

Lecture 9b

Probabilistic Models (cont.)

Outline

Modeling Sequential Data

- ▶ Motivations (recap)
- ▶ Tikhonov regularization (recap)
- ▶ Probabilistic approach: Gaussian processes

Modeling Clustered Data

- ▶ Motivations (recap)
- ▶ The K-means algorithm (recap)
- ▶ Probabilistic approach: Gaussian mixture model

Part 1

Modeling Sequential Data

Recap: Motivations

- ▶ Sequential data is found in many practical applications.
- ▶ It can be time series, images, geographical maps, etc.
- ▶ Sequential data exhibits correlation between adjacent instances → iid. assumption does not hold, and we need ad-hoc ML models.

Browse Through: 38 Data Sets

Name	Data Types	Default Task	Attribute Type	# Instances	# Attributes	Year
UCI AHJ 2020 Predictive Maintenance Dataset	Multivariate, Time-Series	Classification, Regression, Causal Discovery	Real	10000	14	2020
UCI Air Quality	Multivariate, Time-Series	Regression	Real	9358	15	2016
UCI Air quality	Multivariate, Time-Series	Regression	Real	9358	15	2016
UCI Aspirin sensor readings	Multivariate, Time-Series	Regression	Real	19735	29	2017
UCI Asustek Star Webcam Data	Multivariate, Time-Series	Classification	Categorical, Real	6850	15	1999
UCI Asustek Star Webcam Data (HMD Quality)	Multivariate, Time-Series	Classification	Real	2965	22	2002
UCI Baidu Multi-Site Air-Quality Data	Multivariate, Time-Series	Regression	Integer, Real	420768	18	2019
UCI Baidu PM2.5 Data	Multivariate, Time-Series	Regression	Integer, Real	43824	13	2017
UCI BikeSharingUsageDataset	Multivariate, Time-Series	Classification, Clustering	Integer, Real	2919987	10	2020
UCI MLR-RSS Dataset for Indoor Localization and Navigation	Multivariate, Sequential Time-Series	Classification, Clustering	Integer	6611	15	2019

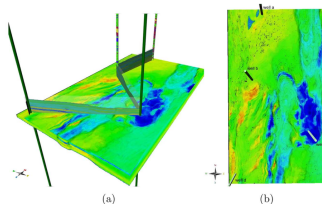
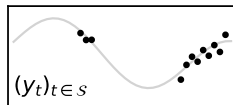


Image source: DOI:10.1016/j.cageo.2017.05.004

Recap: Tikhonov Regularization



Tikhonov regularization:

- ▶ Learn a surrogate model \mathbf{z} of the true sequence \mathbf{y} , and enforce correlation between adjacent elements in the sequence through a **penalty term**.

$$\mathcal{E}(\mathbf{z}) = \|\mathbf{M}\mathbf{z} - \mathbf{y}_S\|^2 + \|\mathbf{\Gamma}\mathbf{z}\|^2$$

The matrix \mathbf{M} selects elements of \mathbf{z} for which we have observations of \mathbf{y} , and $\mathbf{\Gamma}$ is a matrix designed by the user and whose rows are typically of the type:

$$\mathbf{\Gamma}_{i,:} = (0, \dots, 0, -1, 1, 0, \dots, 0) \quad (\text{penalizing slope})$$

$$\mathbf{\Gamma}_{i,:} = (0, \dots, 0, -1, 2, -1, 0, \dots, 0) \quad (\text{penalizing curvature})$$

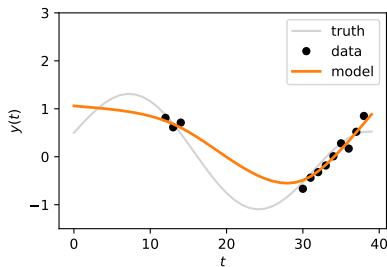
- ▶ A minimum of $\mathcal{E}(\mathbf{z})$ can be obtained in closed form as:

$$\mathbf{z} = \underbrace{(\mathbf{M}^\top \mathbf{M} + \mathbf{\Gamma}^\top \mathbf{\Gamma})^{-1} \mathbf{M}^\top}_{\mathbf{W}} \mathbf{y}_S$$

Recap: Tikhonov Regularization

Example:

- ▶ True sequence, observed data, and sequence model with *Tikhonov regularization*.



