Gaussian Models

Multinomial Models

Markov Chain 00000

< (T) >

GLM 000000

Chapter 11 Unimodal Models

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

August 2020



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

Gaussian Models	Multinomial Models	Markov Chain	GLM

Outline

1 Gaussian Models

2 Multinomial Models

3 Markov Chain Models

4 Generalized Linear Models

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

Gaussian Models	Multinomial Models	Markov Chain	GLM

Unimodal Models

unimodal models:

- $\circ~$ 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





Gaussian Models	Multinomial Models	Markov Chain	GLM
● 000			

Multivariate Gaussian Models (I)

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

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

the log-likelihood function:

$$l(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{i=1}^{N} \ln p_{\boldsymbol{\mu}, \boldsymbol{\Sigma}}(\mathbf{x}_i)$$
$$= C - \frac{N}{2} \ln |\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{i=1}^{N} (\mathbf{x}_i - \boldsymbol{\mu})^{\mathsf{T}} \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu})$$

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

Gaussian Models	Multinomial Models	Markov Chain	GLM
0000			

Multivariate Gaussian Models (II)

$$\frac{\partial l(\boldsymbol{\mu}, \boldsymbol{\Sigma})}{\partial \boldsymbol{\mu}} = 0 \implies \sum_{i=1}^{N} \boldsymbol{\Sigma}^{-1}(\mathbf{x}_{i} - \boldsymbol{\mu}) = 0 \implies \boldsymbol{\mu}_{\text{MLE}} = \frac{1}{N} \sum_{i=1}^{N} \mathbf{x}_{i}$$

$$\frac{\partial l(\boldsymbol{\mu}, \boldsymbol{\Sigma})}{\partial \boldsymbol{\Sigma}} = 0 \implies -\frac{N}{2} (\boldsymbol{\Sigma}^{\mathsf{T}})^{-1} + \frac{1}{2} (\boldsymbol{\Sigma}^{\mathsf{T}})^{-1} \Big[\sum_{i=1}^{N} (\mathbf{x}_{i} - \boldsymbol{\mu}) (\mathbf{x}_{i} - \boldsymbol{\mu})^{\mathsf{T}} \Big] (\boldsymbol{\Sigma}^{\mathsf{T}})^{-1}$$

$$\implies \boldsymbol{\Sigma}_{\mathsf{MLE}} = \frac{1}{N} \sum_{i=1}^{N} (\mathbf{x}_i - \boldsymbol{\mu}_{\mathsf{MLE}}) (\mathbf{x}_i - \boldsymbol{\mu}_{\mathsf{MLE}})^{\mathsf{T}}$$

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

Gaussian Models	Multinomial Models	Markov Chain	GLM
0000			

Gaussian Models for Classification

- assuming K classes $\{\omega_1, \cdots, \omega_K\}$, we collect a training set \mathcal{D}_k for each class ω_k $(k = 1, 2, \cdots, K)$
- if each feature vector is continuous (∈ ℝ^d) and follows a unimodal distribution, we may choose a multivariate Gaussian for each class, i.e. N(x | µ^(k), Σ^(k)) (k = 1, 2, ···, K)

• maximum likelihood estimation: $\mathcal{D}_k \longrightarrow \left\{ \mu_{\scriptscriptstyle \mathsf{MLE}}^{(k)}, \Sigma_{\scriptscriptstyle \mathsf{MLE}}^{(k)}
ight\}$

classify any new x using the plug-in MAP decision rule:

$$g(\mathbf{x}) = \arg\max_{k} \Pr(\omega_{k}) p(\mathbf{x} \,|\, \omega_{k}) = \arg\max_{k} \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_{\mathsf{MLE}}^{(k)}, \boldsymbol{\Sigma}_{\mathsf{MLE}}^{(k)})$$

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

Aultinomial Models

Markov Chain

Separation Boundary of Gaussian Models

quadratic discriminant analysis (QDA)

