# Reasoning (II)

Mingsheng Long

Tsinghua University
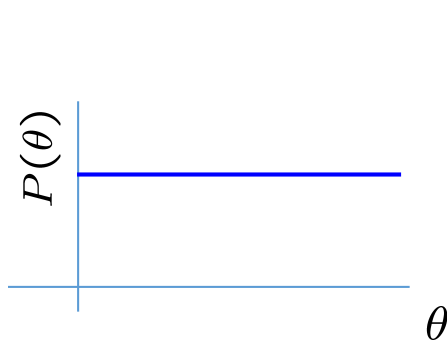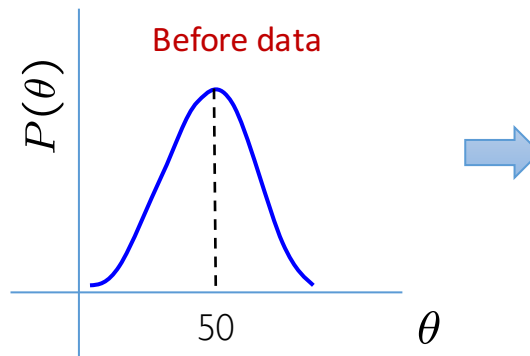
# Outline

- **Probability Basics**

- Discriminative Models

- Generative Models

  - Naïve Bayes Classifier

  - Gaussian Discriminant Analysis

- Mixture Models and EM

  - Gaussian Mixture Model

  - Expectation Maximization
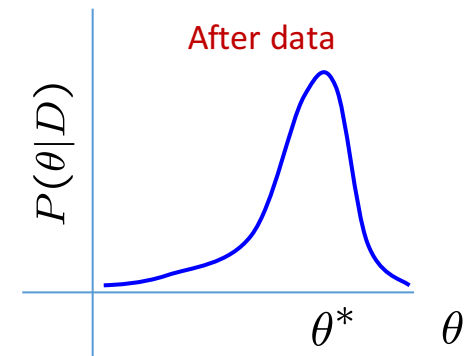
# Bayesian Approach

- Bayesian approaches try to reflect our belief about parameter $\theta$.
  - In this case, we will consider $\theta$ to be a random variable.

- Use the prior information and decide a prior distribution of $\theta$.

- Given the data, estimate a posterior distribution over possible values of $\theta$ with Bayes rule.



Uninformative priors:
uniform distribution

We believe that
$\theta^*$ is around 50

Posterior
distribution

# Bayes Rule

- Bayes rule:

Observed Data      Parameter
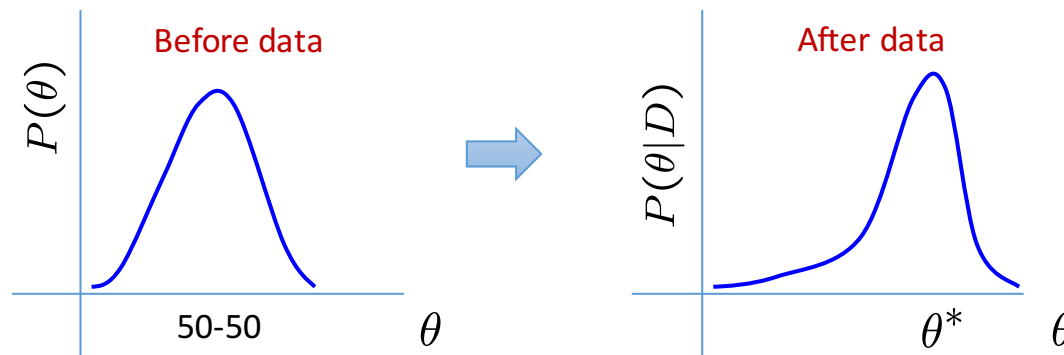
$$p(\theta|D) = \frac{p(D|\theta)p(\theta)}{p(D)}$$

- Equivalently:

$$P(\theta \mid \mathcal{D}) = \frac{P(\mathcal{D} \mid \theta)P(\theta)}{P(\mathcal{D})}$$

$$p(\theta|D) \propto p(D|\theta)p(\theta)$$

Posterior      Likelihood      Prior

$$P(\theta \mid \mathcal{D}) \propto P(\mathcal{D} \mid \theta)P(\theta)$$

Before data

$P(\theta)$

50-50    $\theta$

After data

$P(\theta|D)$

$\theta^*$    $\theta$

Thomas Bayes

Bayes, Thomas. 1763. An essay towards solving a problem in the doctrine of chances. *Philosophical Transactions of the Royal Society of London*, **53:370-418**

# Bernoulli Distribution

- Bernoulli probability mass function (PMF):

$$P(x = 1|q) = q$$
$$P(x = 0|q) = 1 - q$$

- Mean:

$$\mathbb{E}[x] = q$$

- Multinoulli probability mass function (PMF):

$$P(y = l|\boldsymbol{\phi}) = \phi_l$$
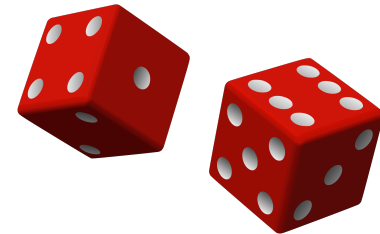
$$\sum_{l=1}^{C} \phi_l = 1$$

- Parameters:

$$\{\boldsymbol{\phi}\}$$

Flipping coin

| | | |
|---|---|---|
| **Sample Space** | $\Omega$ | $\{1,2,3,4,5,6\}$ |
| **Outcome** | $\omega \in \Omega$ | Example: 3 |
| **Event** | $E \subseteq \Omega$ Head | $q$ |
| **Probability** | $P(E)$ Tail | $1 - q$ |

Flipping dice

$$\sum_{l=1}^{6} \phi_l = 1$$

# Gaussian Distribution

- Probability density function (PDF):

$$p(\boldsymbol{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{|2\pi\boldsymbol{\Sigma}|}} \exp(-\frac{1}{2}(\boldsymbol{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{x}-\boldsymbol{\mu}))$$
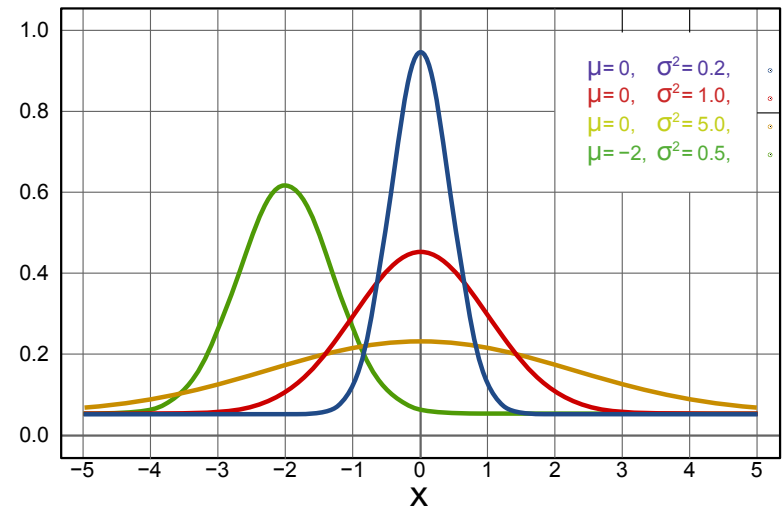
- Mean:

$$\mathbb{E}[\boldsymbol{x}] = \boldsymbol{\mu}$$

- Variance:

$$\mathrm{Var}[\boldsymbol{x}] = \boldsymbol{\Sigma}$$

- Parameters:

$$\{\boldsymbol{\mu}, \boldsymbol{\Sigma}\}$$



μ=0,   σ²= 0.2,
μ=0,   σ²= 1.0,
μ=0,   σ²= 5.0,
μ=−2, σ²= 0.5,

# Likelihood of Parametric Model

- Suppose we have a parametric model $\{p(z;\theta) | \theta \in \Theta\}$ and a sample dataset $\mathcal{D} = (z_1, \dots, z_N)$.