Observations:

- ▶ The model's predicted sequence implements local correlations, and provides plausible interpolations in regions without data.
- ▶ It is hard however to assess the level of uncertainty of the model without seeing the ground-truth.

Question:

- ▶ Can we use the probabilistic modeling approach to obtain an estimate of predictive uncertainty?

Gaussian Processes

Idea:

- ▶ View the predicted sequence as high-dimensional multivariate Gaussian distribution.

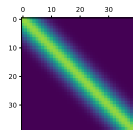
$$\mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$$

with $\boldsymbol{\mu} \in \mathbb{R}^T$ and $\Sigma \in \mathbb{R}^{T \times T}$, and T is the number of time steps in the sequence.

- ▶ The parameter $\boldsymbol{\mu}$ is e.g. a constant sequence set to the mean of \mathbf{y}_S , or something more complex learned from the data.
- ▶ The matrix Σ models the covariance between adjacent elements of the sequence, e.g.

$$\Sigma_{ij} = \alpha \cdot \exp(-\gamma \cdot \|i - j\|^2)$$

$\Sigma =$

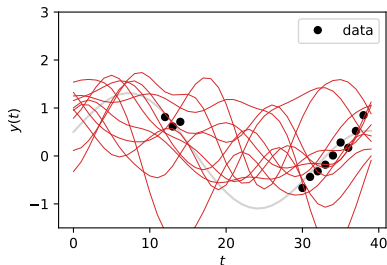


where i and j are positions in the sequence and (α, γ) are hyperparameters.

Gaussian Process

Example:

- Sequences z sampled from $\mathcal{N}(\mu, \Sigma)$, and shown in red:



Note:

- In practice, to avoid overfitting, we still need account for some possible uncorrelated intrinsic or measurement noise. We can achieve this by adding a diagonal term to our model of covariance: $\Sigma \leftarrow \Sigma + \epsilon^2 I$, where ϵ is the standard deviation of the noise.

Gaussian Process: Inferences

- ▶ To 'learn' from the available observations, we condition the Gaussian process to the observed data. We first partition \mathbf{z} in terms of observed and unobserved elements (resp. \mathbf{z}_S and \mathbf{z}_T):

$$\underbrace{\begin{bmatrix} \mathbf{z}_S \\ \mathbf{z}_T \end{bmatrix}}_{\mathbf{z}} \sim \mathcal{N}\left(\underbrace{\begin{bmatrix} \boldsymbol{\mu}_S \\ \boldsymbol{\mu}_T \end{bmatrix}}_{\boldsymbol{\mu}}, \underbrace{\begin{bmatrix} \Sigma_{SS} & \Sigma_{ST} \\ \Sigma_{TS} & \Sigma_{TT} \end{bmatrix}}_{\Sigma}\right)$$

- ▶ We now apply the formula for conditioning Gaussian distribution:

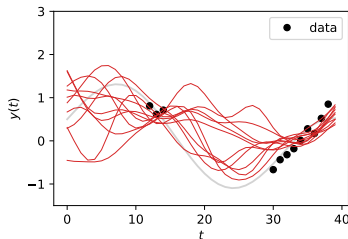
$$\mathbf{z}_T | \mathbf{z}_S \sim \mathcal{N}(\boldsymbol{\mu}', \Sigma')$$

