

Chapter 11

Unimodal Models

supplementary slides to
Machine Learning Fundamentals
© **Hui Jiang 2020**
published by Cambridge University Press

August 2020

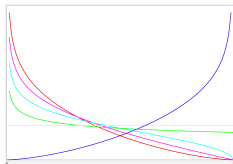
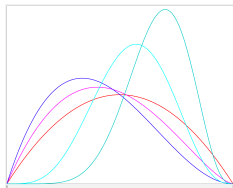


Outline

- 1 Gaussian Models
- 2 Multinomial Models
- 3 Markov Chain Models
- 4 Generalized Linear Models

Unimodal Models

- unimodal models:
 - simple generative models with a single peak extended to include all bounded monotonic functions
 - unimodality for multivariate models: all marginal distributions are unimodal
- unimodal models include almost all common probability distributions, Markov chains, generalized linear models, etc.
- parameter estimation is straightforward
- suitable for simple data distributions where the probability mass is concentrated only in a single region of the space



Multivariate Gaussian Models (I)

- given a training set: $D = \{\mathbf{x}_1; \mathbf{x}_2; \dots; \mathbf{x}_N\}$ ($\mathbf{x}_i \in \mathbb{R}^d$)
- choose a multivariate Gaussian distribution to model D :

$$p_{\mu; \Sigma}(\mathbf{x}) = N(\mathbf{x}; \mu, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} e^{-\frac{(\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu)}{2}}$$

- the log-likelihood function:

$$\begin{aligned} l(\mu; \Sigma) &= \sum_{i=1}^N \ln p_{\mu; \Sigma}(\mathbf{x}_i) \\ &= C - \frac{N}{2} \ln |\Sigma| - \frac{1}{2} \sum_{i=1}^N (\mathbf{x}_i - \mu)^T \Sigma^{-1} (\mathbf{x}_i - \mu) \end{aligned}$$

Multivariate Gaussian Models (II)

$$\frac{\partial l(\boldsymbol{\mu}; \mathbf{X})}{\partial \boldsymbol{\mu}} = 0 \Rightarrow \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu}) = 0 \Rightarrow \boldsymbol{\mu}_{\text{MLE}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$$

$$\begin{aligned} \frac{\partial l(\boldsymbol{\Sigma}; \mathbf{X})}{\partial \boldsymbol{\Sigma}} = 0 &\Rightarrow \frac{N}{2} (\boldsymbol{\Sigma})^{-1} + \frac{1}{2} (\boldsymbol{\Sigma})^{-1} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu}_{\text{MLE}})(\mathbf{x}_i - \boldsymbol{\mu}_{\text{MLE}})^T (\boldsymbol{\Sigma})^{-1} \\ &\Rightarrow \boldsymbol{\Sigma}_{\text{MLE}} = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu}_{\text{MLE}})(\mathbf{x}_i - \boldsymbol{\mu}_{\text{MLE}})^T \end{aligned}$$

Note that $\frac{\partial}{\partial \mathbf{A}} \mathbf{x}^T \mathbf{A}^{-1} \mathbf{y} = -(\mathbf{A}^{-1})^T \mathbf{y} \mathbf{x}^T (\mathbf{A}^{-1})^T$ (square \mathbf{A})

$\frac{\partial}{\partial \mathbf{A}} \ln |j\mathbf{A}| = (\mathbf{A}^{-1})^T = (\mathbf{A}^T)^{-1}$ (square \mathbf{A})

Gaussian Models for Classification

- assuming K classes $f \in \{1; \dots; K\}$, we collect a training set D_k for each class f_k ($k = 1; 2; \dots; K$)
- if each feature vector is continuous ($\in \mathbb{R}^d$) and follows a unimodal distribution, we may choose a multivariate Gaussian for each class, i.e. $N(\mathbf{x}^j | \mu^{(k)}; \Sigma^{(k)})$ ($k = 1; 2; \dots; K$)
- maximum likelihood estimation: $D_k \rightarrow \hat{\mu}^{(k)}_{MLE}; \hat{\Sigma}^{(k)}_{MLE}$
- classify any new \mathbf{x} using the plug-in MAP decision rule:

$$g(\mathbf{x}) = \arg \max_k \Pr(f_k) p(\mathbf{x} | f_k) = \arg \max_k N(\mathbf{x} | \hat{\mu}^{(k)}_{MLE}; \hat{\Sigma}^{(k)}_{MLE})$$

assuming all priors $\Pr(f_k)$ are equiprobable.

