TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Statistical Machine Learning

## Lecture 03: Statistics Refresher

**Kristian Kersting**
**TU Darmstadt**

Summer Term 2020

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Today's Objectives

- Make you remember your sweetest high school dreams: statistics & probabilities.

- This topic is harder than most of remaining chapters, but you will need it to continue!

- Covered Topics:
  - Random Variables: discrete & continuous

  - Distributions: discrete & continuous

- Expected values and moments

- Joint distributions, conditional distributions, independence

**Outline**

**1. Random Variables and Common Distributions**
   Random Variables
   Discrete Distributions
   Continuous Distributions

**2. Basic Rules of Probability**

**3. Expectations, Variance and Moments**

**4. Exponential Family**

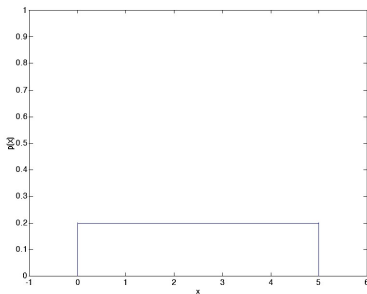**5. Information and Entropy**

**6. Wrap-Up**

# **Outline**

# Random Variables

- What is a random variable?
    - Is a random number determined by chance

    - More formally, drawn according to a probability distribution

    - Typical random variables in statistical learning: input data, output data, noise

- What is a probability distribution?
    - Describes the probability (density) that the random variable will be equal to a certain value.

    - The probability distribution can be given by the physics of an experiment (e.g., throwing dice)

# Random Variables

- Important concept: The data generating model
  - E.g., what is the data generating model for: i) throwing dice, ii) regression, iii) classification, iv) visual perception?

- Problem: On which time scale is a distribution observed?

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Uniform Distribution



- All data is equally probable within a bounded region $R$

$$p(x) = \frac{1}{R}$$

- The uniform distribution plays an important role in entropy methods and information theory.

# Discrete Distributions

- The random variables take on discrete values
  - E.g, when throwing a dice, the possible values are (countably finite set):

    $$x_i \in \{1, 2, 3, 4, 5, 6\}$$

  - E.g., the number of sand grains at the beach (countably infinite set):

    $$x_i \in \mathbb{N}$$

# Discrete Distributions

- The probabilities sum to 1

$$\sum_i p(x_i) = 1$$

- Discrete distributions are particularly important in classification and decision making

- A discrete distribution is described by a probability mass function (or frequency function), which is a normalized histogram

# Bernoulli Distribution

- A Bernoulli random variable only takes on two values, for example 0 and 1

$$
\begin{aligned}
x &\in \{0, 1\} \\
p(x = 1|\mu) &= \mu \\
\text{Bern}(x|\mu) &= \mu^x(1 - \mu)^{1-x} \\
\mathbb{E}[x] &= \mu \\
\text{var}[x] &= \mu(1 - \mu)
\end{aligned}
$$

- The only parameter of a Bernoulli distribution is $\mu$, i.e., it is completely defined using only this parameter

# Bernoulli Distribution

- Bernoulli distributions are often modeled with sigmoidal nonlinearites in statistical learning



Does pattern belong to class C or not?

1   0   1   1   0   0   0   1   1   1

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# **Binomial Distribution**

- Binomial variables are a sequence of *N* repeated Bernoulli variables

- One interpretation is "what is the probability of getting $m \in \mathbb{N}$ heads in *N* trials?"

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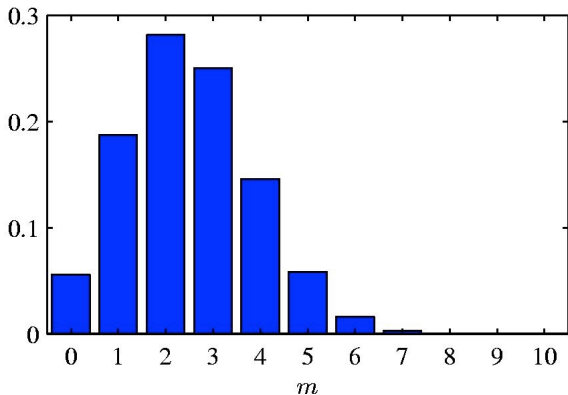\text{var}[m] &= \sum_{m=0}^{N} (m - \mathbb{E}[m])^2 \text{Bin}(m|N, \mu) = N\mu(1 - \mu)
\end{aligned}
$$

# Binomial Distribution

- The Binomial distribution is completely defined with *N* - the number of samples - and $\mu$ - the probability that one sample is equal to 1

- Binomial variables are important for example in density estimation: "What is the probability that *k* out of *n* data points fall into region *R*?"

