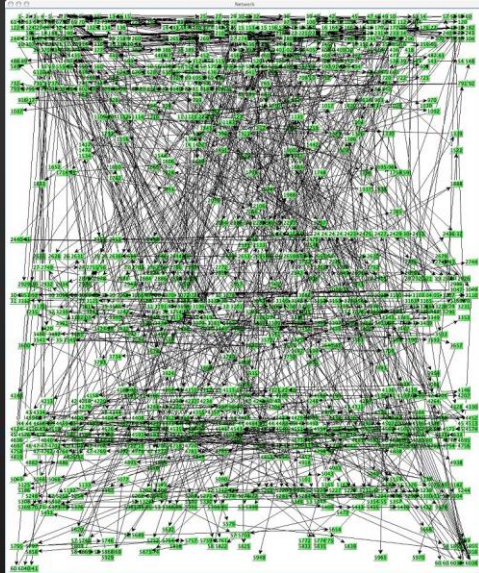


# CSE 575: Statistical Machine Learning

Jingrui He  
CIDSE, ASU

What is ***Machine Learning***?

# Machine Learning



what society thinks I do



what my friends think I do



what my parents think I do

$$L_p = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i y_i (\mathbf{x}_i \cdot \mathbf{w} + b) + \sum_{i=1}^n \alpha_i$$

$$\alpha_i \geq 0, \forall i$$

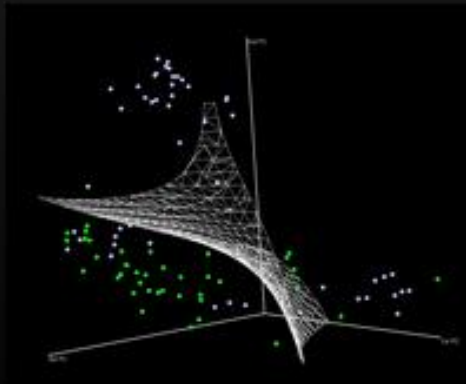
$$\mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i, \sum_{i=1}^n \alpha_i y_i = 0$$

$$\nabla \hat{g}(\theta_t) = \frac{1}{n} \sum_{i=1}^n \nabla \ell(x_i, y_i; \theta_t) + \nabla r(\theta_t)$$

$$\theta_{t+1} = \theta_t - \eta_t \nabla \ell(x_{i(t)}, y_{i(t)}; \theta_t) - \eta_t \cdot \nabla r(\theta_t)$$

$$\mathbb{E}_{i(t)}[\ell(x_{i(t)}, y_{i(t)}; \theta_t)] = \frac{1}{n} \sum_i \ell(x_i, y_i; \theta_t)$$

