

ECE 1513: Introduction to Machine Learning

Lecture 2: Probabilistic Modeling and Density Estimation

Ali Bereyhi

ali.bereyhi@utoronto.ca

Department of Electrical and Computer Engineering
University of Toronto

Winter 2025

Quick Recap: ML General Recipe

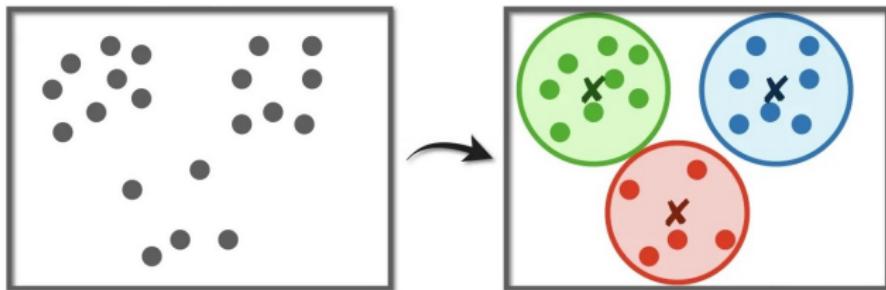
We defined ML as

the set of data-driven approaches that help us understand the environment and its behavior, and generalize it!

Any learning task is accomplished by ML through *three major steps*

- Collect data
- Specify a model *that captures the pattern*
- Develop a learning algorithm

Quick Recap: Clustering



- *Data*
 - ↳ Collection of samples $\mathbb{D} = \{x_n : n = 1, \dots, N\}$
- *Model*
 - ↳ A function mapping x to a cluster, e.g., K -centroid model
- *Learning algorithm*
 - ↳ It takes \mathbb{D} and returns a **good** clustering

Quick Recap: Clustering

K -Means():

- 1: Initiate μ_1, \dots, μ_K
- 2: **while** μ_1, \dots, μ_K changing **do**
- 3: Set $\mathcal{C}_1, \dots, \mathcal{C}_K \leftarrow \text{Cluster_Assignment}(\mu_1, \dots, \mu_K)$
- 4: Update $\mu_1, \dots, \mu_K \leftarrow \text{Centroid_Update}(\mathcal{C}_1, \dots, \mathcal{C}_K)$
- 5: **end while**
- 6: Return μ_1, \dots, μ_K

The **optimal** clustering algorithm should minimize **average risk**

$$\mathcal{J} = \frac{1}{N} \sum_{k=1}^K \sum_{n=1}^N r_{n,k} \|x_n - \mu_k\|^2$$

where $r_{n,k} = 1$ if $x_n \in \mathcal{C}_k$ and zero otherwise \rightsquigarrow the **smaller** \mathcal{J} , the **better**

Today's Agenda: Density Estimation

In today's lecture, we study another *unsupervised learning* problem, i.e.,
density estimation

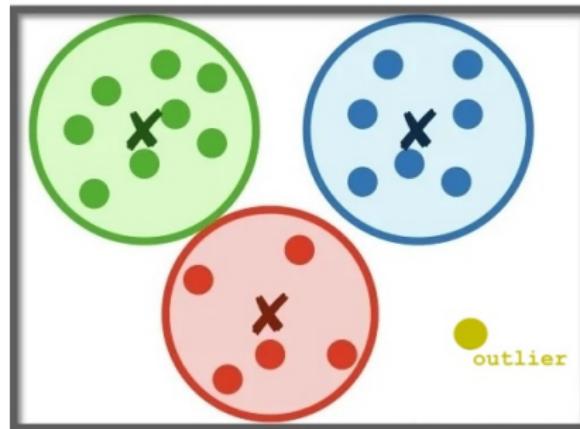
through the following steps

- *Probabilistic modeling of data*
 - ↳ We review key notions from probability theory
- *Maximum likelihood approach to density estimation*
 - ↳ We look at Gaussian maximum likelihood
- *Revisiting the clustering problem*
- *Few notes on how things extend beyond clustering*

We also talk about *validating a model* and *testing its generalization*

Motivation: Outlier Detection

In clustering, we talked about **outliers**



Let's look at **alternative** approach!

Outlier Detection: Example

We collect data from results of a game

$$\mathbb{D} = \{x_n : n = 1, \dots, N\}$$

with x_n being

$$x_n = \begin{cases} 1 & \text{Player } n \text{ wins} \\ 0 & \text{Player } n \text{ loses} \end{cases}$$

For instance, we collect $\mathbb{D} = \{1, 0, 0, 1, 1, 0, 0, 0, 0, 0\}$ in 10 days

Outlier Detection

Is it typical (normal) to see in next 3 days the following observation?

$$\{1, 1, 1\}$$

Outlier Detection: Learning Task

This is a learning task

- *Data*
 - ↳ $\mathbb{D} = \{x_n : n = 1, \dots, N\}$ collected in N days
- *Pattern \rightsquigarrow Model*
 - ↳ There is a *certain winning chance* in *normal* situation
- *Learning Algorithm*
 - ↳ $\mathbb{D} \mapsto$ Is what we see coming from a *normal situation*?

Density Estimation

We can interpret this problem from **stochastic** viewpoint

- We observe some data samples

$$\mathbb{D} = \{\boldsymbol{x}_n : n = 1, \dots, N\}$$

- We know (assume) that they are coming from a random process
 - ↳ The model is the distribution of this random process
- We want to infer from \mathbb{D} what the random process is
 - ↳ This is what the learning algorithm does

Let's recap key notions in Probability Theory

Discrete Random Variables

Say $x \in \mathbb{X} = \{a_1, \dots, a_I\}$; then, $x \sim P(x)$ implies that

- *Probability of $x = a_i$ is $P(a_i)$*
- *$P(x)$ is non-negative for any $x \in \mathbb{X}$ and less than 1*

$$0 \leq P(x) \leq 1$$

- *$P(x)$ adds up to 1*

$$\sum_{x \in \mathbb{X}} P(x) = \sum_{i=1}^I P(a_i) = 1$$

$P(x) \equiv \text{Probability Mass Function} \rightsquigarrow \text{Distribution}$