- The likelihood of estimated parameter $\hat{\theta} \in \Theta$ for sample $\mathcal{D}$ is

$$p(\mathcal{D}; \hat{\theta}) = \prod_{n=1}^{N} p(z_n; \hat{\theta})$$

- Due to numerical instability, we prefer to work with the log-likelihood

$$\log p(\mathcal{D}; \hat{\theta}) = \sum_{n=1}^{N} \log p(z_n; \hat{\theta})$$

# Maximum Likelihood Estimation

- Suppose $\mathcal{D} = (z_1, \ldots, z_N)$ is an <u>i.i.d.</u> sample from some distribution.

- Finding the <span style="color:blue">maximum likelihood estimator</span> (MLE) for parameter $\boldsymbol{\theta}$ in the <span style="color:red">parametric model $\{p(y; \theta) | \theta \in \Theta\}$</span> is an optimization problem:

$$\hat{\theta} \in \underset{\theta \in \Theta}{\operatorname{argmax}} \log p(\mathcal{D}; \theta)$$

$$= \underset{\theta \in \Theta}{\operatorname{argmax}} \sum_{n=1}^{N} \log p(z_n; \theta)$$

- **Note:** MLE of a parametric model leads to a particular loss function.

# MLE for Gaussian Distribution

- Recall that the density of Gaussian distribution is $\boldsymbol{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

$$p(\boldsymbol{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{|2\pi\boldsymbol{\Sigma}|}} \exp(-\frac{1}{2}(\boldsymbol{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{x} - \boldsymbol{\mu}))$$

- The log-density is

$$\log p(\boldsymbol{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = -\frac{1}{2}\log|2\pi\boldsymbol{\Sigma}| - \frac{1}{2}(\boldsymbol{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{x} - \boldsymbol{\mu})$$

- To estimate $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ from an i.i.d. sample $\boldsymbol{x}_1, \dots, \boldsymbol{x}_n \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, we will maximize the log joint density

$$\sum_{i=1}^{n} \log p(\boldsymbol{x}_i|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = -\frac{n}{2}\log|2\pi\boldsymbol{\Sigma}| - \frac{1}{2}\sum_{i=1}^{n}(\boldsymbol{x}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{x}_i - \boldsymbol{\mu})$$

# MLE for Gaussian Distribution

- To estimate $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ from an <u>i.i.d.</u> sample $\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, we will maximize the log joint density

$$\sum_{i=1}^{n} \log p(\boldsymbol{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = -\frac{n}{2}\log|2\pi\boldsymbol{\Sigma}| - \frac{1}{2}\sum_{i=1}^{n}(\boldsymbol{x}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{x}_i - \boldsymbol{\mu})$$

- A solid exercise in vector and matrix differentiation. Find $\widehat{\boldsymbol{\mu}}$ and $\widehat{\boldsymbol{\Sigma}}$ by

$$\nabla_{\boldsymbol{\mu}} J(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = 0 \qquad \nabla_{\boldsymbol{\Sigma}} J(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = 0$$

- We get a closed-form solution:

Check: $\widehat{\boldsymbol{\Sigma}}_{\mathrm{MLE}} = \frac{n-1}{n}\boldsymbol{\Sigma}$

$$\widehat{\boldsymbol{\mu}}_{\mathrm{MLE}} = \frac{1}{n}\sum_{i=1}^{n} \boldsymbol{x}_i \qquad \widehat{\boldsymbol{\Sigma}}_{\mathrm{MLE}} = \frac{1}{n}\sum_{i=1}^{n}(\boldsymbol{x}_i - \widehat{\boldsymbol{\mu}}_{\mathrm{MLE}})(\boldsymbol{x}_i - \widehat{\boldsymbol{\mu}}_{\mathrm{MLE}})^T$$

# Bayes Decision Rule

- Assumption:

  - The learning task $p(X, Y) = p(Y|X)p(X)$ can be sampled from.

- Question:

  - Given instance $\boldsymbol{x}$, how should it be classified to minimize error?

- Bayes Decision Rule:

$$h(\boldsymbol{x}) = \operatorname*{argmax}_{y \in \mathcal{Y}}[p(Y = y | X = \boldsymbol{x})]$$

- How to directly measure $p(Y|X)$ with a parametric model $q(Y|X, \boldsymbol{\theta})$?

  - Discriminative models!

# Outline

- Probability Basics

- **Discriminative Models**

- Generative Models
  - Naïve Bayes Classifier
  - Gaussian Discriminant Analysis

- Mixture Models and EM
  - Gaussian Mixture Model
  - Expectation Maximization

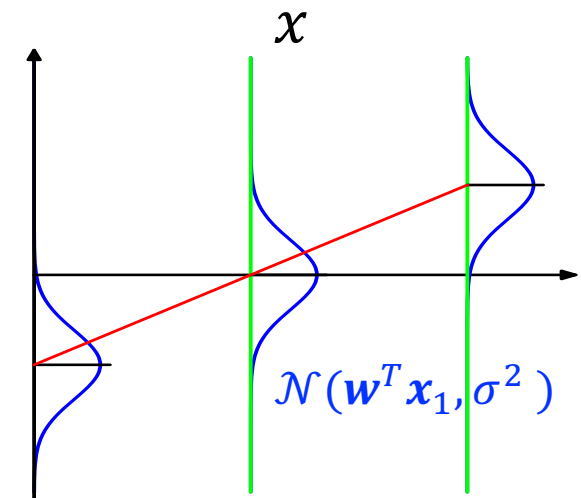# Linear Regression

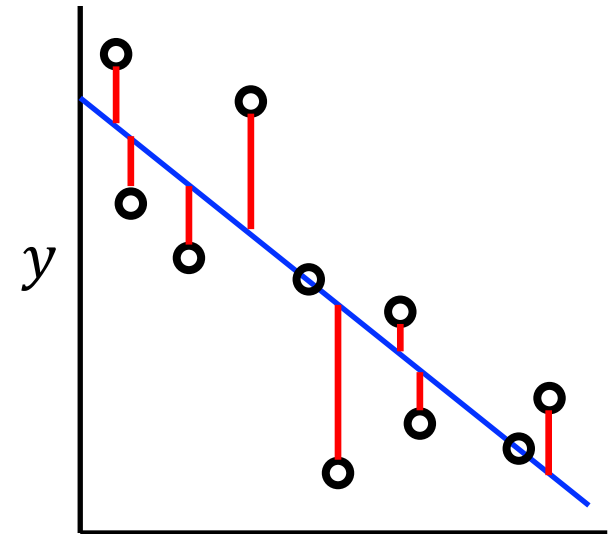- In regression problem, we assume that:

$$y \sim \mathcal{N}(\boldsymbol{w}^T\boldsymbol{x}, \sigma^2)$$

  when $y$ is independent with each other.

- The linear regression model should give expectation of $y$:

$$\mathbb{E}(y|\boldsymbol{x}, \boldsymbol{w}, \sigma^2) = \boldsymbol{w}^T\boldsymbol{x}$$

- Find the best parameter $\boldsymbol{w}$ using i.i.d. sample $\mathcal{D} = \{(x_1, y_1), \dots, (x_N, y_N)\}$.

- $\sigma$ is useless in the final regression model.

$\mathcal{N}(\boldsymbol{w}^T\boldsymbol{x}_1, \sigma^2)$

# Gaussian Linear Regression

- If we assume that $y_n | \boldsymbol{w}, \boldsymbol{x}_n \sim \mathcal{N}(\boldsymbol{w}^T \boldsymbol{x}_n, \sigma^2)$

- Then for point $(\boldsymbol{x}_n, y_n)$

$$p(y_n | \boldsymbol{w}, \boldsymbol{x}_n) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{ -\frac{1}{2\sigma^2}(y_n - \boldsymbol{w}^T \boldsymbol{x}_n)^2 \right\}$$

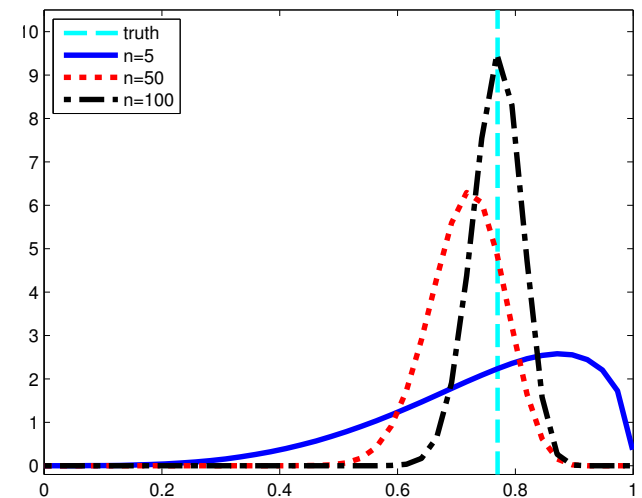- The log-likelihood for linear regression on the whole dataset $\mathcal{D}_n$:

$$\log p(\mathcal{D}_n; \boldsymbol{w}) = \sum_{n=1}^{N} \log p(y_n | \boldsymbol{w}, \boldsymbol{x}_n)$$

$$= \frac{N}{2} \log \frac{1}{2\pi\sigma^2} - \frac{1}{2\sigma^2} \sum_{n=1}^{N} (y_n - \boldsymbol{w}^T \boldsymbol{x}_n)^2$$

