

# Chapter 10

## Overview of Generative Models

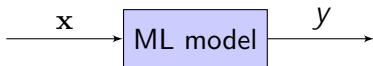
supplementary slides to  
*Machine Learning Fundamentals*  
© **Hui Jiang 2020**  
published by Cambridge University Press

August 2020

# Outline

- 1 Formulation of Generative Models
- 2 Bayesian Decision Theory
- 3 Statistical Data Modelling
- 4 Density Estimation
- 5 Generative Models (in a nutshell)

# Discriminative Models in ML: Review

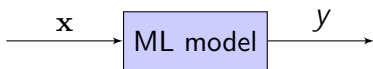


- input  $\mathbf{x}$  is a random vector:  $\mathbf{x} \sim p(\mathbf{x})$
- output  $y$  is generated by an unknown but *deterministic target* function  $y = f(\mathbf{x})$  for each input  $\mathbf{x}$
- our goal: estimate  $f(\cdot)$  in a model space  $\mathcal{H}$
- use a training set:  $D = \{(\mathbf{x}_1; y_1); (\mathbf{x}_2; y_2); \dots; (\mathbf{x}_N; y_N)\}$ , where  $\mathbf{x}_i \sim p(\mathbf{x})$  and  $y_i = f(\mathbf{x}_i)$
- choose a loss function  $l(y; y')$ , and minimize the empirical risk:

$$f^* = \arg \min_{f \in \mathcal{H}} R_{\text{emp}}(f; D) = \arg \min_{f \in \mathcal{H}} \sum_{i=1}^N l(y_i; f(\mathbf{x}_i))$$

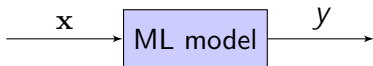
- final performance depends on the generalization bound

# Generative Models in ML



- input  $\mathbf{x}$  and output  $y$  are random variables drawn from an unknown joint distribution, i.e.  $(\mathbf{x}; y) \sim p(\mathbf{x}; y)$
- the relation  $\mathbf{x} \not\perp y$  is stochastic, solely relies on  $p(y|\mathbf{x})$
- our goal: estimate  $p(\mathbf{x}; y)$  using a probabilistic model  $\hat{p}(\mathbf{x}; y)$
- use a training set:  $D = f((\mathbf{x}_1; y_1); (\mathbf{x}_2; y_2); \dots; (\mathbf{x}_N; y_N)g$ , where  $(\mathbf{x}_i; y_i) \sim p(\mathbf{x}; y)$
- the relation  $\mathbf{x} \not\perp y$  may be approximated by  $\hat{p}(y|\mathbf{x})$
- final performance relies on the gap between  $p(\mathbf{x}; y)$  and  $\hat{p}(\mathbf{x}; y)$ , e.g.  $\text{KL}(p(\cdot) \parallel \hat{p}(\cdot))$

# Discriminative vs. Generative Models: Recap



the goal: to estimate an ML model to predict output  $y$  from input  $x$  based on some samples  $D = \{(x_1; y_1); (x_2; y_2); \dots; (x_N; y_N)\}$

## discriminative models

- data generation assumption:  
 $x_i \sim p(\mathbf{x})$  and  $y_i = f(\mathbf{x}_i)$
- $x \neq y$  is deterministic:  
 $y = f(\mathbf{x})$
- use  $D$  to estimate the target function:  $y = f(\mathbf{x})$

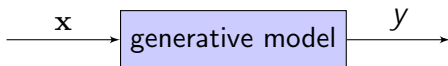
## generative models

- data generation assumption:  
 $(x_i; y_i) \sim p(\mathbf{x}; y)$
- $x \neq y$  is stochastic:  $p(y|\mathbf{x})$
- use  $D$  to estimate the joint distribution:  $p(\mathbf{x}; y)$

# Deterministic vs. Stochastic

- **deterministic**: given the same input  $\mathbf{x}$ , the output  $y$  is always the same as  $y = f(\mathbf{x})$
- **stochastic**: given the same input  $\mathbf{x}$ , the output  $y$  is still a random variable following  $p(y/\mathbf{x})$
- stochasticity may come from noises, parameter variations, etc.
- discriminative models focus on function estimation
- generative models focus on density estimation
- generative models are a more generic setting but also more challenging to learn in general

