

# Probability and Information Theory

Lecture slides for Chapter 3 of *Deep Learning*  
[www.deeplearningbook.org](http://www.deeplearningbook.org)  
Ian Goodfellow  
2016-09-26

adapted by m.n. for CMPS 392

# Probability

- Sample space  $\Omega$ : set of all outcomes of a random experiment
- Set of events  $\mathcal{F}$  : collection of possible outcomes of an experiment.
- Probability measure:  $P: \mathcal{F} \rightarrow \mathbb{R}$ 
  - Axioms of probability
    - $P(A) \geq 0$  for all  $A \in \mathcal{F}$
    - $P(\Omega) = 1$
    - If  $A_1, A_2, \dots$  are disjoint events then

$$P\left(\bigcup_i A_i\right) = \sum_i P(A_i)$$

# Random variable

- Consider an experiment in which we flip 10 coins, and we want to know the number of coins that come up heads.
- Here, the elements of the sample space  $\Omega$  are 10-length sequences of heads and tails.
- For example, we might have

$$w_0 = \langle H, H, T, H, T, H, H, T, T, T \rangle$$

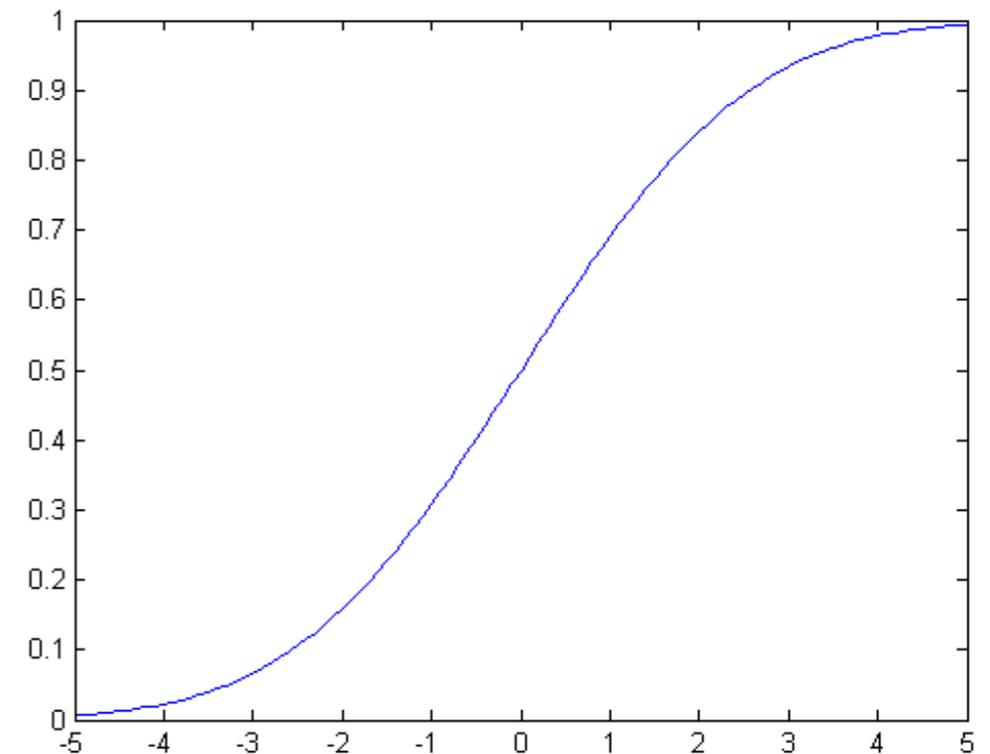
- However, in practice, we usually do not care about the probability of obtaining any particular sequence of heads and tails.
- Instead we usually care about real-valued functions of outcomes, such as
  - the number of heads that appear among our 10 tosses,
  - or the length of the longest run of tails.
- These functions, under some technical conditions, are known as random variables:  $X: \Omega \rightarrow \mathbb{R}$

# Discrete vs. continuous

- Discrete random variable:
  - $P(X = k) = P(\{\omega: X(\omega) = k\})$
- Continuous random variable:
  - $P(a \leq X \leq b) = P(\{\omega: a \leq X(\omega) \leq b\})$

A cumulative distribution function (CDF):

$$P(X \leq x)$$



# Probability Mass Function (discrete variable)

- The domain of  $P$  must be the set of all possible states of  $\mathbf{x}$ .
- $\forall x \in \mathbf{x}, 0 \leq P(x) \leq 1$ . An impossible event has probability 0 and no state can be less probable than that. Likewise, an event that is guaranteed to happen has probability 1, and no state can have a greater chance of occurring.
- $\sum_{x \in \mathbf{x}} P(x) = 1$ . We refer to this property as being **normalized**. Without this property, we could obtain probabilities greater than one by computing the probability of one of many events occurring.

Example: uniform distribution:  $P(\mathbf{x} = x_i) = \frac{1}{k}$

# Probability Density Function (continuous variable)

- The domain of  $p$  must be the set of all possible states of  $\mathbf{x}$ .
- $\forall x \in \mathbf{x}, p(x) \geq 0$ . Note that we do not require  $p(x) \leq 1$ .
- $\int p(x)dx = 1$ .

Example: uniform distribution:  $u(x; a, b) = \frac{1}{b-a}$ .

The pdf at some point  $x$  is not the probability of  $x$ :  $p(x) \neq P(\mathbf{x} = x)$

# Computing Marginal Probability with the Sum Rule

$$\forall x \in \mathbf{x}, P(\mathbf{x} = x) = \sum_y P(\mathbf{x} = x, y = y). \quad (3.3)$$

$$p(x) = \int p(x, y) dy. \quad (3.4)$$

# Conditional Probability

$$P(y = y \mid \mathbf{x} = \mathbf{x}) = \frac{P(y = y, \mathbf{x} = \mathbf{x})}{P(\mathbf{x} = \mathbf{x})}. \quad (3.5)$$

# Chain Rule of Probability

$$P(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}) = P(\mathbf{x}^{(1)}) \prod_{i=2}^n P(\mathbf{x}^{(i)} \mid \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(i-1)}). \quad (3.6)$$