# Binomial Distribution



Bin($m$|10, 0.25)

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Multinoulli Distribution

- Multinoulli variables, also called Categorical variables in some literature, are a generalization of binomial variables to multiple outputs (e.g., multiple classes)

- 1-of-$K$ coding scheme (also called one-hot encoding)

$$\mathbf{x} = (0, 0, 1, 0, 0, 0)^\mathsf{T}$$

$$p(\mathbf{x}|\boldsymbol{\mu}) = \prod_{k=1}^{K} \mu_k^{x_k} \quad \forall k : \mu_k \geq 0 \quad \text{and} \quad \sum_{k=1}^{K} \mu_k = 1$$

$$\mathbb{E}\left[\mathbf{x}|\boldsymbol{\mu}\right] = \sum_{\mathbf{x}} p(\mathbf{x}|\boldsymbol{\mu})\mathbf{x} = (\mu_1, \ldots, \mu_K)^\mathsf{T}$$

$$\sum_{\mathbf{x}} p(\mathbf{x}|\boldsymbol{\mu}) = \sum_{k=1}^{K} u_k = 1$$

# Multinomial Distribution

- *N* independent trials can result in one of *K* types of outcome

- What is the probability that in *N* trials, the frequency of the *K* classes is $m_1, m_2, \ldots, m_K$

$$
\begin{aligned}
\text{Mult}(m_1, m_2, \ldots, m_k | \boldsymbol{\mu}, N) &= \begin{pmatrix} N \\ m_1, m_2, \ldots, m_K \end{pmatrix} \prod_{k=1}^{K} \mu_k^{m_k} \\
\mathbb{E}[m_k] &= N\mu_k \\
\text{var}[m_k] &= N\mu_k(1 - \mu_k) \\
\text{cov}[m_j m_k] &= -N\mu_j \mu_k
\end{aligned}
$$

# Multinomial Distribution

- The multinomial distribution play an important role in multi-class classification ($N = 1$)



To which class does a data vector belong?
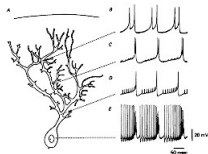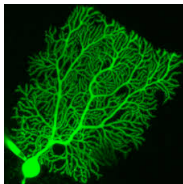
# Poisson Distribution

- The Poisson distribution is the binomial distribution where the number of trials $N$ goes to infinity, and the probability of success on each trial, $\mu$, goes to zero, such that $N\mu = \lambda$ is a constant

$$p(m|\lambda) = \frac{\lambda^m}{m!}e^{-\lambda}$$

- Where the $m$ is the number of "successes"

- For example, Poisson distributions are an important model for t he firing characteristics of biological neurons. They are also used as an approximation to binomial variables with small $p$

# Poisson Distribution

- Example: What is the probability of firing of a *Purkinje* neuron in the cerebellum in a 10ms time interval?
    - We know that the average firing of these neurons is about 40Hz, $\lambda = 40\text{Hz} \times 0.01\text{s}$

    - Note that this approximation only work if the number of spike is low in the given time interval
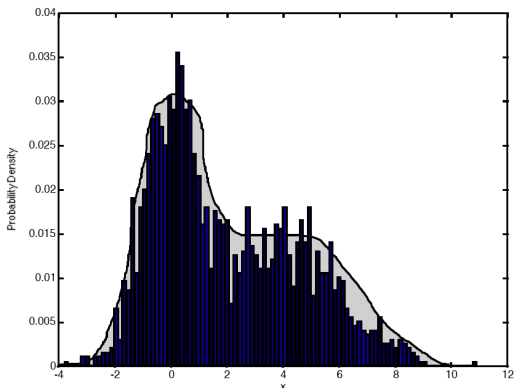
# Continuous Distributions