with

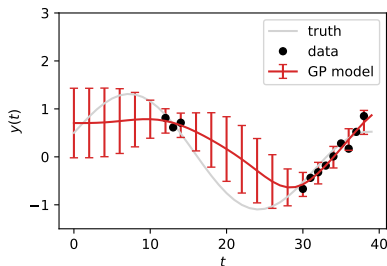
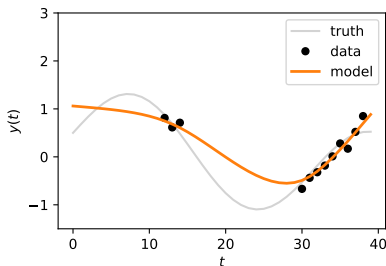
$$\boldsymbol{\mu}' = \boldsymbol{\mu}_T + \Sigma_{TS} \Sigma_{SS}^{-1} (\mathbf{z}_S - \boldsymbol{\mu}_S)$$

$$\Sigma' = \Sigma_{TT} - \Sigma_{TS} \Sigma_{SS}^{-1} \Sigma_{ST}$$

Sequences $\mathbf{z}_T | \mathbf{z}_S$ more closely follow the data, and remain heterogeneous outside the data.



Tikhonov-Regularized Model vs. Gaussian Processes



- ▶ Tikhonov-regularized models require the user to provide a matrix Γ specifying how local variations are penalized.
- ▶ Gaussian processes require the user to provide a matrix Σ (a model of covariance), and a model of mean μ .
- ▶ Gaussian processes provide to the user not only a prediction for missing observations but also a model of predictive uncertainty (shown here as error bars, and where error bars show the standard deviation $\sqrt{\Sigma_{ii}^T}$).

Part 2

Modeling Clustered Data

Recap: Motivations

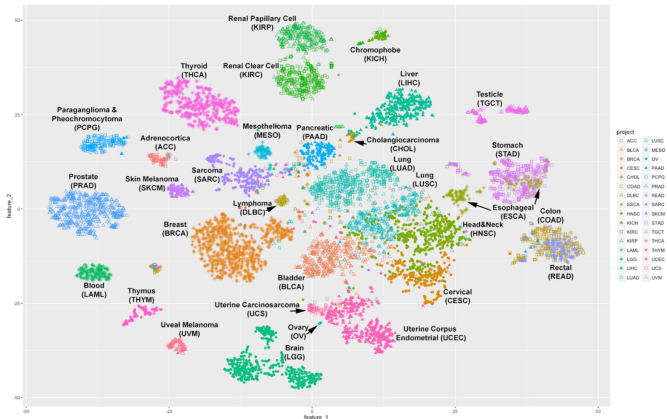


Image source: <https://doi.org/10.1038/s41467-021-21254-9>

- ▶ Methylation profiles of cancer cells organize into clusters, which often correlate with oncology categories.
- ▶ This cluster analysis is useful to verify that existing taxonomies are supported by the data.

Recap: K-means Clustering

Notation:

- ▶ Let $(\mathbf{x}_i)_{i=1}^N$ be our dataset of N instances.
- ▶ Let $(\boldsymbol{\mu}_k)_{k=1}^C$ be the C cluster centroids.
- ▶ Let $c : \{1, \dots, N\} \rightarrow \{1, \dots, K\}$ assign instances to clusters.

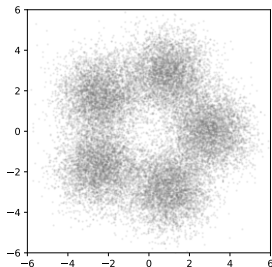
Objective:

- ▶ K-means clustering seeks to find cluster centroids and an assignment onto clusters that minimize the objective

$$J(\boldsymbol{\mu}, c) = \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - \boldsymbol{\mu}_{c(i)}\|^2$$

- ▶ Each data point is assigned to one cluster as a result.

Limits of Clustering



Observations:

- ▶ Instances truly organize into distinct groups but these groups do not form well-separable clusters in \mathbb{R}^d (here \mathbb{R}^2).
- ▶ If one would train a clustering algorithm, it would either find a single cluster, or produce a solution that do not have the desired separability property (cf. e.g. Dunn's index).

Idea:

- ▶ Adopt a probabilistic approach to account for group uncertainty.