Discrete Random Variables: Example

Binary (Bernoulli) random variable $x \in \mathbb{X} = \{0, 1\}$

$$x \sim \text{Ber}(\theta)$$

is distributed as

$$P(x) = \begin{cases} \theta & x = 1 \\ 1 - \theta & x = 0 \end{cases}$$

Continuous Random Variables

Say $x \in \mathbb{R}$; then, $x \sim P(x)$ implies that

- *Probability of $a < x \leq b$ is*

$$\int_a^b P(x) dx$$

- *$P(x)$ is non-negative for any $x \in \mathbb{X}$*

$$P(x) \geq 0$$

- *$P(x)$ adds up to 1*

$$\int_{-\infty}^{+\infty} P(x) dx = 1$$

$P(x) \equiv$ Probability Density Function \rightsquigarrow Distribution

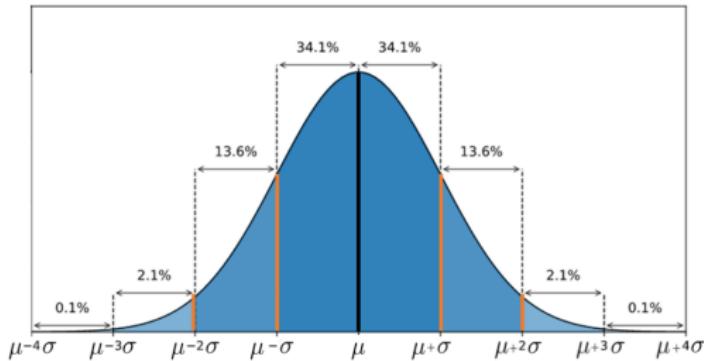
Continuous Random Variables: Example

Gaussian random variable with mean μ and variance σ^2

$$x \sim \mathcal{N}(\mu, \sigma^2)$$

is distributed as

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



Independent Processes

Random variables $x \sim P(x)$ and $y \sim P(y)$ are **independent** if

$$P(x, y) = P(x) P(y)$$

i.i.d. Sequence

x_1, \dots, x_N are *independent and identically distributed (i.i.d.)* if they are **independent** and have the **same distribution** \equiv

independent samples of the same process

Example: i.i.d. Bernoulli Samples

We sample $\text{Ber}(\theta)$ 5 times **independently**; then, the probability of observing

$$x_1, x_2, x_3, x_4, x_5 = 1, 1, 0, 1, 0$$

is given by

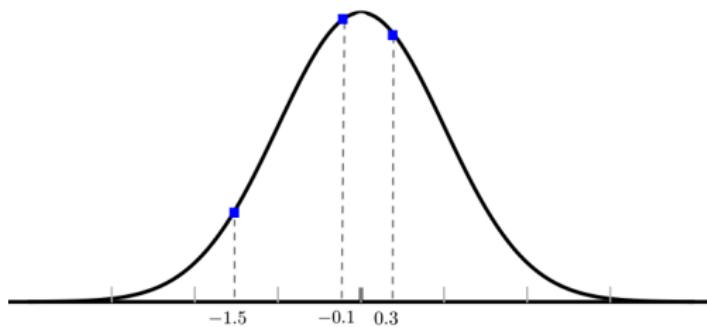
$$\begin{aligned} P(x_1, x_2, x_3, x_4, x_5) &= P(1, 1, 0, 1, 0) \\ &= \theta^3 (1 - \theta)^2 \end{aligned}$$

Expectation: Gaussian Example

We sample $\mathcal{N}(0, 1)$ 3 times **independently**; then, the density at

$$x_1, x_2, x_3 = 0.3, -1.5, -0.1$$

is given by



$$P(x_1, x_2, x_3) = \left(\frac{1}{\sqrt{2\pi}}\right)^3 \exp\{-0.045\} \exp\{-1.125\} \exp\{-0.005\}$$

Expectation: Average in Asymptotic Sense

Say $x \sim P(x)$; then,

$$\mathbb{E}\{f(x)\} = \begin{cases} \sum_{x \in \mathbb{X}} f(x) P(x) & x \text{ is discrete} \\ \int_{-\infty}^{+\infty} f(x) P(x) dx & x \text{ is continuous} \end{cases}$$

Intuitive Meaning

Say x_1, \dots, x_N are i.i.d. samples from $P(x)$; then,

$$\frac{1}{N} \sum_{n=1}^N f(x_n) \approx \mathbb{E}\{f(x)\}$$

when N is large

Expectation: Bernoulli Example

Say $x \sim \text{Ber}(\theta)$: this means that

$$\mathbb{E}\{x\} = \theta$$

Intuitive Meaning

Say x_1, \dots, x_N are independent samples of $\text{Ber}(\theta)$; then,

$$\frac{1}{N} \sum_{n=1}^N x_n \approx \theta$$

when N is large

Expectation: Gaussian Example

Say $x \sim \mathcal{N}(\mu, \sigma^2)$: this means that

$$\mathbb{E}\{x\} = \mu$$

$$\text{Var}\{x\} = \mathbb{E}\{(x - \mu)^2\} = \sigma^2$$

Intuitive Meaning

Say x_1, \dots, x_N are independent Gaussian samples from $\mathcal{N}(\mu, \sigma^2)$; then,

$$\frac{1}{N} \sum_{n=1}^N x_n \approx \mu$$

$$\frac{1}{N} \sum_{n=1}^N (x_n - \mu)^2 \approx \sigma^2$$

when N is large

Random Variables \rightsquigarrow Random Objects

We can extend all these notions to higher dimensions

- x is a random vector whose entries are random variables
- X is a random matrix whose entries are random variables
- ...

We then write

$$x \sim P(x)$$

and mean by $P(x)$ the joint distribution of all random entries

Density Estimation: Formulation

We can interpret this problem from stochastic viewpoint

- We observe some data samples

$$\mathbb{D} = \{\mathbf{x}_n : n = 1, \dots, N\}$$

- We know (assume) that they are coming from a random process
 - ↳ The model is the distribution of this random process
- We want to infer from \mathbb{D} what the random process is
 - ↳ This is what the learning algorithm does

Modeling Distributions

We observe some data samples

$$\mathbb{D} = \{\mathbf{x}_n : n = 1, \dots, N\}$$

and may see a pattern in the data, e.g.,

- ↳ *Samples are binary*
- ↳ *Histogram of samples looks roughly Gaussian*

?

How to model this pattern?

!