- The random variables take on continuous values

- Continuous distributions are discrete distributions where the number of discrete values goes to infinity, while the probability of each value goes to zero

- A continuous distribution is described by a probability density function, which integrates to 1
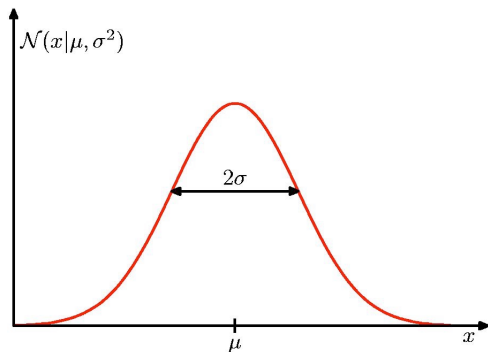
$$\int_{-\infty}^{+\infty} p(x)dx = 1$$

- Continuous distributions are particularly important in regression and unsupervised learning

- A lot of Machine Learning is centered around how to better model a density function

# Example of a probability density function $p(x)$



$$P(a < x < b) = \int_a^b p(x)dx$$
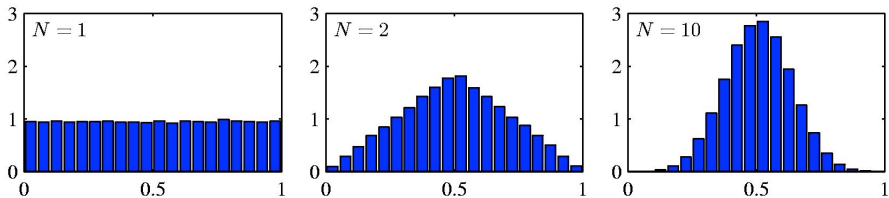
# The Gaussian Distribution



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

# Central Limit Theorem

- Why are Gaussians SO important?

- The distribution of the sum of $N$ i.i.d. (independent and identically distributed) random variables becomes increasingly Gaussian as $N$ grows
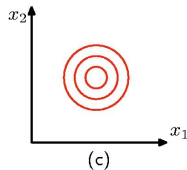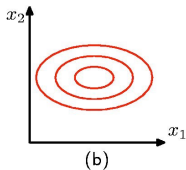
# Central Limit Theorem

- Example: *N* uniform [0,1] random variables



- Gaussians are often a *good* model of data

- Working with Gaussians leads to analytic solutions for complex operations

# The Multivariate Gaussian Distribution



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

# The Multivariate Gaussian Distribution

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

- **To clear some confusion**: for a chosen vector **x**, $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ is a real number with the probability density of **x** (which can be greater than 1, only the integral of the probability density function needs to be 1). The mean $\boldsymbol{\mu}$ is just a specific vector amongst all the possible vectors. The covariance matrix $\boldsymbol{\Sigma}$ tells us how two dimensions of a vector are related to each other.

# Geometry of the Multivariate Gaussian

$$\Delta^2 = (\mathbf{x} - \boldsymbol{\mu})^\intercal \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})$$

$$\boldsymbol{\Sigma}^{-1} = \sum_{i=1}^{D} \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^\intercal$$

$$\Delta^2 = \sum_{i=1}^{D} \frac{y_i^2}{\lambda_i}$$

$$y_i = \mathbf{u}_i^\intercal (\mathbf{x} - \boldsymbol{\mu})$$

$\Delta^2$ is the Mahalanobis distance.

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# **Outline**

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# 2. Basic Rules of Probability

■ Joint Distribution

$$p(x, y)$$

■ Marginal Distribution

$$p(y) = \int p(x, y) dx$$

■ Conditional Distribution

$$p(y|x) = \frac{p(x, y)}{p(x)}$$

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# 2. Basic Rules of Probability

■ Probabilistic Independence

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

■ Chain Rule of Probabilities

$$p(x_1, \ldots, x_n) = p(x_1|x_2, \ldots, x_n)p(x_2, \ldots, x_n)$$
$$= p(x_1|x_2, \ldots, x_n)p(x_2|x_3, \ldots, x_n) \ldots p(x_{n-1}|x_n)p(x_n)$$

