

# 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

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

# History of Entropy: thermodynamics to information theory



- Boltzmann (1884-1906)
  - Created Statistical Mechanics
  - Explain why not all energy is available to do useful work
  - Relate macro state to statistical behavior of microstate
- Shannon (1916-2001)

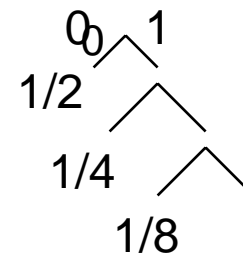
# Useful Property of Entropy: Relationship to Code Length

- Random variable  $x$  has 8 possible states, each equally likely
  - We would need 3 bits to transmit
  - Also,  $H(x) = -8 \times (1/8) \log_2(1/8) = 3$  bits
- If  $x$  has 8 states with probabilities  $(1/2, 1/4, 1/8, 1/16, 1/64, 1/64, 1/64, 1/64)$ 

$H(x) = 2$  bits

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

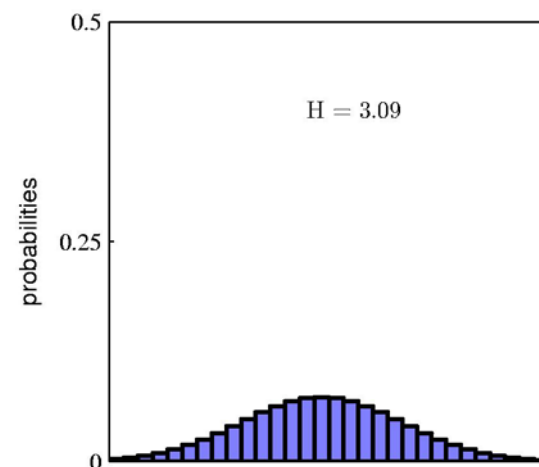
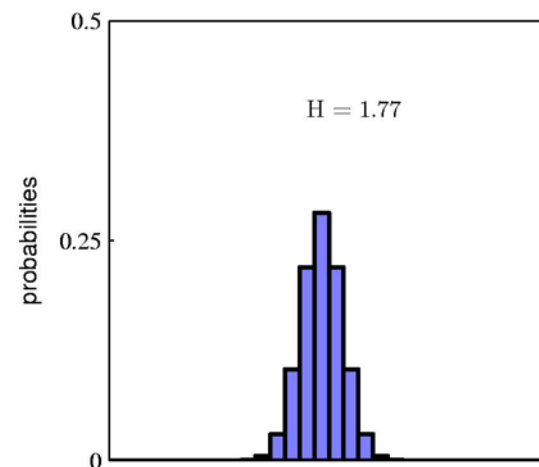
average code length =  $(1/2)1 + (1/4)2 + (1/8)3 + (1/16)4 + 4(1/64)6 = 2$  bits
- Noiseless coding theorem of Shannon
  - Entropy is a lower bound on the no of bits needed to transmit a random variable
- Natural logarithms are used in relationship to other topics
  - Nats instead of bits



# 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  $\lim_{p \rightarrow 0} p \ln p = 0$
- Maximum value of entropy is  $\ln M$ 
  - when all  $M$  values are equally likely

30 bins, higher value for broader distribution



# Entropy with Continuous Variable

- Divide  $x$  into bins of width  $\Delta$
- Gives a discrete distribution with probabilities  $p(x_i)\Delta$
- Entropy  $H_{\Delta} = -\sum_i p(x_i)\Delta \ln p(x_i) - \ln \Delta$
- 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

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

# Entropy as a Functional

- 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

# 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 wrt  $p(x)$ :

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

- Leads to 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

# Conditional Entropy

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

$$H[\mathbf{y} | \mathbf{x}] = -\int \int p(\mathbf{y} | \mathbf{x}) \ln p(\mathbf{y} | \mathbf{x}) d\mathbf{y}d\mathbf{x}$$

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

# Kullback-Leibler (KL) 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
- 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 K-L divergence

$$KL(p//q)$$

# K-L Definition

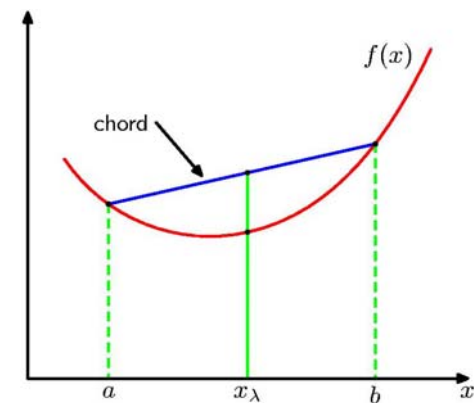
- Measure of dissimilarity between distributions  $p(\mathbf{x})$  and  $q(\mathbf{x})$

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

- 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

# Proof of Positivity of K-L

- 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  $f(\lambda a + (1-\lambda)b) \leq \lambda f(a) + (1-\lambda)f(b)$
- By induction leads to Jensen's inequality
- Applying to KL definition yields desired result



# 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
$$I[\mathbf{x}, \mathbf{y}] = KL(p(\mathbf{x}, \mathbf{y}) \parallel p(\mathbf{x})p(\mathbf{y}))$$
$$= \int \int p(\mathbf{x}, \mathbf{y}) \ln \left( \frac{p(\mathbf{x})p(\mathbf{y})}{p(\mathbf{x}, \mathbf{y})} \right) d\mathbf{x}d\mathbf{y}$$
- 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  $\mathbf{x}$  is observed