

# Information Theory

Sargur N. Srihari

# Topics

## 1. Entropy as an Information Measure

### 1. Discrete variable definition

Relationship to Code Length

### 2. Continuous Variable

Differential Entropy

## 2. Maximum Entropy

## 3. Conditional Entropy

## 4. Kullback-Leibler Divergence (Relative Entropy)

## 5. Mutual Information

# Information Measure

- How much information is received when we observe a specific value for a discrete random variable  $x$ ?
- Amount of information is degree of surprise
  - Certain means no information
  - More information when event is unlikely
- Depends on probability distribution  $p(x)$ , a quantity  $h(x)$
- If there are two unrelated events  $x$  and  $y$  we want  $h(x,y) = h(x) + h(y)$
- Thus we choose  $h(x) = -\log_2 p(x)$ 
  - Negative assures that information measure is positive
- Average amount of information transmitted is the expectation wrt  $p(x)$  referred to as entropy

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

# Usefulness of Entropy

- Uniform Distribution
  - Random variable  $x$  has 8 possible states, each equally likely
    - We would need 3 bits to transmit
    - Also,  $H(x) = -8 \times (1/8) \log_2(1/8) = 3 \text{ bits}$
- Non-uniform Distribution
  - If  $x$  has 8 states with probabilities  $(1/2, 1/4, 1/8, 1/16, 1/64, 1/64, 1/64, 1/64)$   
 $H(x) = 2 \text{ bits}$
- Non-uniform distribution has smaller entropy than uniform
- Has an interpretation of in terms of disorder

# Relationship of Entropy to Code Length

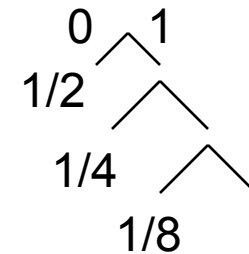
- Take advantage of non-uniform distribution to use shorter codes for more probable events

- If  $x$  has 8 states  $(a,b,c,d,e,f,g,h)$  with probabilities  $(1/2, 1/4, 1/8, 1/16, 1/64, 1/64, 1/64, 1/64)$

Can use codes  $0, 10, 110, 1110, 111100, 111110, 111111$

$$\text{average code length} = (1/2)1 + (1/4)2 + (1/8)3 + (1/16)4 + 4(1/64)6 = 2 \text{ bits}$$

- Same as entropy of the random variable
- Shorter code string is not possible due to need to disambiguate string into component parts
- $11001110$  is uniquely decoded as sequence  $cad$



# Relationship between Entropy and Shortest Coding Length

- Noiseless coding theorem of Shannon
  - Entropy is a lower bound on number of bits needed to transmit a random variable
- Natural logarithms are used in relationship to other topics
  - Nats instead of bits

# History of Entropy: thermodynamics to information theory

- Entropy is average amount of information needed to specify state of a random variable
- Concept had much earlier origin in physics
  - Context of equilibrium thermodynamics
  - Later given deeper interpretation as measure of disorder (developments in statistical mechanics)

# Entropy Persona

## World of Atoms



- Ludwig Eduard Boltzmann (1844-1906)
  - Created Statistical Mechanics
    - First law: conservation of energy
      - Energy not destroyed but converted from one form to other
    - Second law: principle of decay in nature— entropy increases
      - Explains why not all energy is available to do useful work
    - Relate macro state to statistical behavior of microstate
- Claude Shannon (1916-2001)
- Stephen Hawking (Gravitational Entropy)