# Bayes Rule

$$p(y|x) = \frac{p(x|y)p(y)}{p(x)}$$

posterior $\propto$ likelihood $\times$ prior

- **posterior**: $p(y|x)$

- **likelihood:** $p(x|y)$

- **prior:** $p(y)$

- $p(x) = \int p(x,y)\mathrm{d}y = \int p(x|y)p(y)\mathrm{d}y$

# Partitioned Gaussian Distributions

$$p(\mathbf{x}) = \mathcal{N}\left(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}\right)$$

$$\mathbf{x} = \left( \begin{array}{c} \mathbf{x}_a \\ \mathbf{x}_b \end{array} \right) \qquad \boldsymbol{\mu} = \left( \begin{array}{c} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{array} \right) \qquad \boldsymbol{\Sigma} = \left( \begin{array}{cc} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{array} \right)$$

$$\boldsymbol{\Lambda} \equiv \boldsymbol{\Sigma}^{-1} \qquad \boldsymbol{\Lambda} = \left( \begin{array}{cc} \boldsymbol{\Lambda}_{aa} & \boldsymbol{\Lambda}_{ab} \\ \boldsymbol{\Lambda}_{ba} & \boldsymbol{\Lambda}_{bb} \end{array} \right)$$

$\boldsymbol{\Lambda}$ is the precision matrix.

# Partitioned Conditionals and Marginals

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Partitioned Conditionals and Marginals

$$p(\mathbf{x}_a|\mathbf{x}_b) = \mathcal{N}\left(\mathbf{x}_a|\boldsymbol{\mu}_{a|b}, \boldsymbol{\Sigma}_{a|b}\right)$$

$$
\begin{aligned}
\boldsymbol{\Sigma}_{a|b} &= \boldsymbol{\Lambda}_{aa}^{-1} = \boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}\boldsymbol{\Sigma}_{ba} \\
&= \boldsymbol{\Sigma}_{a|b}\left\{\boldsymbol{\Lambda}_{aa}\boldsymbol{\mu} - \boldsymbol{\Lambda}_{ab}\left(\mathbf{x}_b - \boldsymbol{\mu}\right)\right\} \\
&= \boldsymbol{\mu}_a + \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}\left(\mathbf{x}_b - \boldsymbol{\mu}_b\right)
\end{aligned}
$$

$$
\begin{aligned}
p\left(\mathbf{x}_a\right) &= \int p\left(\mathbf{x}_a, \mathbf{x}_b\right) d\mathbf{x}_b \\
&= \mathcal{N}\left(\mathbf{x}_a|\boldsymbol{\mu}_a, \boldsymbol{\Sigma}_{aa}\right)
\end{aligned}
$$

- Important result: If the joint distribution $p(\mathbf{x}_a, \mathbf{x}_b)$ is Gaussian, then the conditional distributions $p(\mathbf{x}_a|\mathbf{x}_b)$ and $p(\mathbf{x}_b|\mathbf{x}_a)$ are also Gaussians. Moreover, the marginal distributions $p(\mathbf{x}_a)$ and $p(\mathbf{x}_b)$ are also Gaussians

# Outline

# Expectations

- Expectation

$$\mathbb{E}_{x \sim p(x)}[f(x)] = \mathbb{E}_x[f] = \mathbb{E}[f] = \begin{cases} \sum_x p(x)f(x) & \text{discrete case} \\ \int p(x)f(x)\mathrm{d}x & \text{continuous case} \end{cases}$$

- Conditional Expectation

$$\mathbb{E}_{x \sim p(x|y)}[f(x)] = \mathbb{E}_x[f|y] = \begin{cases} \sum_x p(x|y)f(x) & \text{discrete case} \\ \int p(x|y)f(x)\mathrm{d}x & \text{continuous case} \end{cases}$$

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Expectations

■ Approximate Expectation

$$\mathbb{E}\left[f\right] = \int f(x)p(x)\mathrm{d}x \approx \frac{1}{N}\sum_{n=1}^{N}f(x_n)$$

■ We sample $N$ points from the distribution $p(x)$ and compute the function at those points. The probability of computing $f(x_n)$ for a certain point $x_n$ is given by the probability of sampling $p(x_n)$

■ This result is very important! When there is no analytical solution, we can use this to approximate integrals by sampling!

