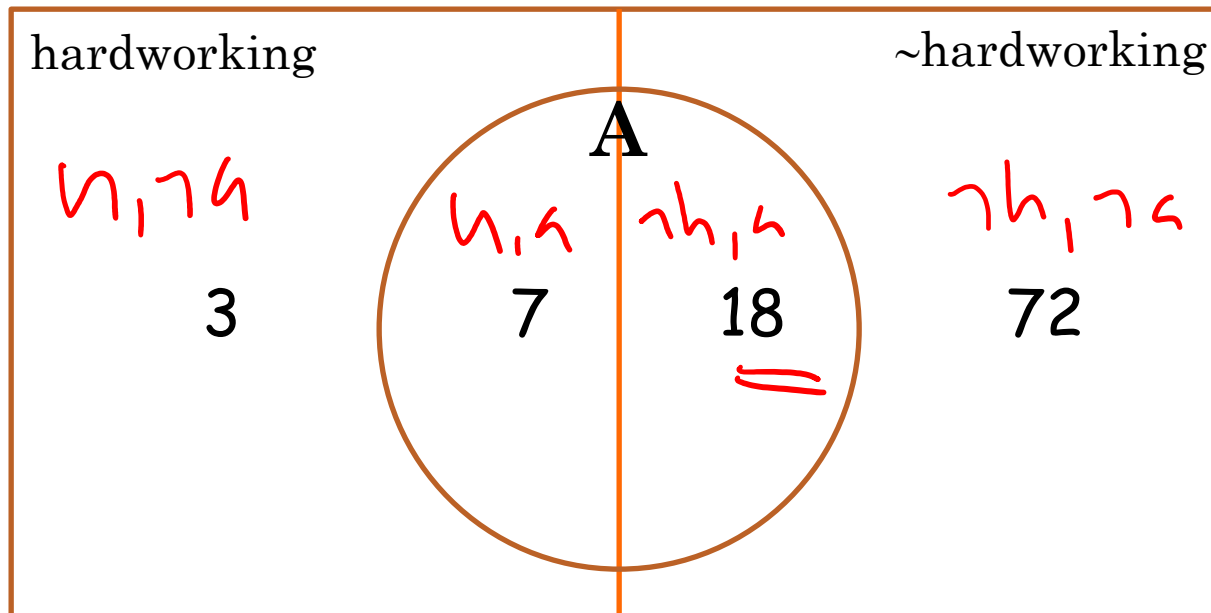


$$P(h | a) = ?$$

$$\frac{h, a}{a} = \frac{7}{7} = \underline{\underline{1}}$$

# MEDIUM HARD CLASS

$$p(h) = 0.1$$

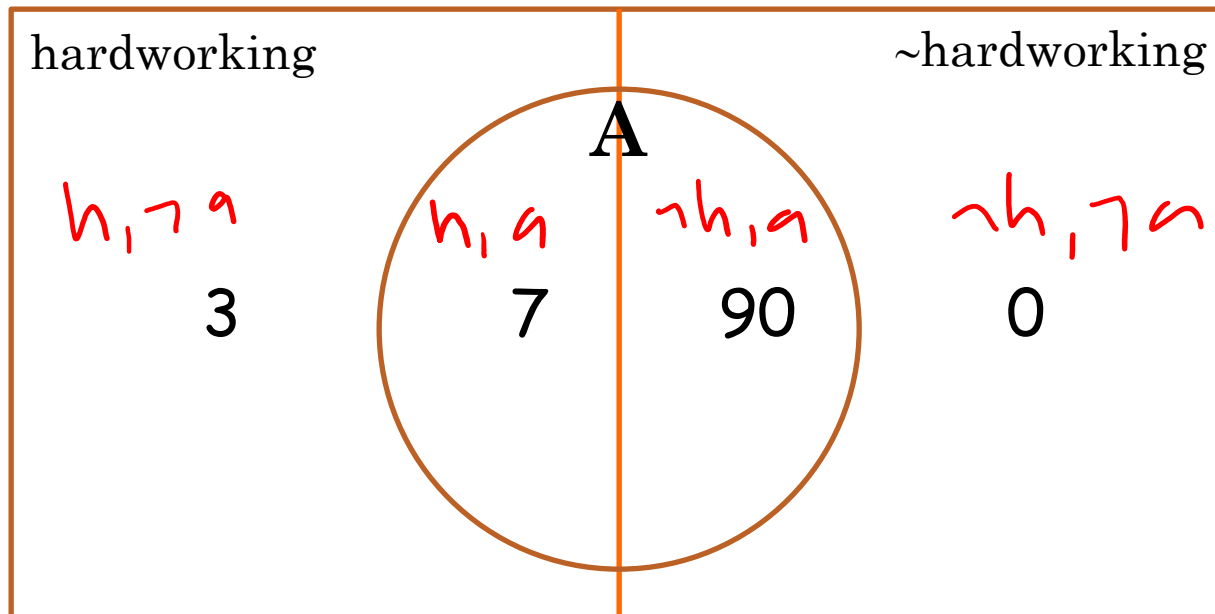


$$P(h | a) = ?$$

$$\frac{h, a}{a} = \frac{7}{25} = 0.28$$

# WEIRD CLASS

$$P(h) = 0.10$$



$$P(h | a) = ?$$

$$\frac{h, a}{a} = \frac{7}{97} \approx 0.07 \dots$$

# CHAIN RULE

- $P(X_1, X_2, X_3, \dots, X_k) =$ 
  - $P(X_1) P(X_2 | X_1) P(X_3 | X_1, X_2) \dots P(X_k | X_1, X_2, X_3, \dots, X_{k-1})$ 
    - or
  - $P(X_2) P(X_1 | X_2) P(X_3 | X_1, X_2) \dots P(X_k | X_1, X_2, X_3, \dots, X_{k-1})$ 
    - or
  - $P(X_2) P(X_3 | X_2) P(X_1 | X_3, X_2) \dots P(X_k | X_1, X_2, X_3, \dots, X_{k-1})$ 
    - or
  - Pick an order, then
    - $P(\text{first})P(\text{second} | \text{first})P(\text{third} | \text{first}, \text{second}) \dots P(\text{last} | \text{all\_previous})$

$$P(A, B, C, D) = P(C)P(A|C)P(D|C, A)P(B|C, A, D)$$

*C, A, D, B*

$$P(A, B, C, D) = \underbrace{P(B) P(D|B)}_{B, D, C, A} P(C|B, D) P(A|B, D, C)$$

$$= P(D, B) P(C|B, D) P(A|B, D, C)$$

$$= P(C, B, D) P(A|B, D, C)$$

$$= P(A, B, C, D)$$

# MARGINAL INDEPENDENCE

- An event  $\alpha$  is **independent** of event  $\beta$  in  $P$ , denoted as  $P \models \alpha \perp \beta$ , if