what other programmers think I do



what I think I do

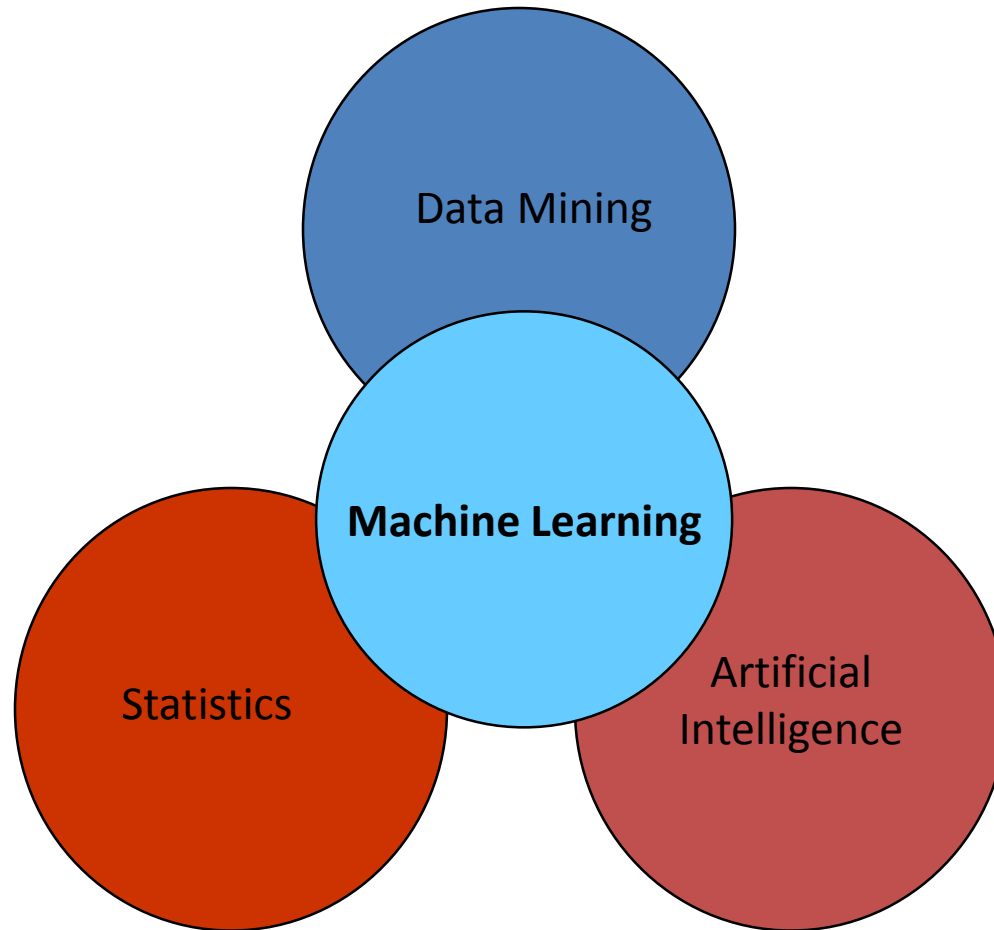
```
>>> from scipy import svm
```

what I really do

# Machine Learning

- Prof. Tom Mitchell@CMU &
- Prof. Carlos Guestrin@UW
  - ‘Study of algorithms that improve their performance, at some task, with experience’
- Prof. Andrew Ng@Stanford
  - ‘Machine learning is the science of getting computers to act without being explicitly programmed’

# Related Fields



# Useful Resources

- The discipline of machine learning:  
<http://www.cs.cmu.edu/~tom/pubs/MachineLearning.pdf>
- Coursera: <https://www.coursera.org/course/ml>
- Andrew Moore's tutorials: <http://www.autonlab.org/tutorials/>
- Alex Smola@CMU's machine learning lectures:  
[https://www.youtube.com/playlist?list=PLZSO\\_6-bSqHQmMKwWVvYwKreGu4b4kMU9](https://www.youtube.com/playlist?list=PLZSO_6-bSqHQmMKwWVvYwKreGu4b4kMU9)
- Mathworks Matlab tutorials:  
[http://www.mathworks.com/academia/student\\_center/tutorials/launchpad.html](http://www.mathworks.com/academia/student_center/tutorials/launchpad.html)
- Ben Taskar@UW's Matlab tutorial:  
<https://alliance.seas.upenn.edu/~cis520/wiki/index.php?n=Recitations.MatlabTutorial>
- Probability review by David Blei@Columbia:  
[http://www.cs.princeton.edu/courses/archive/spring07/cos424/scribe\\_notes/0208.pdf](http://www.cs.princeton.edu/courses/archive/spring07/cos424/scribe_notes/0208.pdf)

# Becoming Famous in ML?

- Winning Netflix prize?  
[http://en.wikipedia.org/wiki/Netflix\\_Prize](http://en.wikipedia.org/wiki/Netflix_Prize)
- Building Watson to win Jeopardy!?  
[http://en.wikipedia.org/wiki/Watson\\_\(computer\)](http://en.wikipedia.org/wiki/Watson_(computer))
- Helping President Obama win the election?  
<http://www.rayidghani.com/>
- Predicting the stock market with Twitter feed?  
<http://arxiv.org/pdf/1010.3003&>
- IEEE/ACM/AAAI fellow?