# Expectations

- Example: What is the expectation of the following distribution?

# Expectations

- Some rules of expectation

  $\mathbb{E}\left[a\mathbf{x}\right] = a\mathbb{E}\left[\mathbf{x}\right]$

  $\mathbb{E}\left[\mathbf{x} + \mathbf{y}\right] = \mathbb{E}\left[\mathbf{x}\right] + \mathbb{E}\left[\mathbf{y}\right]$

  $\mathbb{E}\left[\mathbf{xy}\right] = \mathbb{E}\left[\mathbf{x}\right]\mathbb{E}\left[\mathbf{y}\right]$  only if **x** and **y** are statistically independent!

  $\mathbb{E}\left[\sum_i a_i x_i\right] = \sum_i a_i \mathbb{E}\left[x_i\right]$

- Expectation of functions

  $\mathbb{E}\left[g(\mathbf{x})\right] = \int g(\mathbf{x})p(\mathbf{x})\mathrm{d}\mathbf{x}$

  In general $\mathbb{E}\left[g(\mathbf{x})\right] \neq g\left(\mathbb{E}\left[\mathbf{x}\right]\right)$

# Variance and Covariance

- Variances give a measure of dispersion - the expected spread of the variable in relation to its mean

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

# Variance and Covariance

- Covariances give a measure of correlation - how much two variables change together

$$\text{cov}\,[x, y] = \mathbb{E}_{x,y}\left[(x - \mathbb{E}\,[x])\,(y - \mathbb{E}\,[y])\right]$$
$$= \mathbb{E}_{x,y}\,[xy] - \mathbb{E}_x[x]\mathbb{E}_y\,[y]$$

$$\text{cov}\,[\mathbf{x}, \mathbf{y}] = \mathbb{E}_{\mathbf{x},\mathbf{y}}\left[(\mathbf{x} - \mathbb{E}\,[\mathbf{x}])\,(\mathbf{y} - \mathbb{E}\,[\mathbf{y}])^\mathsf{T}\right]$$
$$= \mathbb{E}_{\mathbf{x},\mathbf{y}}\left[(\mathbf{x} - \mathbb{E}\,[\mathbf{x}])\,(\mathbf{y}^\mathsf{T} - \mathbb{E}\,[\mathbf{y}^\mathsf{T}])\right]$$
$$= \mathbb{E}_{\mathbf{x},\mathbf{y}}\left[\mathbf{xy}^\mathsf{T}\right] - \mathbb{E}_{\mathbf{x}}[\mathbf{x}]\mathbb{E}_{\mathbf{y}}\left[\mathbf{y}^\mathsf{T}\right]$$

# Variance and Covariance

- Note the very important rule

$$\mathbb{E}\left[\mathbf{x}\mathbf{x}^\mathsf{T}\right] = \mathbb{E}_\mathbf{x}[\mathbf{x}]\mathbb{E}_\mathbf{x}\left[\mathbf{x}^\mathsf{T}\right] + \text{cov}\left[\mathbf{x},\mathbf{x}\right]$$
$$= \boldsymbol{\mu}\boldsymbol{\mu}^\mathsf{T} + \boldsymbol{\Sigma}$$

# Moments of Random Variables

- Definition of a Moment
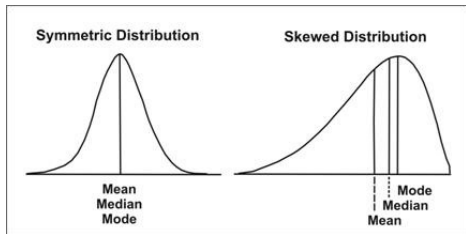  $m_n = \mathbb{E}[x^n]$
- Definition of a Central Moment
  $cm_n = \mathbb{E}\left[(x - \mu)^n\right]$
- $cm_2$: variance
- $cm_3$: skewness (measure of asymmetry)
- $cm_4$: kurtosis (measure of heavy tailed-ness and light tailed-ness)

