# **Speech and Language Processing** Lecture 1 Basics of probability distributions

Information and Communications Engineering Course Takahiro Shinozaki Manabu Okumura

# Basics of probability distributions

# **Probability Space**

- Sample space ( $\Omega$ )
  - Set of all possible outcomes of an experiment
- Probability function (f(x))
  - A function that maps each outcome to a probability  $f(x) \in [0, 1]$  for all  $x \in \Omega$

$$\sum_{x\in\Omega}f(x)=1$$

- Event (E)
  - Subset of the sample space Probability of an event *E* is :

$$P(E) = \sum_{x \in E} f(x)$$

# Random Variable

- A function that maps an outcome of an experiment to a value
  - Notation:

"P[X=x] = p" means "the probability of a random variable X takes a value x is p"

Example



X: The value of a die

Y: Whether the value of a die is odd number or not

Z: Whether the value of a die is larger than 2 or not

Random Variable	1	2	3	4	5	6
Х	1	2	3	4	5	6
Y	1	0	1	0	1	0
Z	0	0	1	1	1	1

# Joint Probability

• Probability that more than one events jointly occur

Example



P(X = i, Y = j): Probability that the value of X is i and the value of Y is j

Note: P(X=i, Y=j) = P(Y=j, X=i)

# **Conditional Probability**

 Probability of an event given that another event has occurred

Example

Randomly picks up two balls sequentially from a box containing 4 blue and 6 green balls



$$P("second \ ball \ is \ blue"|"first \ ball \ is \ geen") = \frac{4}{9}$$
$$P("second \ ball \ is \ blue"|"first \ ball \ is \ blue") = \frac{3}{9}$$

# **Two Principal Rules**

- Sum rule
  - Summing joint probability P(X, Y) for all possible values of Y gives probability of P(X)
  - P(X) is called the marginal probability

$$P(X = x_i) = \sum_j P(X = x_i, Y = y_j)$$

- Product rule
  - Product of probability P(X) and conditional probability P(Y|X) is equal to joint probability P(X,Y)

$$P(X = x_i, Y = y_j) = P(Y = y_j | X = x_i)P(X = x_i)$$

# Bayes' Theorem

• From the product rule, we obtain:

$$P(Y = y_j | X = x_i) = \frac{P(X = x_i | Y = y_j)P(Y = y_j)}{P(X = x_i)} \quad for \; \forall x_i, y_i$$

If we simplify the notation, we have:

$$P(Y | X) = \frac{P(X | Y)P(Y)}{P(X)}$$
 (Bayes' theorem)

Using the sum rule, P(Y|X) is obtained from joint probability as:

$$P(Y \mid X) = \frac{P(X \mid Y)P(Y)}{P(X)} = \frac{P(X,Y)}{\sum_{Y} P(X,Y)}$$

# Independence

• If the joint distribution of two variables X and Y factorizes into the product of the marginals, then X and Y are said to be "independent"

$$P(X,Y) = P(X | Y)P(Y) = P(X)P(Y)$$
  $\Longrightarrow$  X and Y are independent

$$P(X,Y) = P(X | Y)P(Y) \neq P(X)P(Y) \quad \triangleleft$$

X and Y are not independent

**T** 7

#### Exercises 1.1,1.2

Joint probability of random variables A and B are given in the following table. According to the table, for example,

P(A = 0, B = 0) = 0.2

P(A,B)	B=0	B=1	B=2
A=0	0.2	0.1	0.1
A=1	0.3	0.2	0.1

- 1. Obtain P(A = 0)
- 2. Obtain P(A = 0|B = 0)

# Probability Densities

• If the probability of a real-valued variable x falling in the interval  $(x, x+\delta x)$  is given by  $p(x)\delta x$  when  $\delta x \rightarrow 0$ , p(x) is called the probability density of x





#### The Sum and The Product Rules For Continuous Variable

$$P(X = x_i) = \sum_j P(X = x_i, Y = y_j)$$

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

$$P(X = x_i, Y = y_j) = P(Y = y_j | X = x_i)P(X = x_i)$$
  

$$p(x, y) = p(y | x)p(x)$$

#### Expectation

• Expectation of a function f(x) under a probability distribution p(x) is denoted by E[f]

$$E[f] = \sum_{x} p(x) f(x) \qquad (x \text{ is discrete})$$

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

(x is continuous)