# Generative Models for Classification



- input  $\mathbf{x}$ : feature vectors (continuous or discrete)
- output  $y = \{1; 2; \dots; K\}$ : **discrete**, called class labels
- the joint distribution  $p(\mathbf{x}; y) = p(y)p(\mathbf{x}|y)$  breaks down to:

prior probabilities:  $p(y = !_k) = \text{Pr}(!_k)$  ( $!_k = 1; 2; \dots; K$ )

class-conditional distributions:  $p(\mathbf{x}|y = !_k) = p(\mathbf{x}|!_k)$   
( $!_k = 1; 2; \dots; K$ )

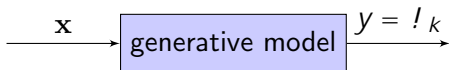
- probabilistic distribution constraints:

priors satisfy  $\sum_{k=1}^K \text{Pr}(!_k) = 1$

if  $\mathbf{x}$  is continuous,  $\int_{\mathbf{x}} p(\mathbf{x}|!_k) d\mathbf{x} = 1$  ( $!_k = 1; 2; \dots; K$ )

if  $\mathbf{x}$  is discrete,  $\sum_{\mathbf{x}} p(\mathbf{x}|!_k) = 1$  ( $!_k = 1; 2; \dots; K$ )

# Bayesian Decision Theory (I): Classification



- given any  $\mathbf{x}$ , determine the best class in  $\{!_1, \dots, !_K\}$
- any decision rule:  $\mathbf{x} \mapsto g(\mathbf{x}) \in \{!_1, \dots, !_K\}$
- Bayesian decision theory: the best decision is

$$\begin{aligned}
 g^*(\mathbf{x}) &= \arg \max_k p(!_k | \mathbf{x}) = \arg \max_k \frac{\Pr(!_k) p(\mathbf{x} | !_k)}{p(\mathbf{x})} \\
 &= \arg \max_k \Pr(!_k) p(\mathbf{x} | !_k)
 \end{aligned}$$

a.k.a. the **maximum a posteriori (MAP) rule** or Bayes decision rule.

- why is this optimal?



# Optimality of the MAP rule (I)

## Theorem 1

Assume  $p(\mathbf{x}; !)$  is known, when  $\mathbf{x}$  is used to predict  $!$ , the MAP rule leads to the lowest expected risk (using 0-1 loss).

- the 0-1 loss function:  $l(!; !') = \begin{cases} 0 & \text{when } ! = !' \\ 1 & \text{otherwise} \end{cases}$
- the expected risk of any rule  $\mathbf{x} \mapsto g(\mathbf{x}) \in \{!_1, \dots, !_K\}$ :

$$\begin{aligned}
 R(g) &= \mathbb{E}_{p(\mathbf{x}; !)} [l(!; g(\mathbf{x}))] = \int_{\mathbf{x}} \sum_{k=1}^K l(!_k; g(\mathbf{x})) p(\mathbf{x}; !_k) d\mathbf{x} \\
 &= \int_{\mathbf{x}} \underbrace{\left[ \sum_{k=1}^K l(!_k; g(\mathbf{x})) p(!_k | \mathbf{x}) \right]}_{\sum_{!_k \in g(\mathbf{x})} p(!_k | \mathbf{x})} p(\mathbf{x}) d\mathbf{x}
 \end{aligned}$$

## Optimality of the MAP rule (II)

- due to  $\sum_{k=1}^K p(!_k | \mathbf{x}) = 1$ , we have

$$\sum_{!_k \neq g(\mathbf{x})} p(!_k | \mathbf{x}) = 1 - p(g(\mathbf{x}) | \mathbf{x})$$

- we have

$$R(g) \stackrel{\#}{=} \int_{\mathcal{X}} \left[ 1 - p(g(\mathbf{x}) | \mathbf{x}) \right] p(\mathbf{x}) d\mathbf{x} \stackrel{\#}{=} \int_{\mathcal{X}} p(g(\mathbf{x}) | \mathbf{x}) p(\mathbf{x}) d\mathbf{x}$$

- since  $g(\mathbf{x}) \in \{!_1, \dots, !_K\}$ , we choose:

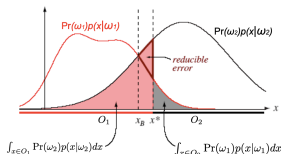
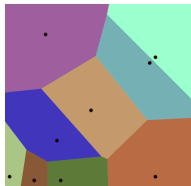
$$g^*(\mathbf{x}) = \arg \max_k p(!_k | \mathbf{x})$$

# Classification Error Probability

- any rule  $x \mapsto g(x) \in \{1, \dots, K\}$  partitions input space into  $K$  regions, i.e.  $O_1; O_2; \dots; O_K$   
 $x \in O_k \Rightarrow g(x) = k$
- the expected risk is the probability of classification error:

$$\begin{aligned}
 R(g) &= \Pr(\text{error}) = 1 - \Pr(\text{correct}) \\
 &= 1 - \sum_{k=1}^K \Pr(x \in O_k; k) \\
 &= 1 - \sum_{k=1}^K \int_{x \in O_k} p(x; k) dx
 \end{aligned}$$

- the Bayes error:  $R(g^*)$  of the MAP rule (the lowest possible error)



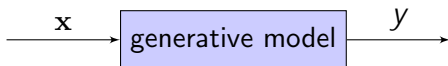
## Example: the MAP rule for independent binary features

- 2-class ( $!_1$  and  $!_2$ ) classification:  $\Pr(!_1)$  and  $\Pr(!_2)$
- use  $d$  independent binary features  $\mathbf{x} = [x_1; x_2; \dots; x_d]^T$ , where  $x_i \in \{0, 1\}$   $\forall i = 1; 2; \dots; d$
- denote  $\theta_i = \Pr(x_i = 1 | !_1)$  and  $\phi_i = \Pr(x_i = 1 | !_2)$ , we have:  
 $p(\mathbf{x} | !_1) = \prod_{i=1}^d \theta_i^{x_i} (1 - \theta_i)^{1 - x_i}$      $p(\mathbf{x} | !_2) = \prod_{i=1}^d \phi_i^{x_i} (1 - \phi_i)^{1 - x_i}$
- the MAP rule: given  $\mathbf{x}$ , classify as  $!_1$  if  $\Pr(!_1) p(\mathbf{x} | !_1) > \Pr(!_2) p(\mathbf{x} | !_2)$ , otherwise  $!_2$ .
- take logarithm to derive a **linear** decision boundary:

$$g(\mathbf{x}) = \sum_{i=1}^d \lambda_i x_i + \lambda_0 = \begin{cases} > 0 & \Rightarrow !_1 \\ < 0 & \Rightarrow !_2 \end{cases}$$

$$\text{where } \lambda_i = \ln \frac{\theta_i(1 - \phi_i)}{\phi_i(1 - \theta_i)} \text{ and } \lambda_0 = \sum_{i=1}^d \ln \frac{1 - \theta_i}{1 - \phi_i} + \ln \frac{\Pr(!_1)}{\Pr(!_2)}$$

# Generative Models for Regression



- input:  $n$ -dimensional vector  $\mathbf{x} \in \mathbb{R}^n$ ; output:  $y \in \mathbb{R}$
- the joint distribution  $p(\mathbf{x}; y)$  is known
- $\mathbf{x}$  is used to predict  $y$  as  $y = g(\mathbf{x})$
- what is the best decision rule  $g(\mathbf{x})$  for  $\mathbf{x} \not\sim y$ ?
- Bayesian decision theory suggests the best rule as:

$$g^*(\mathbf{x}) = E(y|\mathbf{x}) = \int_{\mathcal{Y}} y p(y|\mathbf{x}) dy$$

## Theorem 2

*Assume  $p(\mathbf{x}; y)$  is known, the conditional mean  $E(y|\mathbf{x})$  leads to the lowest expected risk (using mean square loss).*