# Independence

$$\forall x \in \mathbf{x}, y \in \mathbf{y}, p(\mathbf{x} = x, \mathbf{y} = y) = p(\mathbf{x} = x)p(\mathbf{y} = y). \quad (3.7)$$

# Conditional Independence

$$\forall x \in \mathbf{x}, y \in \mathbf{y}, z \in \mathbf{z}, p(\mathbf{x} = x, \mathbf{y} = y \mid \mathbf{z} = z) = p(\mathbf{x} = x \mid \mathbf{z} = z)p(\mathbf{y} = y \mid \mathbf{z} = z). \quad (3.8)$$

# Expectation

$$\mathbb{E}_{\mathbf{x} \sim P}[f(x)] = \sum_x P(x) f(x), \quad (3.9)$$

$$\mathbb{E}_{\mathbf{x} \sim p}[f(x)] = \int p(x) f(x) dx. \quad (3.10)$$

linearity of expectations:

$$\mathbb{E}_{\mathbf{x}}[\alpha f(x) + \beta g(x)] = \alpha \mathbb{E}_{\mathbf{x}}[f(x)] + \beta \mathbb{E}_{\mathbf{x}}[g(x)], \quad (3.11)$$

# Variance and Covariance

$$\text{Var}(f(x)) = \mathbb{E} \left[ (f(x) - \mathbb{E}[f(x)])^2 \right]. \quad (3.12)$$

$$\text{Cov}(f(x), g(y)) = \mathbb{E} [(f(x) - \mathbb{E}[f(x)]) (g(y) - \mathbb{E}[g(y)])]. \quad (3.13)$$

Covariance matrix:

$$\text{Cov}(\mathbf{x})_{i,j} = \text{Cov}(x_i, x_j). \quad (3.14)$$

# Bernoulli Distribution

$$P(\mathbf{x} = 1) = \phi \quad (3.16)$$

$$P(\mathbf{x} = 0) = 1 - \phi \quad (3.17)$$

$$P(\mathbf{x} = x) = \phi^x (1 - \phi)^{1-x} \quad (3.18)$$

$$\mathbb{E}_{\mathbf{x}}[\mathbf{x}] = \phi \quad (3.19)$$

$$\text{Var}_{\mathbf{x}}(\mathbf{x}) = \phi(1 - \phi) \quad (3.20)$$

# Gaussian Distribution

Parametrized by variance:

$$\mathcal{N}(x; \mu, \sigma^2) = \sqrt{\frac{1}{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right). \quad (3.21)$$

Parametrized by precision:

$$\mathcal{N}(x; \mu, \beta^{-1}) = \sqrt{\frac{\beta}{2\pi}} \exp\left(-\frac{1}{2}\beta(x - \mu)^2\right). \quad (3.22)$$

# Gaussian Distribution

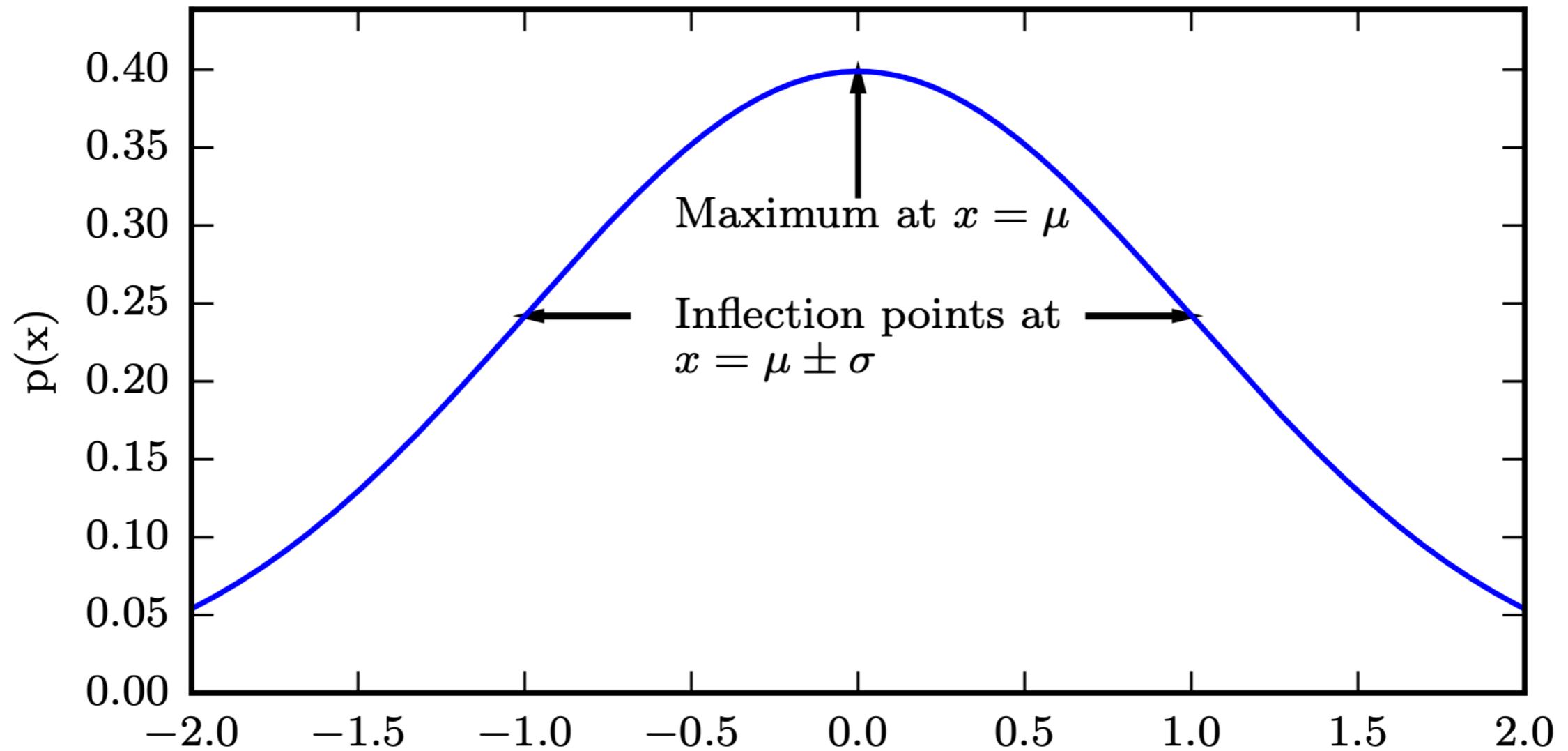


Figure 3.1

# Multivariate Gaussian

Parametrized by covariance matrix:

$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sqrt{\frac{1}{(2\pi)^n \det(\boldsymbol{\Sigma})}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right). \quad (3.23)$$

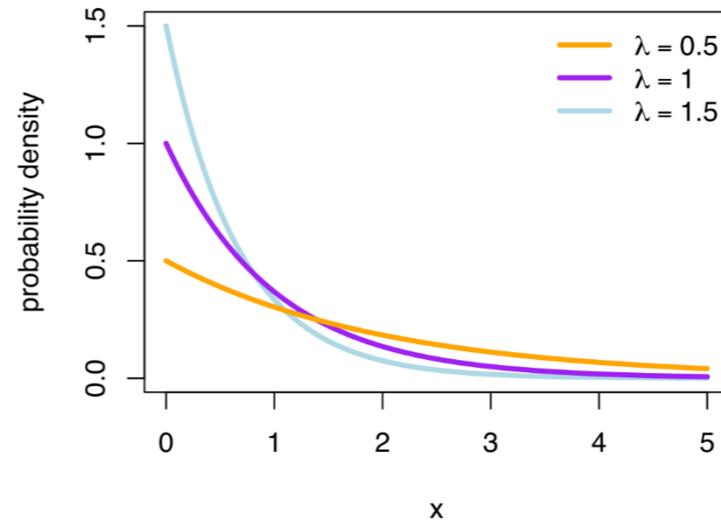
Parametrized by precision matrix:

$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\beta}^{-1}) = \sqrt{\frac{\det(\boldsymbol{\beta})}{(2\pi)^n}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\beta}(\mathbf{x} - \boldsymbol{\mu})\right). \quad (3.24)$$

# More Distributions

Exponential:

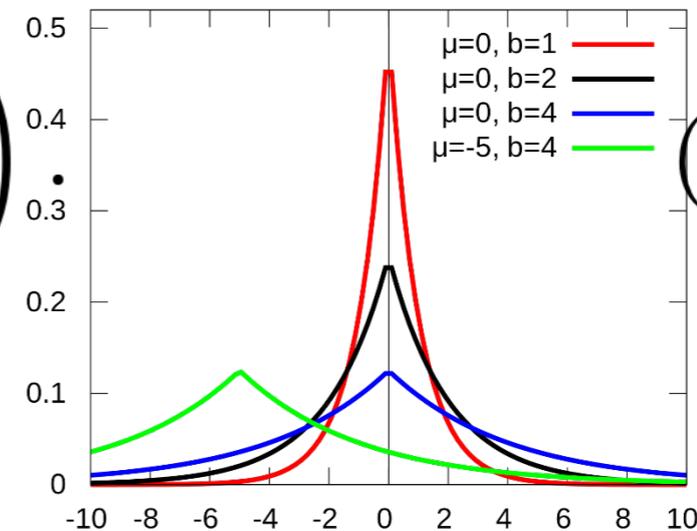
$$p(x; \lambda) = \lambda \mathbf{1}_{x \geq 0} \exp(-\lambda x)$$



$$(3.25)$$

Laplace:

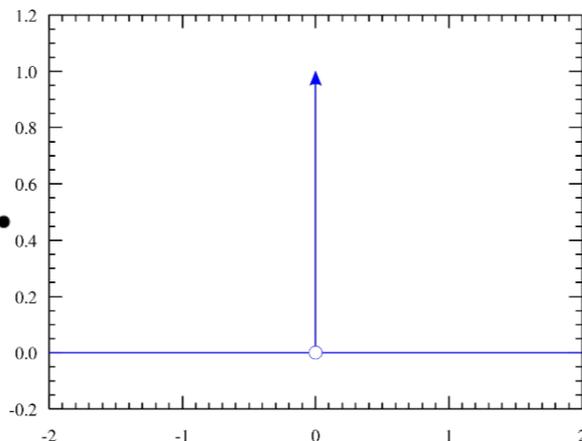
$$\text{Laplace}(x; \mu, \gamma) = \frac{1}{2\gamma} \exp\left(-\frac{|x - \mu|}{\gamma}\right)$$



$$(3.26)$$

Dirac:

$$p(x) = \delta(x - \mu)$$



$$(3.27)$$

# Empirical Distribution

$$\hat{p}(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m \delta(\mathbf{x} - \mathbf{x}^{(i)}) \quad (3.28)$$