Separation Boundary of Gaussian Models

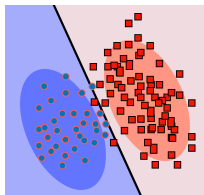
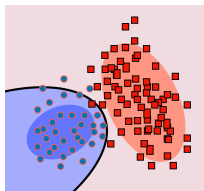
- quadratic discriminant analysis (QDA)
 - Gaussian models lead to quadratic classifiers
 - separation boundaries between classes are a parabola-like quadratic surface
- linear discriminant analysis (LDA)
 - use a common covariance matrix

$$D_1; D_2; \dots; D_K \quad ! \quad \text{MLE}$$

$$D_k \quad ! \quad \mu_{\text{MLE}}^{(k)} \quad (k = 1; \dots; K)$$

the plug-in MAP rule leads to a linear separation boundary:

$$g(\mathbf{x}) = \arg \max_k N(\mathbf{x} | \mu_{\text{MLE}}^{(k)}; \Sigma_{\text{MLE}})$$



Multinomial Models

- discrete data consists of some independent observations, e.g. $X = (x_1; x_2; \dots; x_T)$, each of which is a distinct symbol
- assume there are M distinct symbols in total
- denote p_i as the probability of observing the i th symbol for all $i = 1; 2; \dots; M$
- the sum-to-one constraint: $\sum_{i=1}^M p_i = 1$
- use a multinomial model for X :

$$\Pr(X | p_1; p_2; \dots; p_M) = \frac{(r_1 + r_2 + \dots + r_M)!}{r_1! r_2! \dots r_M!} p_1^{r_1} p_2^{r_2} \dots p_M^{r_M}$$

where r_i ($i = 1; 2; \dots; M$) denotes the frequency of the i th symbol appearing in X

Maximum Likelihood Estimation of Multinomial Models

- estimate all parameters of a multinomial model, i.e. $\{p_1; p_2; \dots; p_M\}$, from some observations
- maximize the log-likelihood function:

$$\arg \max_{p_1; p_2; \dots; p_M} \ln \Pr(X \mid p_1; p_2; \dots; p_M)$$

subject to $\sum_{i=1}^M p_i = 1$

- construct the Lagrangian function:

$$L(p_1; p_2; \dots; p_M; \lambda) = C + \sum_{i=1}^M r_i \ln p_i - \lambda \left(\sum_{i=1}^M p_i - 1 \right)$$

$$\frac{\partial L(p_1; p_2; \dots; p_M; \lambda)}{\partial p_i} = 0 \Rightarrow p_i = \frac{r_i}{\sum_{i=1}^M r_i}$$

- maximum likelihood estimation:

$$p_i^{(\text{MLE})} = \frac{r_i}{\sum_{i=1}^M r_i} \quad (i = 1; 2; \dots; M)$$

Example: Multinomial Models for DNA Sequences

X = GAATTCTTCAAAGAGTTCCAGATATCCACAGGCAGATTCTA
GCACACATCTCAATGAAGTTCCTGAGAAAGCTTCTGTCTAGTTTT
GAAAATATTTCTTTTCCATCATGGGCCTCAAAGCGCTCAAATG
TTGCAGATACTAGAGAAAGACTGTTT

- assume all nucleotides in a sequence are independent
- p_1 denotes the probability of observing G at any location, p_2 for A, p_3 for C, p_4 for T

$$\Pr(X \mid p_1; p_2; p_3; p_4) = \frac{(r_1 + r_2 + r_3 + r_4)!}{r_1! r_2! r_3! r_4!} \prod_{i=1}^4 p_i^{r_i}$$

- maximum likelihood estimation:

$$p_i^{(\text{MLE})} = \frac{r_i}{\sum_{i=1}^4 r_i} \quad (i = 1; 2; 3; 4)$$

Markov Chain Models

- multinomial models treat sequence as a bag of symbols by ignoring the order information
- how to model the sequential information of sequences?

$$X = \overset{n}{x_1 x_2 x_3 \dots x_{t-1} x_t x_{t+1} \dots x_T} \overset{0}$$

- consider the product rule:

$$\Pr(X) = p_{x_1} p_{x_2|x_1} p_{x_3|x_1x_2} \dots p_{x_T|x_1 \dots x_{T-1}}$$

- uncontrollable complexity: as a sequence gets longer and longer, it requires some conditional distributions involving more and more parameters.