Gaussian Mixture Model

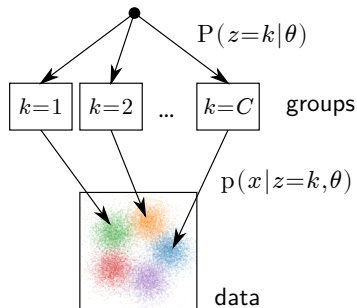
Data is generated via a two-step process:

step 1:

$$P(z|\theta) = \begin{cases} \alpha_1 & \text{if } z = 1 \\ \alpha_2 & \text{if } z = 2 \\ \vdots & \\ \alpha_C & \text{if } z = C \end{cases}$$

step 2:

$$p(\mathbf{x} | z = k, \theta) \sim \mathcal{N}(\boldsymbol{\mu}_k, \Sigma_k)$$



where $\theta = (\alpha_k, \boldsymbol{\mu}_k, \Sigma_k)_{k=1}^C$ are the parameters of the model.

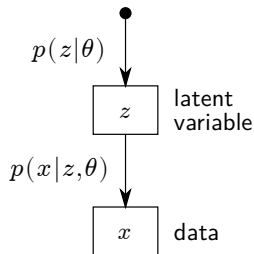
Gaussian Mixture Model

The marginal distribution $p(\mathbf{x}|\theta)$ of the Gaussian mixture model can be written using the law of total probabilities as:

$$p(\mathbf{x}|\theta) = \sum_{k=1}^C \underbrace{P(z=k|\theta)}_{\alpha_k} \cdot \underbrace{p(\mathbf{x}|z=k,\theta)}_{\sim \mathcal{N}(\boldsymbol{\mu}_k, \Sigma_k)}$$

More generally, the marginal distribution of a model with some discrete latent variable $z \in \mathcal{Z}$ can be written as:

$$p(\mathbf{x}|\theta) = \sum_{z \in \mathcal{Z}} p(z|\theta) \cdot p(\mathbf{x}|z,\theta)$$



Learning a Mixture Model

Assuming a dataset \mathcal{D} , and considering examples in the dataset to be iid, the log-likelihood of the data according to the model can be written as:

$$\begin{aligned}\log P(\mathcal{D}|\theta) &= \log \prod_{\mathbf{x} \in \mathcal{D}} p(\mathbf{x}|\theta) \\ &= \sum_{\mathbf{x} \in \mathcal{D}} \log p(\mathbf{x}|\theta)\end{aligned}$$

and we wish to maximize that quantity. Injecting the latent variable model into the objective, we get:

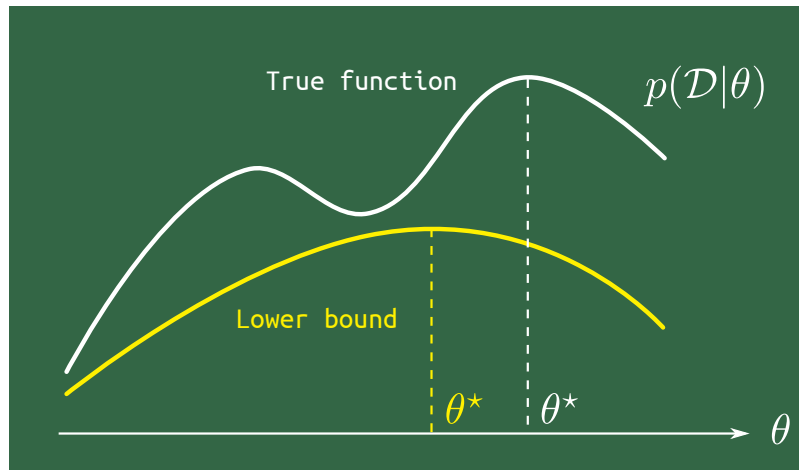
$$= \sum_{\mathbf{x} \in \mathcal{D}} \log \sum_{z \in \mathcal{Z}} p(z|\theta) \cdot p(\mathbf{x}|z, \theta)$$

Problem:

- ▶ The maximum of $\log p(\mathbf{x}|\theta)$ cannot be found analytically.

Learning a Latent Variable Model

Strategy: If the function $p(\mathcal{D}|\theta)$ cannot be easily optimized, find a lower-bound of that function that is easier to optimize.