# Moments of the Multivariate Gaussian

$$\mathbb{E}\left[\mathbf{x}\right] = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\mathbf{\Sigma}|^{1/2}} \int \exp\left\{-\frac{1}{2}\left(\mathbf{x}-\boldsymbol{\mu}\right)^{\mathsf{T}}\mathbf{\Sigma}^{-1}\left(\mathbf{x}-\boldsymbol{\mu}\right)\right\}\mathbf{x}\mathrm{d}\mathbf{x}$$

$$= \frac{1}{(2\pi)^{D/2}} \frac{1}{|\mathbf{\Sigma}|^{1/2}} \int \exp\left\{-\frac{1}{2}\mathbf{z}^{\mathsf{T}}\mathbf{\Sigma}^{-1}\mathbf{z}\right\}\left(\mathbf{z}+\boldsymbol{\mu}\right)\mathrm{d}\mathbf{z}$$

Thanks to the asymmetry of $\mathbf{z}$, $\mathbb{E}\left[\mathbf{x}\right] = \boldsymbol{\mu}$

# Moments of the Multivariate Gaussian

$$\mathbb{E}\left[\mathbf{x}\mathbf{x}^\mathsf{T}\right] = \boldsymbol{\mu}\boldsymbol{\mu}^\mathsf{T} + \boldsymbol{\Sigma}$$

$$\mathrm{cov}\left[\mathbf{x}\right] = \mathrm{cov}\left[\mathbf{x}, \mathbf{x}\right] = \mathbb{E}\left[\left(\mathbf{x} - \mathbb{E}\left[\mathbf{x}\right]\right)\left(\mathbf{x} - \mathbb{E}\left[\mathbf{x}\right]\right)^\mathsf{T}\right] = \boldsymbol{\Sigma}$$



(a)        (b)        (c)

# **Outline**

# 4. Exponential Family

- The exponential family are a large class of distributions that are all analytically appealing, because taking the $\log$ of them decomposes them into simple terms

- All distributions from this family are uni-modal

$$p\left(\mathbf{x}|\boldsymbol{\eta}\right) = h\left(\mathbf{x}\right)g\left(\boldsymbol{\eta}\right)\exp\left\{\boldsymbol{\eta}^{\mathsf{T}}\mathbf{u}(\mathbf{x})\right\}$$

where $\boldsymbol{\eta}$ is the natural parameter and

$$g\left(\boldsymbol{\eta}\right)\int h\left(\mathbf{x}\right)\exp\left\{\boldsymbol{\eta}^{\mathsf{T}}\mathbf{u}(\mathbf{x})\right\}\mathrm{d}\mathbf{x} = 1$$

hence $g$ can be interpreted as a normalization coefficient

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Exponential Family - Bernoulli Distribution

■ The Bernoulli Distribution

$$
\begin{aligned}
p\left(x|\mu\right) = \mathsf{Bern}(x|\mu) &= \mu^{x}\left(1-\mu\right)^{1-x} \\
&= \exp\left\{x\ln\mu + (1-x)\ln\left(1-\mu\right)\right\} \\
&= (1-\mu)\exp\left\{\ln\left(\frac{\mu}{1-\mu}\right)x\right\}
\end{aligned}
$$

■ Comparing with the general form we see that

$$
\eta = \ln\left(\frac{\mu}{1-\mu}\right), \quad \mu = \underbrace{\sigma\left(\eta\right) = \frac{1}{1+\exp\left(-\eta\right)}}_{\text{Logistic sigmoid}}
$$

# Exponential Family - Bernoulli Distribution

- Hence, the Bernoulli Distribution can be written as

$$p\left(x|\mu\right) = \sigma(-\eta)\exp(\eta x)$$

where

$$u(x) = x, \quad h(x) = 1, \quad g\left(\eta\right) = 1 - \sigma(\eta) = \sigma(-\eta)$$

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Exponential Family - Multinoulli Distribution