Markov Chain Models: Markov Assumption

- Markov assumption: every random variable in a sequence only depends on its most recent history, independent from the others given the most recent history
- 1st-order Markov assumption:

$$p(x_t | x_1, \dots, x_{t-1}) = p(x_t | x_{t-1})$$

- 2nd-order Markov assumption:

$$p(x_t | x_1, \dots, x_{t-1}) = p(x_t | x_{t-2}, x_{t-1})$$

- Markov chain models:

$$\Pr(X) = p(x_1) \prod_{t=2}^T p(x_t | x_{t-1})$$

$$\Pr(X) = p(x_1) p(x_2 | x_1) \prod_{t=3}^T p(x_t | x_{t-2}, x_{t-1})$$

Markov Chain Models: Two More Assumptions

- stationary assumption :

conditional distributions do not change over time

$$p(x_t | x_{t-1}) = p(x_t | x_{t-2}, x_{t-1})$$

allow to use one conditional distribution for all time instances

- discrete observation assumption: all observations in a sequence are the same random variable out of a finite set of distinct symbols, i.e. $\{s_1, s_2, \dots, s_M\}$

represent the conditional distribution as a transition matrix:

$$A = \begin{matrix} & \begin{matrix} h \\ i \end{matrix} \\ \begin{matrix} h \\ i \end{matrix} & a_{ij} \end{matrix} \quad \begin{matrix} M \\ M \end{matrix}$$

with $a_{ij} = \Pr(x_t = s_j | x_{t-1} = s_i)$

represent Markov chain models as directed graphs

Example 1: Markov Chain Models for DNA Sequences

- 1st-order Markov chain model for DNA sequence

- MLE formula: $a_{ij}^{(\text{MLE})} = \frac{r(!_i!_j)}{r(!_i)}$ for all $(1 \leq i, j \leq M)$

- for any new sequence:

$$\Pr(\text{GAATC}) = p(\text{G}|\text{begin})p(\text{A}|\text{G})p(\text{A}|\text{A})p(\text{T}|\text{A})p(\text{C}|\text{T})p(\text{end}|\text{C})$$

$$= 0:25 \quad 0:16 \quad 0:18 \quad 0:12 \quad 0:35 \quad 0:01$$

Example 2: N-gram Language Models

- n-gram language models: use Markov chain models for languages
each word only depends on its previous word(s)
a set of word conditional probabilities

- given any sentence:

S = I would like to fly from Toronto to Chicago this Friday

- 1st-order Markov chain (bi-gram model) N^2 parameters

$$\Pr(S) = p(\text{I}|\text{begin}) p(\text{would}|\text{I}) p(\text{like}|\text{would}) \dots p(\text{end}|\text{Friday})$$

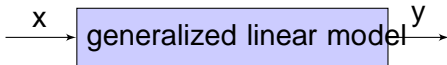
- 2nd-order Markov chain (tri-gram model) N^3 parameters

$$\Pr(S) = p(\text{I}|\text{begin}) p(\text{would}|\text{begin}; \text{I}) p(\text{like}|\text{I}; \text{would}) \dots p(\text{end}|\text{this}; \text{Friday})$$

- zero-order Markov chain (uni-gram model) N parameters

$$\Pr(S) = p(\text{I}) p(\text{would}) p(\text{like}) \dots p(\text{Friday})$$

Generalized Linear Models



generalized linear models (GLMs): extend linear regression to deal with non-Gaussian distributions

- 1 select a unimodal probability distribution for output, i.e. $\phi(y)$
e.g. binomial, multinomial, Poisson
- 2 choose a link function $g(\cdot)$ to associate the mean of y to a linear predictor of input x :

$$E y = g(w^T x)$$

e.g. $\exp(\cdot)$, sigmoid, softmax

the range of the link function must match the domain of mean

Generalized Linear Models: Some Examples

GLM	y	distribution	$g(\cdot)$
linear regression	\mathbb{R}	Gaussian	identity
logistic regression	binary	binomial	sigmoid
probit regression	binary	binomial	probit
Poisson regression	count	Poisson	$\exp(\cdot)$
log-linear model	categorical	multinomial	softmax

Table: Some popular generalized linear models (GLMs)

