

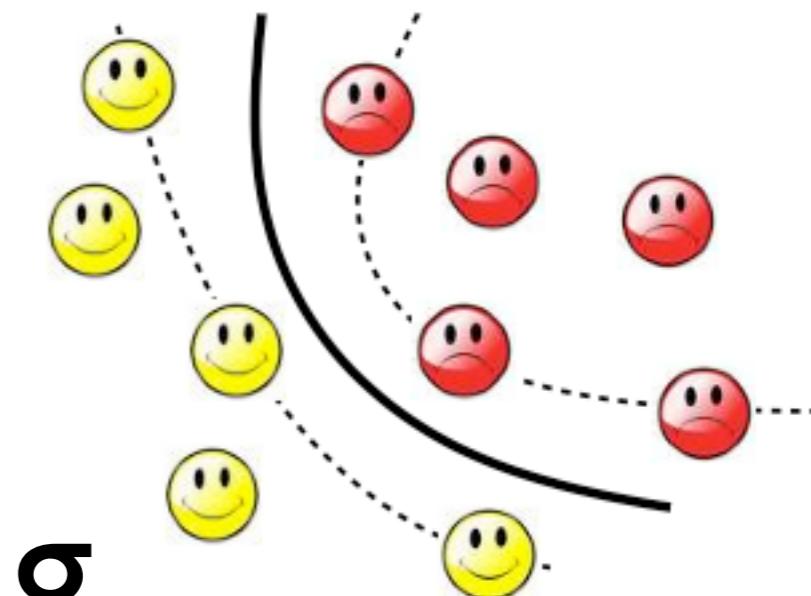


THE UNIVERSITY OF
SYDNEY

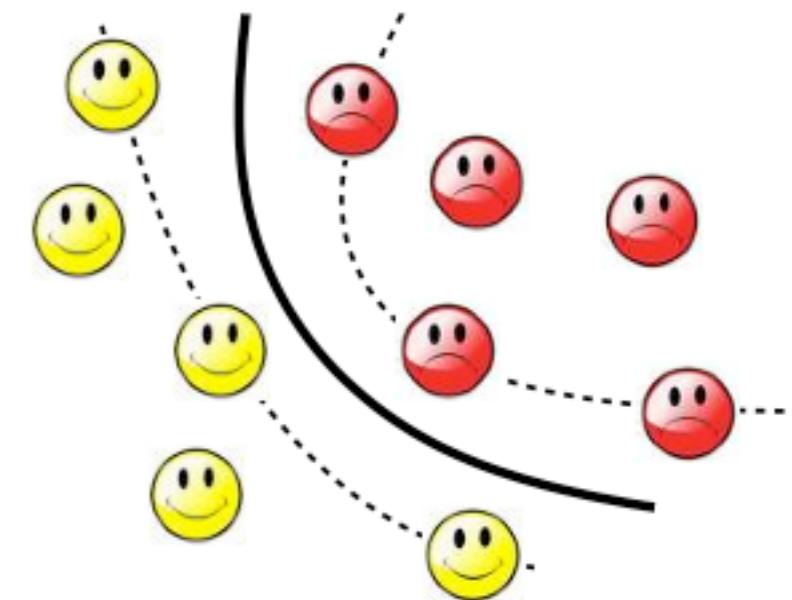
Machine Learning and Data Mining

(COMP 5318)

Basics of probability theory and Bayes' rule



Tongliang Liu



Review



Basics I

$$S = \begin{bmatrix} 3 & 0 & 1 \\ 0 & 2 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

- This is a 3×3 matrix.
- In general $m \times n$.
 - m rows and n columns
 - Square matrix when $m = n$
- Each row or column could represent one object. If rows are objects then columns are features/attributes/components



Basics II

- Identity matrix I

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

- If A is a square matrix, $AI = IA = A$
- I is an example of a **diagonal** matrix.
- If $A = [a_1, \dots, a_m]$ is matrix where a_i are the columns, then
 - A is orthogonal if $a_i \cdot a_j = 0$ for $i \neq j$
 - A is orthonormal if above and $a_i \cdot a_i = 1$



Basics III

- Every vector can be written as a linear combination of some finitely many “special” vectors.
- These are called **basis-vectors**.

$$S = \begin{bmatrix} 3 \\ 2 \\ 2 \end{bmatrix} = 3 \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} + 2 \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} + 2 \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$



Linear independence

- Intuitively, a set of vectors is linearly independent if any element of the set cannot be expressed as a linear combination of the others.
- The columns are not linearly independent:

$$S = \begin{bmatrix} 3 & 0 & 1 \\ 0 & 2 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$



Rank of a matrix I

- Given a matrix X , the **rank** of a matrix is the maximum number of linearly independent columns.
- A rank 2 matrix:

$$S = \begin{bmatrix} 3 & 0 & 1 \\ 0 & 2 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$



Eigen Decomposition

- For any square matrix A we say that λ is an eigenvalue and \mathbf{u} is its eigenvector if

$$A\mathbf{u} = \lambda\mathbf{u}, \quad \mathbf{u} \neq 0.$$

- Stacking up all eigenvectors/values gives

$$AU = U\Lambda = \begin{bmatrix} & & & \\ | & | & & | \\ \mathbf{u}_1 & \mathbf{u}_2 & \dots & \mathbf{u}_n \\ | & | & & | \end{bmatrix} \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_n \end{bmatrix}$$



Singular Value Decomposition

- Given **any** real matrix X of size (m,n) , it can be expressed as:

$$X = U\Sigma V^T$$

The diagram illustrates the dimensions of the matrices in the SVD formula. It shows three boxes with dimensions: $m \times r$ (left), $r \times r$ (middle), and $r \times n$ (right). Arrows point from each dimension box to its corresponding term in the formula: the $m \times r$ box points to U , the $r \times r$ box points to Σ , and the $r \times n$ box points to V^T .

- r is the rank of matrix X
- U is a (m,r) column-orthonormal matrix
- V is a (n, r) column-orthonormal matrix
- Σ is diagonal $r \times r$ matrix



Example: compression of X

- Now a compact way of writing the spectral representation is:

$$A = U\Lambda U^T = \sum_{i=1}^r \lambda_i \mathbf{u}_i \times \mathbf{u}_i^T$$

$$X = U\Sigma V^T = \sum_{i=1}^r \lambda_i \mathbf{u}_i \times \mathbf{v}_i^T$$

- However, can approximate it as:

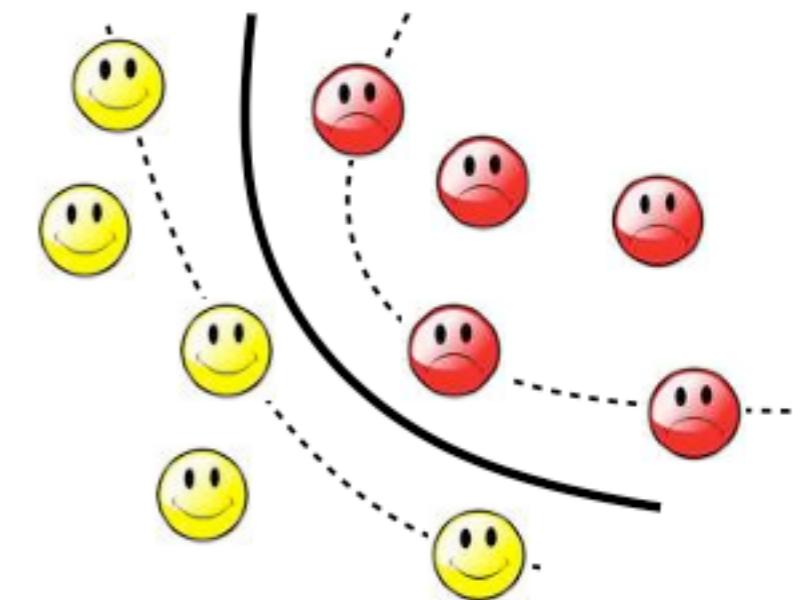
$$\hat{X} = \sum_{i=1}^k \lambda_i \mathbf{u}_i \times \mathbf{v}_i^T$$

- This new compression ratio is:

$$\frac{mk + k + nk}{mn} = \frac{k(m+1+n)}{mn} \approx \frac{km}{mn} = \frac{k}{n} \leqslant \frac{r}{n} \approx \frac{mr + r + nr}{mn}$$



THE UNIVERSITY OF
SYDNEY



Probability Theory



THE UNIVERSITY OF
SYDNEY

Why Probabilities?

- As stated by Laplace, “Probability is common sense reduced to calculation”.
- Probability theory is useful in understanding, studying, and analysis complex real world systems



THE UNIVERSITY OF
SYDNEY

Understanding uncertainty



- Aleatory: chance, no ability to predict outcome
- Epistemic: encoding knowledge, ability to predict outcome
- Sensing: ability to encode noisy measurements

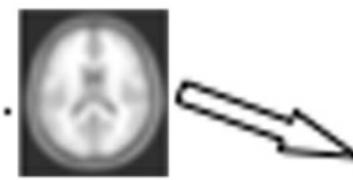
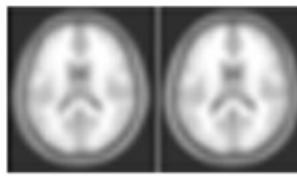
It is better to be imprecisely right than precisely wrong!