# Seriously

- Do a great job in CSE 575!
- Read many many many ... many papers
- Publish many many many ... many papers
  - ICML: <http://icml.cc/2017/>
  - NIPS: <http://nips.cc/Conferences/2017/>
  - UAI: <http://auai.org/uai2017/>
  - IJCAI: <http://ijcai17.org/>
  - AAI: <http://www.aaai.org/Conferences/AAAI/aaai17.php>
  - ACM KDD: <http://www.kdd.org/kdd2017/>
  - ICDM: <http://icdm2017.bigke.org/>
  - SDM: <http://www.siam.org/meetings/sdm17/>
  - Journal of Machine Learning Research: <http://jmlr.org/>
  - IEEE Transactions on Knowledge and Data Engineering: <http://www.computer.org/portal/web/tkde>



# Who Wants Machine Learning People?

- IT
  - Outlier/fraud detection
  - Web image search
  - Recommendation
  - Information filtering
  - Community detection
  - Ad placement
  - Sentiment analysis
  - ...
- Companies
  - Facebook, Google, LinkedIn, Twitter, Microsoft, IBM, AT&T, Apple, Amazon, Siemens, Foursquare, Yelp, Walmart Lab, NEC, Generic Electric, Baidu, Samsung, ...

# Who Wants Machine Learning People?

- Finance
  - Stock market prediction
  - Algorithmic trading
  - Return forecasting
  - ...
- Companies
  - Goldman Sachs, Morgan Stanley, American Express, Citadel LLC, Barclays Capital, Rotella Capital Management, Citi Bank, Pequot Capital, Zestfinance, Federal Reserve Board, WorldQuant LLC, ...

# Who Wants Machine Learning People?

- Speech recognition, natural language processing
- Computer vision
- Healthcare
- Robot control
- Computational biology
- Sensor networks
- ...

# Basics on Probability

# Coin Flips

- You flip a coin
  - Head with probability 0.5
- You flip 100 coins
  - How many heads would you expect

# Coin Flips cont.

- You flip a coin
  - Head with probability  $p$
  - Binary random variable
  - Bernoulli trial with success probability  $p$
- You flip  $k$  coins
  - How many heads would you expect
  - Number of heads  $X$ : discrete random variable
  - Binomial distribution with parameters  $k$  and  $p$

# Discrete Random Variables

- Random variables (RVs) which may take on only a **countable** number of **distinct** values
  - E.g., the total number of heads  $X$  you get if you flip 100 coins
- $X$  is a RV with arity  $k$  if it can take on exactly one value out of  $\{x_1, \dots, x_k\}$ ,
  - E.g., the possible values that  $X$  can take are 0, 1, 2,..., 100

# Probability of Discrete RV

- Probability mass function (pmf):  $P(X = x_i)$
- Easy facts about pmf
  - $\sum_i P(X = x_i) = 1$
  - $P(X = x_i \cap X = x_j) = 0$  if  $i \neq j$
  - $P(X = x_i \cup X = x_j) = P(X = x_i) + P(X = x_j)$  if  $i \neq j$
  - $P(X = x_1 \cup X = x_2 \cup \dots \cup X = x_k) = 1$



# Common Distributions

- Uniform  $X \sim U[1, \dots, N]$ 
  - $X$  takes values  $1, 2, \dots, N$
  - $P(X = i) = 1/N$
  - E.g., picking balls of different colors from a box
- Binomial  $X \sim \text{Bin}(n, p)$ 
  - $X$  takes values  $0, 1, \dots, n$
  - $P(X = i) = \binom{n}{i} p^i (1-p)^{n-i}$
  - E.g., coin flips

# Coin Flips of Two Persons

- Your friend and you both flip coins
  - Head with probability 0.5
  - You flip 50 times; your friend flip 100 times
  - How many heads will both of you get

# Joint Distribution

- Given two discrete RVs  $X$  and  $Y$ , their **joint distribution** is the distribution of  $X$  and  $Y$  together