We assume \mathbf{x}_n 's are sampled from a distribution with some unknowns

$$\mathbf{x}_n \sim P_{\theta}(\mathbf{x}_n)$$

Modeling Distributions: Binary Example

We observe the results of a game: *samples are*

$$\mathbb{D} = \{x_n \in \{0, 1\} : n = 1, \dots, N\}$$

If we know each sample is given by an **independent** player, we could **assume**

$$x_n \sim P_{\theta}(x_n) \equiv \text{Ber}(\theta)$$

and hence the probability of whole dataset is

$$\begin{aligned}\mathbb{D} \sim P_{\theta}(x_1, \dots, x_N) &= \prod_{n=1}^N P_{\theta}(x_n) \\ &= \theta^{\# \text{ of } 1's} (1 - \theta)^{\# \text{ of } 0's}\end{aligned}$$

Modeling Distributions: Gaussian Example

We observe the a Gaussian random generator: *samples are*

$$\mathbb{D} = \{x_n : n = 1, \dots, N\}$$

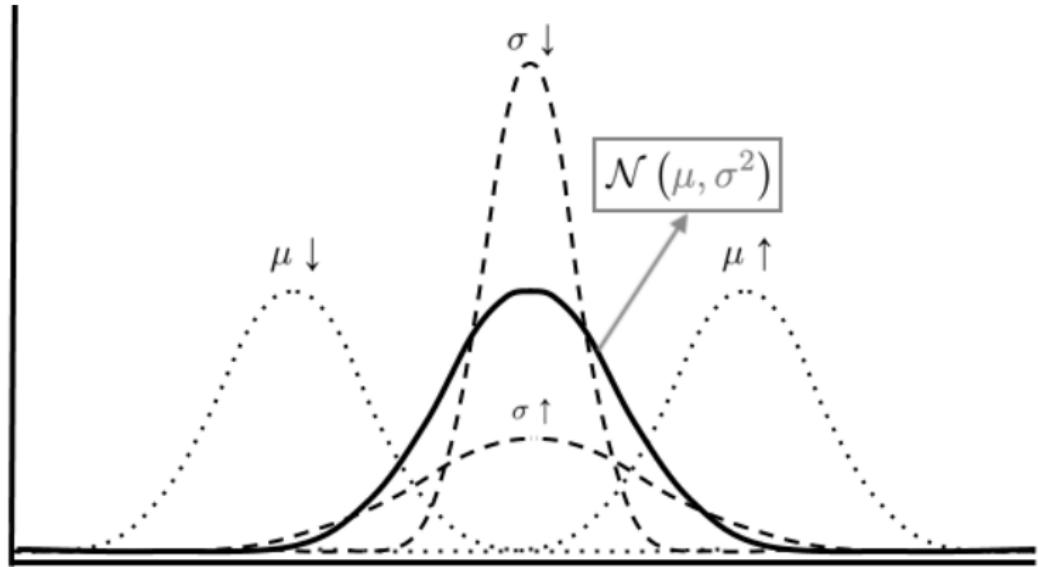
If we know *samples are generated independently*, we know

$$x_n \sim P_{\mu, \sigma}(x_n) \equiv \mathcal{N}(\mu, \sigma^2)$$

and hence the probability of whole dataset is

$$\begin{aligned}\mathbb{D} \sim P_{\mu, \sigma}(x_1, \dots, x_N) &= \prod_{n=1}^N P_{\mu, \sigma}(x_n) \\ &= \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right)^N \exp \left\{ -\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 \right\}\end{aligned}$$

Modeling Distributions: Gaussian Example



Likelihood: Key to Modeling Distributions

We assume a **model** for the distribution as

$$\mathbf{x}_n \sim P_{\theta}(\mathbf{x}_n)$$

Likelihood

The likelihood of the dataset is

$$\mathcal{L}_{\mathbb{D}}(\boldsymbol{\theta}) = P_{\theta}(\mathbf{x}_1, \dots, \mathbf{x}_N)$$

Attention!

Likelihood is only a function of the **unknown θ**

Likelihood: Binary Example

We observe 5 sample results of a game

$$\mathbb{D} = \{1, 1, 0, 1, 0\}$$

Assuming them to be samples of $\text{Ber}(\theta)$, the likelihood is

$$\begin{aligned}\mathcal{L}_{\mathbb{D}}(\theta) &= P_{\theta}(1, 1, 0, 1, 0) \\ &= \theta^3 (1 - \theta)^2\end{aligned}$$

Attention!

Likelihood is only a function of the unknown θ

Likelihood: Gaussian Example

We observe three samples of a Gaussian random generator

$$\mathbb{D} = \{0.3, -1.5, -0.1\}$$

The likelihood is

$$\begin{aligned}\mathcal{L}_{\mathbb{D}}(\mu, \sigma) &= P_{\mu, \sigma}(0.3, -1.5, -0.1) \\ &= \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right)^3 \exp \left\{ -\frac{(\mu + 0.1)^2 + (\mu - 0.5)^2 + (\mu - 1.5)^2}{2\sigma^2} \right\}\end{aligned}$$

Attention!

Likelihood is only a function of the unknowns μ and σ

Optimal Model \equiv Maximum Likelihood

- ? How can we find right choice for the **unknown θ**
- ! We find the **unknown θ** that maximizes the likelihood of dataset

Notion of **Optimality** \equiv Likelihood

This does the same job that **distortion** did for us in clustering

Maximum Likelihood Estimation

The optimal model is given by

$$\begin{aligned}\theta^* &= \operatorname{argmax}_{\theta} \mathcal{L}_{\mathbb{D}} (\theta) \\ &= \operatorname{argmax}_{\theta} P_{\theta} (x_1, \dots, x_N)\end{aligned}$$

Classic case is that we assume samples in dataset are i.i.d.

$$\begin{aligned}\theta^* &= \operatorname{argmax}_{\theta} \mathcal{L}_{\mathbb{D}} (\theta) \\ &= \operatorname{argmax}_{\theta} \prod_{n=1}^N P_{\theta} (x_n) = \operatorname{argmax}_{\theta} \prod_{n=1}^N \mathcal{L}_n (\theta)\end{aligned}$$

Maximum Likelihood Estimation

The optimal model is given by

$$\boldsymbol{\theta}^* = \operatorname{argmax}_{\boldsymbol{\theta}} \prod_{n=1}^N \mathcal{L}_n(\boldsymbol{\theta})$$

Log Trick

Assume $f(\boldsymbol{\theta}) \geq 0$; then, we can write

$$\operatorname{argmax}_{\boldsymbol{\theta}} f(\boldsymbol{\theta}) = \operatorname{argmax}_{\boldsymbol{\theta}} \log f(\boldsymbol{\theta})$$

Using log trick we have

$$\boldsymbol{\theta}^* = \operatorname{argmax}_{\boldsymbol{\theta}} \sum_{n=1}^N \log \mathcal{L}_n(\boldsymbol{\theta}) = \operatorname{argmax}_{\boldsymbol{\theta}} \sum_{n=1}^N \mathcal{S}_n(\boldsymbol{\theta})$$