# Mixture Distributions

$$P(\mathbf{x}) = \sum_i P(c = i)P(\mathbf{x} | c = i) \quad (3.29)$$

Gaussian mixture  
with three  
components

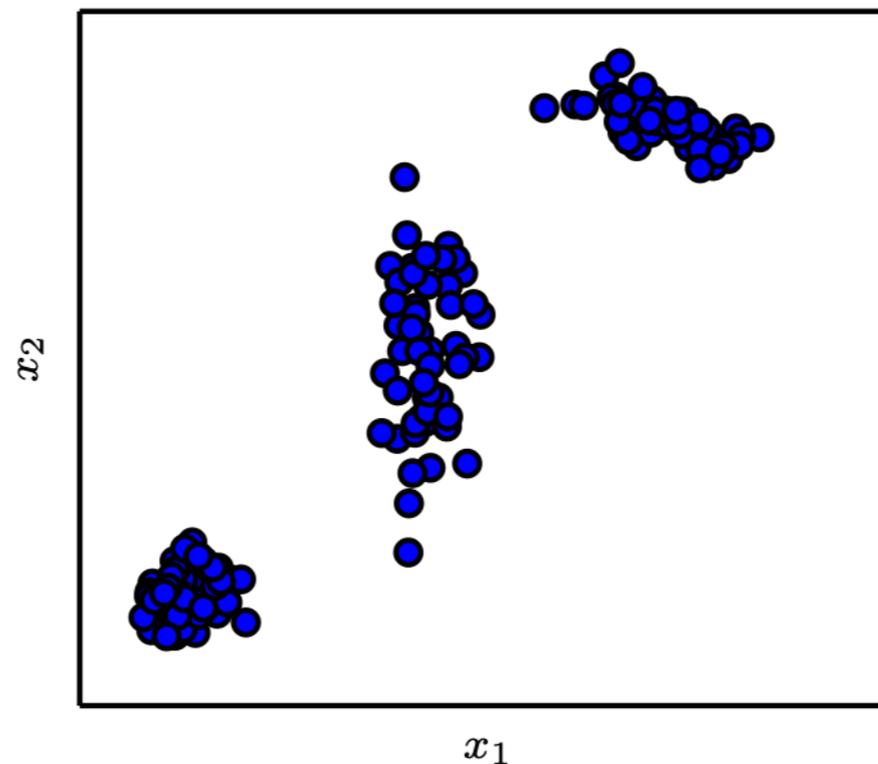


Figure 3.2

# Logistic Sigmoid

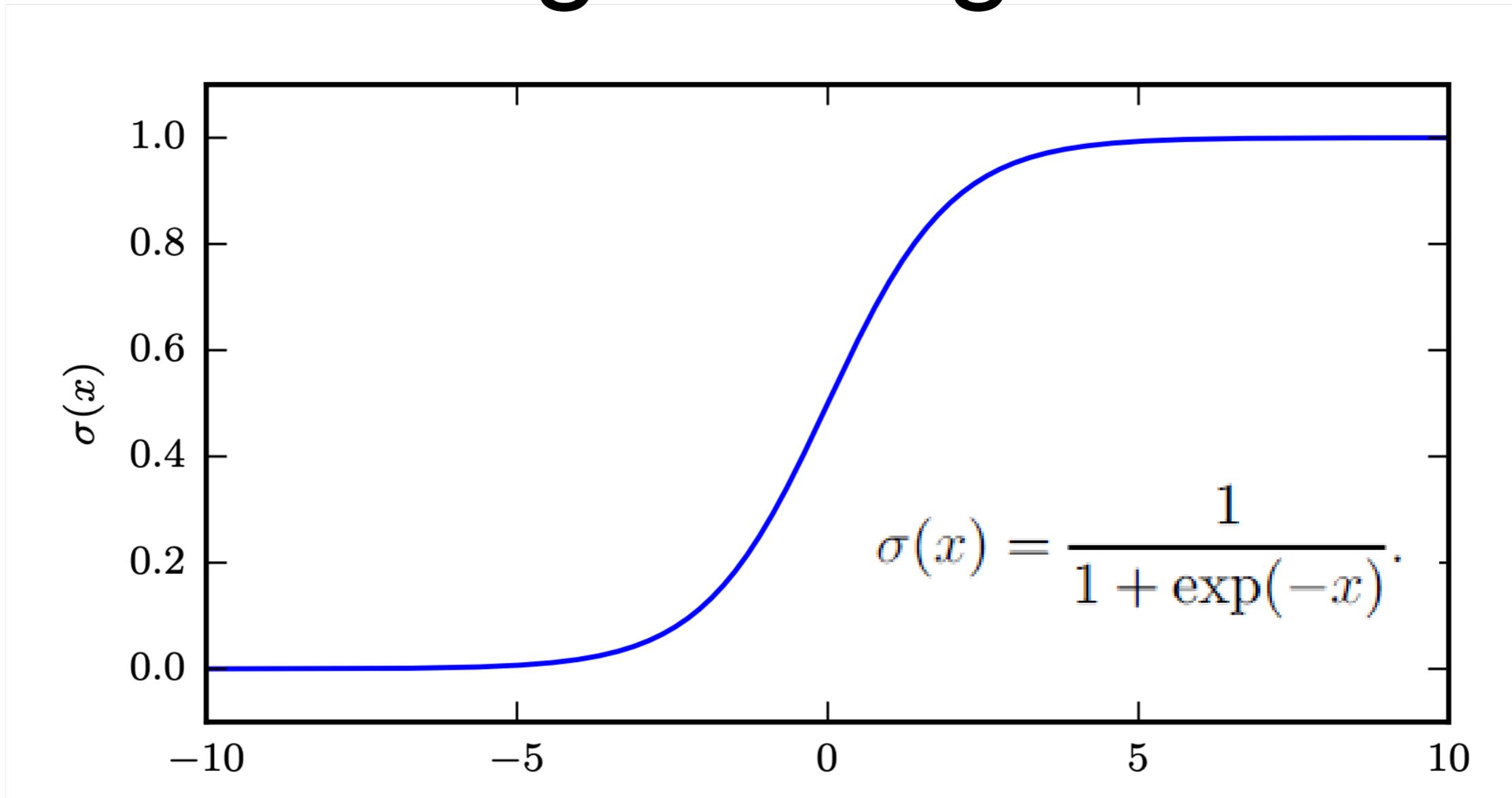
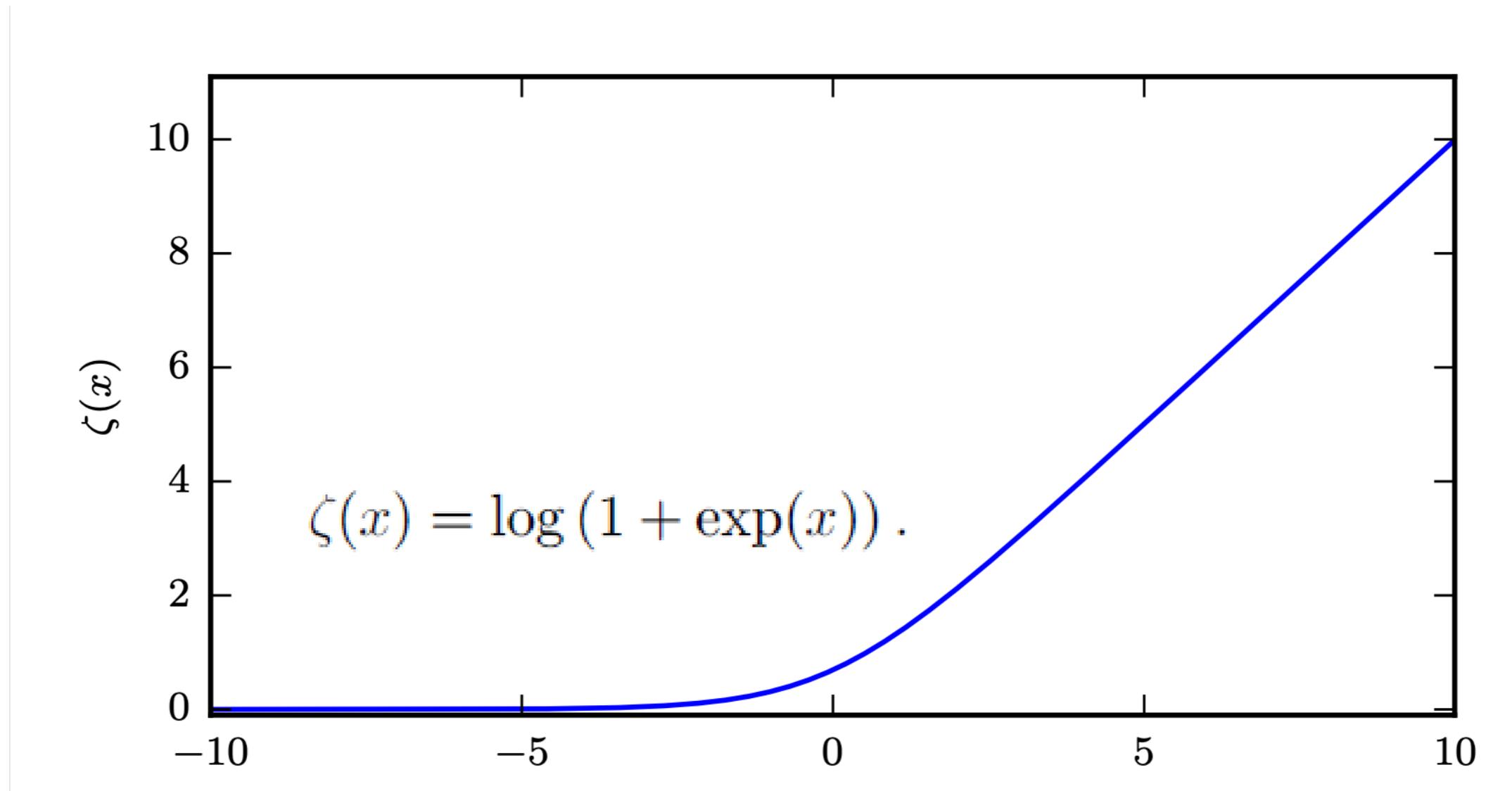