- o 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

$$\{\mathcal{D}_1, \mathcal{D}_2, \cdots, \mathcal{D}_K\} \longrightarrow \Sigma_{\mathsf{MLE}}$$

$$\mathcal{D}_k \longrightarrow \boldsymbol{\mu}_{\mathsf{MLE}}^{(k)} \quad (k = 1, \cdots, K)$$

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

$$g(\mathbf{x}) = \arg \max_{k} \ \mathcal{N}(\mathbf{x} \,|\, \boldsymbol{\mu}_{\mathsf{MLE}}^{(k)}, \boldsymbol{\Sigma}_{\mathsf{MLE}})$$





Gaussian Models	Multinomial Models	Markov Chain	GLM
	000		

Multinomial Models

- discrete data consists of some *independent* observations, e.g. $\mathbf{X} = \{x_1, x_2, \cdots, x_T\}$, 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 $i = 1, 2, \cdots, M$
- the sum-to-one constraint: $\sum_{i=1}^{M} p_i = 1$

use a multinomial model for X:

$$\Pr(\mathbf{X} \mid p_1, p_2, \cdots p_M) = \frac{(r_1 + r_2 + \cdots + r_M)!}{r_1! r_2! \cdots r_M!} p_1^{r_1} p_2^{r_2} \cdots p_M^{r_M}$$

where $r_i \ (i = 1, 2, \cdots, M)$ denotes the frequency of the *i*-th symbol appearing in ${f X}$

Gaussian Models	Multinomial Models	Markov Chain	GLM
	000		

Maximum Likelihood Estimation of Multinomial Models

- estimate all parameters of a multinomial model, i.e. $\{p_1, p_2, \cdots p_M\}$, from some observation ${f X}$
- maximize the log-likelihood function:

 $\arg \max_{p_1, p_2, \cdots p_M} \ln \Pr(\mathbf{X} \mid p_1, p_2, \cdots p_M)$

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

construct the Lagrangian function:

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

$$\frac{\partial}{\partial p_i} \mathcal{L}(p_1, p_2, \cdots, p_M, \lambda) = 0 \implies p_i = \frac{r_i}{\lambda}$$

maximum likelihood estimation:

$$p_i^{(\mathrm{MLE})} = \frac{r_i}{\sum_{i=1}^M r_i} \quad (i = 1, 2, \cdots, M)$$

supplementary slides to Machine Learning Fundamentals ^(C) Hui Jiang 2020 published by Cambridge University Press