Maximum Likelihood: Binary Example

We observe 5 sample results of a game

$$\mathbb{D} = \{1, 1, 0, 1, 0\}$$

Assuming them to be samples of $\text{Ber}(\theta)$, the likelihood is

$$\mathcal{L}_{\mathbb{D}}(\theta) = \theta^3 (1 - \theta)^2$$

Maximum likelihood estimate is

$$\theta^* = \operatorname{argmax}_{\theta} [3 \log \theta + 2 \log (1 - \theta)] = 0.6$$

Maximum Likelihood: Gaussian Example

We observe three samples of a Gaussian random generator $\mathcal{N}(\mu, 1)$

$$\mathbb{D} = \{0.3, -1.5, -0.1\}$$

The likelihood is

$$\mathcal{L}_{\mathbb{D}}(\mu) = \left(\frac{1}{\sqrt{2\pi}} \right)^3 \exp \left\{ -\frac{(\mu + 0.1)^2 + (\mu - 0.5)^2 + (\mu - 1.5)^2}{2} \right\}$$

Maximum likelihood estimate is

$$\begin{aligned}\theta^* &= \underset{\theta}{\operatorname{argmax}} \left[-\frac{3}{2} \log 2\pi - \frac{(\mu + 0.1)^2 + (\mu - 0.5)^2 + (\mu - 1.5)^2}{2} \right] \\ &= 0.63\end{aligned}$$

Multivariate Gaussian Distribution

A basic approach to model the data is to fit it into a **Gaussian** distribution

Multivariate Gaussian Distribution

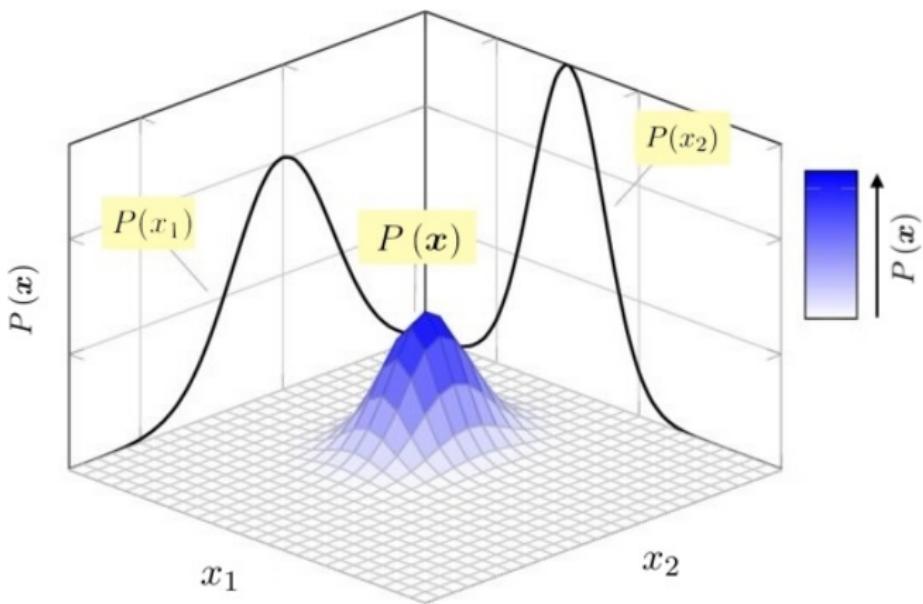
Random vector $\boldsymbol{x} \in \mathbb{R}^d$ is Gaussian with mean $\boldsymbol{\mu}$ and covariance $\boldsymbol{\Sigma}$ and noted as $\boldsymbol{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ if

$$P(\boldsymbol{x}) = \frac{1}{\sqrt{(2\pi)^d \det(\boldsymbol{\Sigma})}} \exp \left\{ -\frac{1}{2} (\boldsymbol{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\boldsymbol{x} - \boldsymbol{\mu}) \right\}$$

Note that

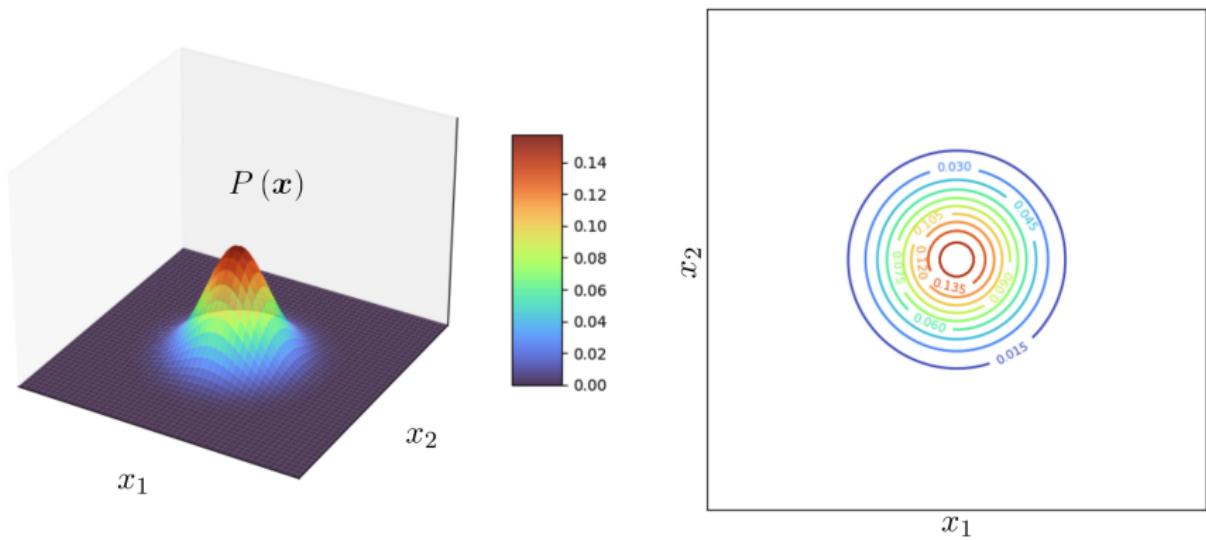
$$\boldsymbol{\mu} = \mathbb{E}\{\boldsymbol{x}\} \quad \boldsymbol{\Sigma} = \mathbb{E}\{(\boldsymbol{x} - \boldsymbol{\mu})(\boldsymbol{x} - \boldsymbol{\mu})^\top\}$$