Figure 3.3: The logistic sigmoid function.

Commonly used to parametrize Bernoulli distributions

# Softplus Function



A  
smoothed  
version of  $x^+ = \max(0, x)$ .

Figure 3.4: The softplus function.

# Useful Properties

- $\sigma(x) = \frac{\exp(x)}{\exp(x) + \exp(0)}$
- $\frac{d}{dx} \sigma(x) = \sigma(x)(1 - \sigma(x))$
- $1 - \sigma(x) = \sigma(-x)$
- $\log(\sigma(x)) = -\zeta(-x)$
- $\frac{d}{dx} \zeta(x) = \sigma(x)$
- $\forall x \in (0, 1), \sigma^{-1}(x) = \log\left(\frac{x}{1-x}\right)$
- $\forall x > 0, \zeta^{-1}(x) = \log(\exp(x) - 1)$
- $\zeta(x) = \int_{-\infty}^x \sigma(y) dy$
- $\zeta(x) - \zeta(-x) = x$

# Bayes' Rule

$$P(\mathbf{x} | y) = \frac{P(\mathbf{x})P(y | \mathbf{x})}{P(y)}. \quad (3.42)$$

# Bayes Rule

Posterior

Likelihood

Prior

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

Marginal likelihood

The diagram illustrates the components of Bayes' Rule. The word 'Posterior' has an arrow pointing to the left side of the equation,  $P(A|B)$ . The word 'Likelihood' has an arrow pointing to the numerator's first term,  $P(B|A)$ . The word 'Prior' has an arrow pointing to the numerator's second term,  $P(A)$ . The words 'Marginal likelihood' have an arrow pointing to the denominator,  $P(B)$ .

computed by the total probability rule:

$$P(B) = \sum_a P(B|A = a) P(A = a)$$

# Bag-of-words Naïve Bayes:

- Features:  $W_i$  is the word at position  $i$
- Called “bag-of-words” because model is insensitive to word order or reordering:
  - In a bag-of-words model, each position is identically distributed
- $$P(\text{Spam}|W_1, W_2, \dots, W_N) \propto P(W_1|\text{Spam})P(W_2|\text{Spam}) \dots P(W_N|\text{Spam})$$
- Start with a bunch of probabilities:
  - **Prior** distribution  $P(\text{Spam}), P(\text{Ham})$
  - and the **likelihood** probabilities (The CPT tables)  $P(W_i|Y)$
- Use standard inference to compute the **posterior** probabilities  $P(Y|W_1 \dots W_n)$
- We can use the normalization trick:
$$P(\text{Ham}|W_1 \dots W_n) + P(\text{Spam}|W_1 \dots W_n) = 1$$
- Computing the **log posterior** (instead of the posterior) prevents numerical errors

# Example

$P(Y)$

```
ham : 0.66
spam: 0.33
```

$P(W|spam)$

```
the : 0.0156
to : 0.0153
and : 0.0115
of : 0.0095
you : 0.0093
a : 0.0086
with: 0.0080
from: 0.0075
...
```

$P(W|ham)$

```
the : 0.0210
to : 0.0133
of : 0.0119
2002: 0.0110
with: 0.0108
from: 0.0107
and : 0.0105
a : 0.0100
...
```

Word	$P(w   spam)$	$P(w   ham)$	Total spam (log)	Total ham (log)
(prior)	0.333	0.666	-1.1	-0.4
The	0.005	0.013	-5.27	-4.27
Year	...	...	...	...
2002				
...				