  - $P(\alpha \mid \beta) = P(\alpha)$ , or
  - $P(\beta) = 0$
- Proposition: A distribution  $P$  satisfies  $\alpha \perp \beta$  if and only if
  - $P(\alpha, \beta) = P(\alpha) P(\beta)$
  - *Can you prove it?*
- Corollary:  $\alpha \perp \beta$  implies  $\beta \perp \alpha$

# MARGINAL INDEPENDENCE

X	Y	P(X, Y)
t	t	0.18
t	f	0.42
f	t	0.12
f	f	0.28

Is  $X \perp Y$ ?



# CONDITIONAL INDEPENDENCE

- Two events are independent given another event
- An event  $\alpha$  is **independent** of event  $\beta$  given event  $\gamma$  in  $P$ , denoted as  $P \models (\alpha \perp \beta \mid \gamma)$ , if
  - $P(\alpha \mid \beta, \gamma) = P(\alpha \mid \gamma)$ , or
  - $P(\beta, \gamma) = 0$
- Proposition: A distribution  $P$  satisfies  $\alpha \perp \beta \mid \gamma$  if and only if
  - $P(\alpha, \beta \mid \gamma) = P(\alpha \mid \gamma) P(\beta \mid \gamma)$

Height & Knowledge

$H \perp K$  ? No

$H \perp K \mid A$  Yes

# NUMBER OF PARAMETERS

- Assuming everything is binary
- $P(X_1)$  requires
  - 1 independent parameter
- $P(X_1, X_2, \dots, X_n)$  requires
  - $2^n - 1$  independent parameters
- $P(X_1 | X_2)$  requires
  - 2 independent parameters
- $P(X_1, X_2, \dots, X_n | X_{n+1}, X_{n+2}, \dots, X_{n+m})$  requires
  - $2^m \times (2^n - 1)$  independent parameters