Example: Bivariate Gaussian Distribution



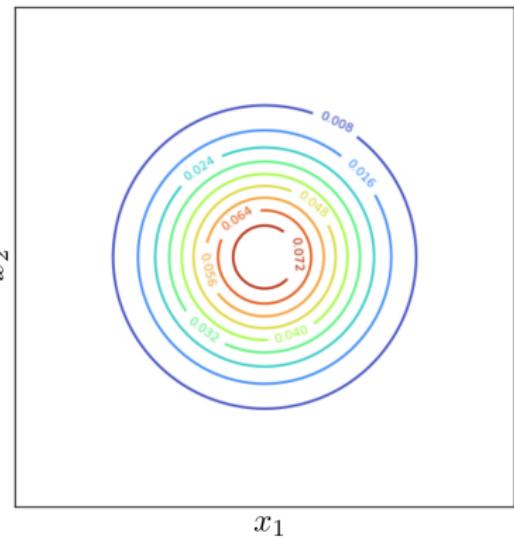
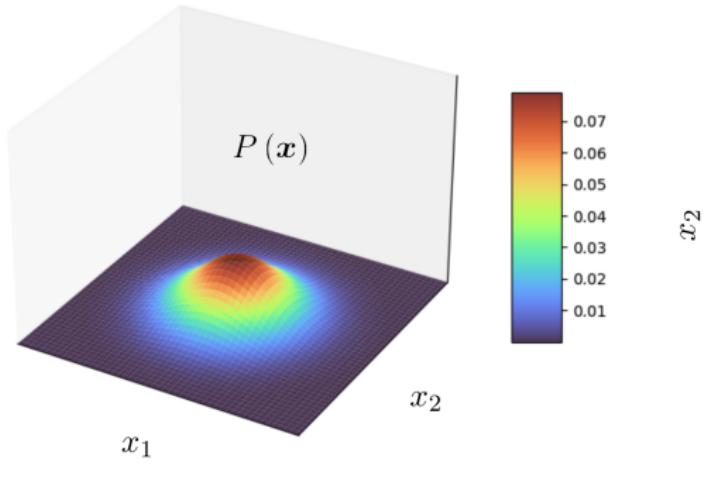
Example: Bivariate Gaussian Distribution

$$\Sigma = \mathbf{I} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



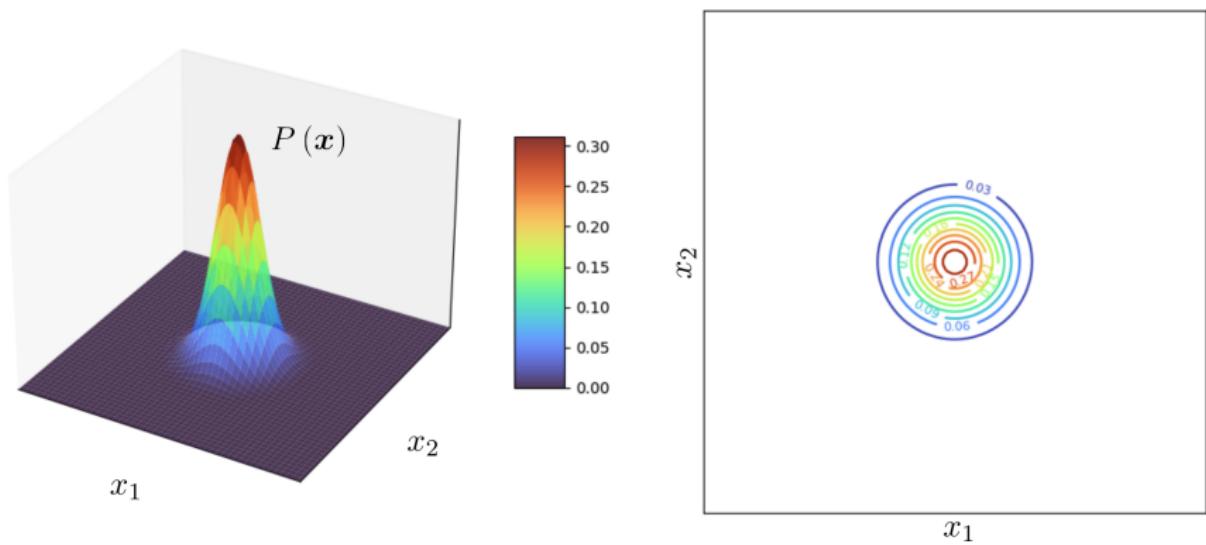
Example: Bivariate Gaussian Distribution

$$\Sigma = 2\mathbf{I} = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$



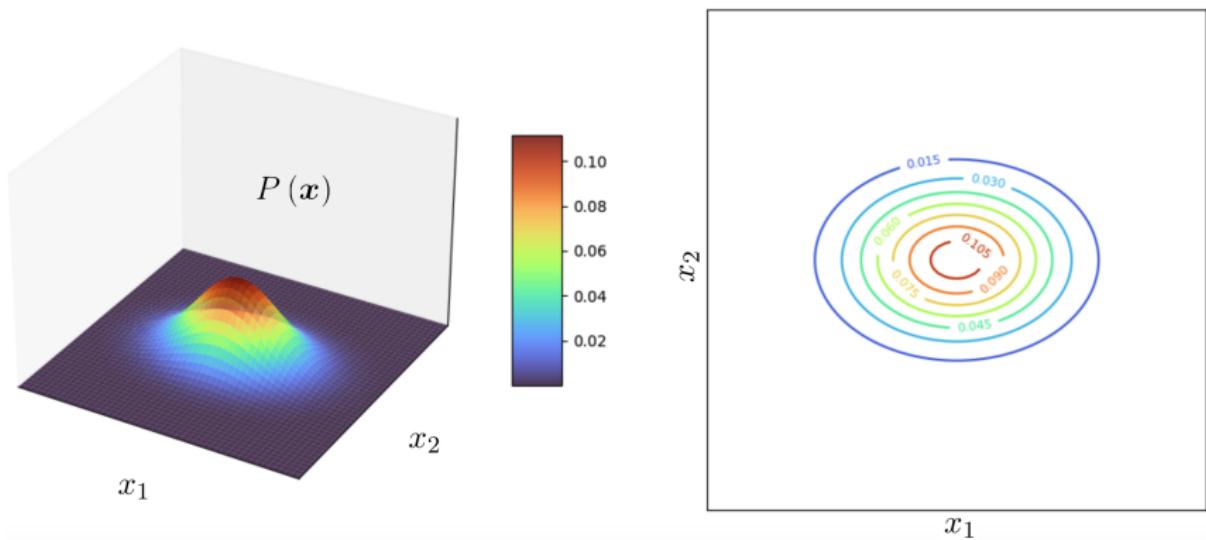
Example: Bivariate Gaussian Distribution

$$\Sigma = 0.5\mathbf{I} = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}$$