## World of Bits





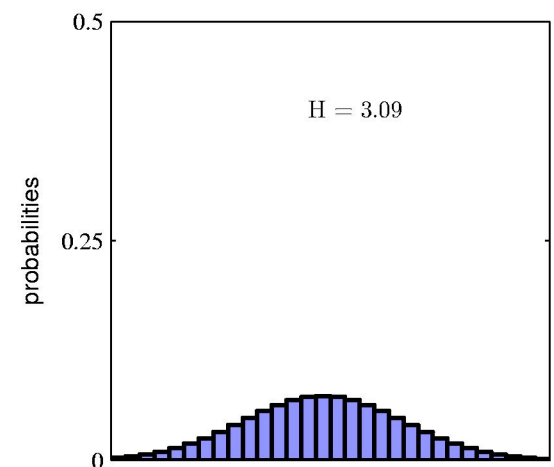
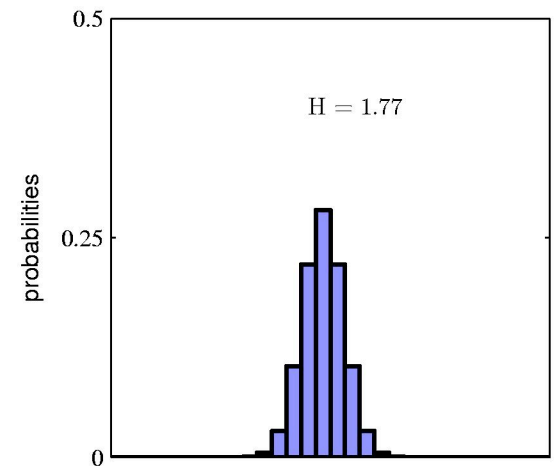
# Physics view of Entropy

- $N$  objects into bins so that  $n_i$  are in  $i^{\text{th}}$  bin where  $\sum_i n_i = N$
- No of different ways of allocating objects to bins
  - $N$  ways to choose first,  $N-1$  ways for second leads to  $N.(N-1) \dots 2.1 = N!$
  - We don't distinguish between rearrangements within each bin
    - In  $i^{\text{th}}$  bin there are  $n_i!$  ways of reordering objects
  - Total no of ways of allocating  $N$  objects to bins is  $W = \frac{N!}{\prod_i n_i!}$ 
    - Called Multiplicity (also weight of macrostate)
- Entropy: scaled log of multiplicity  $H = \frac{1}{N} \ln W = \frac{1}{N} \ln N! - \frac{1}{N} \sum_i \ln n_i!$ 
  - Sterlings approx as  $N \rightarrow \infty \ln N! \approx N \ln N - N$
  - Which gives  $H = - \lim_{N \rightarrow \infty} \sum_i \left( \frac{n_i}{N} \right) \ln \left( \frac{n_i}{N} \right) = - \sum_i p_i \ln p_i$
- Overall distribution, as ratios  $n_i/N$ , called *macrostate*
- Specific arrangement of objects in bin is *microstate*

# Entropy and Histograms

- If  $X$  can take one of  $M$  values (bins, states) and  $p(X=x_i)=p_i$  then 
$$H(p) = -\sum_i p_i \ln p_i$$
- Minimum value of entropy is 0 when one of the  $p_i=1$  and other  $p_i$  are 0
  - noting that  $\lim_{p \rightarrow 0} p \ln p = 0$
- Sharply peaked distribution has low entropy
- Distribution spread more evenly will have higher entropy

30 bins, higher value for broader distribution



# Maximum Entropy Configuration

- Found by maximizing  $H$  using Lagrange multiplier to enforce constraint of probabilities
- Maximize 
$$\tilde{H} = -\sum_i p(x_i) \ln p(x_i) + \lambda \left( \sum_i p(x_i) - 1 \right)$$
- Solution: all  $p(x_j)$  are equal or  $p(x_j) = 1/M$   $M = \text{no of states}$
- Maximum value of entropy is  $\ln M$
- To verify it is a maximum, evaluate second derivative of entropy 
$$\frac{\partial^2 \tilde{H}}{\partial p(x_i) \partial p(x_j)} = -I_{ij} \frac{1}{p_i}$$
  - where  $I_{ij}$  are elements of identity matrix