GLM Example (1): Logistic & Probit Regression

- GLMs for binary classification $y \in \{0, 1\}$
- choose binomial distribution for $p(y)$:

$$y \sim B(y | N = 1; p) = p^y (1 - p)^{1 - y}$$

$$\text{where } E[y] = p \in (0, 1)$$

- logistic regression uses the sigmoid function as the link function: $p = \text{I}(w^T x)$

$$\hat{p}_w(y | x) = \text{I}(w^T x)^y (1 - \text{I}(w^T x))^{1 - y}$$

- probit regression uses the probit function as the link function: $p = \Phi(w^T x)$

$$\hat{p}_w(y | x) = \Phi(w^T x)^y (1 - \Phi(w^T x))^{1 - y}$$

- sigmoid

$$\text{I}(x) = \frac{1}{1 + e^{-x}}$$

- probit

$$\Phi(x) = \frac{1}{2} (1 + \text{erf}(x))$$

GLM Example (2): Poisson Regression

- A GLM for count data: $y = 0; 1; 2; 3;$
- choose Poisson distribution for $p(y)$:

$$p(y) = \frac{e^{-\mu} \mu^y}{y!} \quad y = 0; 1; 2; 3;$$

where $\mu = \exp(w^T x) > 0$

- Poisson regression uses the exponential function as the link function: $\mu = \exp(w^T x)$

$$p_w(y|x) = \frac{1}{y!} \exp(-\exp(w^T x)) \exp(y w^T x) \quad y = 0; 1; 2; 3;$$

- derive the maximum likelihood estimation \hat{w}_{MLE}

GLM Example (3): Log-linear Model

- A GLM for K -class pattern classification, i.e. $y \in \{1, 2, \dots, K\}$
- use the 1-of- K representation to encode y as a K -dimension one-hot vector $\mathbf{y} = [y_1, y_2, \dots, y_K]^T$

- choose multinomial distribution for \mathbf{y} :

$$\mathbf{y} \sim \text{Mult}(\mathbf{y}; N=1; p_1, \dots, p_K) = \prod_{k=1}^K p_k^{y_k}$$

$$\text{where } \mathbf{E} \mathbf{y} = [p_1, p_2, \dots, p_K]^T$$

- choose the softmax function as the link function:

$$\mathbf{E} \mathbf{y} = \text{softmax}(\mathbf{x}) = \left[\frac{e^{w_1^T \mathbf{x}}}{\sum_{k=1}^K e^{w_k^T \mathbf{x}}}, \frac{e^{w_2^T \mathbf{x}}}{\sum_{k=1}^K e^{w_k^T \mathbf{x}}}, \dots, \frac{e^{w_K^T \mathbf{x}}}{\sum_{k=1}^K e^{w_k^T \mathbf{x}}} \right]^T$$

- log-linear model:

$$p_{w_1; \dots; w_K}(\mathbf{y} | \mathbf{x}) = \prod_{k=1}^K \frac{e^{w_k^T \mathbf{x}}}{\sum_{k=1}^K e^{w_k^T \mathbf{x}}}^{y_k}$$

Log-linear Models for Text Categorization

- text categorization: automatically classify text documents into different categories
- feature engineering: use some pre-defined rules to extract a fixed-size feature vector \mathbf{x} to represent each text document
- use log-linear models as the classifier
- given a training set as $D = (\mathbf{x}^{(i)}; \mathbf{y}^{(i)}) \quad j \quad i = 1; 2; \dots; N$
- MLE: to maximize log-likelihood function

$$l(\mathbf{w}_1; \dots; \mathbf{w}_K) = \sum_{i=1}^N \sum_{k=1}^K y_k^{(i)} \ln \Pr \left(\frac{e^{\mathbf{w}_k | \mathbf{x}^{(i)}}}{\sum_{k=1}^K e^{\mathbf{w}_k | \mathbf{x}^{(i)}}} \right) !$$

- classify any new document \mathbf{x} :

$$\begin{aligned} \hat{k} &= \arg \max_k \Pr(k | \mathbf{x}) = \arg \max_k \Pr \left(\frac{e^{(\mathbf{w}_k^{(\text{MLE})}) | \mathbf{x}}}{\sum_{k=1}^K e^{(\mathbf{w}_k^{(\text{MLE})}) | \mathbf{x}}} \right) \\ &= \arg \max_k \mathbf{x} | \mathbf{w}_k^{(\text{MLE})} \end{aligned}$$