# NUMBER OF PARAMETERS

- Assuming everything is binary

- $P(X_1)$  requires

- 1 independent parameter

- $P(X_1, X_2, \dots, X_n)$  requires

- $2^n - 1$  independent parameters

- $P(X_1 | X_2)$  requires

- 2 independent parameters

- $P(X_1, X_2, \dots, X_n | X_{n+1}, X_{n+2}, \dots, X_{n+m})$  requires

- $2^m \times (2^n - 1)$  independent parameters

$$\begin{array}{c|c} x_1 & P(x_1) \\ \hline T & \bar{0} \\ F & \bar{1} \end{array}$$

$$\begin{array}{c|c} x_1 & x_1 | x_2 = T \\ \hline T & \bar{0} \\ F & \bar{1} \end{array}$$

$$\begin{array}{c|c} x_1 & x_1 | x_2 = F \\ \hline T & \bar{0} \\ F & \bar{1} \end{array}$$

# of entries  $(2^n - 1) \times 2^m$  # of tables

$X: T, F$   
 $Y: R, G, B$

$X$	$Y$	
$T$	$R$	$P_1$
$T$	$B$	$P_2$
$T$	$G$	$P_3$
$F$	$R$	$P_4$
$F$	$B$	$P_5$
$F$	$G$	$P_6$

$$\sum P_i = 1$$

ind params  
 $2 \times 3 - 1$

$P(Y|X)$

$Y$	$P(Y X=T)$
$R$	$-$
$G$	$-$
$B$	$-$

$\pm$

1

2 + 2 ind

$Y$	$P(Y X=F)$
$R$	$-$
$G$	$-$
$B$	$-$

$\pm$

1

$X$	$P(X Y=R)$
$T$	$-$
$F$	$-$

$\pm$

1

$X$	$P(X Y=G)$
$T$	$-$
$F$	$-$

$\pm$

1

$X$	$P(X Y=B)$
$T$	$-$
$F$	$-$

$\pm$

1

1 + 1 + 1 = 3 ind

# CONTINUOUS SPACES

- Assume  $X$  is continuous and  $\text{Val}(X) = [0,1]$
- If you would like to assign the same probability to all real numbers in  $[0, 1]$ , what is, for e.g.,  $P(X=0.5) = ?$
- Answer:  $P(X=0.5) = 0$ .

# PROBABILITY DENSITY FUNCTION

- We define **probability density function**,  $p(x)$ , a non-negative integrable function, such that  $\int_{\text{Val}(X)} p(x)dx = 1$

$$P(X \leq a) = \int_{-\infty}^a p(x)dx$$

$$P(a \leq X \leq b) = \int_a^b p(x)dx$$

# CONDITIONAL PROBABILITY

- We want  $P(Y | X=x)$  where  $X$  is continuous,  $Y$  is discrete
- $P(Y | X=x) = P(Y, X=x) / P(X=x)$ 
  - What's wrong with this expression?
- Instead, we use the following expression

$$P(Y | X = x) = \lim_{\varepsilon \rightarrow 0} P(Y | x - \varepsilon \leq X \leq x + \varepsilon)$$

# CONDITIONAL PROBABILITY

- We want  $p(Y|X)$  where  $X$  is discrete,  $Y$  is continuous
- How would you represent it?

$$p(y|X=T) - \text{pdf}$$

$$p(y|X=\bar{T}) - \text{pdf}$$



# EXPECTATION

$$E_P[X] = \sum_x xP(x)$$

$$E_P[X] = \int_x xp(x)dx$$

$$E_P[aX + b] = aE_P[X] + b$$

$$E_P[X + Y] = E_P[X] + E_P[Y]$$

$$E_P[X | y] = \sum_x xP(x | y)$$

What about  $E[X*Y]$ ?

## VARIANCE

$$\text{Var}_P[X] = E_P \left[ \left( X - E_P[X] \right)^2 \right]$$

$$\text{Var}_P[X] = E_P[X^2] - \left( E_P[X] \right)^2$$

Can you derive the second expression using the first expression?

$$\text{Var}_P[aX + b] = a^2 \text{Var}_P[X]$$

What is  $\text{Var}[X+Y]$ ?