# Maximum A Posteriori Estimation

$$\log p(\theta \mid x) = \arg\max_{\theta} \log p(x \mid \theta) + \log p(\theta)$$

$$\hat{\theta} = \arg\max_{\theta} \log p(\theta \mid D) = \arg\max_{\theta} \{\log p(D \mid \theta) + \log p(\theta)\}$$

Posterior   Likelihood   Prior

- MAP: Maximum a posteriori estimation of parameters $\theta$

- We can view MLE as MAP

  - with a uniform prior distribution.

$$x_n$$
$$N$$

- As amount of data becomes large, posterior variance becomes small, and MAP behaves like other estimators such as MLE.

# Bayesian Linear Regression

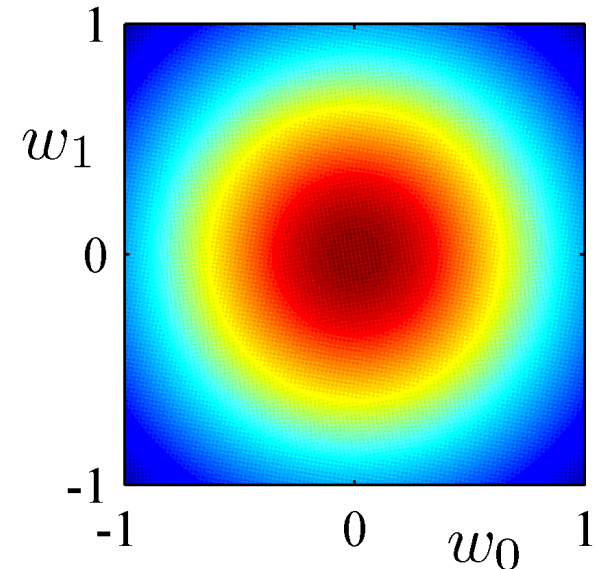- Recall the Gaussian noise assumption for linear regression:

$$y_n = \boldsymbol{w}^T \boldsymbol{x}_n + \epsilon_n, \qquad \epsilon_n \sim \mathcal{N}(0, \sigma^2)$$

- Then we can get:

$$y_n | \boldsymbol{w}, \boldsymbol{x}_n \sim \mathcal{N}(\boldsymbol{w}^T \boldsymbol{x}_n, \sigma^2)$$

- A common choice for the prior:

$$p(\boldsymbol{w}) = \mathcal{N}(\boldsymbol{w} | 0, \alpha^{-1} \boldsymbol{I})$$

# MAP and Regularization

- MAP estimation of $\boldsymbol{w}$:

$$\hat{\boldsymbol{w}} = \arg\min_{\boldsymbol{w}} -\log p(\boldsymbol{w}|D_n) = \arg\min_{\boldsymbol{w}} \{-\log p(\boldsymbol{y}|\boldsymbol{w}, \boldsymbol{X}) - \log p(\boldsymbol{w})\}$$

$$-\log p(\boldsymbol{y}|\boldsymbol{w}) = -\sum_{i=1}^{n} \log p(y_i|\boldsymbol{w}, \boldsymbol{x}_i) = -\frac{n}{2}\log\frac{1}{2\pi\sigma^2} + \frac{1}{2\sigma^2}\sum_{i=1}^{n}(y_i - \boldsymbol{w}^T\boldsymbol{x}_i)^2$$

$$-\log p(\boldsymbol{w}) = -\frac{d}{2}\log\frac{1}{2\pi\sigma^2} + \frac{\alpha}{2}\sum_{j=1}^{d} w_j^2$$

$$\Rightarrow \hat{\boldsymbol{w}} = \arg\min_{\boldsymbol{w}} \frac{1}{2\sigma^2}\sum_{i=1}^{n}(y_i - \boldsymbol{w}^T\boldsymbol{x}_i)^2 + \frac{\alpha}{2}\sum_{j=1}^{d} w_j^2$$

$$\Rightarrow \hat{\boldsymbol{w}} = \arg\min_{\boldsymbol{w}} \|\boldsymbol{X}\boldsymbol{w} - y\|_2^2 + \frac{\alpha}{\beta}\|\boldsymbol{w}\|_2^2 \qquad \text{denote } \beta = \frac{1}{\sigma^2}$$

# Bayesian Model Averaging

- The posterior predictive distribution ($\tilde{y}$ is the prediction):

$$p(\tilde{y}|y) = \int p(\tilde{y}|\theta)p(\theta|y)d\theta$$

- Assume there is a true model $p(y|\theta)$

- Account for the uncertainty in $\theta$.

- To account for model uncertainty among some models $M_1, \ldots, M_h$, we use **Bayesian model averaged (BMA)** posterior predictive distribution

$$p(\tilde{y}|y) = \sum_{h=1}^{H} p(\tilde{y}|M_h, y)p(M_h|y)$$

predictive distribution under model $M_h$     posterior model probability

# Outline

- Probability Basics

- Discriminative Models

- **Generative Models**

  - **Naïve Bayes Classifier**

  - Gaussian Discriminant Analysis

- Mixture Models and EM

  - Gaussian Mixture Model

  - Expectation Maximization

# Bayes Rule

- Alternative idea:

  - It is possible to switch conditioning according to Bayes rule.

  - Given any two random variables $X$ and $Y$, it holds that:

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

- We try to model $p(X = x | Y = y)$ and $p(Y = y)$ in this problem.

  - Generative models!

- We also write $p(Y = y | X = x) \propto p(X = x | Y = y) p(Y = y)$

  - $\propto$ is commonly used to avoid non-necessary normalization term.

# Naïve Bayes Classifier

- Model distribution of high-dimensional data $p(X = \boldsymbol{x}|Y = y)$ is hard:

  - Because we need to find a proper description of data distribution.

  - Especially the dependency between dimensions.

- Naïve Bayes Classifier (NBC) assumes conditional-independence:

  - Each dimension is independent given label $y$:

$$p(X = \boldsymbol{x}|Y = y) = \prod_{j=1}^{d} p(x_{\cdot j}|y)$$

Naïve!

$j$th dimension of $\boldsymbol{x}$

  - Thus $p(y)$ and $p(x_i|y)$ can be computed directly from dataset.

# Naïve Bayes Classifier

- Naïve Bayes Classifier is used when the features are all discrete.

- For binary feature, model $p(x._j|y)$ and $p(y)$ as **Bernoulli distribution**:

$x._j \in \{0,1\}$

$$p(x._j = 1|y = +1) \sim \text{Bernoulli}(\phi_j^+),$$
$$p(x._j = 1|y = -1) \sim \text{Bernoulli}(\phi_j^-),$$
$$p(y = +1) \sim \text{Bernoulli}(\phi)$$

- So we can estimate pareameters $\phi_j^+$ and $\phi$ as ($\phi_j^-$ is similar):

$$\phi_j^+ = \frac{\sum_{i=1}^n \mathbf{1}\{x_{ij} = 1 \wedge y_i = +1\}}{\sum_{i=1}^n \mathbf{1}\{y_i = +1\}}$$

$$\phi = \frac{\sum_{i=1}^n \mathbf{1}\{y_i = +1\}}{n}$$

# Naïve Bayes Classifier

- For feature with many values, use **multinoulli distribution** instead.

| D1 | Sunny | Hot | High | Weak | No |
| D2 | Sunny | Hot | High | Strong | No |
| D3 | Overcast | Hot | High | Weak | Yes |
| D4 | Rain | Mild | High | Weak | Yes |
| D5 | Rain | Cool | Normal | Weak | Yes |
| D6 | Rain | Cool | Normal | Strong | No |
| D7 | Overcast | Cool | Normal | Strong | Yes |

- What will Naïve Bayes Classifier predict for (Rain, Cool, High, Weak)?

$$p(\text{Rain}|\text{Yes}) = 2/4, p(\text{Cool}|\text{Yes}) = 2/4, p(\text{High}|\text{Yes})$$
$$= 2/4, p(\text{Weak}|\text{Yes}) = 3/4, p(\text{Yes}) = 4/7$$

- $p(\text{Yes}|(\text{Rain}, \text{Cool}, \text{High}, \text{Weak})) \propto 3/56$

- $p(\text{No}|(\text{Rain}, \text{Cool}, \text{High}, \text{Weak})) \propto 2/189$

# Naïve Bayes Classifier

- What if $p(x_{\cdot j} = r_j | Y = +1) = \frac{\sum_{i=1}^{n} 1\{x_{\cdot j} = r_j \wedge y_i = +1\}}{\sum_{i=1}^{n} 1\{y_i = +1\}} = 0$?