# Optimality of Conditional Mean for Regression

## Proof:

- The expected risk of any rule  $\mathbf{x} \mapsto g(\mathbf{x}) \in \mathbb{R}$ :

$$\begin{aligned} R(g) &= \mathbb{E}_{p(\mathbf{x};y)} [l(\cdot; g(\mathbf{x}))] = \int_{\mathbf{x}} \int_y (y - g(\mathbf{x}))^2 p(\mathbf{x}; y) d\mathbf{x} dy \\ &= \int_{\mathbf{x}} \underbrace{\left[ \int_y (y - g(\mathbf{x}))^2 p(y|\mathbf{x}) dy \right]}_{Q(g|\mathbf{x})} p(\mathbf{x}) d\mathbf{x} \end{aligned}$$

- Take partial derivative w.r.t.  $g$  as:

$$\begin{aligned} \frac{\partial Q(g|\mathbf{x})}{\partial g(\cdot)} = 0 &\Rightarrow \int_y (g(\mathbf{x}) - y) p(y|\mathbf{x}) dy = 0 \\ \Rightarrow g^*(\mathbf{x}) &= \int_y y p(y|\mathbf{x}) dy = \mathbb{E}(y|\mathbf{x}) \end{aligned}$$

## Plug-in MAP Decision Rule for classification

- since the true distributions  $\Pr(!_k)$  and  $p(\mathbf{x}j!_k)$  are unknown, the optimal MAP decision rule is not feasible in practice
- given training data:  $D = \{(\mathbf{x}_1; y_1); (\mathbf{x}_2; y_2); \dots; (\mathbf{x}_N; y_N)\}$
- choose two probabilistic models:
  - $\hat{p}(!_k)$  to approximate  $\Pr(!_k)$
  - $\hat{p}_k(\mathbf{x})$  to approximate  $p(\mathbf{x}j!_k)$  ( $8k = 1; 2; \dots; K$ )
- parameter estimation: estimate  $f; 1; \dots; K g$  using  $D$
- the optimal MAP rule in theory:

$$!^* = \arg \max_k \Pr(!_k) p(\mathbf{x}j!_k)$$

- the plug-in MAP decision rule in practice:

$$!^* = \arg \max_k \hat{p}(!_k) \hat{p}_k(\mathbf{x})$$

# Plug-in MAP Decision Rule





# Statistical Data Modeling

Assume we have collected some training samples:

$$D = \{(\mathbf{x}_1; y_1); \dots; (\mathbf{x}_N; y_N)\}$$

where each  $(\mathbf{x}_i; y_i) \sim p(\mathbf{x}; y)$  ( $i = 1; 2; \dots; N$ ).

- 1 choose some probabilistic models:

$$\Pr(!_k) = \hat{p}_\lambda(!_k)$$

$$p(\mathbf{x}; !_k) = \hat{p}_{\theta_k}(\mathbf{x}) \quad (k = 1; 2; \dots; K)$$

- 2 estimate the model parameters:

$$D \sim \{\lambda; \theta_1; \dots; \theta_K\}$$

- 3 apply the plug-in MAP rule:

$$\hat{g}(\mathbf{x}) = \arg \max_k \hat{p}_\lambda(!_k) \hat{p}_{\theta_k}(\mathbf{x})$$

# Maximum Likelihood Estimation (I)

- generative models for classification  $f: \mathcal{X} \rightarrow \{1, \dots, K\}$ :  
 prior probabilities:  $\Pr(\mathcal{I}_k)$  ( $k = 1, \dots, K$ )  
 class-conditional distribution:  $p(\mathbf{x} | \mathcal{I}_k)$  ( $k = 1, \dots, K$ )
- collect training data for each class:  $D_k \sim p(\mathbf{x} | \mathcal{I}_k)$
- **density estimation**: estimate the probability distribution from a finite number of samples
- select probabilistic models:  $\hat{p}_k(\mathbf{x}) \approx p(\mathbf{x} | \mathcal{I}_k)$
- **maximum likelihood estimation (MLE)**: learn  $\hat{p}_k(\mathbf{x})$  to maximize the probability of observing the training data  $D_k$

$$\hat{p}_k^* = \arg \max_k \hat{p}_k(D_k) \quad (k = 1, \dots, K)$$

- MLE: fit data best; best interpret the observed data

## Maximum Likelihood Estimation (II)

- drop index  $k$  and  $\hat{\rho}(\cdot) \rightarrow \rho(\cdot)$ , MLE turns to be:

$$\text{MLE} = \arg \max p(D) = \arg \max p(\mathbf{x}_1; \mathbf{x}_2; \dots; \mathbf{x}_N)$$

where  $D = \{ \mathbf{x}_1; \mathbf{x}_2; \dots; \mathbf{x}_N \}$

- assume all data are *i.i.d.* (independent and identically distributed), i.e., all samples are drawn independently from the same distribution:

$$p(\mathbf{x}_1; \mathbf{x}_2; \dots; \mathbf{x}_N) = \prod_{i=1}^N p(\mathbf{x}_i)$$

- why called maximum *likelihood* (not probability)?