# SOME DISTRIBUTIONS

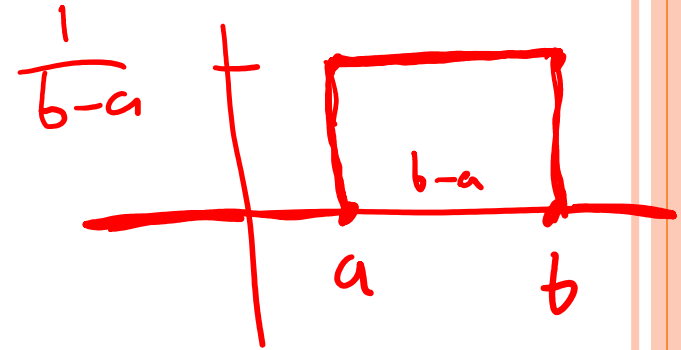
# BINOMIAL DISTRIBUTION

- Two parameters
  - $n$  – number of independent experiments each measuring a binary outcome (e.g., Yes/No, Heads/Tails, Positive/Negative, ...)
  - $p$  – “success” probability for each individual experiment (e.g., Yes, Heads, Positive, ...)
- Probability of exactly  $k$  successes
  - $P(k; n, p) = \binom{n}{k} p^k (1 - p)^{n-k}$ , where  $\binom{n}{k} = \frac{n!}{k!(n-k)!}$
- Expected value:  $np$
- Variance:  $np(1 - p)$
- An important use for binomial distribution: estimate a binary measure using a sample
  - Given  $k$  success in  $n$  experiments, estimate  $p$  and a confidence interval for  $p$

# UNIFORM DISTRIBUTION

- A variable  $X$  has a uniform distribution over  $[a, b]$  if it has the PDF

$$p(x) = \begin{cases} \frac{1}{b-a} & a \leq x \leq b \\ 0 & \text{otherwise} \end{cases}$$



Check and make sure that  $p(x)$  integrates to 1.  
What are the mean and variances of this distribution?

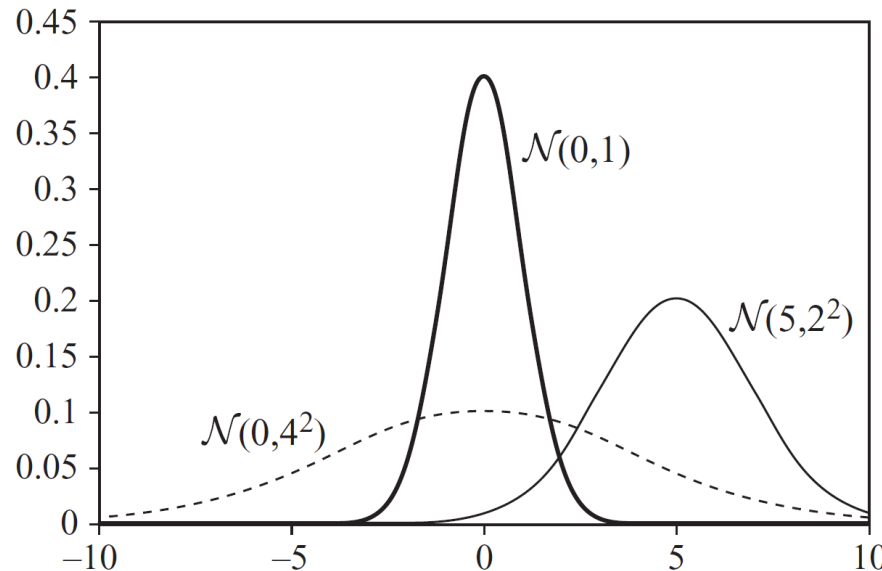
$$\int_a^b x p(x) dx$$

# GAUSSIAN DISTRIBUTION

- A variable  $X$  has a Gaussian distribution with mean  $\mu$  and variance  $\sigma^2$ , if it has the PDF

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

$$\int \underline{p(x)} x \, dx = 1$$



Can  $p(x)$  be ever greater than 1?

# OTHER TOPICS

- Information theory
- Parameter estimation
- Decision-making under uncertainty