# Entropy with Continuous Variable

- Divide  $x$  into bins of width  $\Delta$
- For each bin there must exist a value  $x_i$  such that  $\int_{i\Delta}^{(i+1)\Delta} p(x)d(x) = p(x_i)\Delta$
- Gives a discrete distribution with probabilities  $p(x_i)\Delta$
- Entropy  $H_\Delta = -\sum_i p(x_i)\Delta \ln(p(x_i)\Delta) = -\sum_i p(x_i)\Delta \ln p(x_i) - \ln \Delta$
- Omit the second term and consider the limit  $\Delta \rightarrow 0$

$$H_\Delta = -\int p(x) \ln p(x) dx$$

- Known as Differential Entropy
- *Discrete and Continuous forms of entropy differ by quantity  $\ln \Delta$  which diverges*
  - Reflects to specify continuous variable very precisely requires a large no of bits

# Entropy with Multiple Continuous Variables

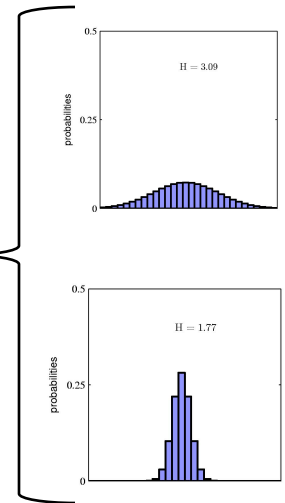
- Differential Entropy for multiple continuous variables

$$H(\mathbf{x}) = -\int p(\mathbf{x}) \ln p(\mathbf{x}) d\mathbf{x}$$

- For what distribution is differential entropy maximized?
  - For discrete distribution, it is uniform
  - For continuous, it is Gaussian
    - as we shall see

# Entropy as a Functional

- Ordinary calculus deals with functions
- A *functional* is an operator that takes a function as input and returns a scalar
- A widely used functional in machine learning is entropy  $H[p(x)]$  which is a scalar quantity
- We are interested in the maxima and minima of functionals analogous to those for functions
  - Called calculus of variations



# Maximizing a Functional

- Functional: mapping from set of functions to real value
- For what function is it maximized?
- Finding shortest curve length between two points on a sphere (geodesic)
  - With no constraints it is a straight line
  - When constrained to lie on a surface solution is less obvious— may be several
- Constraints incorporated using Lagrangian

# Maximizing Differential Entropy

- Assuming constraints on first and second moments of  $p(x)$  as well as normalization

$$\int p(x)dx = 1 \quad \int xp(x)dx = \mu \quad \int (x - \mu)^2 p(x)dx = \sigma^2$$

- Constrained maximization is performed using Lagrangian multipliers. Maximize following functional

$$\begin{aligned} \text{wrt } p(x): \quad & - \int p(x) \ln p(x) dx + \lambda_1 (\int p(x) dx - 1) \\ & + \lambda_2 (\int xp(x) dx - \mu) + \lambda_3 (\int (x - \mu)^2 p(x) dx - \sigma^2) \end{aligned}$$

- Using the calculus of variations derivative of functional is set to zero giving  $p(x) = \exp\{-1 + \lambda_1 + \lambda_2 x + \lambda_3 (x - \mu)^2\}$
- Backsubstituting into three constraint equations leads to the result that distribution that maximizes differential entropy is Gaussian!



# Differential Entropy of Gaussian

- Distribution that maximizes Differential Entropy is Gaussian

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

- Value of maximum entropy is

$$H(x) = \frac{1}{2} \{1 + \ln(2\pi\sigma^2)\}$$

- Entropy increases as variance increases
- Differential entropy, unlike discrete entropy, can be negative for  $\sigma^2 < 1/2\pi e$

# Conditional Entropy

- If we have joint distribution  $p(x,y)$ 
  - We draw pairs of values of  $x$  and  $y$
  - If value of  $x$  is already known, additional information to specify corresponding value of  $y$  is  $-\ln p(y|x)$
- Average additional information needed to specify  $y$  is the conditional entropy

$$H[y | x] = -\int \int p(y | x) \ln p(y | x) dy dx$$

- By product rule  $H[x,y] = H[y|x] + H[x]$ 
  - where  $H[x,y]$  is differential entropy of  $p(x,y)$
  - $H[x]$  is differential entropy of  $p(x)$
  - Information needed to describe  $x$  and  $y$  is given by information needed to describe  $x$  plus additional information needed to specify  $y$  given  $x$