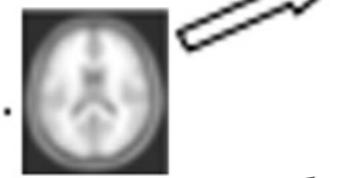
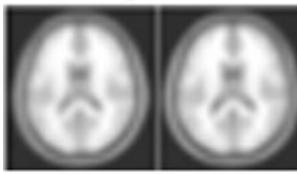
Elements of Machine Learning

Input training data

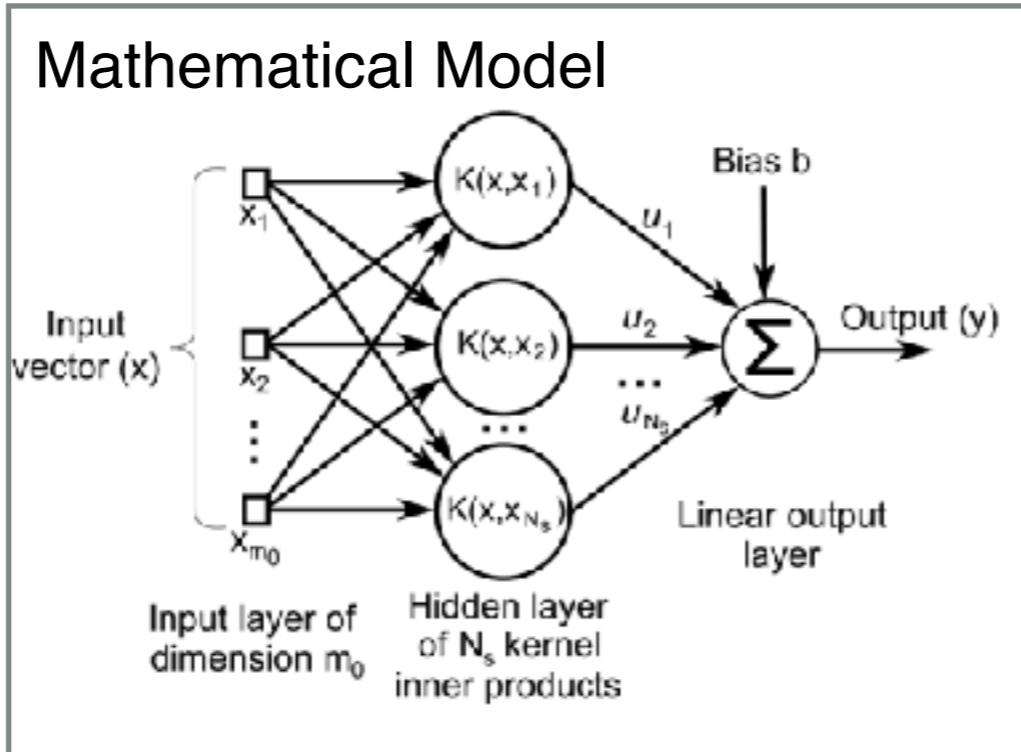
Group 1



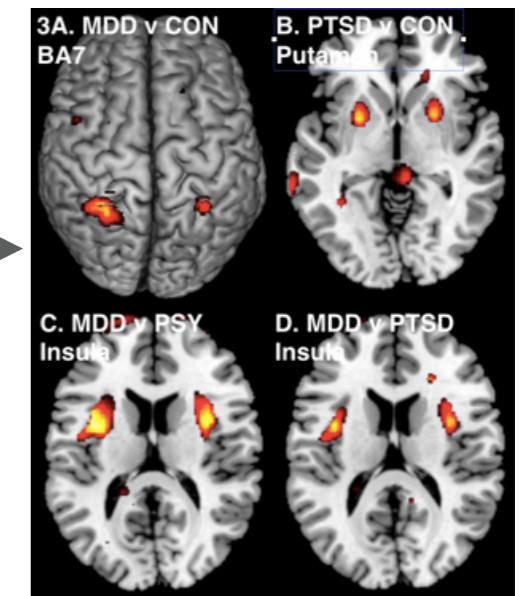
Group 2



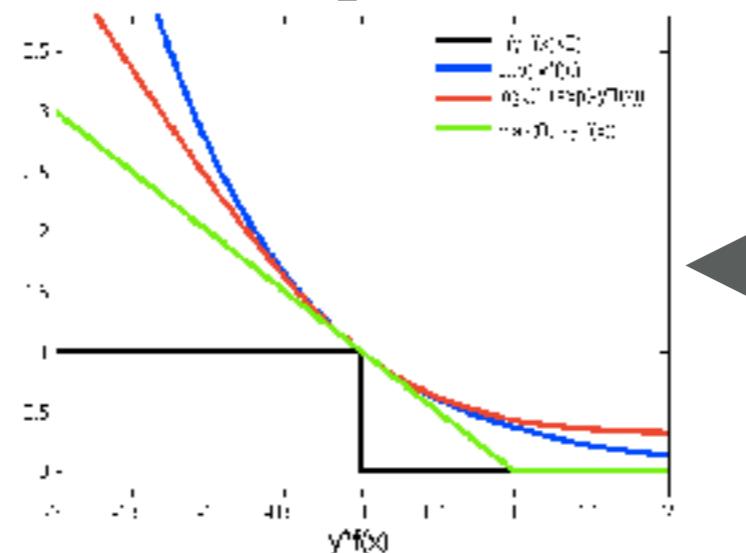
Data



Output hypothesis
Predictions/Patterns



Input predefined
function class



Objective function

Optimisation
method



THE UNIVERSITY OF
SYDNEY

How good is the output h_S ?

- The hypothesis is learned from training data.
How good is it for the test data?
- Given a previously unseen data, how will h_S perform on it?



Predictions and Probabilities

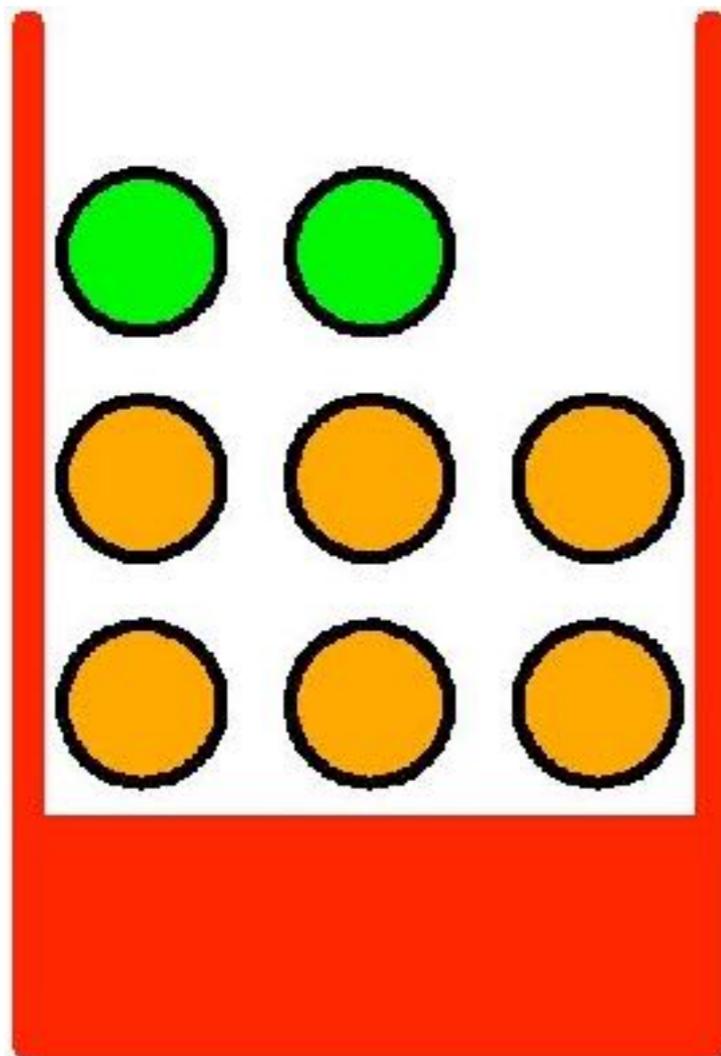
- When we make predictions we should assign “probabilities” with the prediction.
- Examples:
 - 20% chance it will rain tomorrow.
 - 50% chance that the tumour is malignant.
 - 60% chance that the stock market will fall by the end of the week.
 - 30% that the next president of the United States will be a Democrat.
 - 0.1% chance that the user will click on a banner-ad.
- How do we assign probabilities to complex events... using smart data algorithms... and counting.