Building a Lower-Bound

Jensen's inequality

For any element λ of the d -dimensional simplex (i.e. $\lambda \succeq \mathbf{0}$ and $\lambda^\top \mathbf{1} = 1$, and any concave function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, we have

$$f\left(\sum_{i=1}^d \lambda_i a_i\right) \geq \sum_{i=1}^d \lambda_i f(a_i)$$

Application to the latent variable model:

- ▶ Let λ be some probability distribution $(q(z|\mathbf{x}))_{z \in \mathcal{Z}}$ independent from θ .
- ▶ Because our objective

$$\sum_{\mathbf{x} \in \mathcal{D}} \log \sum_{z \in \mathcal{Z}} p(z|\theta) \cdot p(\mathbf{x}|z, \theta)$$

does not contain such terms $q(z|\mathbf{x})$ independent from θ , create them!

Building a Lower-Bound

Step 1: Add $q(z|\mathbf{x})$ both on the numerator and denominator:

$$\log p(\mathcal{D}|\theta) = \sum_{\mathbf{x} \in \mathcal{D}} \log \sum_z \frac{p(z|\theta) \cdot p(\mathbf{x}|z, \theta)}{q(z|\mathbf{x})} \cdot q(z|\mathbf{x})$$

Step 2: Applying the Jensen inequality

$$\geq \sum_{\mathbf{x} \in \mathcal{D}} \sum_z \log \left(\frac{p(z|\theta) \cdot p(\mathbf{x}|z, \theta)}{q(z|\mathbf{x})} \right) \cdot q(z|\mathbf{x})$$

... and verify the bound:

$$\begin{aligned} &= \sum_{\mathbf{x} \in \mathcal{D}} \sum_z \log \left(\frac{p(\mathbf{x}|\theta) \cdot p(z|\mathbf{x}, \theta)}{q(z|\mathbf{x})} \right) \cdot q(z|\mathbf{x}) \\ &= \underbrace{\sum_{\mathbf{x} \in \mathcal{D}} \sum_z \log \left(p(\mathbf{x}|\theta) \right) \cdot q(z|\mathbf{x})}_{\log P(\mathcal{D}|\theta)} - \underbrace{\sum_{\mathbf{x} \in \mathcal{D}} \sum_z \log \left(\frac{q(z|\mathbf{x})}{p(z|\mathbf{x}, \theta)} \right) \cdot q(z|\mathbf{x})}_{\text{KL}(q(z|\mathbf{x}) \parallel p(z|\mathbf{x}, \theta)) \geq 0} \end{aligned}$$

Building a Lower-Bound

Question: How to ensure that

$$\text{KL}(q(z|\mathbf{x}) \parallel p(z|\mathbf{x}, \theta))$$

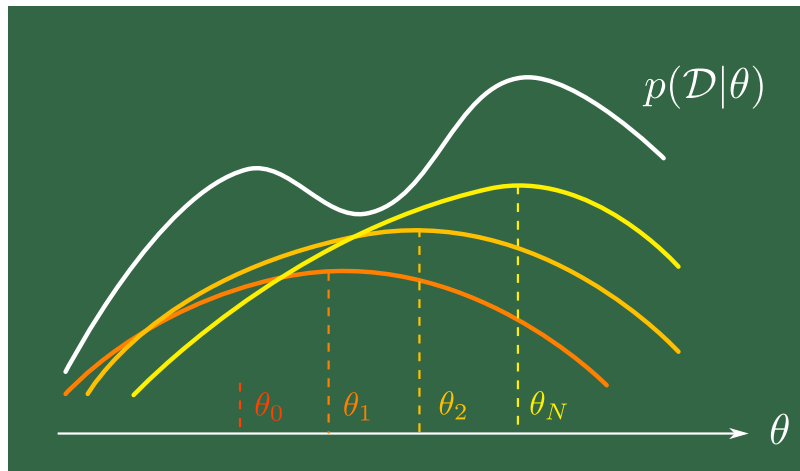
is small, so that maximizing the lower-bound with θ gives a good approximation of the true log-likelihood?