# Mean and Variance

- Mean
  - Synonym of the expectation *E*[*f*(*x*)]
- Variance
  - A measure of how much variability there is in f(x) around its mean value E[f(x)]

$$\operatorname{var}[f] = E[(f(x) - E[f(x)])^2] = E[f(x)^2] - E[f(x)]^2$$

• In particular, the variance of the variable x itself is:

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

#### Covariance

- Covariance
  - The extent to which x and y vary together

$$cov[x, y] \equiv E_{p(x,y)}[(x - E[x])(y - E[y])]$$
$$= E_{p(x,y)}[xy] - E[x]E[y]$$

Expectation with respect to joint probability of *x* and *y* 

# Entropy

• Amount of randomness in the random variable

$$H[p] = E\left[-\log(p(x))\right]$$

Example

x	0	1
p(x)	0.5	0.5

 $H[x] = -\sum_{x} p(x) \log p(x)$ = -0.5 log(0.5) - 0.5 log(0.5) = 0.693

X	0	1
p(x)	0.1	0.9

$$H[x] = -\sum_{x} p(x) \log p(x)$$
  
= -0.1 log(0.1) - 0.9 log(0.9)  
= 0.325

# Relative Entropy

- A measure of dissimilarity of two distributions p and q
  - Also called as kullback-Leibler (KL) divergence

$$KL[p||q] = E_p\left[\log\left(\frac{p(x)}{q(x)}\right)\right] = -\int p(x)\log\left(\frac{q(x)}{p(x)}\right)dx$$

- KL[p||q] is nonnegative. KL[p||q] = 0 if and only if p(x) = q(x)
- In general,  $KL[p||q] \neq KL[q||p]$

# Cross Entropy

 The cross-entropy between p and q over the same underlying set of events measures the average number of bits needed to identify an event drawn from the set if a coding scheme used for the set is optimized for an estimated probability distribution q rather than the true distribution p.

$$CE[p,q] = E_p[-\log q(x)] = -\int p(x)\log q(x) \, dx$$

• Relationship with entropy and KL divergence

CE[p,q] = H[p] + KL[p||q]

# Categorical distribution

- A probability model for a categorical random variable
- The distribution is represented by a table
- An example is the probability distribution of a skewed die



Number	1	2	3	4	5	6
Probability	0.3	0.1	0.1	0.1	0.3	0.1

# 1-of-K Representation

• The same probability as the table description can be expressed as an equation by using 1-of-K representation

Value	1-of-K representation	Probability
V	$X_{v}=(x_{v,1}, x_{v,2}, x_{v,3}, x_{v,4}, x_{v,5})$	$\boldsymbol{\alpha} = (\alpha_{1,} \alpha_{2,} \alpha_{3,} \alpha_{4,} \alpha_{5})$
1 (あ)	1,0,0,0,0	$Pr(X=1)=\alpha_1=0.3$
2 (l)	0,1,0,0,0	$Pr(X=2)=\alpha_2=0.1$
3 (う)	0,0,1,0,0	Pr(X=3)=α <sub>3</sub> =0.2
4 (え)	0,0,0,1,0	Pr(X=4)=α <sub>4</sub> =0.1
5(お)	0,0,0,0,1	Pr(X=5)=α <sub>5</sub> =0.3

$$p(v|\boldsymbol{\alpha}) = p(X_v|\boldsymbol{\alpha}) = \prod_{k=1}^{K} \alpha_k^{x_{v,k}}$$

#### Exercise 1.3

• When  $P(v|\alpha)$  is given as follows, obtain  $P("i"|\alpha)$ 

	Value v	1-of-K representation $X_v = (x_{v,1}, x_{v,2}, x_{v,3}, x_{v,4}, x_{v,5})$
1	(a)	1,0,0,0,0
2	(i)	0,1,0,0,0
3	(u)	0,0,1,0,0
4	(e)	0,0,0,1,0
5	(0)	0,0,0,0,1

$$P(v|\boldsymbol{\alpha}) = \prod_{k=1}^{5} \alpha_{k}^{x_{v,k}}$$

 $\alpha = (0.3, 0.2, 0.1, 0.1, 0.3)$ 

### Gaussian Distribution

• Defined by two parameters mean  $\mu$  and standard deviation  $\sigma$  ( $\sigma^2$  is variance)