– E.g.,  $P(\text{You get 21 heads AND you friend get 70 heads})$

- $$\sum_x \sum_y P(X = x \cap Y = y) = 1$$

– E.g.,

$$\sum_{i=0}^{50} \sum_{j=0}^{100} P(\text{You get } i \text{ heads AND your friend get } j \text{ heads}) = 1$$

# Conditional Probability

- $P(X = x | Y = y)$  is the probability of  $X = x$ , given the occurrence of  $Y = y$ 
  - E.g., you get 0 heads, given that your friend gets 61 heads
- $$P(X = x | Y = y) = \frac{P(X = x \cap Y = y)}{P(Y = y)}$$

# Law of Total Probability

- Given two discrete RVs  $X$  and  $Y$ , which take values in  $\{x_1, \dots, x_m\}$  and  $\{y_1, \dots, y_n\}$ , we have

$$\begin{aligned} P(X = x_i) &= \sum_j P(X = x_i \cap Y = y_j) \\ &= \sum_j P(X = x_i | Y = y_j) P(Y = y_j) \end{aligned}$$

# Marginalization

Marginal Probability

Joint Probability

$$\begin{aligned} P(X = x_i) &= \sum_j P(X = x_i \cap Y = y_j) \\ &= \sum_j P(X = x_i | Y = y_j) P(Y = y_j) \end{aligned}$$

Conditional Probability

Marginal Probability

# Bayes Rule

- X and Y are discrete RVs...

$$P(X = x | Y = y) = \frac{P(X = x \cap Y = y)}{P(Y = y)}$$



$$P(X = x_i | Y = y_j) = \frac{P(Y = y_j | X = x_i) P(X = x_i)}{\sum_k P(Y = y_j | X = x_k) P(X = x_k)}$$

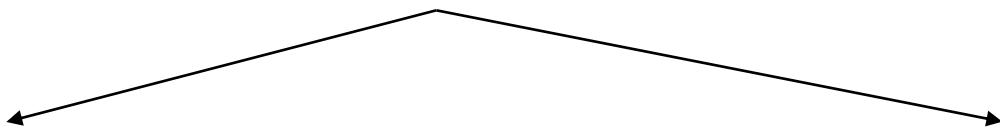
# Independent RVs

- Intuition:  $X$  and  $Y$  are independent means that  $X = x$  **neither** makes it **more or less** probable that  $Y = y$
- Definition:  $X$  and  $Y$  are independent iff
$$P(X = x \cap Y = y) = P(X = x)P(Y = y)$$



# More on Independence

- $P(X = x \cap Y = y) = P(X = x)P(Y = y)$



A diagram consisting of a horizontal line with a central point. From this point, two arrows branch out downwards and outwards, pointing towards the two conditional probability terms in the equation below.

$$P(X = x|Y = y) = P(X = x) \quad P(Y = y|X = x) = P(Y = y)$$

- **E.g.**, no matter how many heads you get, your friend will not be affected, and vice versa

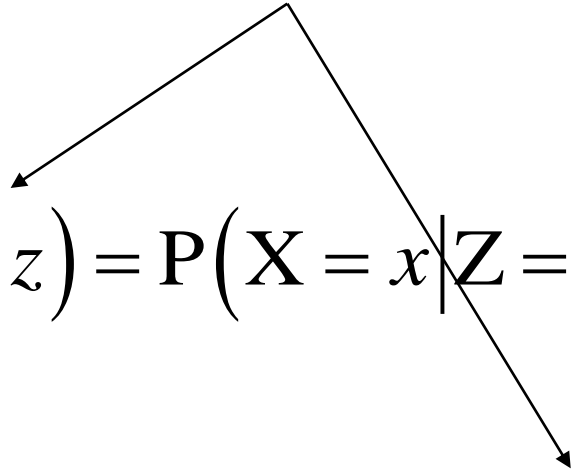
# Conditionally Independent RVs

- Intuition:  $X$  and  $Y$  are conditionally independent given  $Z$  means that once  $Z$  is **known**, the value of  $X$  does not add any **additional** information about  $Y$
- Definition:  $X$  and  $Y$  are conditionally independent given  $Z$  iff