THE UNIVERSITY OF
SYDNEY

Probability Theory

Apples and Oranges



$$P(\text{apples}) = 2/8 = 0.25$$

$$P(\text{apples}) = 3/4 = 0.75$$



Probability Basics

- Probability is a deep topic....but for most cases the rules are straightforward to apply.
- Terminology
 - Experiment
 - Sample Space
 - Events
 - Probability
 - Rules of probability
 - Conditional probability – Bayes' Rule



Experiments and Sample Space

- Consider an experiment and let S be the space of possible outcomes.
- Example:
 - Experiment is tossing a coin; $S=\{h,t\}$
 - Experiment is rolling a pair of dice: $S=\{(1,1),(1,2),\dots,(6,6)\}$
 - Experiment is a race consisting of three cars: 1,2 and 3. The sample space is $\{(1,2,3),(1,3,2),(2,1,3),(2,3,1),(3,1,2),(3,2,1)\}$



Probability

- Let Sample Space $S = \{1, 2, \dots, m\}$
- Consider numbers $p_i \geq 0, i = 1, 2, \dots, m; \sum_i p_i = 1$
- p_i is the probability that the outcome of the experiment is i .
- Suppose we toss a fair coin. Sample space is $S = \{h, t\}$. Then $p_h = 0.5$ and $p_t = 0.5$.



Assigning probabilities

- Experiment: Will it rain or not in Sydney:
 $S = \{\text{rain, no-rain}\}$
– $P_{\text{rain}} = 138/365 = 0.38$; $P_{\text{no-rain}} = 227/365 = 0.62$
- Assigning (or rather how to obtain) probabilities is a deep philosophical problem.
– What is the probability that the “green object standing outside my house is a burglar dressed in green?”



Events

- An *Event A* is a set of possible outcomes of the experiment. Thus A is a subset of S.
- Let A be the event of getting a seven when we roll a pair of dice.
 - $A = \{(1,6), (6,1), (2,5), (5,2), (4,3), (3,4)\}$
 - $P(A) = 6/36 = 1/6$
- In general
$$P(A) = \sum_{i \in A} p_i$$



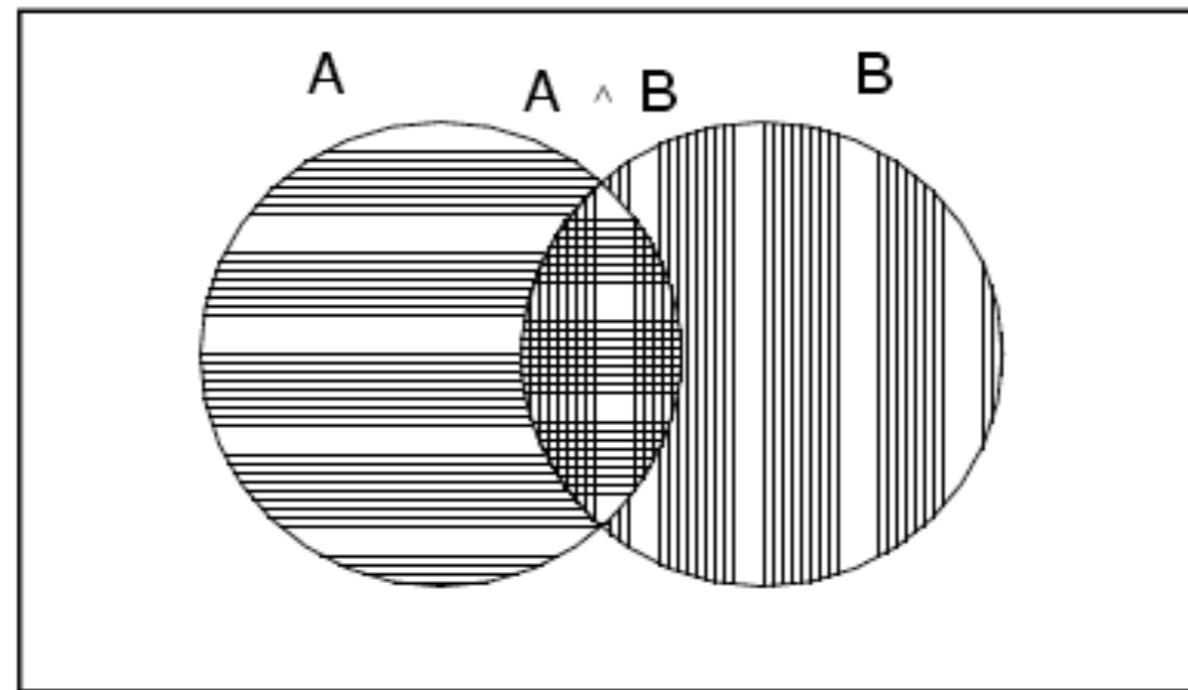
Events and Sample Space

- The sample space S and events are “sets”.
- $P(S) = 1$; probability of everything
- $P(\Phi) = 0$; probability of “null”
- Addition:
$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$
 - Often
$$P(A \cap B) \equiv P(AB) \equiv P(A, B)$$
- Complement:
$$P(A^c) = 1 - P(A)$$



Axioms of probability

True



$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

$$P(A \cap B) \equiv P(AB) \equiv P(A, B)$$

$$P(A^c) = 1 - P(A)$$



Example

- Suppose the probability of raining today is 0.4 and tomorrow is also 0.4 and on both days is 0.1.What is the probability it does not rain on either day?



Example

- Suppose the probability of raining today is 0.4 and tomorrow is also 0.4 and on both days is 0.1.What is the probability it does not rain on either day?
- $S=\{(R,N), (R,R), (N,N), (N,R)\}$
- Let A be the event that it will rain today and B it will rain tomorrow. Then
 $A=\{(R,N), (R,R)\} ; B=\{(N,R),(R,R)\}$
- Rain at least today or tomorrow:
$$P(A \cup B) = 0.4 + 0.4 - 0.1 = 0.7$$
- Will not rain on either day: $1 - 0.7 = 0.3$



Discrete Random Variables

- Events like “ASX is up” are binary events.
- We can extend this: by defining a **discrete random variable**.

$P(X = x)$ the probability that event $X = x$

- Two properties need to be satisfied

$$0 \leq P(X = x) \leq 1$$

$$\sum_{x \in X} P(X = x) = 1 \quad P(X=x) \leq 1 \text{ only for discrete variables}$$



Continuous Random Variables

- Random variables can also be continuous:
Height, rainfall, salary, chemical concentration...
- We can talk about the average (mean) and standard deviation or variance.
e.g., the average height of students in COMP5318 is 175 cm with a standard deviation of 15 cm.