  - This will cause $p(\boldsymbol{x}|Y = +1)$ to zero, no matter how large other $p(x_{\cdot l} = r_l | Y = +1)$.

  - We can add prior to solve this problem (Laplacian smoothing):

  $$p(x_{\cdot j} = r_j | Y = +1) = \frac{\sum_{i=1}^{n} 1\{x_{\cdot j} = r_j \wedge y_i = +1\} + 1}{\sum_{i=1}^{n} 1\{y_i = +1\} + k_j}$$

- Continuous variables:

  - We can discretize the variable.

  How many values does this feature have?

  - Or we can use another model based on a different assumption.
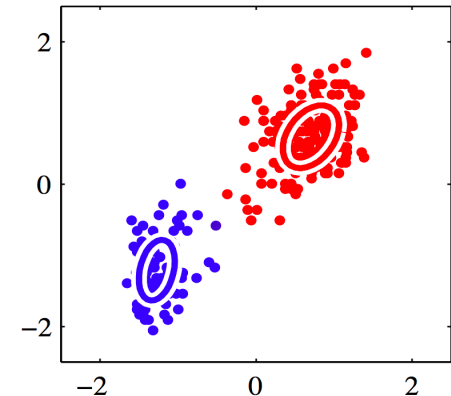
# Outline

- Probability Basics

- Discriminative Models

- **Generative Models**

  - Naïve Bayes Classifier

  - **Gaussian Discriminant Analysis**

- Mixture Models and EM

  - Gaussian Mixture Model

  - Expectation Maximization

# Gaussian Discriminant Analysis

- Alternative methods for dataset with all continuous features:

  - Using parametric distribution to represent $p(X = \boldsymbol{x}|Y = y)$.

- A **common assumption** in classification:

  - We always assume that data

    points in a class is a cluster.

- So we can model $p(X = \boldsymbol{x}|Y = y)$ by **Gaussian distribution**:

$$p(X = \boldsymbol{x}|Y = +1) \propto \exp\left(-\frac{1}{2}(\boldsymbol{x} - \boldsymbol{\mu}_+)^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{x} - \boldsymbol{\mu}_+)\right)$$

$$p(X = \boldsymbol{x}|Y = -1) \propto \exp\left(-\frac{1}{2}(\boldsymbol{x} - \boldsymbol{\mu}_-)^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{x} - \boldsymbol{\mu}_-)\right)$$

Usually share $\boldsymbol{\Sigma}$.

# Gaussian Discriminant Analysis

- We still model $p(Y = y)$ as Bernoulli distribution: Bernoulli($\phi$).

- Now we use MLE to find the best parameter estimation:

$$\ell(\phi, \boldsymbol{\mu}_+, \boldsymbol{\mu}_-, \boldsymbol{\Sigma}) = \log \prod_{i=1}^{n} p(\boldsymbol{x}_i, y_i; \phi, \boldsymbol{\mu}_+, \boldsymbol{\mu}_-, \boldsymbol{\Sigma})$$

$$= \boxed{\log \prod_{i=1}^{n} p(\boldsymbol{x}_i | y_i; \boldsymbol{\mu}_+, \boldsymbol{\mu}_-, \boldsymbol{\Sigma})} + \boxed{\log \prod_{i=1}^{n} p(y_i | \phi)}$$

- The computing process is very similar to the process of Gaussian.

  - The main difference is that $\boldsymbol{\mu}_+, \boldsymbol{\mu}_-$ are different.

$$\phi = \frac{\sum_{i=1}^{n} \mathbf{1}\{y_i = +1\}}{n}, \boldsymbol{\mu}_+ = \frac{\sum_{i=1}^{n} \mathbf{1}\{y_i = +1\} \boldsymbol{x}_i}{\sum_{i=1}^{n} \mathbf{1}\{y_i = +1\}}, \boldsymbol{\mu}_- = \frac{\sum_{i=1}^{n} \mathbf{1}\{y_i = -1\} \boldsymbol{x}_i}{\sum_{i=1}^{n} \mathbf{1}\{y_i = -1\}}$$

$$\boldsymbol{\Sigma} = \frac{1}{n} \sum_{i=1}^{n} (\boldsymbol{x}_i - \boldsymbol{\mu}_{y_i})(\boldsymbol{x}_i - \boldsymbol{\mu}_{y_i})^T$$

# Gaussian Discriminant Analysis

- On a test data $\boldsymbol{x}$, Gaussian Discriminant Analysis (GDA) outputs label:

$$\underset{y \in \{+1, -1\}}{\text{argmax}} \left[ p(\boldsymbol{x}|y)p(y) \right]$$



Shared $\boldsymbol{\Sigma}$ leads to "large margin model"

# Discriminative vs. Generative



- <span style="color:blue">Discriminative models:</span>
    - Concentrate on the prediction of label or certain variables.
    - Usually simpler and more efficient on general data.
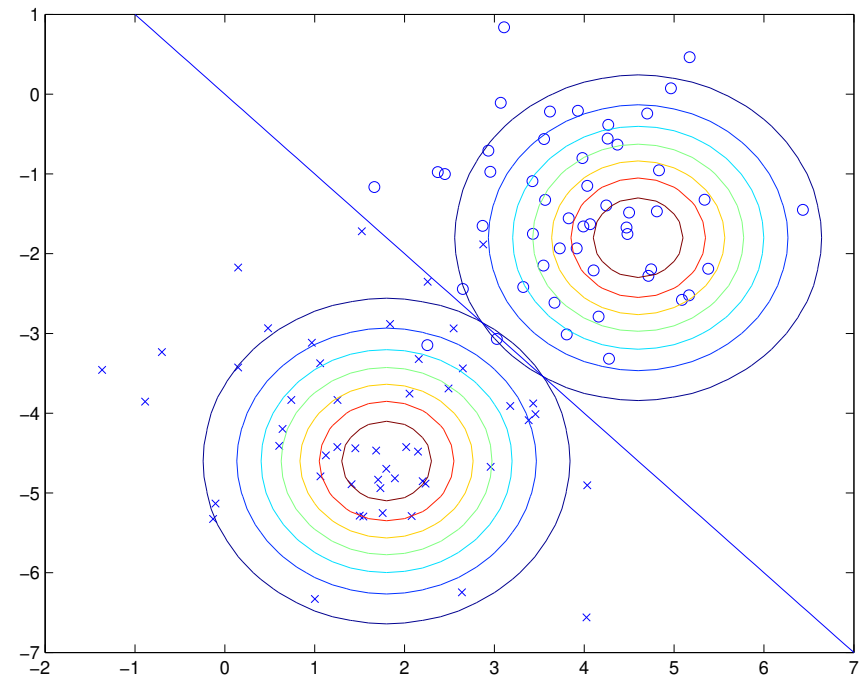
- <span style="color:red">Generative models:</span>
    - Usually stronger assumption, better results on smaller data.
    - Can capture structure of data distribution and generate new data.

# Outline

- Probability Basics

- Discriminative Models

- Generative Models

  - Naïve Bayes Classifier

  - Gaussian Discriminant Analysis

- **Mixture Models and EM**

  - **Gaussian Mixture Model**

  - Expectation Maximization

# Gaussian Mixture Model

- The generating process of the GDA Model:

  - Choose $y \in \{+1, -1\}$ with $p(+1) = p(-1) = \frac{1}{2}$.