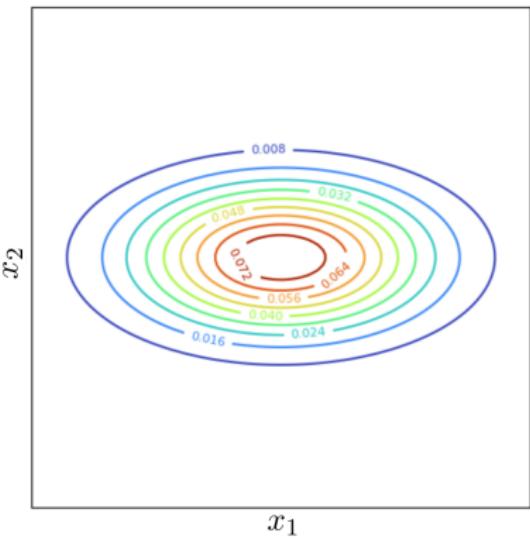
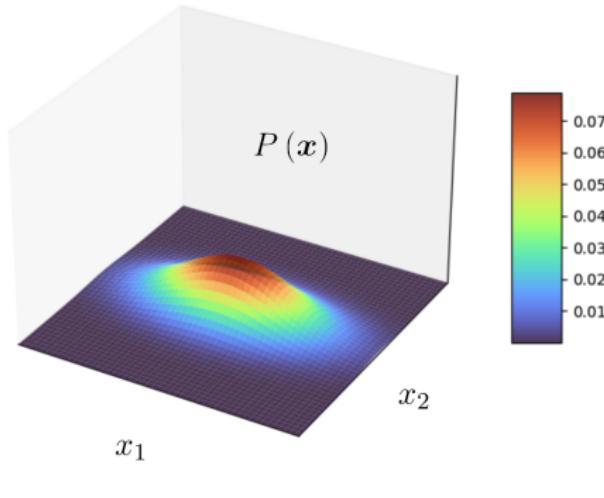
Example: Bivariate Gaussian Distribution

$$\Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}$$



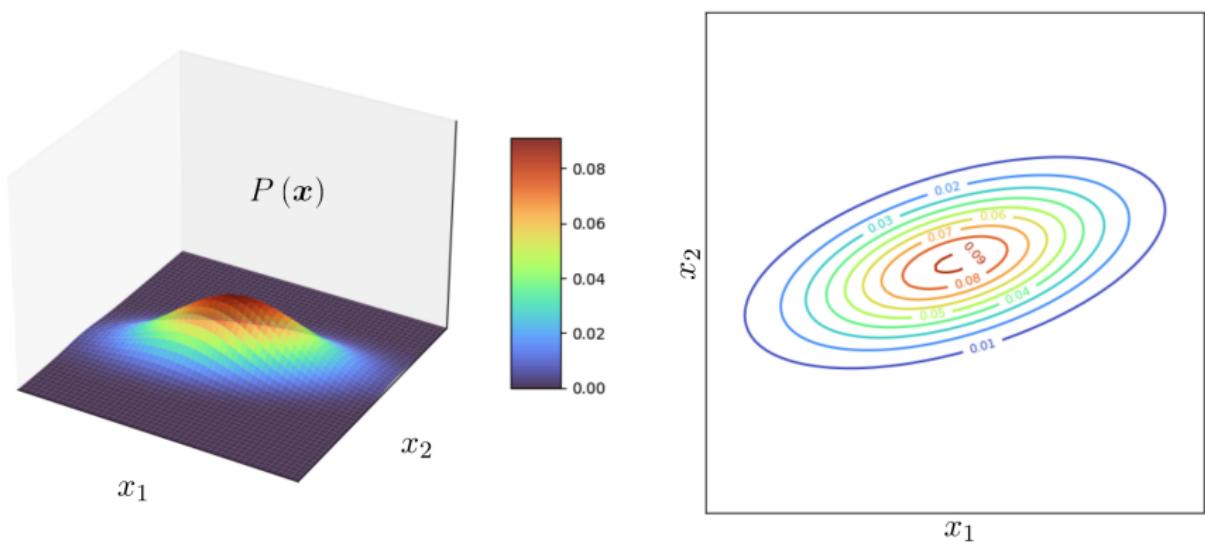
Example: Bivariate Gaussian Distribution

$$\Sigma = \begin{bmatrix} 4 & 0 \\ 0 & 1 \end{bmatrix}$$



Example: Bivariate Gaussian Distribution

$$\Sigma = \begin{bmatrix} 4 & 1 \\ 1 & 1 \end{bmatrix}$$



Multivariate Gaussian Distribution

A Gaussian model assumes the data is sampled from $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

Gaussian Model

Points $\mathbf{x}_n \in \mathbb{D}$ are sampled from $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$; thus,

$$P_{\boldsymbol{\mu}, \boldsymbol{\Sigma}}(\mathbf{x}_n) = \frac{1}{\sqrt{(2\pi)^d \det(\boldsymbol{\Sigma})}} \exp \left\{ -\frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \right\}$$

Log-likelihood is then

$$\mathcal{S}_n(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = -\frac{d}{2} \log 2\pi - \frac{1}{2} \log \det(\boldsymbol{\Sigma}) - \frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu})$$

Why Gaussian?

Many processes in nature are approximately Gaussian

Central Limit Theorem – *Informal*

Sum of large enough independent random objects tends to a Gaussian object

Superposition of Gaussians can describe a large class of distributions

$$\mathcal{Q} \equiv \sum_{m=1}^M \pi_m \mathcal{N}(\boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m)$$

↳ We mention it again later

Gaussian Maximum Likelihood

We assume

$$\mathbb{D} = \{\mathbf{x}_n : n = 1, \dots, N\}$$

contains samples of $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

No harm if we further *normalize*

$$\begin{aligned}\mathcal{S}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) &= \frac{1}{N} \sum_{n=1}^N \mathcal{S}_n(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \\ &= -\frac{d}{2} \log 2\pi - \frac{1}{2} \log \det(\boldsymbol{\Sigma}) - \frac{1}{2N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu})\end{aligned}$$

Gaussian Maximum Likelihood: Identity Covariance

Assume we know that $\Sigma = \mathbf{I}$

$$\mathcal{S}(\boldsymbol{\mu}) = -\frac{d}{2} \log 2\pi - \frac{1}{2N} \sum_{n=1}^N \|\mathbf{x}_n - \boldsymbol{\mu}\|^2$$

Maximum likelihood estimate is at the center of cluster

$$\boldsymbol{\mu}^* = \operatorname{argmax}_{\boldsymbol{\mu}} \mathcal{S}(\boldsymbol{\mu}) = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$$

! That makes sense!