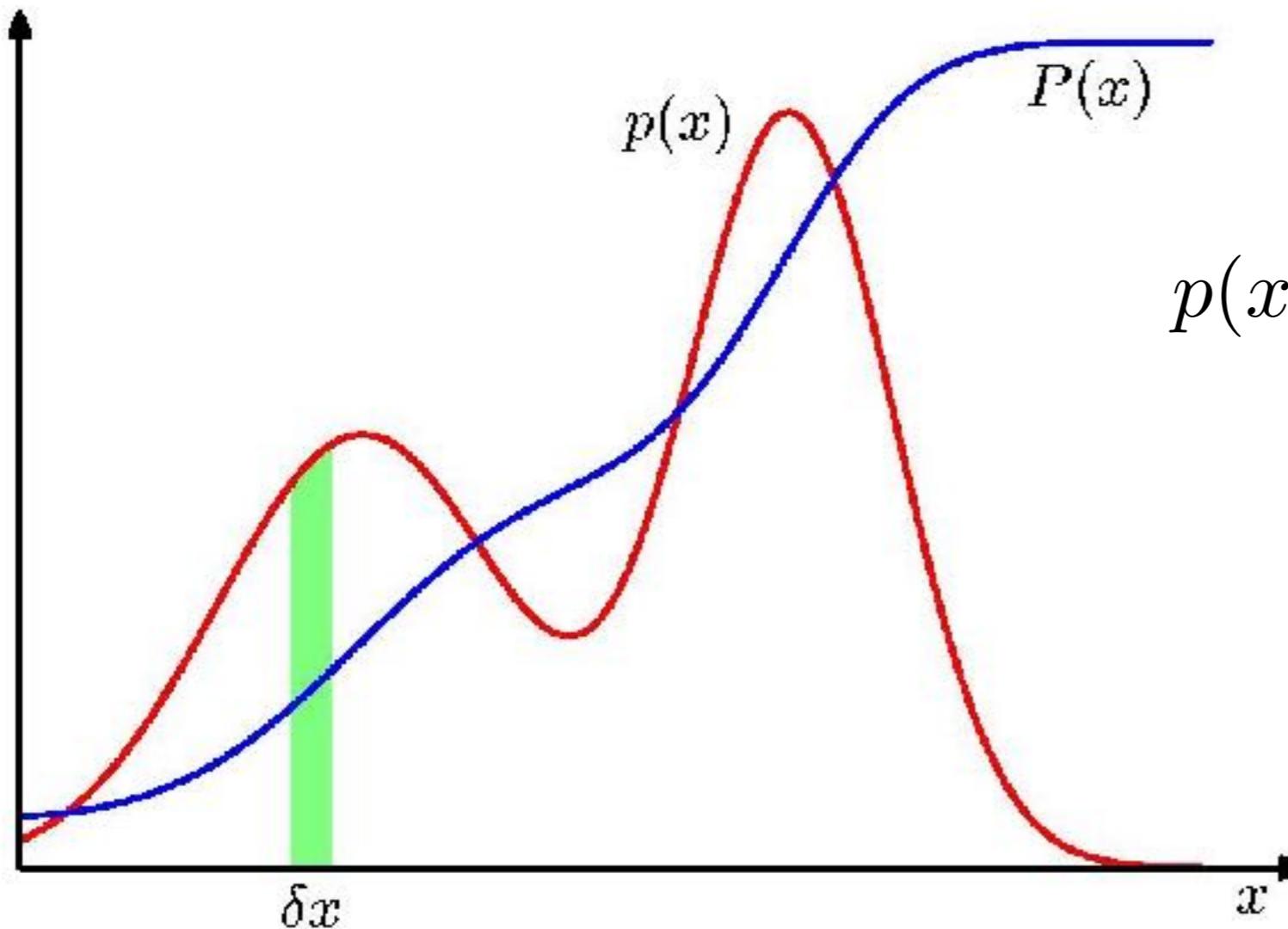
Probability Densities

- Random variables (both continuous and discrete) are associated with distributions.
- Common examples of discrete distributions are: Bernoulli, binomial, multinomial, Poisson.
- Common examples of continuous distributions are: Gaussian (Normal), Laplacian, Exponential, Gamma.
- Associated with distributions are parameters...
- One of the key problems in Statistics is to learn the parameters of a distribution from data.

This is **like summarising data**.



Probability Densities



$$p(x) \geq 0$$

$$\int_{-\infty}^{\infty} p(x) dx = 1$$

probability density
function (pdf)

$$p(x \in (a, b)) = \int_a^b p(x) dx$$

$$P(z) = \int_{-\infty}^z p(x) dx$$

Cumulative distribution
function (cdf)



Expectations

$$\mathbb{E}[f] = \sum_x p(x)f(x)$$

$$\mathbb{E}[f] = \int p(x)f(x)dx$$

$$\mathbb{E}_x[f|y] = \sum_x p(x|y)f(x)$$



Conditional Expectation
(discrete)

$$\mathbb{E}[f] \simeq \frac{1}{N} \sum_{n=1}^N f(x_n)$$

Approximate Expectation
(discrete and continuous)



Variance and Covariance

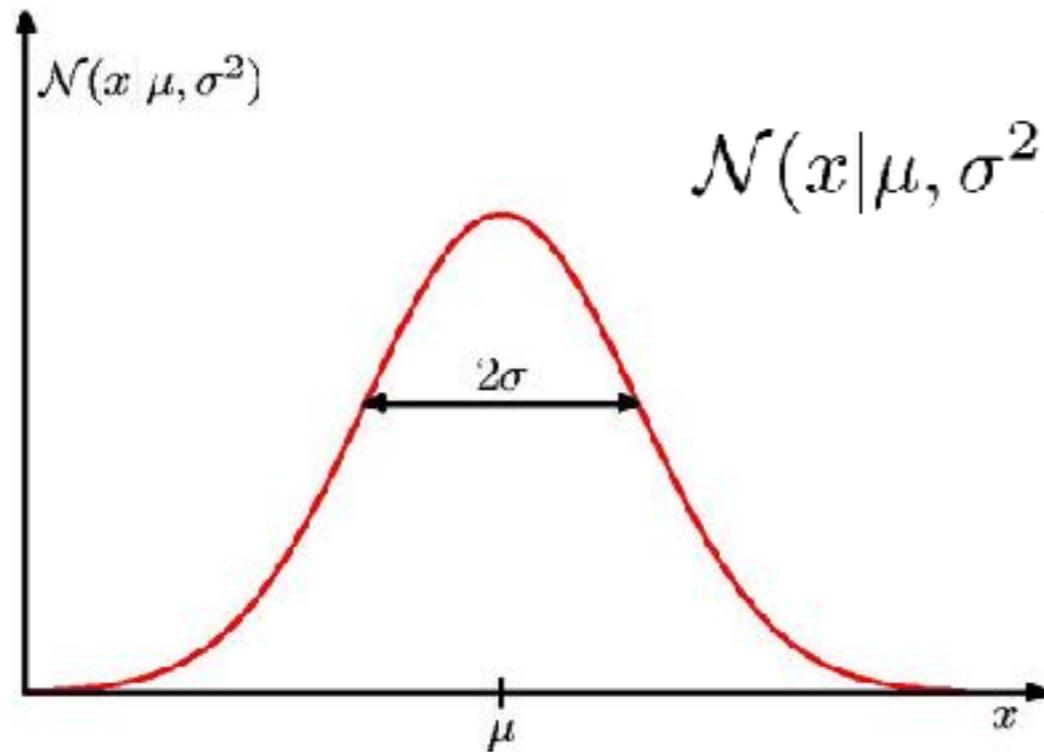
$$\text{var}[f] = \mathbb{E} \left[(f(x) - \mathbb{E}[f(x)])^2 \right] = \mathbb{E}[f(x)^2] - \mathbb{E}[f(x)]^2$$

$$\begin{aligned}\text{cov}[x, y] &= \mathbb{E}_{x,y} [\{x - \mathbb{E}[x]\} \{y - \mathbb{E}[y]\}] \\ &= \mathbb{E}_{x,y}[xy] - \mathbb{E}[x]\mathbb{E}[y]\end{aligned}$$

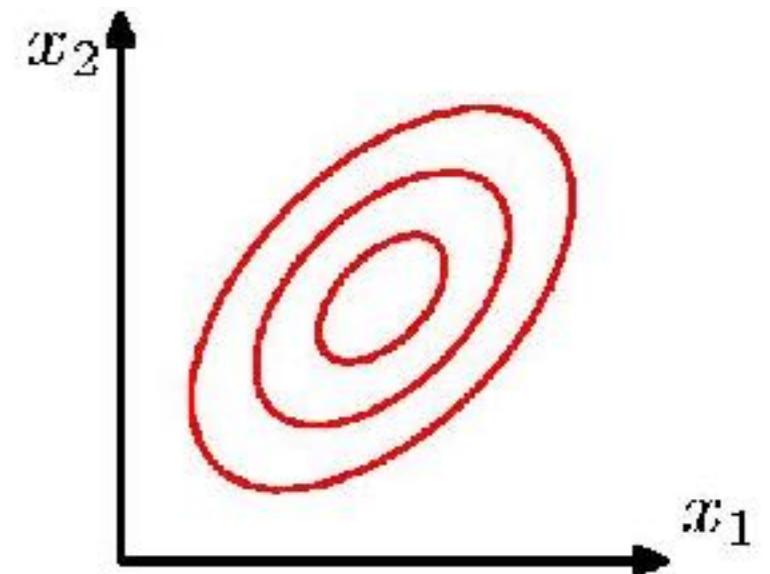
$$\begin{aligned}\text{cov}[\mathbf{x}, \mathbf{y}] &= \mathbb{E}_{\mathbf{x},\mathbf{y}} [\{\mathbf{x} - \mathbb{E}[\mathbf{x}]\}\{\mathbf{y}^T - \mathbb{E}[\mathbf{y}^T]\}] \\ &= \mathbb{E}_{\mathbf{x},\mathbf{y}}[\mathbf{x}\mathbf{y}^T] - \mathbb{E}[\mathbf{x}]\mathbb{E}[\mathbf{y}^T]\end{aligned}$$



The Gaussian Distribution



$$N(x|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp \left\{ -\frac{1}{2\sigma^2}(x - \mu)^2 \right\}$$



$$N(\mathbf{x}|\mu, \Sigma) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

Most entropic distribution given a mean and variance



THE UNIVERSITY OF
SYDNEY

Gaussian Mean and Variance

$$\mathbb{E}[x] = \int_{-\infty}^{\infty} \mathcal{N}(x|\mu, \sigma^2) x \, dx = \mu$$

$$\mathbb{E}[x^2] = \int_{-\infty}^{\infty} \mathcal{N}(x|\mu, \sigma^2) x^2 \, dx = \mu^2 + \sigma^2$$

$$\text{var}[x] = \mathbb{E}[x^2] - \mathbb{E}[x]^2 = \sigma^2$$



Binary Variables

Coin flipping: heads=1, tails=0

$$p(x = 1|\mu) = \mu$$

Bernoulli Distribution

$$\text{Bern}(x|\mu) = \mu^x(1 - \mu)^{1-x}$$

$$\mathbb{E}[x] = \mu$$

$$\text{var}[x] = \mu(1 - \mu)$$



Binary Variables

- N coin flips:

$$p(m \text{ heads} | N, \mu)$$

- Binomial Distribution

$$\text{Bin}(m|N, \mu) = \binom{N}{m} \mu^m (1 - \mu)^{N-m}$$

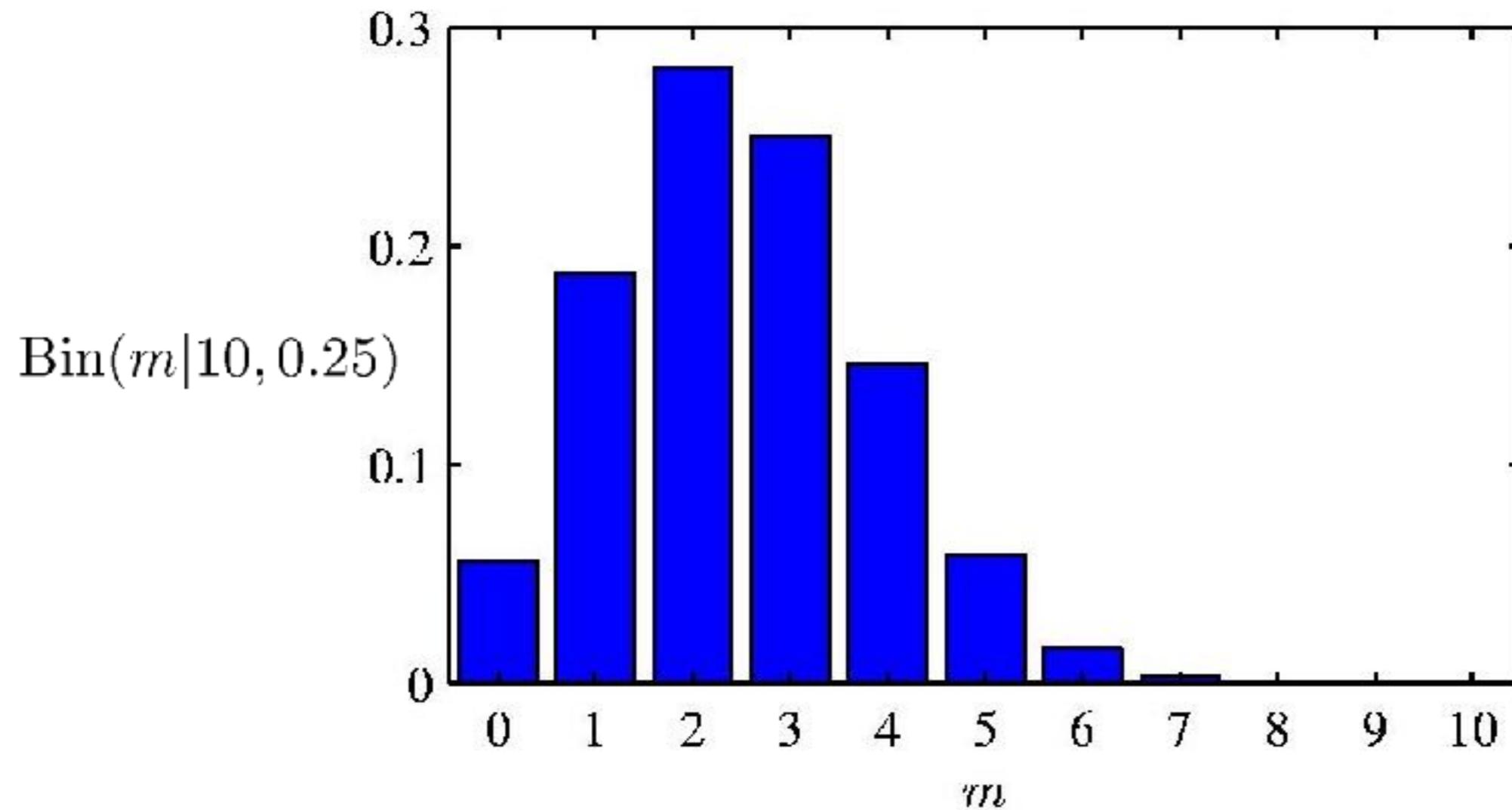
$$\mathbb{E}[m] \equiv \sum_{m=0}^N m \text{Bin}(m|N, \mu) = N\mu$$

$$\text{var}[m] \equiv \sum_{m=0}^N (m - \mathbb{E}[m])^2 \text{Bin}(m|N, \mu) = N\mu(1 - \mu)$$



Binomial Distribution

THE UNIVERSITY OF
SYDNEY



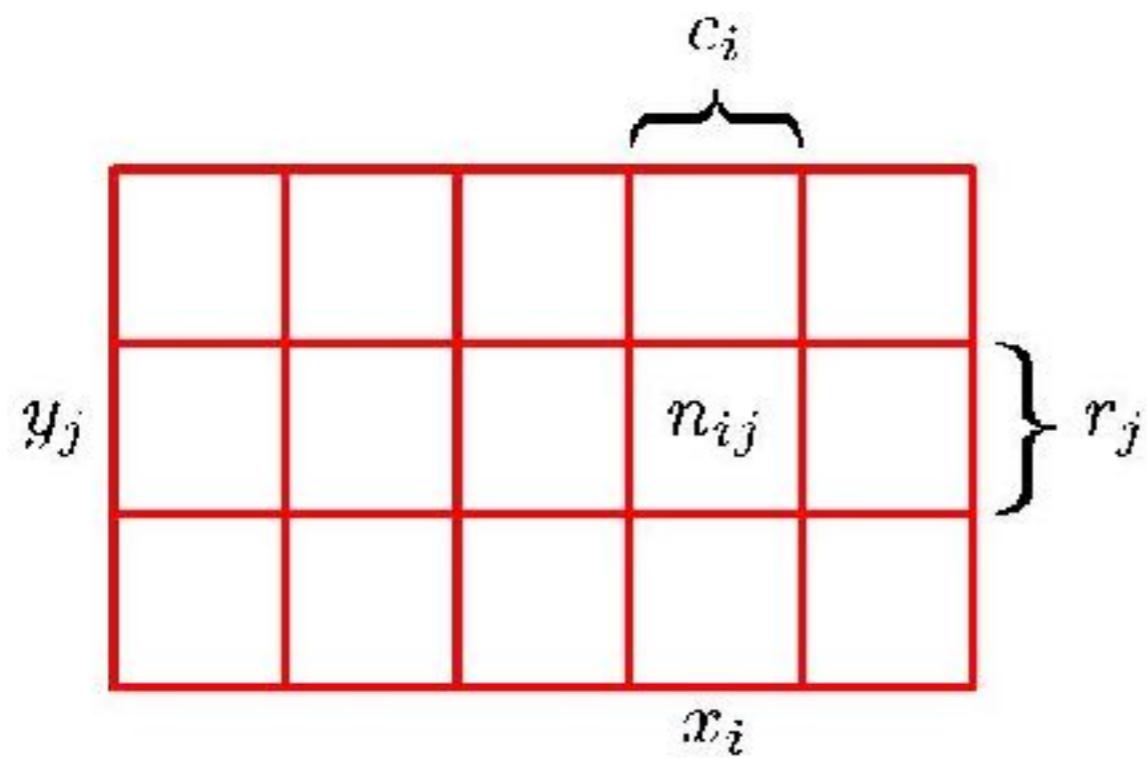


Conditional Probabilities

- One of the most important concepts in all of Machine Learning
- $P(A | B) = P(A, B)/P(B)$...assuming $P(B)$ not equal 0.
 - Conditional probability of A given B has occurred.
- Probability it will rain tomorrow given it has rained today.
 - $P(A | B) = P(A, B)/P(B) = 0.1/0.4 = 1/4 = 0.25$
 - In general $P(A | B)$ is not equal to $P(B | A)$