  - Choose $\boldsymbol{x}|y \sim \mathcal{N}(\boldsymbol{x}|\boldsymbol{\mu}_y, \boldsymbol{\Sigma})$.

- We can compute $p(\boldsymbol{x})$:

$$p(\boldsymbol{x}) = \frac{1}{2} p(\boldsymbol{x}|\boldsymbol{\mu}_{+1}, \boldsymbol{\Sigma})$$
$$+ \frac{1}{2} p(\boldsymbol{x}|\boldsymbol{\mu}_{-1}, \boldsymbol{\Sigma})$$

- This is not a GDA Model, but a Gaussian Mixture Model (GMM).

# Gaussian Mixture Model

- Parameters of Gaussian Mixture Model (GMM):

  - Cluster probabilities:  $\pi = (\pi_1, \dots, \pi_k)$.

  - Cluster means: $\mu = (\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k)$.

  - Cluster covariance matrices:  $\Sigma = (\boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_k)$.

- Generating process of GMM:

  - First generate cluster index:

    - $z \sim (\pi_1, \dots, \pi_k)$

  - Then generate data:

    - $x \sim \mathcal{N}(\boldsymbol{x}|\boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z)$.

- Density: $p(\boldsymbol{x}) = \sum_{z=1}^{k} \pi_z \mathcal{N}(\boldsymbol{x}|\boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z)$



Mixture of Three Gaussians

$N(\mu_1, \Sigma_1)$

$N(\mu_2, \Sigma_2)$

$N(\mu_3, \Sigma_3)$

# Mixture Distribution

- A probability density $p(x)$ represents a mixture distribution
  - if we can write it as a convex combination of probability densities:

$$p(x) = \sum_{i=1}^{k} w_i p_i(x)$$

  - where $w_i \geq 0, \sum_{i=1}^{k} w_i = 1$, and each $p_i$ is a probability density.

- Gaussian mixture model (GMM): $p(\boldsymbol{x}) = \sum_{z=1}^{k} \pi_z \mathcal{N}(\boldsymbol{x}|\boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z)$.

- More constructively, let $S$ be a set of probability distributions:
  - Choose a distribution randomly from $S$.
  - Sample $\boldsymbol{x}$ from the chosen distribution.
    - Then $\boldsymbol{x}$ has a mixture distribution.

# Generative Models for Clustering

- What do we model in unsupervised learning setting?

  - There is no longer $p(x, y)$. We can sample from $p(x)$ only.

- Consider a **clustering** problem.

  - Suppose there are $k$ clusters.

  - We have a distribution assumption (Gaussian) for each cluster.

- And the whole dataset can be generated as follows:

  - Choose a random cluster $z \in \{1, 2, \dots, k\}$.

  - Choose a point $x$ from the distribution for cluster $z$.

- We can see that GMM is very suitable for modeling $p(x)$.

  - **Difficulty:** We do not know which $x$ is sampled from which $z$!

# Gaussian Mixture Model: Learning

- Can we compute MLE of GMM directly?

- The log-likelihood for $\mathcal{D} = (\boldsymbol{x}_1, \dots, \boldsymbol{x}_n)$ sampled <u>i.i.d.</u> from a GMM is

$$\ell(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \log \prod_{i=1}^{n} \sum_{z=1}^{k} \pi_z \mathcal{N}(\boldsymbol{x}_i | \boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z) = \sum_{i=1}^{n} \log \left[ \sum_{z=1}^{k} \pi_z \mathcal{N}(\boldsymbol{x}_i | \boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z) \right]$$

- Plug the Gaussian density in it:

$$\ell(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{i=1}^{n} \log \left[ \sum_{z=1}^{k} \frac{\pi_z}{\sqrt{|2\pi\boldsymbol{\Sigma}_z|}} \exp(-\frac{1}{2}(\boldsymbol{x} - \boldsymbol{\mu}_z)^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{x} - \boldsymbol{\mu}_z)) \right]$$

- The sum inside the log is intractable:

  - A general challenge in mixture models, need approximate methods.

# Gaussian Mixture Model: Learning

- In GDA, we know which $\boldsymbol{x}$ is sampled from which cluster $\boldsymbol{z}$.

  - So the solution of MLE is easy to find.

  - Without $\boldsymbol{z}$, there will be computational difficulties.

  - $\boldsymbol{z}$ is called latent variable.

- An iterative idea that solves one set of variables by fixing the others:

- Iterate between

  - **Step I:** Known each $(\boldsymbol{x}, \boldsymbol{z})$, find best $(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$.

  - **Step II:** Known $(\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$, find $\boldsymbol{z}$ for each $\boldsymbol{x}$.

- We have a general method with strict theoretical foundation here!

  - Expectation-Maximization (EM)!

# Outline

- Probability Basics

- Discriminative Models

- Generative Models

  - Naïve Bayes Classifier

  - Gaussian Discriminant Analysis

- **Mixture Models and EM**

  - Gaussian Mixture Model

  - **Expectation Maximization**

# Latent Variable Model

- Two (**abstract**) sets of random variables: $z$ and $x$

  - $z$ consists of latent variables.

  - $x$ consists of observed variables.

- Joint probability model parameterized by $\theta \in \Theta$:

$$p(x, z | \theta)$$

- A **latent variable model** is a probabilistic model for which certain variables are never observed.

  - The Gaussian mixture model is a latent variable model.

- An observation of $x$ is called an **incomplete** dataset.

  - An observation of $(x, z)$ is called a **complete** dataset.

Unobserved

Observed

# Objectives

- Learning problem:
  - Given incomplete dataset $\mathcal{D} = (x_1, \ldots x_n)$, find MLE
  $$\hat{\theta} = \underset{\theta}{\text{argmax}}\, p(\mathcal{D}|\theta).$$

- Inference problem:
  - Given $x$, find conditional distribution over latent variable $z$:
  $$p(z|x, \theta)$$

- Expectation-Maximization (EM) for both problems!

- For Gaussian mixture model, learning is hard, inference is easy.

- For more complicated models (next lectures), inference can be hard.

# Expectation-Maximization (EM): Key Idea

- Marginal log-likelihood is **hard** to optimize:

$$\max_{\theta} \log p(x|\theta)$$

Objective!

- Typically, the **complete** data log-likelihood is **easy** to optimize:

$$\max_{\theta} \log p(x, z|\theta)$$

  - What if we had a distribution $q(z)$ for the latent variables $z$?

- Then maximize the expected **complete** data log-likelihood:

$$\max_{\theta} \sum_{z} q(z) \log p(x, z|\theta)$$

  - **Assumption:** EM assumes this maximization is relatively easy.

# Evidence Lower Bound (ELBO)

- Let $q(z)$ be any probability function on $\mathcal{Z}$, the support of $z$:

$$\log p(x|\theta) = \log\left[\sum_z p(x, z|\theta)\right]$$

Objective!

$$= \log\left[\sum_z q(z)\left(\frac{p(x, z|\theta)}{q(z)}\right)\right]$$

$\longrightarrow$ log of an expectation

Jenson's inequality
$$\log(\mathbb{E}[X]) \geq \mathbb{E}[\log(X)]$$

$$\geq \underbrace{\sum_z q(z)\log\left(\frac{p(x, z|\theta)}{q(z)}\right)}_{\mathcal{L}(q, \theta)}$$

$\longrightarrow$ expectation of log

Evidence lower bound (ELBO)
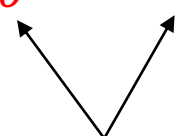
# MLE, EM and ELBO

- For any probability function $q(z)$, we have a lower bound on the marginal log-likelihood

$$\log p(x|\theta) \geq \mathcal{L}(q, \theta).$$

- The MLE is defined as a maximum over $\theta$:

$$\hat{\theta}_{\mathrm{MLE}} = \underset{\theta}{\mathrm{argmax}} \log p(x|\theta)$$

- The EM algorithm maximizes the ELBO over $\theta$ and $q$:

$$\hat{\theta}_{\mathrm{EM}} = \underset{\theta}{\mathrm{argmax}}[\underset{q}{\max} \mathcal{L}(q, \theta)]$$

Lead to an **Iterative** Algorithm!