Markov Chain

Example: Multinomal Models for DNA Sequences

- $\mathbf{X} = \mathbf{G} \mathbf{A} \mathbf{A} \mathbf{T} \mathbf{C} \mathbf{T} \mathbf{C} \mathbf{A} \mathbf{A} \mathbf{G} \mathbf{A} \mathbf{G} \mathbf{G} \mathbf{C} \mathbf{A} \mathbf{C} \mathbf{A} \mathbf{G} \mathbf{G} \mathbf{C} \mathbf{A} \mathbf{C} \mathbf{A} \mathbf{C} \mathbf{A} \mathbf{A} \mathbf{A} \mathbf{G} \mathbf{G} \mathbf{G} \mathbf{C} \mathbf{A} \mathbf{C} \mathbf{A} \mathbf{T} \mathbf{C} \mathbf{T} \mathbf{C} \mathbf{A} \mathbf{A} \mathbf{A} \mathbf{G} \mathbf{G} \mathbf{C} \mathbf{T} \mathbf{C} \mathbf{T} \mathbf{G} \mathbf{A} \mathbf{A} \mathbf{A} \mathbf{T} \mathbf{A} \mathbf{T} \mathbf{T} \mathbf{C} \mathbf{C} \mathbf{T} \mathbf{T} \mathbf{C} \mathbf{C} \mathbf{T} \mathbf{G} \mathbf{G} \mathbf{G} \mathbf{G} \mathbf{C} \mathbf{T} \mathbf{C} \mathbf{A} \mathbf{A} \mathbf{A} \mathbf{G} \mathbf{G} \mathbf{C} \mathbf{T} \mathbf{C} \mathbf{A} \mathbf{A} \mathbf{A} \mathbf{G} \mathbf{G} \mathbf{C} \mathbf{T} \mathbf{C} \mathbf{C} \mathbf{A} \mathbf{A} \mathbf{A} \mathbf{G} \mathbf{G} \mathbf{C} \mathbf{T} \mathbf{C} \mathbf{A} \mathbf{A} \mathbf{A} \mathbf{G} \mathbf{G} \mathbf{C} \mathbf{T} \mathbf{C} \mathbf{A} \mathbf{A} \mathbf{A} \mathbf{A} \mathbf{G} \mathbf{G} \mathbf{C} \mathbf{T} \mathbf{C} \mathbf{A} \mathbf{A} \mathbf{A} \mathbf{C} \mathbf{A} \mathbf{A} \mathbf{A} \mathbf{C} \mathbf{A} \mathbf{A} \mathbf{A} \mathbf{C} \mathbf{A} \mathbf{A} \mathbf{A} \mathbf{C} \mathbf{C} \mathbf{A} \mathbf{A} \mathbf{A} \mathbf{C} \mathbf{C} \mathbf{A} \mathbf{A} \mathbf{A} \mathbf{C} \mathbf{A} \mathbf{A} \mathbf{A} \mathbf{C} \mathbf{C} \mathbf{A} \mathbf{A} \mathbf{A} \mathbf{C} \mathbf{A} \mathbf{A} \mathbf{A} \mathbf{C} \mathbf{C} \mathbf{A} \mathbf{A} \mathbf{A} \mathbf{C} \mathbf{A} \mathbf{A} \mathbf{C} \mathbf{A} \mathbf{A} \mathbf{A} \mathbf{C} \mathbf{A} \mathbf{A} \mathbf{C} \mathbf{A} \mathbf{A} \mathbf{A} \mathbf{C} \mathbf{A} \mathbf{A} \mathbf{A} \mathbf{C} \mathbf{A} \mathbf{A} \mathbf{A} \mathbf{C} \mathbf{A} \mathbf{A} \mathbf{C} \mathbf{A} \mathbf{A} \mathbf{C} \mathbf{A} \mathbf{A} \mathbf{A} \mathbf{C} \mathbf{A} \mathbf{C} \mathbf{A} \mathbf{A} \mathbf{C} \mathbf{A} \mathbf{C} \mathbf{A} \mathbf{C} \mathbf{A} \mathbf{A} \mathbf{C} \mathbf{A} \mathbf$
 - assume all nucleotides in a sequence are independent
 - $\blacksquare \ p_1$ denotes the probability of observing G at any location, p_2 for A, p_3 for C, p_4 for T

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

maximum likelihood estimation:

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

Gaussian Models	Multinomial Models	Markov Chain	GLM
		00000	

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?

$$\mathbf{X} = \left\{ x_1 \ x_2 \ x_3 \cdots x_{t-1} \ x_t \ x_{t+1} \cdots x_T \right\}$$

consider the product rule:

$$\Pr(\mathbf{X}) = p(x_1)p(x_2|x_1)p(x_3|x_1x_2)\cdots$$
$$p(x_t|x_1\cdots x_{t-1})\cdots p(x_T|x_1\cdots x_{T-1})$$

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

Gaussian Models	Multinomial Models	Markov Chain	GLM
		0000	

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:

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

2nd-order Markov assumption:

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

Markov chain models:

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

Gaussian Models	Multinomial Models	Markov Chain	GLM
		00000	

Markov Chain Models: Two More Assumptions

stationary assumption:

conditional distributions do not change over time

$$p(x_t | x_{t-1}) = p(x_{t'} | x_{t'-1})$$

o 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 M distinct symbols, i.e. { $\omega_1, \omega_2, \cdots, \omega_M$ }
 - represent the conditional distribution as a transition matrix:

$$\mathbf{A} = \begin{bmatrix} a_{ij} \end{bmatrix}_{M \times M}$$

with $a_{ij} = \Pr(x_t = \omega_j | x_{t-1} = \omega_i)$ • represent Markov chain models as directed graphs