Probability Theory



Joint Probability

$$p(X = x_i, Y = y_i) = \frac{n_{ij}}{N}$$

Marginal Probability

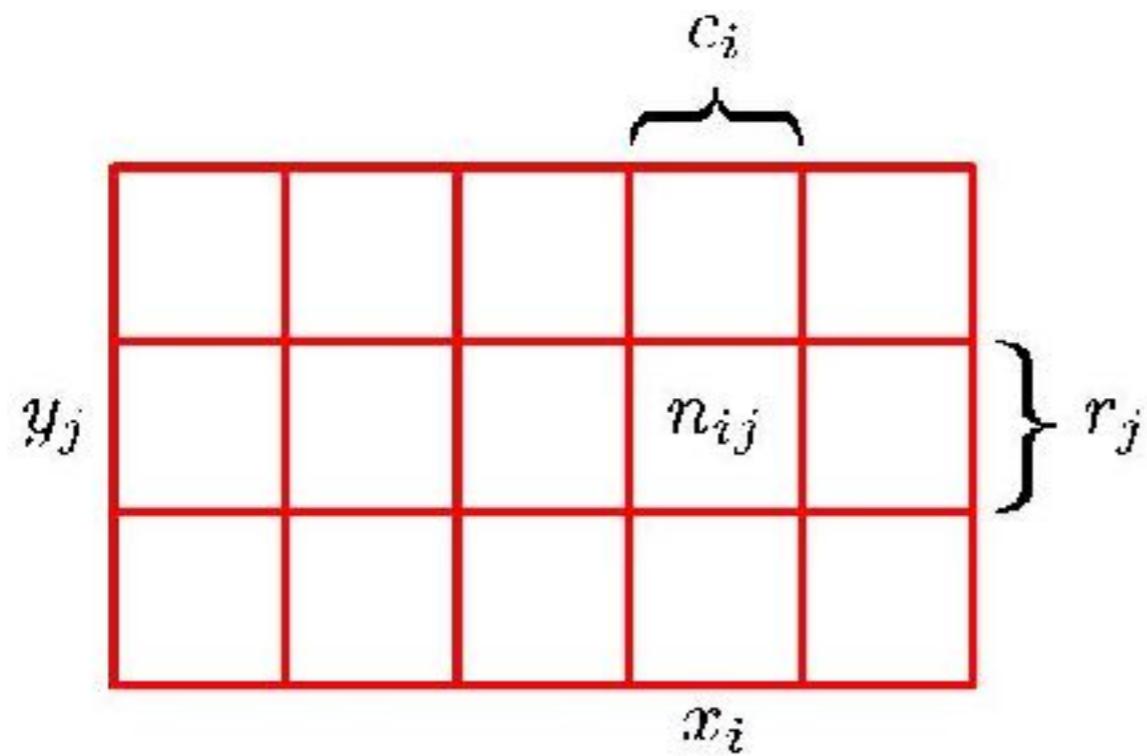
$$p(X = x_i) = \frac{c_i}{N}$$

Conditional Probability

$$p(Y = y_i | X = x_i) = \frac{n_{ij}}{c_i}$$



Probability Theory



Sum Rule

$$\begin{aligned} p(X = x_i) &= \frac{c_i}{N} = \frac{1}{N} \sum_{j=1}^L n_{ij} \\ &= \sum_{j=1}^L p(X = x_i, Y = y_j) \end{aligned}$$

Product Rule

$$\begin{aligned} p(X = x_i, Y = y_i) &= \frac{n_{ij}}{N} = \frac{n_{ij}}{c_i} \cdot \frac{c_i}{N} \\ &= p(Y = y_i | X = x_i) p(X = x_i) \end{aligned}$$



Rules of Probability

- Sum Rule

$$p(X) = \sum_Y p(X, Y)$$

- Product Rule

$$p(X, Y) = p(Y|X)p(X)$$



Bayes' Rule

- $P(A | B) = P(A, B) / P(B); P(B | A) = P(B, A) / P(A)$
- Now $P(A, B) = P(B, A)$
- Thus $P(A | B) P(B) = P(B | A) P(A)$
- Thus $P(A | B) = [P(B | A)P(A)] / [P(B)]$
 - This is called Bayes' Rule
 - Basis of almost all prediction
 - Latest theories hypothesise that human memory and action is Bayes' rule in action.



Bayes' Rule

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Diagram illustrating Bayes' Rule:

- Posterior**: Points to the term $P(A|B)$.
- Likelihood**: Points to the term $P(B|A)$.
- Prior**: Points to the term $P(A)$.
- Normaliser**: Points to the term $P(B)$.

$$P(hypothesis|data) = \frac{P(data|hypothesis)P(hypothesis)}{P(data)}$$



Example

The ASX market goes up 60% of the days of a year. 40% of the time it stays the same or goes down. The day the ASX is up, there is a 50% chance that the Shanghai Index is up. On other days there is 30% chance that Shanghai goes up. Suppose the Shanghai market is up. What is the probability that ASX was up?



Example cont.

- We want to calculate $P(A_1 | S_1)$?
- $P(A_1) = 0.6; P(A_2) = 0.4;$
 $P(S_1 | A_1) = 0.5; P(S_1 | A_2) = 0.3$
 $P(S_2 | A_1) = 1 - P(S_1 | A_1) = 0.5;$
 $P(S_2 | A_2) = 1 - P(S_1 | A_2) = 0.7;$
- $P(A_1 | S_1) = P(S_1 | A_1)P(A_1) / (P(S_1))$
- How do we calculate $P(S_1)$?



Example cont.

- $P(S_1) = P(S_1|A_1) + P(S_1|A_2)$ [Key Step]
 $= P(S_1|A_1)P(A_1) + P(S_1|A_2)P(A_2)$
 $= 0.5 \times 0.6 + 0.3 \times 0.4$
 $= 0.42$
- Finally,
 $P(A_1|S_1) = P(S_1|A_1)P(A_1) / P(S_1)$
 $= (0.5 \times 0.6) / 0.42$
 $= 0.71$



Independence

- Two events A and B are independent if

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

- Example: Toss a coin twice. Then what is the probability of two heads?
- The outcome of the two tosses are not dependent on each other

$$P(H,H) = P(H)P(H) = 0.5 \times 0.5 = 0.25$$

- If A and B are independent then

$$P(A | B) = P(A,B) / P(B) = P(A)P(B) / P(B) = P(A) !$$

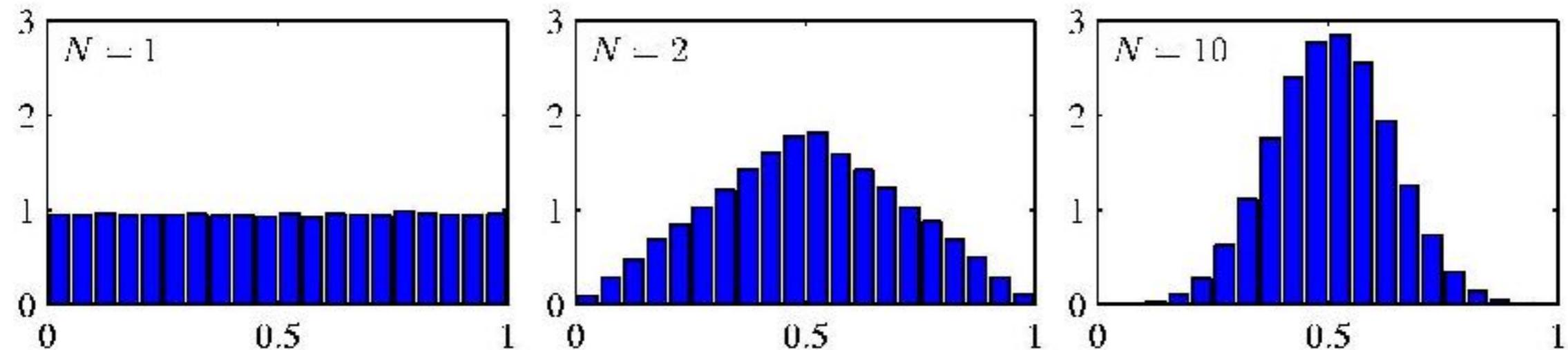


THE UNIVERSITY OF
SYDNEY

Central Limit Theorem

The distribution of the sum of N i.i.d. random variables becomes increasingly Gaussian as N grows.

Example: N uniform $[0,1]$ random variables.





Expectations

$$\mathbb{E}[f] = \sum_x p(x)f(x)$$

$$\mathbb{E}[f] = \int p(x)f(x)dx$$

$$\mathbb{E}_x[f|y] = \sum_x p(x|y)f(x)$$



Conditional Expectation
(discrete)

$$\mathbb{E}[f] \simeq \frac{1}{N} \sum_{n=1}^N f(x_n)$$

Approximate Expectation
(discrete and continuous)



THE UNIVERSITY OF
SYDNEY

The law of large numbers

LLN describes the result of performing the same experiment a large number of times.

The average of the results obtained from a large number of independent trials should converge to the expected value.

$$\frac{\sum_{I=1}^n \mathbf{1}_{\{x_i = \text{"head"}\}}}{n} \rightarrow \int P(x = \text{"head"}) \mathbf{1}_{\{x = \text{"head"}\}} dx \\ = P(x = \text{"head"})$$