# EM: Iterative Optimization

- Choose sequence of $q$'s and $\theta$'s by **coordinate ascent.**

- EM Algorithm (high level):
  - Choose initial $\theta^{\mathbf{old}}$
  - Let $q^* = \underset{q}{\mathrm{argmax}}\, \mathcal{L}\big(q, \theta^{\mathbf{old}}\big)$
  - Let $\theta^{\mathbf{new}} = \underset{\theta}{\mathrm{argmax}}\, \mathcal{L}(q^*, \theta)$
  - Go to Step 2, until converged.



- Will show: $p(x|\theta^{\mathbf{new}}) \geq p(x|\theta^{\mathbf{old}})$
  - Get sequence of $\theta$'s with monotonically increasing likelihood.

- What left: What are $\underset{q}{\mathrm{argmax}}\, \mathcal{L}\big(q, \theta^{\mathbf{old}}\big)$ and $\underset{\theta}{\mathrm{argmax}}\, \mathcal{L}(q^*, \theta^{\mathbf{old}})$?

# ELBO via KL Divergence

- Investigate the evidence lower bound:

$$\mathcal{L}(q, \theta) = \sum_z q(z) \log\left(\frac{p(x, z|\theta)}{q(z)}\right)$$

**KL-Divergence**

$$\sum_z q(z) \log\left(\frac{q(z)}{p(z)}\right)$$
$$= \text{KL}[q(z)\|p(z)]$$
Properties:
$$\text{KL}(p\|q) \geq 0,$$
$$\text{KL}(p\|p) = 0.$$

$$= \sum_z q(z) \log\left(\frac{p(z|x, \theta)p(x|\theta)}{q(z)}\right)$$

$$= \sum_z q(z) \log\left(\frac{p(z|x, \theta)}{q(z)}\right) + \sum_z q(z) \log\left(p(x|\theta)\right)$$

$$= -\text{KL}[q(z)\|p(z|x, \theta)] + \log p(x|\theta)$$

- **Amazing!** We get back an equality for the marginal likelihood:

$$\log p(x|\theta) = \mathcal{L}(q, \theta) + \text{KL}[q(z)\|p(z|x, \theta)]$$

# E-Step: Maximizing Over $q$ for Fixed $\theta = \theta^{\text{old}}$

- Find $q$ maximizing

$$\mathcal{L}(q, \theta^{\text{old}}) = -\text{KL}[q(z), p(z|x, \theta^{\text{old}})] + \underbrace{\log p(x|\theta^{\text{old}})}_{\text{no } q \text{ here}}$$

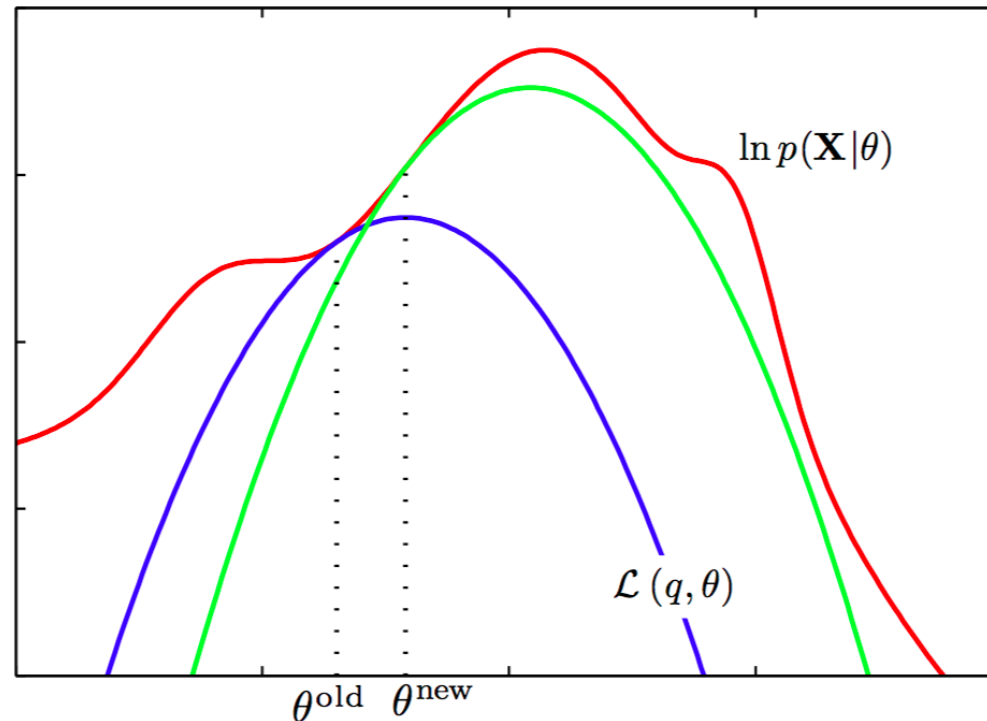  - Recall $\text{KL}(p||q) \geq 0$, and $\text{KL}(p||p) = 0$

- The best $q$ is $\boxed{q^*(z) = p(z|x, \theta^{\text{old}})}$

$$\mathcal{L}(q^*, \theta^{\text{old}}) = \underbrace{-\text{KL}[p(z|x, \theta^{\text{old}}), p(z|x, \theta^{\text{old}})]}_{=0} + \log p(x|\theta^{\text{old}})$$

- Summary:

$$\log p(x|\theta^{\text{old}}) = \mathcal{L}(q^*, \theta^{\text{old}}) \qquad (\text{Tangent at } \theta^{\text{old}})$$

$$\log p(x|\theta) \geq \mathcal{L}(q^*, \theta) \quad \forall \theta$$

# Tight Lower Bound for Any Chosen $\theta$



- For $\theta^{\text{old}}$, take $q(z) = p(z|x, \theta^{\text{old}})$. Then

  $-\log p(x|\theta) \geq \mathcal{L}(q, \theta) \quad \forall \theta.$ [Global lower bound]

  $-\log p(x|\theta^{\text{old}}) = \mathcal{L}(q, \theta^{\text{old}}).$ [Lower bound is tight at $\theta^{\text{old}}$]

# M-Step: Maximizing Over $\theta$ for Fixed $q$

- Consider maximizing the evidence lower bound (EBLO) $\mathcal{L}(q, \theta)$:

$$\mathcal{L}(q, \theta) = \sum_z q(z) \log \left( \frac{p(x, z | \theta)}{q(z)} \right)$$

$$= \underbrace{\sum_z q(z) \log \left( p(x, z | \theta) \right)}_{\mathbb{E}[\text{complete data log-likelihood}]} - \underbrace{\sum_z q(z) \log \left( q(z) \right)}_{\text{no } \theta \text{ here}}$$

- For fixed $q$, maximizing $\mathcal{L}(q, \theta)$ by $\theta$ is equivalent to maximizing

$$\mathbb{E}[\text{complete data log-likelihood}]$$

# Expectation-Maximization (EM): Algorithm

Donald Rubin

- Choose initial $\theta^{\text{old}}$.

- Expectation Step

  - Let $q^*(z) = p(z|x, \theta^{\text{old}})$. [$q^*$ gives best lower bound at $\theta^{\text{old}}$]

  - Let
  $$J(\theta) := \mathcal{L}(q^*, \theta) = \underbrace{\sum_z q^*(z)\log\left(\frac{p(x, z|\theta)}{q^*(z)}\right)}_{\text{Expectation w.r.t. } z \sim q^*(z)}$$

- Maximization Step

$$\theta^{\text{new}} = \underset{\theta}{\text{argmax}}\, J(\theta) \quad \boxed{\text{You can use SGD}}$$

[Equivalent to maximizing the expected complete log-likelihood.]

- Go to the Expectation Step, until converged.

# EM for MAP

- Suppose we have a prior $p(\theta)$.

- Want to find MAP estimate: $\hat{\theta}_{\mathbf{MAP}} = \underset{\theta}{\mathrm{argmax}}\, p(\theta|x)$

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

$$\log p(\theta|x) = \log p(x|\theta) + \log p(\theta) - \log p(x)$$

- Still can use our evidence lower bound on $\log p(x, \theta)$.

$$J(\theta) := \mathcal{L}(q^*, \theta) = \sum_z q^*(z)\log\left(\frac{p(x, z|\theta)}{q^*(z)}\right)$$

- Maximization step becomes $\theta^{\mathbf{new}} = \underset{\theta}{\mathrm{argmax}}\, [J(\theta) + \log p(\theta)]$

# GMM: E-Step

- Denote probability (responsibility) that $\boldsymbol{x_i}$ comes from cluster $j$ by

$$\gamma_i^j = p(z = j | \boldsymbol{x} = \boldsymbol{x_i}, \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$\boxed{q^*(z) = p(z|x, \theta^{\text{old}})}$$

  - The vector $(\gamma_i^1, \ldots, \gamma_i^k)$ is exactly the soft assignment for $\boldsymbol{x_i}$.