$$P(X = x \cap Y = y | Z = z) = P(X = x | Z = z)P(Y = y | Z = z)$$

# More on Conditional Independence

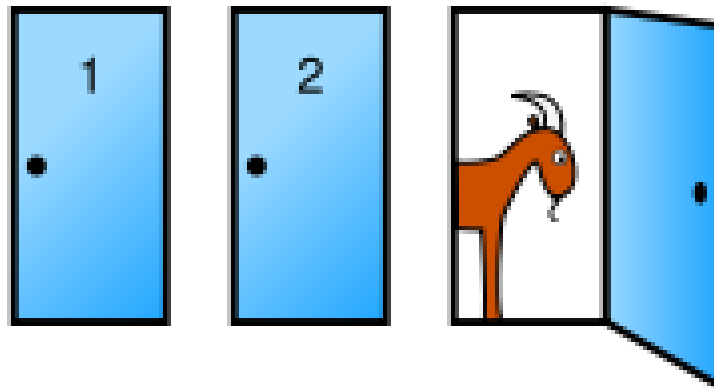
$$P(X = x \cap Y = y | Z = z) = P(X = x | Z = z) P(Y = y | Z = z)$$

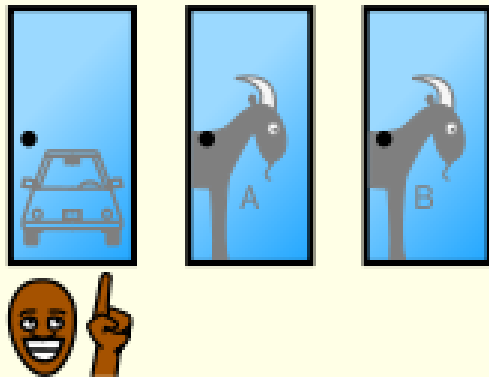

$$P(X = x | Y = y, Z = z) = P(X = x | Z = z)$$

$$P(Y = y | X = x, Z = z) = P(Y = y | Z = z)$$

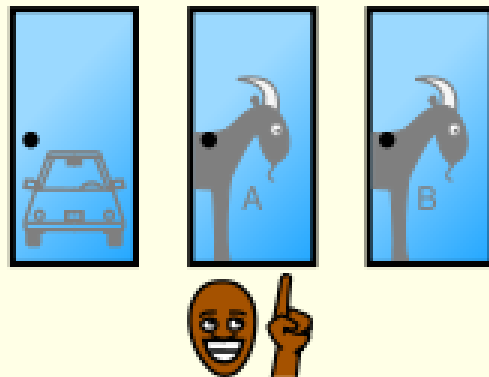
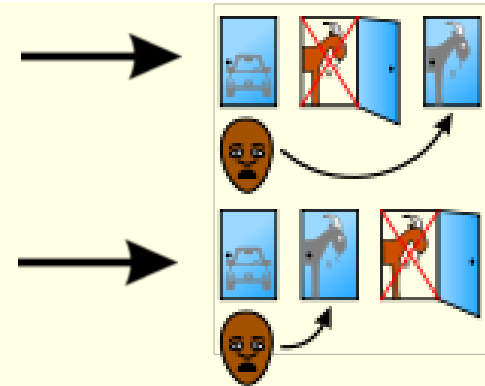
# Monty Hall Problem

- You're given the choice of three doors: Behind one door is a car; behind the others, goats.
- You pick a door, say No. 1
- The host, who knows what's behind the doors, opens another door, say No. 3, which has a goat.
- Do you want to pick door No. 2 instead?