Entropy

$$H[x] = - \sum_x p(x) \log_2 p(x)$$

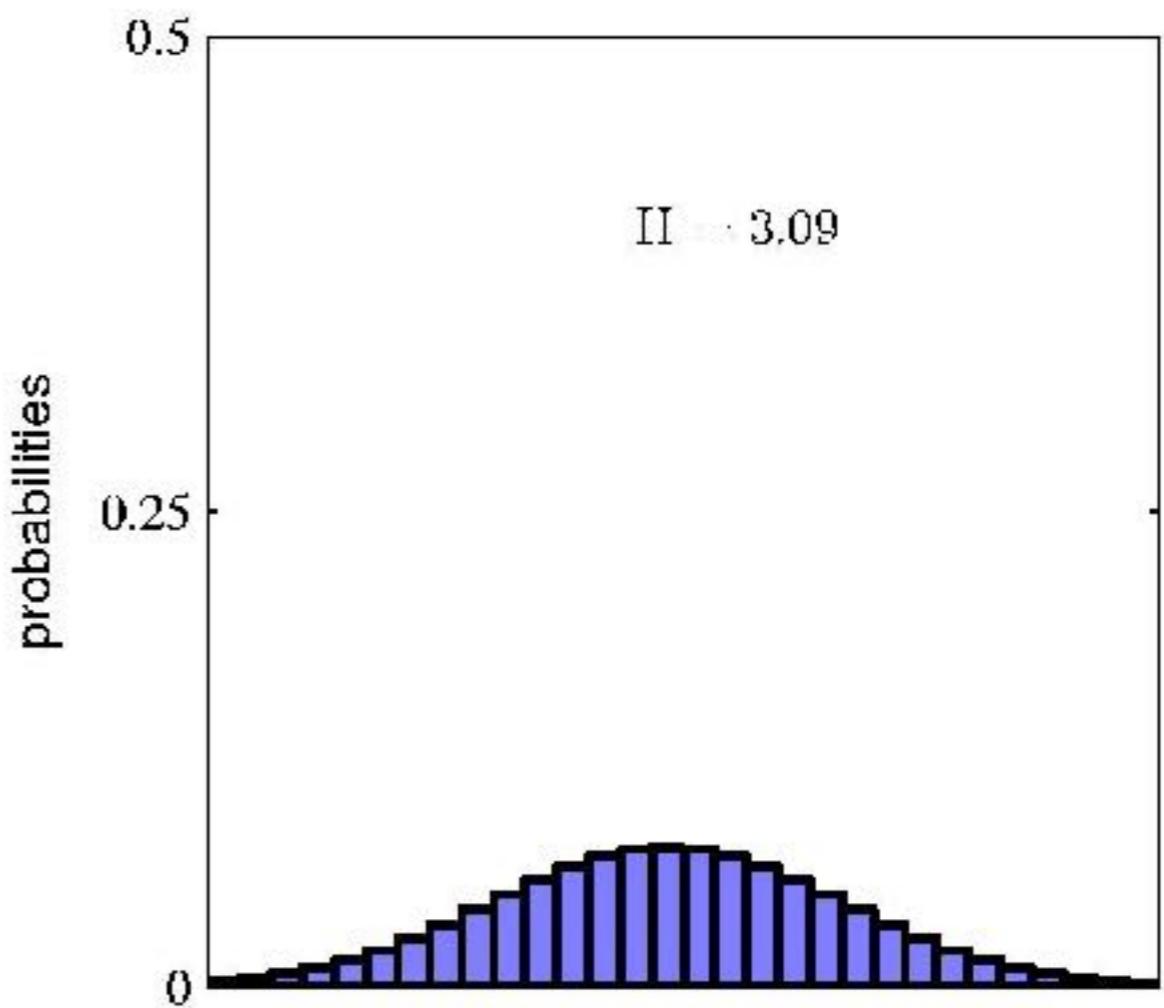
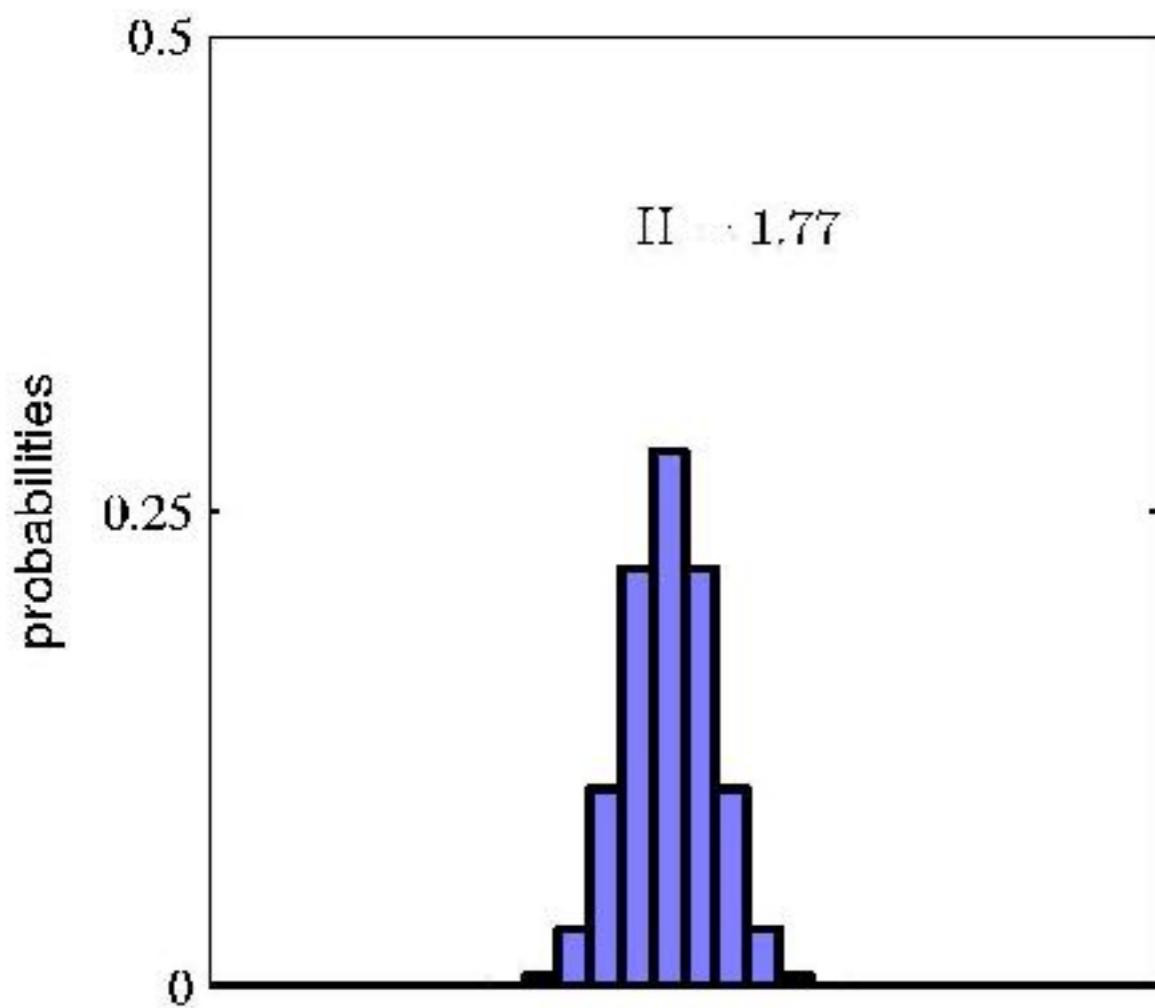
Important quantity in

- coding theory
- statistical physics
- machine learning



Entropy

THE UNIVERSITY OF
SYDNEY





Entropy: Example I

Coding theory: x discrete with 8 possible states; how many bits to transmit the state of x ?

All states equally likely

$$H[x] = -8 \times \frac{1}{8} \log_2 \frac{1}{8} = 3 \text{ bits.}$$



Entropy: Example 2

x	a	b	c	d	e	f	g	h
$p(x)$	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{16}$	$\frac{1}{64}$	$\frac{1}{64}$	$\frac{1}{64}$	$\frac{1}{64}$
code	0	10	110	1110	111100	111101	111110	111111

$$\begin{aligned} H[x] &= -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{4} \log_2 \frac{1}{4} - \frac{1}{8} \log_2 \frac{1}{8} - \frac{1}{16} \log_2 \frac{1}{16} - \frac{4}{64} \log_2 \frac{1}{64} \\ &= 2 \text{ bits} \end{aligned}$$

$$\begin{aligned} \text{average code length} &= \frac{1}{2} \times 1 + \frac{1}{4} \times 2 + \frac{1}{8} \times 3 + \frac{1}{16} \times 4 + 4 \times \frac{1}{64} \times 6 \\ &= 2 \text{ bits} \end{aligned}$$



Mutual Information

$$\begin{aligned} I[\mathbf{x}, \mathbf{y}] &\equiv \text{KL}(p(\mathbf{x}, \mathbf{y}) || p(\mathbf{x})p(\mathbf{y})) \\ &= - \iint p(\mathbf{x}, \mathbf{y}) \ln \left(\frac{p(\mathbf{x})p(\mathbf{y})}{p(\mathbf{x}, \mathbf{y})} \right) d\mathbf{x} d\mathbf{y} \end{aligned}$$

$$I[\mathbf{x}, \mathbf{y}] = H[\mathbf{x}] - H[\mathbf{x}|\mathbf{y}] = H[\mathbf{y}] - H[\mathbf{y}|\mathbf{x}]$$