- From probabilistic computation:

$$\gamma_i^j = p(z = j | \boldsymbol{x_i}) = \frac{p(z = j, \boldsymbol{x_i})}{p(\boldsymbol{x_i})} = \frac{\pi_j \mathcal{N}(\boldsymbol{x_i} | \boldsymbol{\mu_j}, \boldsymbol{\Sigma_j})}{\sum_{c=1}^{k} \pi_c \mathcal{N}(\boldsymbol{x_i} | \boldsymbol{\mu_c}, \boldsymbol{\Sigma_c})}$$

- If we know $\boldsymbol{\mu_j}, \Sigma_j, \pi_j$ for all clusters $j = 1, \ldots k$, then easy to compute:

$$\gamma_i^j = \frac{\pi_j \mathcal{N}(\boldsymbol{x_i} | \boldsymbol{\mu_j}, \boldsymbol{\Sigma_j})}{\sum_{c=1}^{k} \pi_c \mathcal{N}(\boldsymbol{x_i} | \boldsymbol{\mu_c}, \boldsymbol{\Sigma_c})}$$

# GMM: M-Step

$$\underset{\theta}{\operatorname{argmax}} \mathcal{L}(q^*, \theta) = \underset{\theta}{\operatorname{argmax}} \sum_z q^*(z) \log\left(\frac{p(x, z|\theta)}{q^*(z)}\right)$$

$$\boxed{\theta = \{\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}\}} \qquad\qquad = \underset{\theta}{\operatorname{argmax}} \sum_z {\color{red} p(z|x, \theta^{\text{old}}) \log\left(p(x, z|\theta)\right)}$$

- So we have the loss function for Gaussian Mixture Model parameters:

$$\boxed{\color{red} \text{By MLE}} \qquad \underset{\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}}{\operatorname{argmax}} \sum_{i=1}^{n} \sum_{j=1}^{k} {\color{red} \gamma_i^j \log\left[\pi_j \mathcal{N}\left(\boldsymbol{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j\right)\right]}$$

- Let $n_c = \sum_{i=1}^{n} \gamma_i^c$ be the number of points soft-assigned to cluster $c$.

$$\pi_c^{\text{new}} = \frac{n_c}{n}, \boldsymbol{\mu}_c^{\text{new}} = \frac{1}{n_c} \sum_{i=1}^{n} \gamma_i^c \boldsymbol{x}_i, \boldsymbol{\Sigma}_c^{\text{new}} = \frac{1}{n_c} \sum_{i=1}^{n} \gamma_i^c (\boldsymbol{x}_i - \boldsymbol{\mu}_c^{\text{new}})(\boldsymbol{x}_i - \boldsymbol{\mu}_c^{\text{new}})^T$$

# EM for GMM: Overview

- Initialize parameters $\boldsymbol{\pi}, \boldsymbol{\Sigma}, \boldsymbol{\mu}$.

- **E-step.** Evaluate all responsibilities using current parameters:

$$\gamma_i^j = \frac{\pi_j \mathcal{N}(\boldsymbol{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\sum_{c=1}^{k} \pi_c \mathcal{N}(\boldsymbol{x}_i | \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)}$$

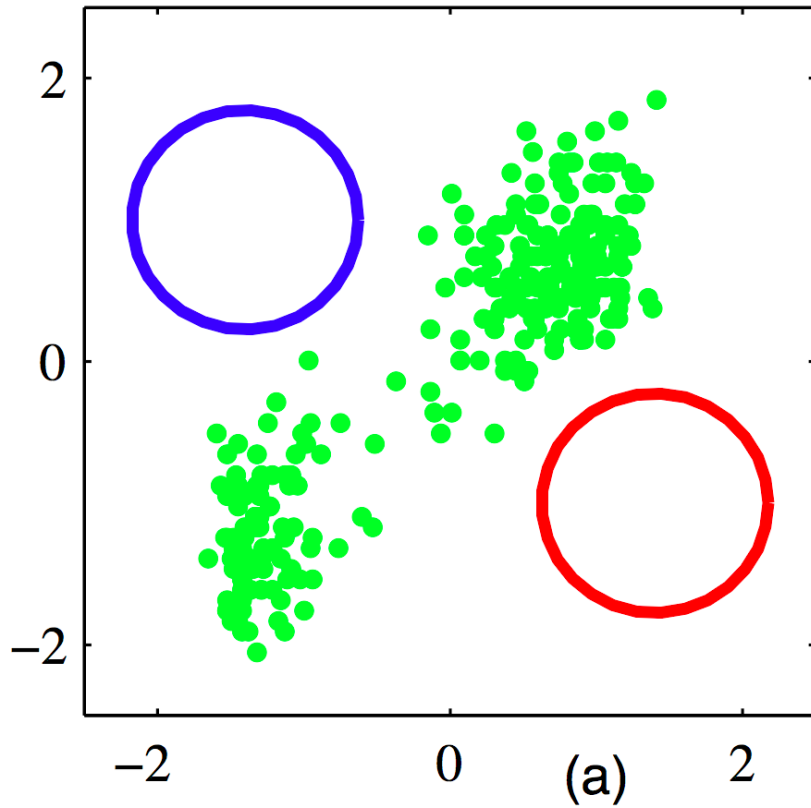- **M-step.** Re-estimate the parameters using the responsibilities:

$$\pi_c^{\text{new}} = \frac{n_c}{n}$$

$$\boldsymbol{\mu}_c^{\text{new}} = \frac{1}{n_c} \sum_{i=1}^{n} \gamma_i^c \boldsymbol{x}_i, \boldsymbol{\Sigma}_c^{\text{new}} = \frac{1}{n_c} \sum_{i=1}^{n} \gamma_i^c (\boldsymbol{x}_i - \boldsymbol{\mu}_c^{\text{new}})(\boldsymbol{x}_i - \boldsymbol{\mu}_c^{\text{new}})^T$$

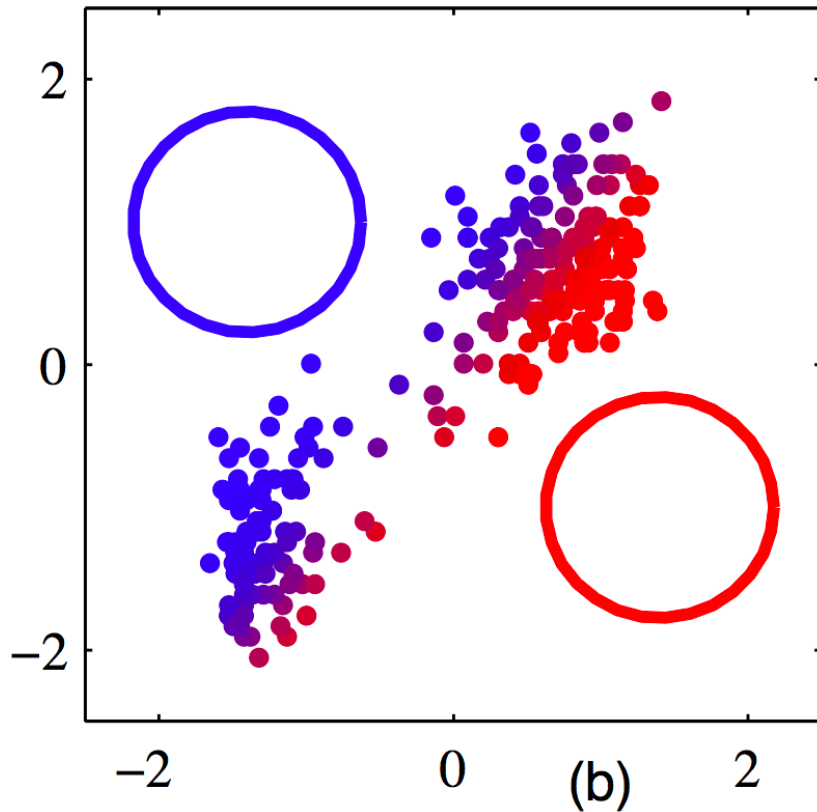- Repeat E-step and M-step until log-likelihood converges.