# Bag of words exercise

- Spam messages:
  - Offer is secret
  - Click secret link
  - Secret sports link
- Ham messages:
  - Play sports today
  - Went play sports
  - Secret sports event
  - Sports is today
  - Sports costs money
- Size of vocabulary = ?
- $P(\text{SPAM}) = ?$
- $P_{ML}(\text{"Secret"} | \text{SPAM}) = ?$
- $P_{ML}(\text{"Secret"} | \text{HAM}) = ?$
- how many parameters to represent the Naïve Bayes Network?
- $P(\text{SPAM} | \text{"Sports"}) = ?$
- $P(\text{SPAM} | \text{"Secret is secret"}) = ?$
- $P(\text{SPAM} | \text{"Today is secret"}) = ?$

# Change of Variables

Example of common mistake:

$$y = \frac{x}{2} \text{ and } x \sim U(0,1)$$

$$p_y(y) = p_x\left(\frac{x}{2}\right) \Rightarrow \begin{cases} y = 1 \text{ if } x \in [0, 1/2] \\ y = 0 \text{ elsewhere} \end{cases}$$

$$\int p(y) dy = \frac{1}{2} !!$$

The right thing is:  $|p_y(g(x))dy| = |p_x(x)dx|$

$$p_x(x) = p_y(g(x)) \left| \frac{\partial y}{\partial x} \right|$$

In higher dimensions:

$$p_x(\mathbf{x}) = p_y(g(\mathbf{x})) \left| \det \left( \frac{\partial g(\mathbf{x})}{\partial \mathbf{x}} \right) \right|.$$

# Information theory

- Learning that an unlikely event has occurred is more informative than learning that a likely event has occurred!
- Which statement has more information?
  - “The sun rose this morning”
  - “There was a solar eclipse this morning”
- Independent events should have additive information:
  - Finding out that a tossed coin has come up heads twice has two times more information than finding out that a tossed coin has come up heads one time!

# Self-Information

$$I(x) = -\log P(x). \quad (3.48)$$

Log base  $e \Rightarrow$  unit is **nats**

Log base 2  $\Rightarrow$  unit is **bits** or **shannons**

# Entropy

Entropy:

$$H(\mathbf{x}) = \mathbb{E}_{\mathbf{x} \sim P}[I(x)] = -\mathbb{E}_{\mathbf{x} \sim P}[\log P(x)]. \quad (3.49)$$

- Entropy is a lower bound on the number of bits needed on average to encode symbols drawn from a distribution  $P$ .
- Distributions that are nearly deterministic have low entropy
- Distributions that are nearly uniform have high entropy

# Entropy of a Bernoulli Variable

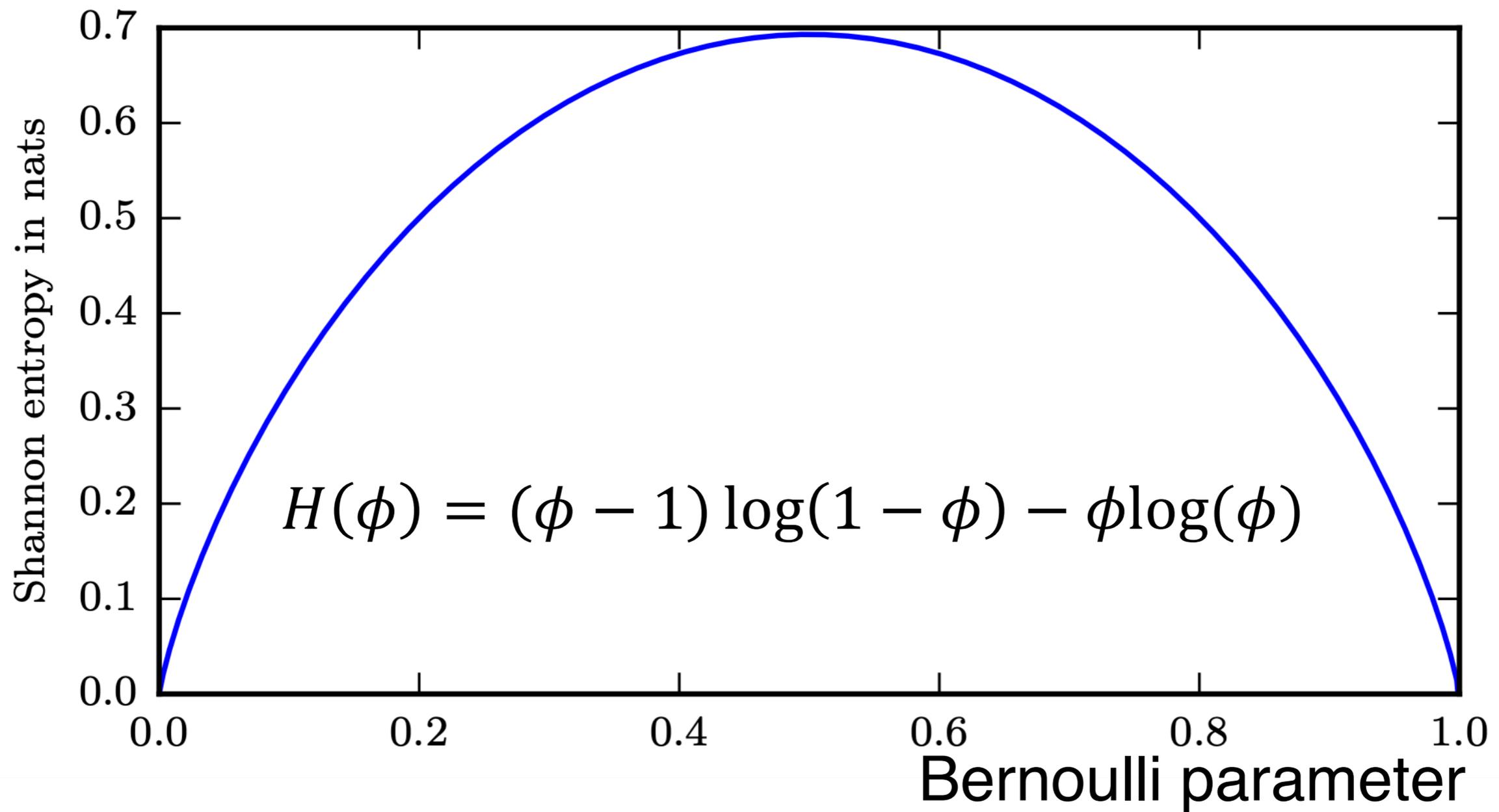


Figure 3.5

# Kullback-Leibler Divergence

KL divergence:

$$D_{\text{KL}}(P\|Q) = \mathbb{E}_{x \sim P} \left[ \log \frac{P(x)}{Q(x)} \right] = \mathbb{E}_{x \sim P} [\log P(x) - \log Q(x)]. \quad (3.50)$$