■ The Multinoulli Distribution also belongs to the exponential family

$$p\left(\mathbf{x}|\boldsymbol{\mu}\right) = \prod_{k=1}^{M} \mu_k^{x_k} = \exp\left\{\sum_{k=1}^{M} x_k \ln \mu_k\right\} = h(\mathbf{x})g(\boldsymbol{\eta})\exp\left\{\boldsymbol{\eta}^{\mathsf{T}}\mathbf{u}(\mathbf{x})\right\}$$

where

$$\mathbf{x} = (x_1, \ldots, x_M)^{\mathsf{T}}, \quad \boldsymbol{\eta} = (\eta_1, \ldots, \eta_M)^{\mathsf{T}}, \quad \eta_k = \ln u_k$$
$$\mathbf{u}(\mathbf{x}) = \mathbf{x}, \quad h(\mathbf{x}) = 1, \quad g(\boldsymbol{\eta}) = 1$$

■ Note that the parameters $\eta_k$ have to be chosen in a way to guarantee that $p\left(\mathbf{x}|\boldsymbol{\mu}\right)$ is a valid probability distribution. Particularly, they must satisfy

$$\sum_{\mathbf{x}} p\left(\mathbf{x}|\boldsymbol{\mu}\right) = 1 \implies \sum_{k=1}^{M} \mu_k = 1$$

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Exponential Family - Multinoulli Distribution

- Let $\mu_M = 1 - \sum_{k=1}^{M-1} \mu_k$, which ensures that the distribution is well defined. We can rewrite $p(\mathbf{x}|\boldsymbol{\mu})$ and observe that

$$\eta_k = \ln \left( \frac{\mu_k}{1 - \sum_{j=1}^{M-1} \mu_j} \right), \quad \mu_k = \underbrace{\frac{\exp(\eta_k)}{1 + \sum_{j=1}^{M-1} \exp(\eta_j)}}_{\text{Softmax}}$$

- Here the parameters $\eta_k$ can be chosen independently, since

$$0 \leq \mu_k \leq 1, \quad \sum_{k=1}^{M-1} \mu_k \leq 1$$

# Exponential Family - Multinoulli Distribution

TECHNISCHE
UNIVERSITÄT
DARMSTADT

■ The Multinoulli Distribution can then be written as

$$p\left(\mathbf{x}|\boldsymbol{\mu}\right) = h\left(\mathbf{x}\right) g\left(\boldsymbol{\eta}\right) \exp\left\{\boldsymbol{\eta}^{\mathsf{T}} \mathbf{u}\left(\mathbf{x}\right)\right\}$$

where

$$\boldsymbol{\eta} = (\eta_1, \ldots, \eta_{M-1}, 0)^{\mathsf{T}}, \quad \mathbf{u}(\mathbf{x}) = \mathbf{x}, \quad h(\mathbf{x}) = 1$$

$$g(\boldsymbol{\eta}) = \left(1 + \sum_{k=1}^{M-1} \exp\left(\eta_k\right)\right)^{-1}$$

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Exponential Family - Gaussian Distribution

■ The Gaussian Distribution can be rewritten as

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

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

$$= h(x)g(\boldsymbol{\eta}) \exp\left\{\boldsymbol{\eta}^\mathsf{T}\mathbf{u}(x)\right\}$$

where

$$\boldsymbol{\eta} = \left(-\frac{1}{2\sigma^2}, \frac{\mu}{\sigma^2}\right)^\mathsf{T}, \quad \mathbf{u}(x) = (x^2, x)^\mathsf{T}, \quad h(\mathbf{x}) = 1$$

$$g(\boldsymbol{\eta}) = \sqrt{\frac{-\eta_1}{\pi}} \exp\left(\frac{\eta_2^2}{4\eta_1}\right)$$

# **Outline**

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Information Theory - Core Questions

- Classical Question: How can we represent information compactly, i.e., using as few bits as possible?
  - Compressing text like with GZIP

  - Compressing pictures like in JPEG, movies like in MPEG

  - Compressing sound using MP3

- Classical Question: How can we transmit or store data reliably?
  - ECC memory

  - Error Correction on CDs

  - Communication with space probes

# Information Theory - Core Questions

- Machine Learning Questions:
    - How can we measure complexity?

    - How can we measure "distances" between probability distributions?