Chicken & egg problem:

- ▶ If we use a simple q , e.g. uniformly distributed, we get a loose lower-bound from which we get a bad θ .
- ▶ To get a tight lower-bound, we need to choose $q(z|\mathbf{x}) \approx p(z|\mathbf{x}, \theta)$ which requires that we know θ .

The Expectation-Maximization (EM) Algorithm

Strategy: Start with some random solution θ and alternately estimate q and θ , until we converge.



The Expectation-Maximization (EM) Algorithm

$$\theta_0 \leftarrow \text{random}()$$

$$q_1(z|\mathbf{x}) \leftarrow p(z|\mathbf{x}, \theta_0) \quad (\text{E-step})$$

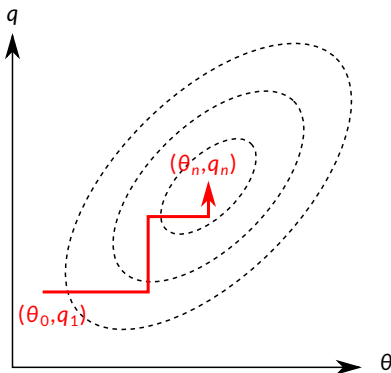
$$\theta_1 \leftarrow \arg \max_{\theta} \sum_{\mathbf{x} \in \mathcal{D}} \sum_z \log \left(\frac{p(\mathbf{x}, z|\theta)}{q_1(z|\mathbf{x})} \right) \cdot q_1(z|\mathbf{x}) \quad (\text{M-step})$$

$$q_2(z|\mathbf{x}) \leftarrow p(z|\mathbf{x}, \theta_1) \quad (\text{E-step})$$

$$\theta_2 \leftarrow \arg \max_{\theta} \sum_{\mathbf{x} \in \mathcal{D}} \sum_z \log \left(\frac{p(\mathbf{x}, z|\theta)}{q_2(z|\mathbf{x})} \right) \cdot q_2(z|\mathbf{x}) \quad (\text{M-step})$$

$$\vdots$$

The Expectation-Maximization (EM) Algorithm



Properties of the algorithm:

- ▶ Block coordinate descent
- ▶ Locally optimal step size
- ▶ The algorithm lands in a local minimum of the function $p(\mathcal{D}|\theta)$.

Advantages of EM compared to gradient descent:

- ▶ no learning rate
- ▶ no need to compute the gradients.

The Gaussian Mixture Model (GMM)

The GMM is formally defined by the two equations:

$$p(z = k | \theta) = \alpha_k$$
$$p(\mathbf{x} | z = k, \theta) = \frac{1}{\sqrt{(2\pi)^d |\Sigma_k|}} \cdot \exp \left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^\top \Sigma_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right)$$

The **parameters** θ of the model to be learned are:

- ▶ The mixing coefficients $(\alpha_k)_{k=1}^C$ subject to $\alpha_k \geq 0$ and $\sum_{k=1}^C \alpha_k = 1$
- ▶ The mean vectors $(\boldsymbol{\mu}_k)_{k=1}^C$.
- ▶ The covariances $(\Sigma_k)_{k=1}^C$ subject to positive semi-definiteness.

Various **simplifications** of the GMM model can be used in practice, e.g. for speed, statistical robustness, or simplicity of implementation.

- ▶ diagonal/isotropic/tied/fixed covariance matrices,
- ▶ fixed mixing coefficients, etc.