Gaussian Maximum Likelihood: Unknown Covariance

For general Σ , we have

$$\mathcal{S}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = -\frac{d}{2} \log 2\pi - \frac{1}{2} \log \det(\boldsymbol{\Sigma}) - \frac{1}{2N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu})$$

Let's first think of one-dimensional case $\boldsymbol{\mu} = \mu$ and $\boldsymbol{\Sigma} = \sigma^2$

$$\mathcal{S}(\mu, \sigma^2) = -\frac{d}{2} \log 2\pi - \frac{1}{2} \log \sigma^2 - \frac{1}{2N\sigma^2} \sum_{n=1}^N (x_n - \mu)^2$$

? How to find the optimal μ and σ^2 ?

! Set the partial derivatives to zero

Gaussian Maximum Likelihood: Unknown Covariance

For general μ and σ^2 , we have

$$\mathcal{S}(\mu, \sigma^2) = -\frac{d}{2} \log 2\pi - \frac{1}{2} \log \sigma^2 - \frac{1}{2N\sigma^2} \sum_{n=1}^N (x_n - \mu)^2$$

Take partial derivative w.r.t. μ

$$\frac{\partial \mathcal{S}}{\partial \mu} = \frac{1}{N\sigma^2} \sum_{n=1}^N (x_n - \mu)$$

and set to zero

$$\frac{\partial \mathcal{S}}{\partial \mu} = 0 \rightsquigarrow \mu^\star = \frac{1}{N} \sum_{n=1}^N x_n$$

Gaussian Maximum Likelihood: Unknown Covariance

For general μ and σ^2 , we have

$$\mathcal{S}(\mu, \sigma^2) = -\frac{d}{2} \log 2\pi - \frac{1}{2} \log \sigma^2 - \frac{1}{2N\sigma^2} \sum_{n=1}^N (x_n - \mu)^2$$

Take partial derivative w.r.t. σ^2

$$\frac{\partial \mathcal{S}}{\partial \sigma^2} = -\frac{1}{2\sigma^2} + \frac{1}{N(\sigma^2)^2} \sum_{n=1}^N (x_n - \mu)^2$$

and set to zero

$$\frac{\partial \mathcal{S}}{\partial \sigma^2} = 0 \rightsquigarrow \sigma^{*2} = \frac{1}{N} \sum_{n=1}^N (x_n - \mu^*)^2$$

Gaussian Maximum Likelihood: Unknown Covariance

In vector form, we do the same with gradients: set gradient w.r.t. μ to zero

$$\nabla_{\mu} \mathcal{S}(\mu, \Sigma) = \mathbf{0} \rightsquigarrow \mu^* = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n = \text{avg}(\mathbb{D})$$

And then gradient w.r.t. Σ to zero

$$\nabla_{\Sigma} \mathcal{S}(\mu, \Sigma) = \mathbf{0} \rightsquigarrow \Sigma^* = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \mu^*)(\mathbf{x}_n - \mu^*)^T = \text{cov}(\mathbb{D})$$

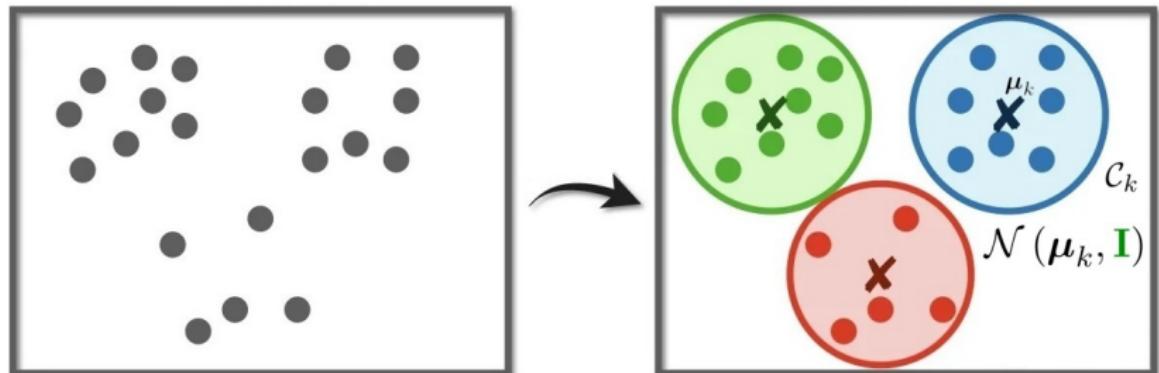
Attention

No worries if it's new for you! We will discuss gradient descent and function optimization in the next weeks!

Gaussian Model for Clustering

For clustering, assume *each cluster is sampled from a Gaussian distribution*

$$\mathbf{x}_n \in \mathcal{C}_k \rightsquigarrow \mathbf{x}_n \sim \mathcal{N}(\boldsymbol{\mu}_k, \mathbf{I})$$



Maximum Likelihood: Estimating Cluster Set

? Given K centroids, to which cluster belongs x_n ?

Log-likelihood of $x_n \in \mathcal{C}_k$ is

$$\mathcal{S}_n(\mathbf{k}) = \log P_{\mathbf{k}}(x_n) = -\frac{d}{2} \log 2\pi - \frac{1}{2} \|x_n - \mu_{\mathbf{k}}\|^2$$

The maximum likelihood estimate is thus given by

$$\begin{aligned} k^* &\leftarrow \operatorname{argmax}_{\mathbf{k}} \mathcal{S}_n(\mathbf{k}) = \operatorname{argmin}_{\mathbf{k}} \|x_n - \mu_{\mathbf{k}}\|^2 \\ &\equiv \text{Cluster_Assignment}() \end{aligned}$$

Maximum Likelihood: Estimating Centroids

? Given samples of a cluster \mathcal{C}_k , what is best μ_k ?