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



# Multivariate Gaussian Distribution

For *D*-dimensional vector *x*, it is defined using a mean vector *μ* and a covariance matrix *S*:

$$N(\boldsymbol{x}|\boldsymbol{\mu}, \boldsymbol{S}) = \frac{1}{\sqrt{(2\pi)^{D}|\boldsymbol{S}|}} \exp\left\{-\frac{1}{2}(\boldsymbol{x}-\boldsymbol{\mu})^{T}\boldsymbol{S}^{-1}(\boldsymbol{x}-\boldsymbol{\mu})\right\}$$
$$|\boldsymbol{S}| \text{ denotes determinant of } \boldsymbol{S}$$

Contour plot of an example of a two dimensional Gaussian distribution



#### Transformation of Continuous Random Variable

• Suppose X and Y are random variables and Y=f(X)



# Transformation of Vector Variable

$$Y = f(X)$$

The volume (area) ratio by a



$$P_{Y}(Y)\Delta V_{Y} = P_{X}(X)\Delta V_{X}$$

$$P_{Y}(Y) = \left|\frac{\Delta V_{X}}{\Delta V_{Y}}\right| P_{X}(X)$$

$$P_{Y}(y_{1}, y_{2}, \dots, y_{N})$$

$$= \left|\frac{\partial(x_{1}, x_{2}, \dots, x_{N})}{\partial(y_{1}, y_{2}, \dots, y_{N})}\right| P_{X}(x_{1}, x_{2}, \dots, x_{N})$$

$$= \left|\frac{\partial x_{1}}{\partial y_{1}} \frac{\partial x_{1}}{\partial y_{2}} \dots \frac{\partial x_{1}}{\partial y_{N}}\right| P_{X}(x_{1}, x_{2}, \dots, x_{N})$$

$$Jacobian determinant of the inverse function$$

### Sampling From a Uniform Distribution

• Samples distribute uniformly over some region

#### Example:

Histogram of samples obtained from a uniform distribution over (0, 1)



1000 samples



#### Sampling From a Gaussian Distribution

Standard normal (Gaussian) distribution has a mean 0.0 and a variance 1.0



#### Exercise 1.4

Assume p(x) and y = f(x) are given as follows

$$p(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right) = N(x|0,1) \quad x \in (-\infty, \infty), \qquad y = 3x + 4$$

Then, q(y) becomes a Gaussian distribution  $N(y|E, 3^2)$ . Obtain E



#### Approximating Expectation with Sampling

When  $x_1, x_2, \dots, x_N$  are samples independently drawn from a distribution p(x)

$$E[f] = \sum_{x} p(x)f(x) \approx \frac{1}{N} \sum_{n=1}^{N} f(x_n) \qquad (x \text{ is discrete})$$

$$E[f] = \int p(x) f(x) dx \approx \frac{1}{N} \sum_{n=1}^{N} f(x_n) \quad (x \text{ is continuous})$$

# Advanced

### Maximum Likelihood (ML) Method

- Assume that we have a set of samples D={x<sub>1</sub>, x<sub>2</sub>, ... x<sub>n</sub>} drawn from a distribution p(x/θ) with unknown parameters θ, and we want to estimate θ
- Maximum likelihood method estimates the parameters by maximizing likelihood  $p(D|\theta)$

$$\hat{\theta} = \arg \max_{\theta} p(D \mid \theta) = \arg \max_{\theta} \prod_{i=1}^{n} p(x_i \mid \theta)$$

Probability of the data set D is decomposed to a product of samples when they are drawn independently

# Bernoulli Distribution

- Probability distribution of a binary random variable which takes value 1 with probability  $\mu$  and value 0 with probability 1- $\mu$ 



Is the result Head or Tail?