# EM for GMM

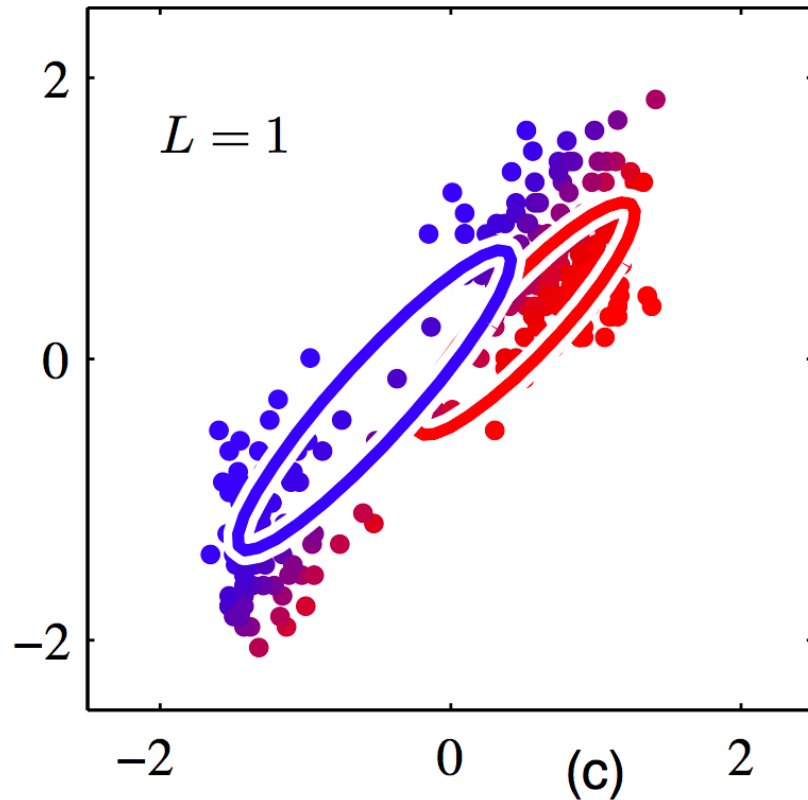- Initialization



(a)

# EM for GMM

- First soft assignment:



responsibilities

$$\gamma_i^j = \frac{\pi_j \mathcal{N}(\boldsymbol{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\sum_{c=1}^{k} \pi_c \mathcal{N}(\boldsymbol{x}_i | \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)}$$
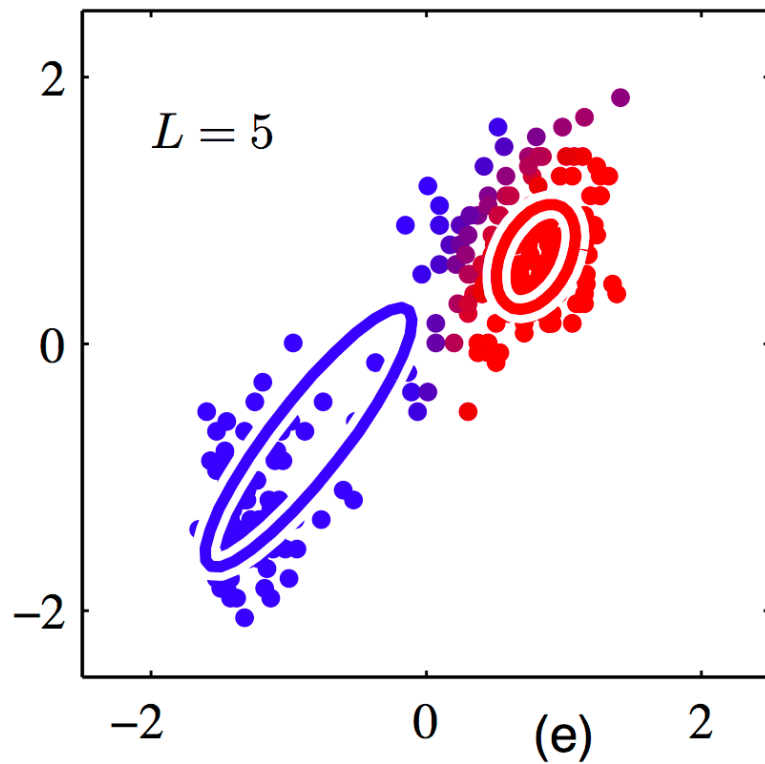
# EM for GMM

- First soft assignment:



parameters

$$\pi_c^{\text{new}} = \frac{n_c}{n}$$

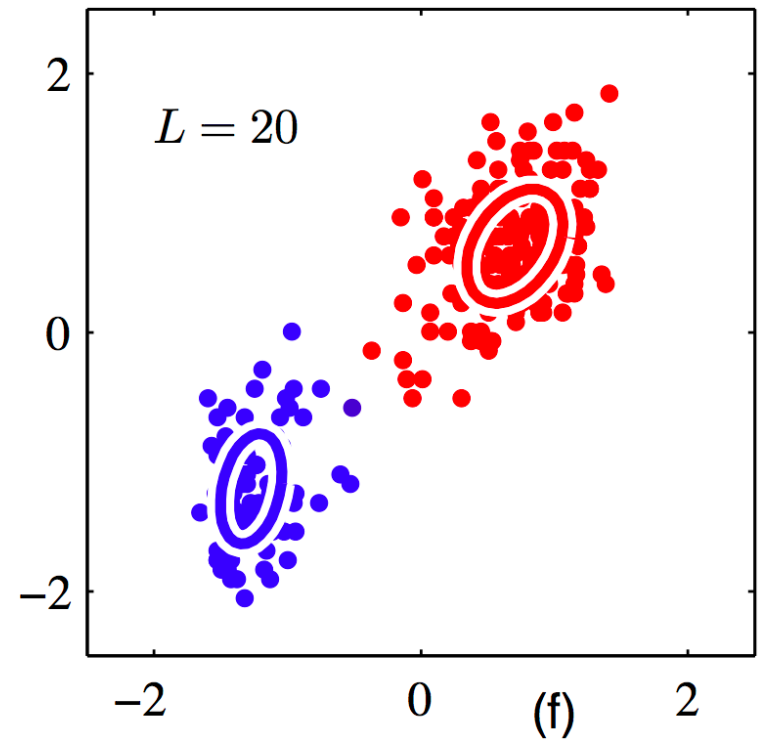$$\boldsymbol{\mu}_c^{\text{new}} = \frac{1}{n_c} \sum_{i=1}^{n} \gamma_i^c \boldsymbol{x}_i$$

$$\boldsymbol{\Sigma}_c^{\text{new}} = \frac{1}{n_c} \sum_{i=1}^{n} \gamma_i^c (\boldsymbol{x}_i - \boldsymbol{\mu}_c^{\text{new}})(\boldsymbol{x}_i - \boldsymbol{\mu}_c^{\text{new}})^T$$

# EM for GMM

- After 5 rounds of EM:



$L = 5$

(e)

- After 20 rounds of EM:



$L = 20$

(f)

# EM and Variational Methods

- When E-step is <span style="color:red">difficult</span>:

  - Hard to take expectation w.r.t. $q^*(z) = p(z|x, \theta^{\mathbf{old}})$

  - For example, hierarchical latent variable models (next lectures).

- **Solution:** Restrict to <span style="color:blue">distributions $Q$ that are easy to work with</span>.

- The evidence lower bound (ELBO) now <span style="color:red">looser</span>:

$$q^* = \underset{q \in Q}{\operatorname{argmin}} \operatorname{KL}\left[q(z), p(z|x, \theta^{\mathbf{old}})\right]$$

- Find an easy-to-work <span style="color:blue">variational distribution $q^*$</span> to approximate the <span style="color:red">inference distribution $p(z|x, \theta^{\mathbf{old}})$</span>.

  - This group of methods are called <span style="color:red">variational methods</span>.

# Thank You

# Questions?

Mingsheng Long
[mingsheng@tsinghua.edu.cn](mailto:mingsheng@tsinghua.edu.cn)
[http://ise.thss.tsinghua.edu.cn/~mlong](http://ise.thss.tsinghua.edu.cn/~mlong)
答疑：东主楼11区413室