Log-likelihood of μ_k is

$$\mathcal{S}_{\mathcal{C}_k}(\mu_k) = \frac{1}{|\mathcal{C}_k|} \log P_{\mu_k}(\mathcal{C}_k) = -\frac{d}{2} \log 2\pi - \frac{1}{2|\mathcal{C}_k|} \sum_{x_n \in \mathcal{C}_k} \|x_n - \mu_k\|^2$$

The maximum likelihood is thus given by

$$\begin{aligned}\mu_k^* &\leftarrow \underset{\mu_k}{\operatorname{argmax}} \mathcal{S}_{\mathcal{C}_k}(\mu_k) = \frac{1}{|\mathcal{C}_k|} \sum_{x_n \in \mathcal{C}_k} \|x_n - \mu_k\|^2 \\ &\equiv \text{Centroid_Update}()\end{aligned}$$

Conclusion: K -Means Maximizes Likelihoods

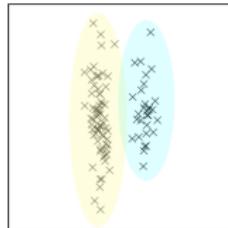
K -Means():

- 1: Initiate μ_1, \dots, μ_K
- 2: **while** μ_1, \dots, μ_K changing **do**
- 3: Set $\mathcal{C}_1, \dots, \mathcal{C}_K \leftarrow \text{Cluster_Assignment}(\mu_1, \dots, \mu_K)$
- 4: Update $\mu_1, \dots, \mu_K \leftarrow \text{Centroid_Update}(\mathcal{C}_1, \dots, \mathcal{C}_K)$
- 5: **end while**
- 6: Return μ_1, \dots, μ_K

K -Means Clustering \equiv Iterative Maximum Likelihood

Possible Extensions

- In clustering, we could also assume Σ to be *unknown*



- In general, we can assume a distribution \mathcal{Q} that is *mixture of Gaussians*

$$\mathcal{Q} \equiv \sum_{m=1}^M \pi_m \mathcal{N}(\boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m) \text{ with } \sum_{m=1}^M \pi_m = 1$$

- ↳ *Gaussian Mixture Model (GMM)*
- ↳ We estimate $\pi_m, \boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m$ for $m = 1, \dots, M$ all together
- ↳ Larger M makes it better, but needs more computation!

Data Generation

Say x collects all pixel values of an image

$$\mathbb{D} = \{x_n = \text{cat image} : n = 1, \dots, N\}$$



We learn the distribution by finding a **good** estimate $P_{\theta}(x)$

$$x_{\text{new}} \sim P_{\theta}(x) \approx P_{\text{data}}(x)$$

x_{new} is a new cat image!



In practice, we need more sophisticated method to estimate $P_{\text{data}}(x)$

Further Read

- Bishop
 - ↳ Chapter 1: Section 1.2
 - ↳ Chapter 2: Sections 2.1 and 2.3
 - ↳ Chapter 9: Section 9.2
 - MacKay
 - ↳ Chapter 22
 - ↳ Chapter 23
 - Goodfellow
 - ↳ Chapter 3
 - ↳ Chapter 5: Section 5.5
 - ESL
 - ↳ Chapter 14: Section 14.3.7
- Probability Review**
Gaussian Maximum Likelihood
Mixture of Gaussians
- Maximum Likelihood vs Clustering**
More on Probability
- Probability Review**
Maximum Likelihood
- Mixture of Gaussians**

Testing A Model

- ? How can we know if our model has learned the pattern?
- ! We can **test** it

Collect a new set of samples

$$\mathbb{T} = \{\mathbf{x}_i : i = 1, \dots, I\}$$

Use the **trained** model θ^* to find the likelihood

$$\mathcal{S}_{\mathbb{T}}(\theta^*) = \log P_{\theta^*}(\mathbb{T})$$

- ↳ If $\mathcal{S}_{\mathbb{T}}(\theta^*) \approx \mathcal{S}_{\mathbb{D}}(\theta^*) \rightsquigarrow$ the model generalizes
- ↳ If $\mathcal{S}_{\mathbb{T}}(\theta^*) \ll \mathcal{S}_{\mathbb{D}}(\theta^*) \rightsquigarrow$ the model does not generalize!
 - ↳ Maybe we need more complicated or simpler model!

Validating the Model

We typically reserve some samples to validate

$$\mathbb{V} = \{\mathbf{x}_j : i = 1, \dots, J\}$$

We train the model θ^* on \mathbb{D} ; then, check

$$\mathcal{S}_{\mathbb{V}}(\theta^*) = \log P_{\theta^*}(\mathbb{V})$$

- ↳ If $\mathcal{S}_{\mathbb{V}}(\theta^*) \approx \mathcal{S}_{\mathbb{D}}(\theta^*) \rightsquigarrow$ the model is valid
 - ↳ We now check generalization on the test set \mathbb{T}
- ↳ If $\mathcal{S}_{\mathbb{V}}(\theta^*) \ll \mathcal{S}_{\mathbb{D}}(\theta^*) \rightsquigarrow$ the model is not valid!
 - ↳ We change some **hyperparameters** in the model and repeat

Test and Validation: Example

Consider clustering problem

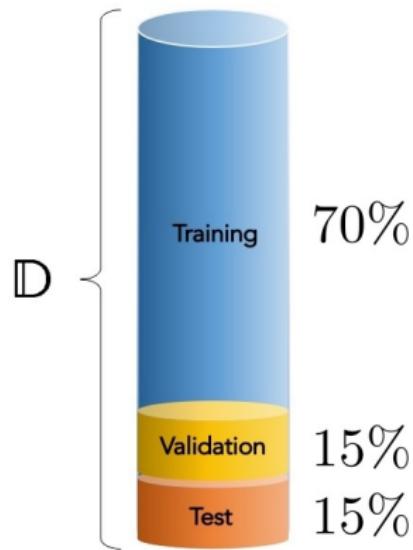
- ① Choose some *hyperparameter* K
- ② Train via K -means clustering on \mathbb{D}
 - ↳ Find μ_1, \dots, μ_K and save the final risk

$$\mathcal{J}_{\text{train}} = \frac{1}{N} \sum_{k=1}^K \sum_{n=1}^N \textcolor{blue}{r}_{n,k} \|\boldsymbol{x}_n - \boldsymbol{\mu}_k\|^2$$

- ③ Compute the risk on the set \mathbb{V} and call it \mathcal{J}_{val}
 - ↳ If $\mathcal{J}_{\text{val}} \approx \mathcal{J}_{\text{train}} \rightsquigarrow$ the model is valid
 - ↳ We now test if $\mathcal{J}_{\text{test}}$ on \mathbb{T} is also close
 - ↳ If $\mathcal{J}_{\text{val}} \gg \mathcal{J}_{\text{train}} \rightsquigarrow$ the model is not valid!
 - ↳ We change K to higher/lower value and get back to ①

Training vs Test Data

In practice, we split our collected data



Further Read

- Bishop
 - ↳ Chapter 1: *Section 1.3* **Validation and Test**
- Goodfellow
 - ↳ Chapter 5: *Sections 5.2 and 5.3* **Generalization**