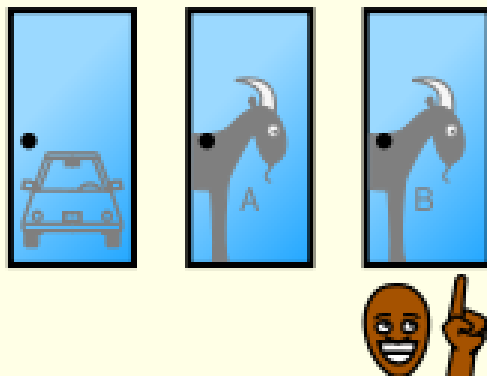
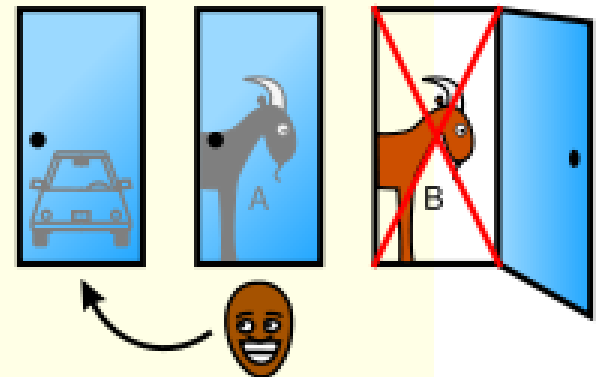




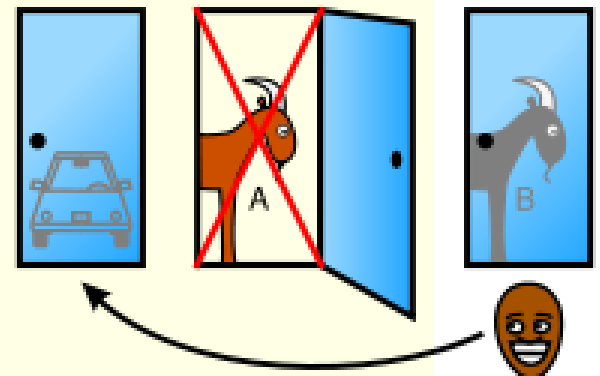
*Host reveals  
Goat A  
or  
Host reveals  
Goat B*



*Host must  
reveal Goat B*



*Host must  
reveal Goat A*



# Monty Hall Problem: Bayes Rule

- $C_k$ : the car is behind door  $k$ ,  $k = 1, 2, 3$
- $P(C_k) = 1/3$
- $H_{ij}$ : the host opens door  $j$  after you pick door  $i$

- $$P(H_{ij} | C_k) = \begin{cases} 0 & i = j \\ 0 & j = k \\ 1/2 & i = k \\ 1 & i \neq k, j \neq k \end{cases}$$

# Monty Hall Problem: Bayes Rule cont.

- WLOG,  $i=1$  (your choice),  $j=3$  (the host's choice)

- $$P(C_1 | H_{13}) = \frac{P(H_{13} | C_1) P(C_1)}{P(H_{13})}$$

- $$P(H_{13} | C_1) P(C_1) = \frac{1}{2} \cdot \frac{1}{3} = \frac{1}{6}$$

# Monty Hall Problem: Bayes Rule cont.

- $$\begin{aligned} P(H_{13}) &= P(H_{13}, C_1) + P(H_{13}, C_2) + P(H_{13}, C_3) \\ &= P(H_{13} | C_1) P(C_1) + P(H_{13} | C_2) P(C_2) \\ &= \frac{1}{6} + 1 \cdot \frac{1}{3} \\ &= \frac{1}{2} \end{aligned}$$
- $$P(C_1 | H_{13}) = \frac{1/6}{1/2} = \frac{1}{3}$$



# Monty Hall Problem: Bayes Rule cont.

- $P(C_1 | H_{13}) = \frac{1/6}{1/2} = \frac{1}{3}$
- $P(C_2 | H_{13}) = 1 - \frac{1}{3} = \frac{2}{3} > P(C_1 | H_{13})$
- *You should switch!*

# Continuous Random Variables

- What if  $X$  is continuous?
- Probability density function (pdf) instead of probability mass function (pmf)
- A pdf is any function  $f(x)$  that describes the probability density in terms of the input variable  $x$ .

# PDF

- Properties of pdf

- $f(x) \geq 0, \forall x$

- $\int_{-\infty}^{+\infty} f(x) = 1$

- $f(x) \leq 1$  ???