# Relative Entropy

- If we have modeled unknown distribution  $p(\mathbf{x})$  by approximating distribution  $q(\mathbf{x})$ 
  - i.e.,  $q(\mathbf{x})$  is used to construct a coding scheme of transmitting values of  $\mathbf{x}$  to a receiver
  - Average additional amount of information required to specify value of  $\mathbf{x}$  as a result of using  $q(\mathbf{x})$  instead of true distribution  $p(\mathbf{x})$  is given by relative entropy or K-L divergence
- Important concept in Bayesian analysis
  - Entropy comes from Information Theory
  - *K-L Divergence*, or *relative entropy*, comes from Pattern Recognition, since it is a distance (dissimilarity) measure

# Relative Entropy or K-L Divergence

- Additional information required as a result of using  $q(\mathbf{x})$  in place of  $p(\mathbf{x})$

$$\begin{aligned} KL(p \parallel q) &= - \int p(x) \ln q(x) dx - \left( \int p(x) \ln p(x) dx \right) \\ &= - \int p(x) \ln \left\{ \frac{p(x)}{q(x)} \right\} dx \end{aligned}$$

- Not a symmetrical quantity:  $KL(p \parallel q) \neq KL(q \parallel p)$
- K-L divergence satisfies  $KL(p \parallel q) > 0$  with equality iff  $p(x) = q(x)$ 
  - Proof involves convex functions

# Convex Function

- A function  $f(x)$  is convex if every chord lies on or above function

- Any value of  $x$  in interval from  $x=a$  to  $x=b$  can be written as  $\lambda a + (1-\lambda)b$  where  $0 \leq \lambda \leq 1$

- Corresponding point on chord is

$$\lambda f(a) + (1-\lambda)f(b)$$

- Convexity implies

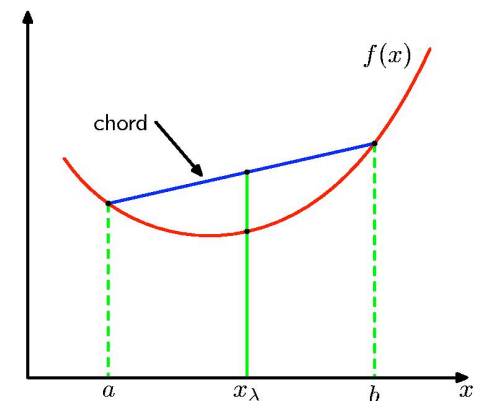
$$f(\lambda a + (1-\lambda)b) \leq \lambda f(a) + (1-\lambda)f(b)$$

Point on curve  $\leq$  Point on chord

- By induction, we get Jensen's inequality

$$f\left(\sum_{i=1}^M \lambda_i x_i\right) \leq \sum_{i=1}^M \lambda_i f(x_i)$$

where  $\lambda_i \geq 0$  and  $\sum_i \lambda_i = 1$



# Proof of positivity of K-L Divergence

- Applying to KL definition yields desired result

$$KL(p \parallel q) = - \int p(x) \ln \left\{ \frac{p(x)}{q(x)} \right\} dx \geq - \ln \int q(x) dx = 0$$
$$\geq 0$$

# Mutual Information

- Given joint distribution of two sets of variables  $p(\mathbf{x}, \mathbf{y})$ 
  - If independent, will factorize as  $p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x})p(\mathbf{y})$
  - If not independent, whether close to independent is given by

- KL divergence between joint and product of marginals

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

- Called Mutual Information between variables  $\mathbf{x}$  and  $\mathbf{y}$

# Mutual Information

- Using Sum and Product Rules

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

- Mutual Information is reduction in uncertainty about  $\mathbf{x}$  given value of  $\mathbf{y}$  (or vice versa)

- Bayesian perspective:

- if  $p(\mathbf{x})$  is prior and  $p(\mathbf{x}|\mathbf{y})$  is posterior, mutual information is reduction in uncertainty after  $\mathbf{y}$  is observed