$p(\mathbf{x})$ : data distribution of various  $\mathbf{x}$  if  $\theta$  is given (fixed)

$\rho(\theta)$ : likelihood function of  $\theta$  if  $\mathbf{x}$  is given (fixed)

## Maximum Likelihood Estimation (III)

- in many cases, it is more convenient to work with the logarithm of the likelihood rather than the likelihood itself
- denote the log-likelihood function  $l(\theta) = \ln p(D, \theta)$ , we have

$$\theta_{ML} = \arg \max_{\theta} l(\theta) = \arg \max_{\theta} \sum_{i=1}^N \ln p(\mathbf{x}_i; \theta)$$

- optimization methods for ML estimation:
  - differential calculus for simple models, e.g., single univariate/multivariate Gaussian, ...
  - Lagrange optimization for models with constraints, e.g., multinomial, markov chain, ...
  - Expectation-Maximization (EM) method for mixture models, e.g., Gaussian mixture models (GMM), hidden Markov models

# MLE Example: Univariate Gaussian Models

- the training set:  $D = \{x_1; x_2; \dots; x_N\}$  ( $x_i \in \mathbb{R}$ )
- choose a univariate Gaussian to approximate the unknown distribution:

$$p(x) = N(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

- the log-likelihood function:

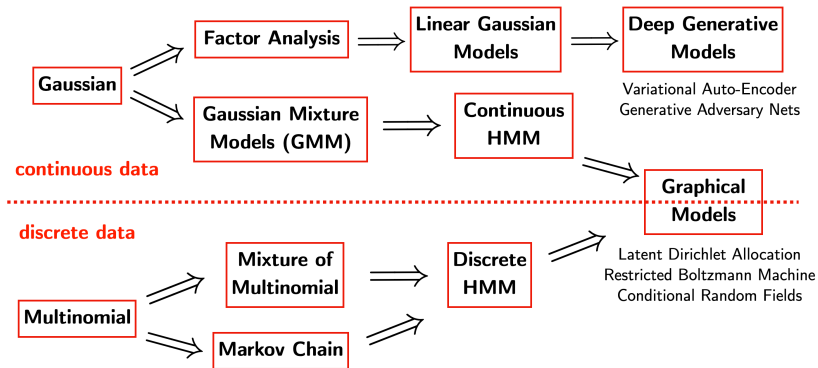
$$l(\mu, \sigma^2) = \sum_{i=1}^N \ln p(x_i) = \sum_{i=1}^N \left[ -\frac{\ln(2\pi\sigma^2)}{2} - \frac{(x_i - \mu)^2}{2\sigma^2} \right]$$

- the MLE of the unknown Gaussian mean and variance:

$$\frac{\partial l(\mu, \sigma^2)}{\partial \mu} = 0 \Rightarrow \mu_{MLE} = \frac{1}{N} \sum_{i=1}^N x_i$$

$$\frac{\partial l(\mu, \sigma^2)}{\partial \sigma^2} = 0 \Rightarrow \sigma_{MLE}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu_{MLE})^2$$

# Roadmap of Generative Models



# Generative Models (in a nutshell)

- Gaussian-derived generative models for continuous data
- multinomial-derived generative models for discrete data
- **unimodal models**: Gaussian, multinomial, Markov chains, generalized linear models, etc.
- **mixture models**: Gaussian mixture models, hidden Markov models, etc.
- **entangled models**: factor analysis, linear Gaussian models, deep generative models (e.g. VAE, GAN)
- **graphical models**: naive Bayes, latent Dirichlet allocation, restricted Boltzmann machine, conditional random fields, etc.
- Bayesian learning: treat model parameters as random variables