Gaussian Models	Multinomial Models	Markov Chain	GLM
		00000	

Example 1: Markov Chain Models for DNA Sequences

Ist-order Markov chain model for DNA sequence



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

 $= 0.25 \times 0.16 \times 0.18 \times 0.12 \times 0.35 \times 0.01$

Gaussian Models	Multinomial Models	Markov Chain	GLM
		00000	

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)
 - o a set of word conditional probabilities
- given any sentence:

 $\mathbf{S}=\mathbf{I}$ would like to fly from Toronto to Chicago this Friday

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

 $\Pr(\mathbf{S}) = p(\mathsf{I}|begin) \, p(\mathsf{would}|\mathsf{I}) \, p(\mathsf{like}|\mathsf{would}) \, \cdots \, p(end|\mathsf{Friday})$

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

 $\Pr(\mathbf{S}) = p(\mathsf{I}|begin) \, p(\mathsf{would}|begin, \mathsf{I}) \, p(\mathsf{like}|\mathsf{I}, \mathsf{would}) \, \cdots \, p(end|\mathsf{this}, \mathsf{Friday})$

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

 $\Pr(\mathbf{S}) = p(\mathsf{I}) p(\mathsf{would}) p(\mathsf{like}) \cdots p(\mathsf{Friday})$

向下 イヨト イヨト ニヨ

Gaussian Models	Multinomial Models	Markov Chain	GLM
0000	000	00000	●00000

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. p(y) \circ e.g. binomial, multinomial, Poisson
 - e.g. billomai, mutilioniai, i oisson
- 2 choose a link function $g(\cdot)$ to associate the mean of y to a linear predictor of input \mathbf{x} :

$$\mathbb{E}\big[y\big] = g(\mathbf{w}^\intercal \mathbf{x})$$

- $\circ~$ e.g. $\exp(\cdot)\text{, sigmoid, softmax}$
- $\circ\;$ the range of the link function must match the domain of mean

Markov Chain 00000

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)

Gaussian Models	Multinomial Models	Markov Chain	GLM
			000000

GLM Example (1): Logistic & Probit Regression

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

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

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

logistic regression uses the sigmoid function as the link function: $p = l(\mathbf{w}^{\mathsf{T}}\mathbf{x})$

$$\hat{p}_{\mathbf{w}}(y|\mathbf{x}) = \left(l(\mathbf{w}^{\mathsf{T}}\mathbf{x})\right)^{y} \left(1 - l(\mathbf{w}^{\mathsf{T}}\mathbf{x})\right)^{1-y}$$

probit regression uses the probit function as the link function: $p = \Phi(\mathbf{w}^{\mathsf{T}}\mathbf{x})$

$$\hat{p}_{\mathbf{w}}(y|\mathbf{x}) = \left(\Phi(\mathbf{w}^{\mathsf{T}}\mathbf{x})\right)^{y} \left(1 - \Phi(\mathbf{w}^{\mathsf{T}}\mathbf{x})\right)^{1-y}$$



sigmoid



probit

くぼ ト く ヨ ト く ヨ ト

$$\Phi(x) = \frac{1}{2} \Bigl(1 \! + \! \mathrm{erf}(x) \Bigr)$$

Gaussian Models	Multinomial Models	Markov Chain	GLM
			000000

GLM Example (2): Poisson Regression

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

$$y \sim p(y \mid \lambda) = \frac{e^{-\lambda} \cdot \lambda^y}{y!} \quad \forall y = 0, 1, 2, 3, \cdots$$

where $\mathbb{E}[y] = \lambda > 0$

Poisson regression uses the exponential function as the link function: $\lambda = \exp(\mathbf{w}^{\mathsf{T}}\mathbf{x})$

$$\hat{p}_{\mathbf{w}}(y|\mathbf{x}) = \frac{1}{y!} \exp\left(-\exp(\mathbf{w}^{\mathsf{T}}\mathbf{x})\right) \exp\left(y\mathbf{w}^{\mathsf{T}}\mathbf{x}\right) \quad y = 0, 1, 2, 3, \cdots$$