- Actual probability can be obtained by taking the integral of pdf

- E.g., the probability of  $X$  being between 0 and 1 is

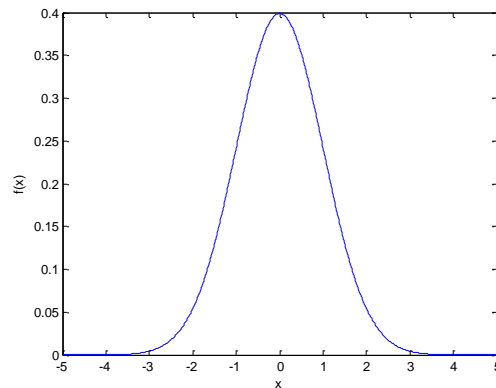
$$P(0 \leq X \leq 1) = \int_0^1 f(x) dx$$

# Cumulative Distribution Function

- $F_X(v) = P(X \leq v)$
- Discrete RVs
  - $F_X(v) = \sum_{v_i} P(X = v_i)$
- Continuous RVs
  - $F_X(v) = \int_{-\infty}^v f(x) dx$
  - $\frac{d}{dx} F_X(x) = f(x)$

# Common Distributions

- Normal  $X \sim N(\mu, \sigma^2)$ 
  - $f(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}, x \in \mathbb{R}$
  - E.g., the height of the entire population



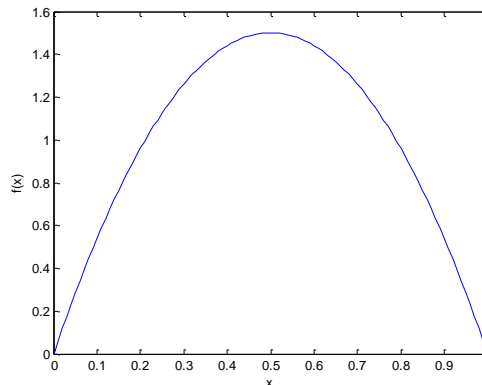
# Common Distributions cont.

- Beta  $X \sim \text{Beta}(\alpha, \beta)$

- $f(x; \alpha, \beta) = \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1}, \quad x \in [0, 1]$