#### ML Estimation for Bernoulli Distribution

- Parameter  $\theta$  in this case is :  $\mu$
- Training sample  $x_i = 0 \text{ or } 1$

$$\hat{\mu} = \arg \max_{\mu} p(D \mid \mu) = \arg \max_{\mu} \prod_{i=1}^{n} \mu^{x_i} (1-\mu)^{1-x_i}$$

$$= \arg \max_{\mu} \log \left( \prod_{i} \mu^{x_i} (1-\mu)^{1-x_i} \right)$$

$$= \arg \max_{\mu} \left\{ \sum_{i} x_i \log(\mu) + \sum_{i} (1-x_i) \log(1-\mu) \right\}$$

$$= \frac{\partial}{\partial \mu} \left( \sum_{i} x_i \log(\mu) + \sum_{i} (1-x_i) \log(1-\mu) \right) = 0$$

$$\mu = \frac{1}{n} \sum_{i} x_i$$
*n*: the

Taking log does not change the problem and makes the equation a bit easier

*n*: the number of samples

# Categorical Distribution

• As a generalization of the Bernoulli Distribution, lets consider a discrete random variable X that takes K values

X	1	2	•••	К
1-of-K	<10 0>	<01 0>		<0.0 1>
< <i>x</i> <sub>1</sub> , <i>x</i> <sub>2</sub> ,, <i>x</i> <sub>K</sub> >	<1,0,,0>	<0,1,,0>	•••	<0,0,,1>
p(X)	$\mu_{1}$	$\mu_2$	•••	$\mu_K$

$$(\mathbf{\mu}) = \prod_{k=1}^{K} \mu_k^{x_k}$$

# ML for Categorical Distribution

- Parameter  $\theta$  in this case is :  $\boldsymbol{\mu} = \{\mu_{1, \mu_{2, \dots, \mu_{K}}}\}$
- Training sample  $x_i$  is a vector of 1-of-K representation. When  $x_i$  represents k-th value,  $x_{i,k}=1$ , and  $x_{i,j}=0$  for  $j \neq k$

$$\hat{\boldsymbol{\mu}} = \arg \max_{\boldsymbol{\mu}} p(\boldsymbol{D} \mid \boldsymbol{\mu}) = \arg \max_{\boldsymbol{\mu}} \prod_{i=1}^{n} \prod_{k=1}^{K} \mu_{k}^{x_{i,k}} = \arg \max_{\boldsymbol{\mu}} \prod_{k=1}^{K} \mu_{k}^{m_{k}}$$

$$\sum_{k=1}^{K} \mu_{k} = 1 \quad \text{Constraint} \quad m_{k} \text{ is the number of the occurrence of } k \text{-th value, } m_{k} = \sum_{i=1}^{n} x_{i,k} \text{ of samples}$$

This is a maximization problem with a constraint



Use the method of Lagrange multiplier

# Method of Lagrange Multiplier



### Solution

$$\hat{\boldsymbol{\mu}} = \arg \max_{\boldsymbol{\mu}} \prod_{k=1}^{K} \mu_{k}^{m_{k}} = \arg \max_{\boldsymbol{\mu}} \sum_{k=1}^{K} m_{k} \log(\mu_{k})$$
$$\sum_{k=1}^{K} \mu_{k} = 1$$
$$\hat{\boldsymbol{\mu}} = \arg \max_{\boldsymbol{\mu}} \left\{ \sum_{k=1}^{K} m_{k} \log \mu_{k} + \lambda \left( \sum_{k=1}^{K} \mu_{k} - 1 \right) \right\}$$

$$\mu_k = \frac{m_k}{n}$$

#### Exercise 1.5

• Show the derivation process of obtaining  $\mu_k = \frac{m_k}{m_k}$ 

for the categorical distribution by maximizing

$$L(\boldsymbol{\mu}, \boldsymbol{\lambda}) = \sum_{k=1}^{K} m_k \log \mu_k + \boldsymbol{\lambda} \left( \sum_{k=1}^{K} \mu_k - 1 \right)$$

where  $\lambda$  is the Lagrange multiplier.

n