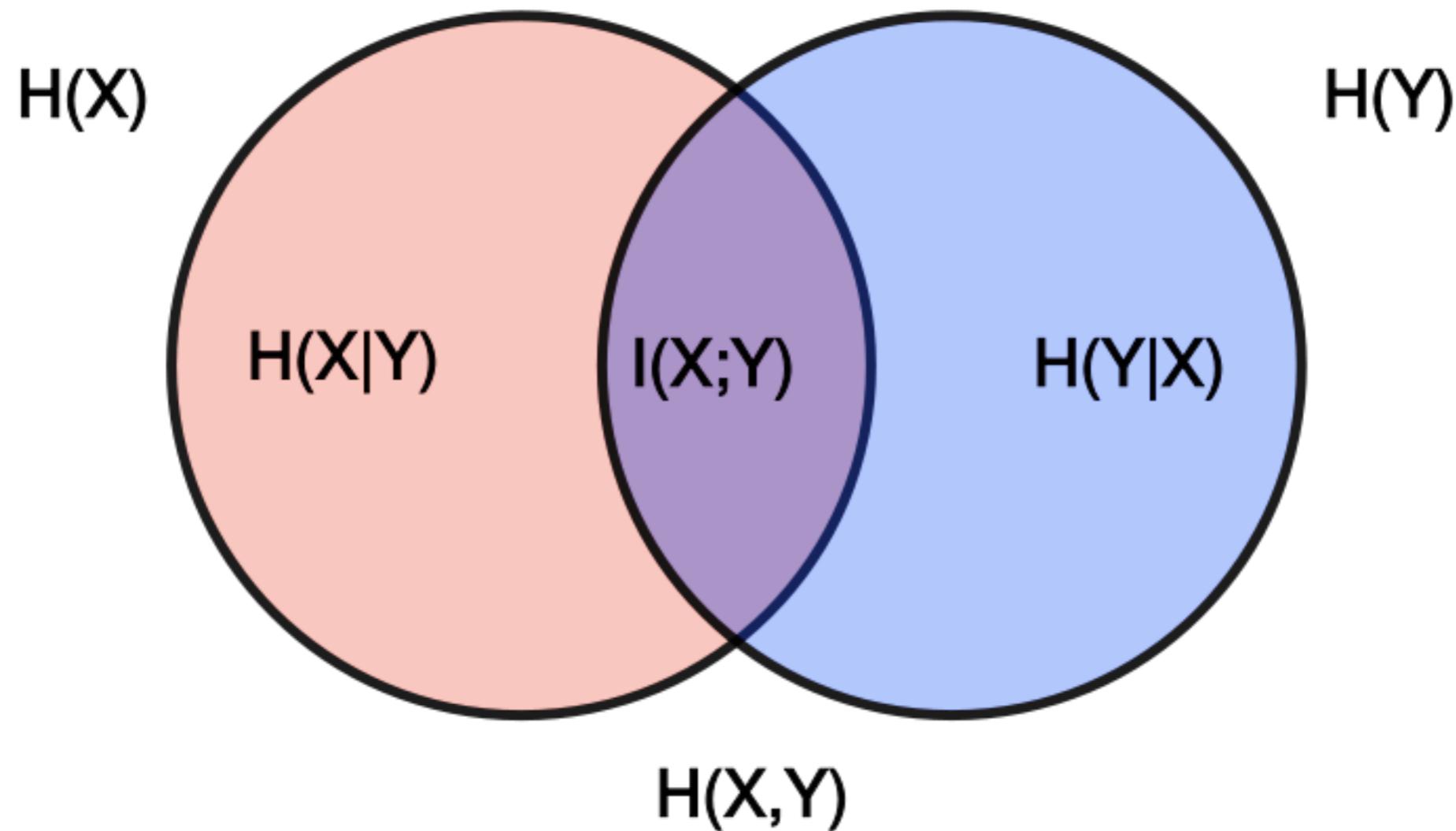
Two random variables \mathbf{x} and \mathbf{y} are independent *if and only if*,

$$I[\mathbf{x}, \mathbf{y}] = 0.$$



THE UNIVERSITY OF
SYDNEY

Mutual Information and Entropy

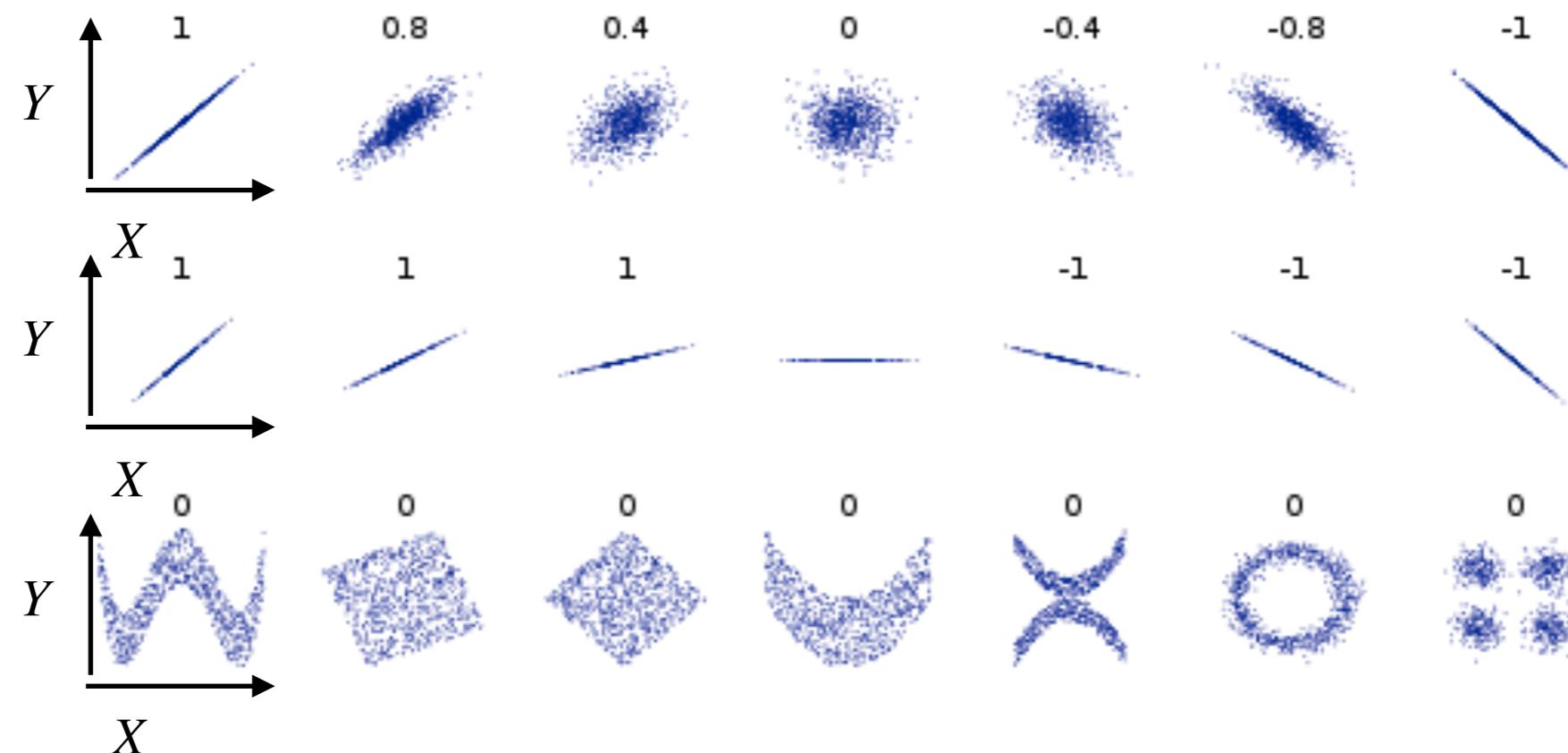




Correlation vs dependence

Correlation coefficient:

$$\rho_{X,Y} = \text{corr}(X, Y) = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{E[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \sigma_Y},$$



Even though the correlation coef is zero they are still dependent!



Exercise (homework)

A diagnostic test has a probability 0.95 of giving a positive result when applied to a person suffering from a certain disease, and a probability 0.10 of giving a (false) positive when applied to a non-sufferer. It is estimated that 0.5 % of the population are sufferers. Suppose that the test is now administered to a person about whom we have no relevant information relating to the disease (apart from the fact that he/she comes from this population). Calculate the following probabilities:

- (a) that the test result will be positive;
- (b) that, given a positive result, the person is a sufferer;
- (c) that, given a negative result, the person is a non-sufferer;
- (d) that the person will be misclassified.



THE UNIVERSITY OF
SYDNEY

Quiz Next Week

- Content: first three lectures and basic programming
- How to study
 - Lecture slides and labs
 - Examples and exercises in class
 - (“Mining of massive datasets” chapter 11)
 - (“Machine learning” chapters 1 and 2)



Quiz Next Week

- Time: 8pm-9pm
- Venue: tutorial labs
- If your tutorials are on Tuesday (or currently haven't been allocated to any tutorial room), please go to Codrington Computer Lab I to sit the quiz.



Quiz Next Week

- Codrington Computer Lab I:



Name:
Codrington Computer Laboratory I
Short Name:
Codrington Computer Lab I
Venue Code:
H69.151



THE UNIVERSITY OF
SYDNEY

Quiz Next Week

- Questions to TAs:



Jiayan Qiu (TA)

jqiu3225@uni.sydney.edu.au



Baosheng Yu (TA)

bayu0826@uni.sydney.edu.au