- $\alpha = \beta = 1$ : uniform distribution between 0 and 1

- E.g., the conjugate prior for the parameter  $p$  in Binomial distribution

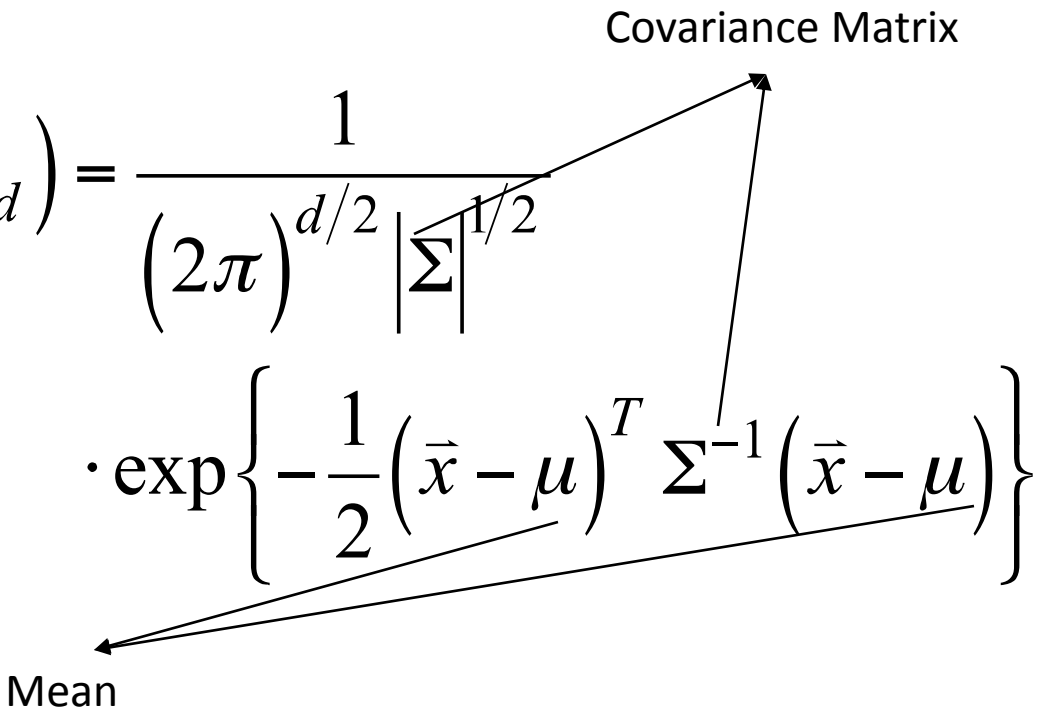


# Joint Distribution

- Given two continuous RVs  $X$  and  $Y$ , the **joint pdf** can be written as  $f_{X,Y}(x, y)$
- $\int_x \int_y f_{X,Y}(x, y) dx dy = 1$

# Multivariate Normal

- Generalization to higher dimensions of the one-dimensional normal

- $$f_{\vec{X}}(x_1, \dots, x_d) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \cdot \exp \left\{ -\frac{1}{2} (\vec{x} - \mu)^T \Sigma^{-1} (\vec{x} - \mu) \right\}$$


The diagram illustrates the components of the Multivariate Normal distribution formula. An arrow points from the label "Covariance Matrix" to the symbol  $\Sigma$  in the denominator. Another arrow points from the label "Mean" to the symbol  $\mu$  in the exponent.



# Moments

- Mean (Expectation):  $\mu = E(X)$ 
  - Discrete RVs:  $E(X) = \sum_{v_i} v_i P(X = v_i)$
  - Continuous RVs:  $E(X) = \int_{-\infty}^{+\infty} x f(x) dx$
- Variance:  $V(X) = E(X - \mu)^2$ 
  - Discrete RVs:  $V(X) = \sum_{v_i} (v_i - \mu)^2 P(X = v_i)$
  - Continuous RVs:  $V(X) = \int_{-\infty}^{+\infty} (x - \mu)^2 f(x) dx$

# Properties of Moments

- Mean

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

- $E(aX) = aE(X)$

- If  $X$  and  $Y$  are independent,  $E(XY) = E(X) \cdot E(Y)$

- Variance

- $V(aX + b) = a^2V(X)$

- If  $X$  and  $Y$  are independent,  $V(X + Y) = V(X) + V(Y)$

# Moments of Common Distributions

- Uniform  $X \sim U[1, \dots, N]$ 
  - Mean  $(1 + N)/2$ ; variance  $(N^2 - 1)/12$
- Binomial  $X \sim \text{Bin}(n, p)$ 
  - Mean  $np$ ; variance  $np^2$
- Normal  $X \sim N(\mu, \sigma^2)$ 
  - Mean  $\mu$ ; variance  $\sigma^2$
- Beta  $X \sim \text{Beta}(\alpha, \beta)$ 
  - Mean  $\alpha/(\alpha + \beta)$ ; variance  $\frac{\alpha\beta}{(\alpha + \beta)^2 (\alpha + \beta + 1)}$

# Probability of Events

- $X$  denotes an event that could possibly happen
  - E.g.,  $X$  = “you will fail in this course”
- $P(X)$  denotes the **likelihood** that  $X$  happens, or  $X$  = true
  - E.g., what’s the probability that you will fail in this course?
- $\Omega$  denotes the entire event set
  - $\Omega = \{X, \bar{X}\}$

# The Axioms of Probabilities

- $0 \leq P(X) \leq 1$
- $P(\Omega) = 1$
- $P(X_1 \cup X_2 \cup \dots) = \sum_i P(X_i)$ , where  $X_i$  are disjoint events
- Useful rules
  - $P(X_1 \cup X_2) = P(X_1) + P(X_2) - P(X_1 \cap X_2)$
  - $P(\bar{X}) = 1 - P(X)$

# Interpreting the Axioms