 \blacksquare derive the maximum likelihood estimation $\hat{\mathbf{w}}_{\text{MLE}}$

Gaussian Models	Multinomial Models	Markov Chain	GLM
			000000

GLM Example (3): Log-linear Model

- A GLM for K-class pattern classification, i.e. $y \in \{\omega_1, \omega_2, \cdots , \omega_K\}$
- use the 1-of-K representation to encode y as a K-dimension one-hot vector $\mathbf{y} \stackrel{\Delta}{=} \begin{bmatrix} y_1 \ y_2 \ \cdots \ y_K \end{bmatrix}^{\mathsf{T}}$
- choose multinomial distribution for $p(\mathbf{y})$:

$$\mathbf{y} \sim \mathsf{Mult}(\mathbf{y} \mid N = 1, p_1, \cdots, p_K) = \prod_{k=1}^K p_k^{y_k}$$
ere $\mathbb{E}[\mathbf{y}] = [p_1 \ p_2 \ \cdots \ p_K]^\mathsf{T}$

choose the softmax function as the link function:

$$\mathbb{E}[\mathbf{y}] = \mathsf{softmax}(\mathbf{x}) = \left[\frac{e^{\mathbf{w}_1^\mathsf{T}\mathbf{x}}}{\sum_{k=1}^{K} e^{\mathbf{w}_k^\mathsf{T}\mathbf{x}}} \frac{e^{\mathbf{w}_2^\mathsf{T}\mathbf{x}}}{\sum_{k=1}^{K} e^{\mathbf{w}_k^\mathsf{T}\mathbf{x}}} \cdots \frac{e^{\mathbf{w}_K^\mathsf{T}\mathbf{x}}}{\sum_{k=1}^{K} e^{\mathbf{w}_k^\mathsf{T}\mathbf{x}}}\right]^\mathsf{T}$$

log-linear model:

wh

$$\hat{p}_{\mathbf{w}_{1},\cdots,\mathbf{w}_{K}}(\mathbf{y} \,|\, \mathbf{x}) = \prod_{k=1}^{K} \left(\frac{e^{\mathbf{w}_{k}^{\mathsf{T}} \mathbf{x}}}{\sum_{k=1}^{K} e^{\mathbf{w}_{k}^{\mathsf{T}} \mathbf{x}}} \right)^{y_{k}}$$

Gaussian Models	Multinomial Models	Markov Chain	GLM
			000000

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 x to represent each text document
- use log-linear models as the classifier
- given a training set as $\mathcal{D} = \{ (\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \mid i = 1, 2, \cdots N \}$
- MLE: to maximize log-likelihood function

$$l(\mathbf{w}_1, \cdots \mathbf{w}_K) = \sum_{i=1}^N \sum_{k=1}^K y_k^{(i)} \ln\left(\frac{e^{\mathbf{w}_k^\mathsf{T} \mathbf{x}^{(i)}}}{\sum_{k=1}^K e^{\mathbf{w}_k^\mathsf{T} \mathbf{x}^{(i)}}}\right)$$

classify any new document x:

$$\hat{k} = \arg \max_{k} \Pr(\omega_{k} | \mathbf{x}) = \arg \max_{k} \frac{e^{(\mathbf{w}_{k}^{(\mathsf{MLE})})^{\mathsf{T}} \mathbf{x}}}{\sum_{k=1}^{K} e^{(\mathbf{w}_{k}^{(\mathsf{MLE})})^{\mathsf{T}} \mathbf{x}}}$$
$$= \arg \max_{k} \mathbf{x}^{\mathsf{T}} \mathbf{w}_{k}^{(\mathsf{MLE})}$$