- KL-divergence is the extra amount of information needed to send a message containing symbols drawn from  $P$ , when we use a code designed to minimize the length of messages containing symbols drawn from  $Q$
- KL-divergence is non-negative
- KL-divergence = 0 if  $P$  and  $Q$  are the same distribution

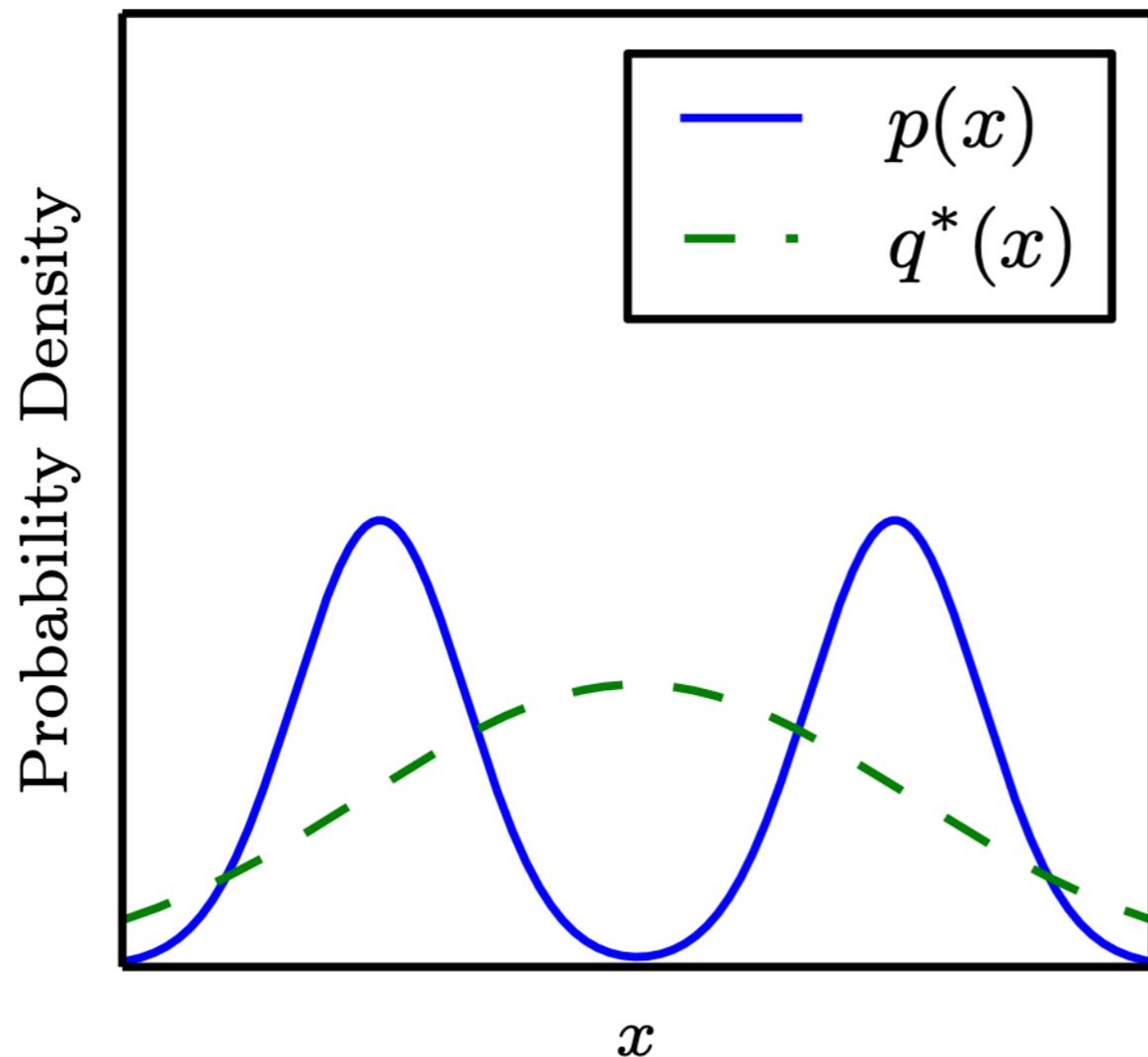
# KL-divergence

- It can be used as a distance measure between distributions
- But it is not a true distance measure since it is not symmetric:
  - $D_{KL}(P||Q) \neq D_{KL}(Q||P)$