EM for the GMM (simplified)

Consider the simplified GMM:

$$p(z = k|\theta) = \frac{1}{K}$$
$$p(\mathbf{x}|z = k, \theta) = \frac{1}{(\pi/\gamma)^{d/2}} \exp(-\gamma \|\mathbf{x} - \boldsymbol{\mu}_k\|^2)$$

E-step: (Apply Bayes rule)

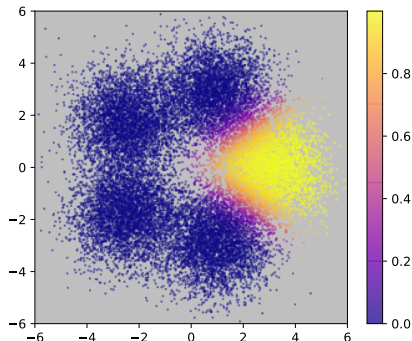
$$q(z = k|\mathbf{x}) = \frac{p(z = k|\theta) \cdot p(\mathbf{x}|z = k, \theta)}{p(\mathbf{x}|\theta)} = \frac{\exp(-\gamma \|\mathbf{x} - \boldsymbol{\mu}_k\|^2)}{\sum_{k=1}^C \exp(-\gamma \|\mathbf{x} - \boldsymbol{\mu}_k\|^2)}$$

M-step: (Set lower-bound gradient to zero)

$$\frac{\partial}{\partial \theta} \sum_{\mathbf{x} \in \mathcal{D}} \sum_{k=1}^C \log(\exp(-\gamma \|\mathbf{x} - \boldsymbol{\mu}_k\|^2)) \cdot q(z = k|\mathbf{x}) = 0$$
$$\Rightarrow \boldsymbol{\mu}_k = \frac{\sum_{\mathbf{x} \in \mathcal{D}} \mathbf{x} \cdot q(z = k|\mathbf{x})}{\sum_{\mathbf{x} \in \mathcal{D}} q(z = k|\mathbf{x})}$$

Gaussian Mixture Model on Toy Data

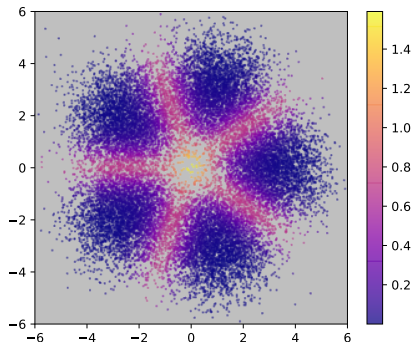
Probability of membership to a given group k , i.e. $q(z = k|\mathbf{x})$, as predicted by a GMM with 5 mixture components:



- Unlike K-means, the GMM does not exhibit an abrupt and unnatural transition between predicted group members and non-group members.

Gaussian Mixture Model on Toy Data

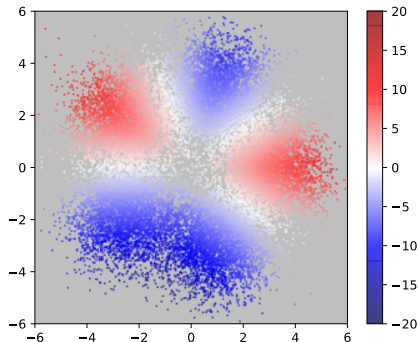
Overall uncertainty measured by the entropy of the vector of probabilities of group memberships, i.e. $-\sum_k q(z = k|\mathbf{x}) \log(q(z = k|\mathbf{x}))$:



- This analysis can be useful to identify prototypical points for each group, or conversely, unresolved instances (that can be of scientific interest).

Gaussian Mixture Model on Toy Data

Log-probability ratio between two sets of groups $\mathcal{K}, \mathcal{K}'$ ($\log \sum_{k \in \mathcal{K}} q(z = k | \mathbf{x}) - \log \sum_{k \in \mathcal{K}'} q(z = k | \mathbf{x})$) as predicted by a GMM with 5 mixture components:



- Interesting nonlinear discriminants functions can be generated from the GMM model.

Summary

Summary

- ▶ There exists a wide variety of probability models for various tasks (regression, discriminant analysis, interpolation, recovering latent structure, etc.)
- ▶ Unlike their 'classical' counterparts, probabilistic models often provide additional functionality (e.g. predictive uncertainty, nonlinearity). This comes at the cost of some extra complexity.
- ▶ Some probabilistic models admit closed form solutions whereas other models (such as the Gaussian Mixture Model) require other forms of optimization (e.g. expectation-maximization).