    - How can we reconstruct data?
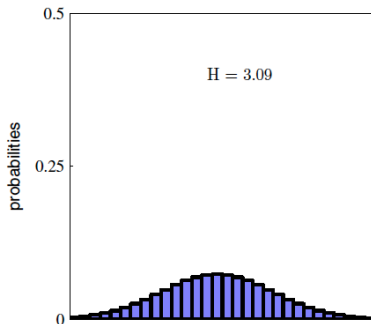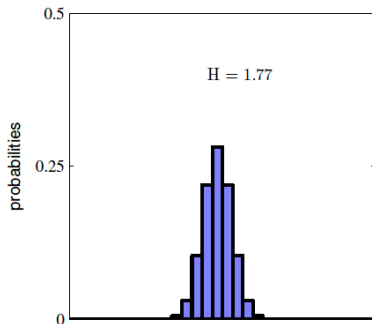
- We are not covering all questions here... :)

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# What is Information?

| $i$ | $a_i$ | $p_i$ |
|---|---|---|
| 1 | a | .0575 |
| 2 | b | .0128 |
| 3 | c | .0263 |
| 4 | d | .0285 |
| 5 | e | .0913 |
| 6 | f | .0173 |
| 7 | g | .0133 |
| 8 | h | .0313 |
| 9 | i | .0599 |
| 10 | j | .0006 |
| 11 | k | .0084 |
| 12 | l | .0335 |
| 13 | m | .0235 |
| 14 | n | .0596 |
| 15 | o | .0689 |
| 16 | p | .0192 |
| 17 | q | .0008 |
| 18 | r | .0508 |
| 19 | s | .0567 |
| 20 | t | .0706 |
| 21 | u | .0334 |
| 22 | v | .0069 |
| 23 | w | .0119 |
| 24 | x | .0073 |
| 25 | y | .0164 |
| 26 | z | .0007 |
| 27 | – | .1928 |

■ All letters in the English alphabet have a very different probability $p_i$ of occurring

■ What is the number of bits you need to represent 27 characters? $\lceil \log_2 27 \rceil \approx \lceil 4.75 \rceil = 5$ bits

■ How can we measure the information in a single character? $h(p_i) = -\log_2 p_i$. Events with a low probability correspond to high information content

■ So, what is the average information in a character in an English text?

　■ $H(p) = \mathbb{E}[h(.)] = \sum_i p_i h(p_i) = -\sum_i p_i \log_2 p_i \approx 4.1$

This quantity is called the entropy. On average, with the right encoding, we can represent each letter with 4.1 bits instead of 4.7

# Entropy of Distributions



What is the "difference" between these distributions?

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Kullback-Leibler Divergence

- The Kullback-Leibler Divergence - KL Divergence - is a similarity measure between two distributions, and is defined as

$$
\text{KL}\,(p||q) = - \int p(x) \ln q(x) \mathrm{d}x - \left( - \int p(x) \ln p(x) \mathrm{d}x \right)
$$
$$
= - \int p(x) \ln \frac{q(x)}{p(x)} \mathrm{d}x
$$

- It represents the average additional amount of extra bits required to specify a symbol $x$, given that its underlying probability distribution is the estimated $q(x)$ and not the true one $p(x)$

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Kullback-Leibler Divergence

- Some properties
  - It is not a distance: $\text{KL}\,(p||q) \neq \text{KL}\,(q||p)$

  - It is non-negative: $\text{KL}\,(p||q) \geq 0$

  - If $\forall x\ p(x) = q(x)$: $\text{KL}\,(p||q) = 0$

- There are other metrics of similarity, but as we will see further in the course, the KL Divergence is deeply connected with maximum likelihood estimation

# **Outline**

TECHNISCHE
UNIVERSITÄT
DARMSTADT

## 6. Wrap-Up

You know now:

- What random variables are (both continuous and discrete)
- What probability distributions are
- Some basic rules of probability theory
- What expectation and variance are
- What a Gaussian distribution is and why it is so important
- What information and entropy are
- How to measure the similarity between two probability distributions

TECHNISCHE
UNIVERSITÄT
DARMSTADT

# Self-Test Questions

- What is a random variable?
- What is a distribution?
- What is a Binomial distribution?
- How does a Poisson distribution relate to Binomial distributions?
- What is a Gaussian distribution?
- What is an expectation?
- What is a joint distribution?
- What is a conditional distribution?
- What is a distribution with a lot of information?
- How to measure the difference between distributions?

# **Homework**

- Reading Assignment for next lecture
  - Bishop appendix E