# The KL Divergence is Asymmetric

Mixture of two Gaussians for P, One Gaussian for Q

$$q^* = \operatorname{argmin}_q D_{\text{KL}}(p||q)$$



$$q^* = \operatorname{argmin}_q D_{\text{KL}}(q||p)$$

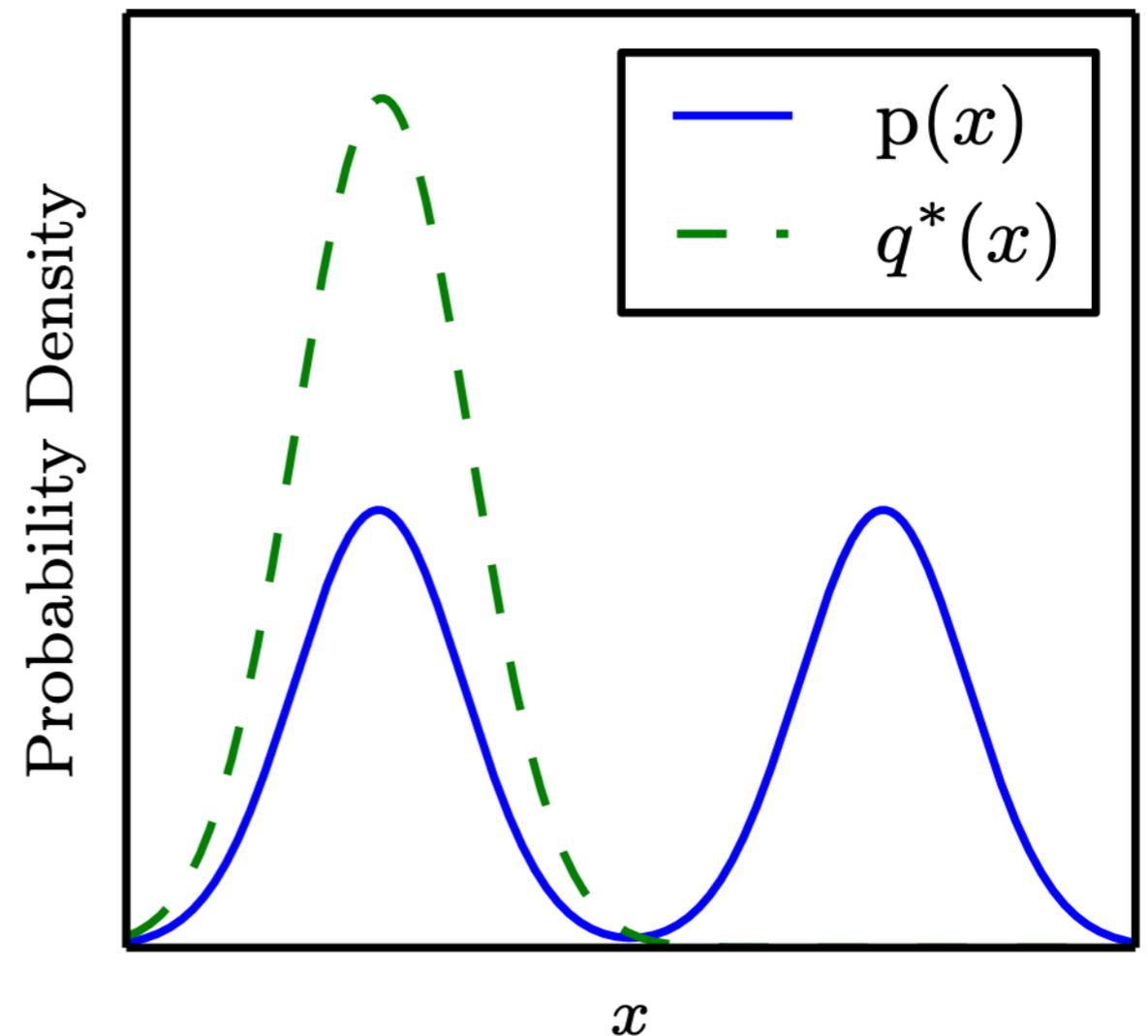


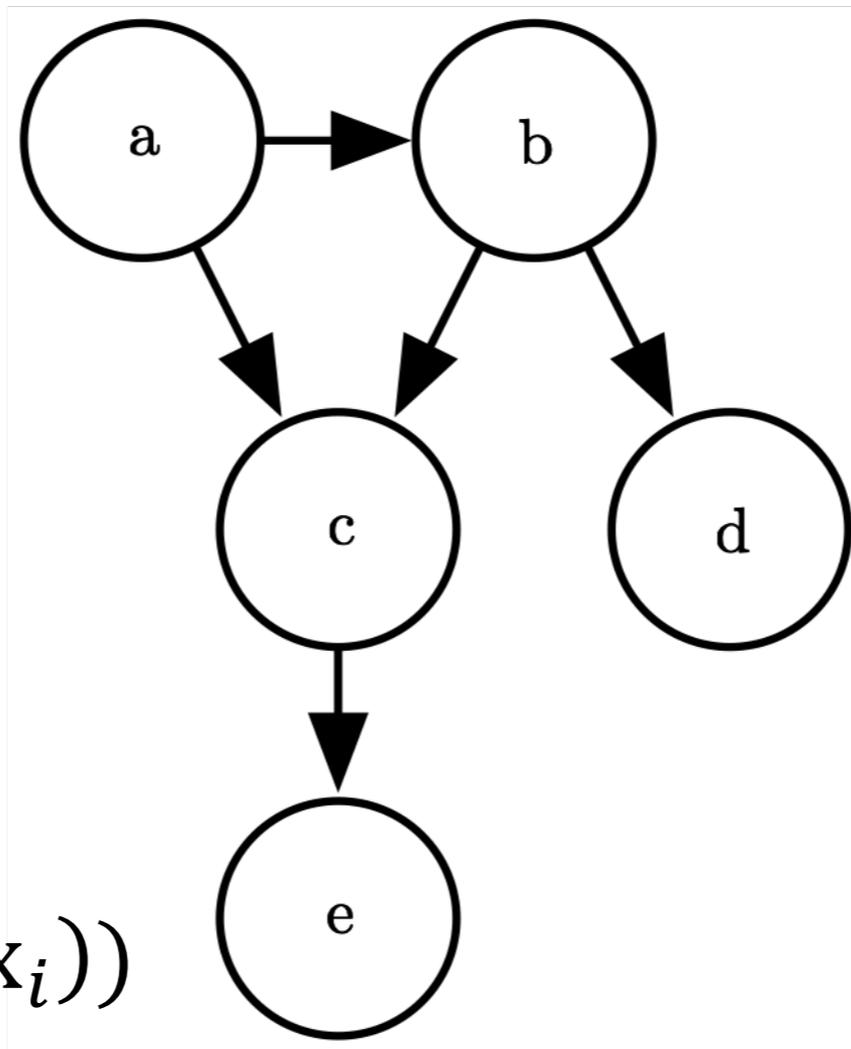
Figure 3.6

# Cross-entropy

- $H(P, Q) = H(P) + D_{KL}(P||Q)$   
 $= -\mathbb{E}_{x \sim P}(\log P(x)) + \mathbb{E}_{x \sim P}(\log P(x)) - \mathbb{E}_{x \sim P}(\log Q(x))$   
 $= -\mathbb{E}_{x \sim P}(\log Q(x))$
- Minimizing the cross entropy with respect to Q is equivalent to minimize the KL divergence!
- Remark: usually we consider  $0 \log 0 = 0$

# Directed Model

Figure 3.7



$$p(\mathbf{x}) = \prod_i p(x_i | Pa_G(x_i))$$

$$p(a, b, c, d, e) = p(a)p(b | a)p(c | a, b)p(d | b)p(e | c). \quad (3.54)$$