

# Parameter Estimation for Generative Models

Ziping Zhao

School of Information Science and Technology  
ShanghaiTech University, Shanghai, China

CS182: Introduction to Machine Learning (Fall 2022)  
<http://cs182.sist.shanghaitech.edu.cn>

Ch. 4 & Ch. 5 of I2ML

# Outline

## Introduction

## Univariate Data

- Maximum Likelihood Estimation

- Bayesian Estimation

- Parametric Classification

- Regression

- Model Selection

## Multivariate Data

- Parameter Estimation

- Multivariate Normal Distribution

- Parametric Classification

- Discrete Features

- Regression

# Outline

## Introduction

### Univariate Data

- Maximum Likelihood Estimation
- Bayesian Estimation
- Parametric Classification
- Regression
- Model Selection

### Multivariate Data

- Parameter Estimation
- Multivariate Normal Distribution
- Parametric Classification
- Discrete Features
- Regression

## Different Approaches to Supervised Learning

- ▶ Three major approaches corresponding to different model assumptions:
  - Parametric approach
  - Nonparametric approach
  - Semiparametric approach
- ▶ Current topic: parametric approach

## Parametric Approach

- ▶ **Assumption**: data follows a distribution that obeys a **parametric model**, e.g., Gaussian.
- ▶ **Advantage of the parametric approach**: the model is fully specified by a small number of **parameters**  $\theta$  as **sufficient statistics** of the distribution
- ▶ The sample  $\mathcal{X} \in \{\mathbf{x}^t\}$  is assumed to be drawn (usually i.i.d.) from an underlying distribution, i.e.,  $\mathbf{x}^t \sim p(\mathbf{x})$
- ▶ The number of parameters, i.e.,  $\dim(\theta)$ , is **independent** of the sample size  $|\mathcal{X}|$ .
- ▶ **Parameter estimation** (or **density estimation**): assuming some parametric form, as a kind of inductive bias, for  $p(\mathbf{x} \mid \theta)$  (or  $p(\mathbf{x}; \theta)$ ),  $\theta$  is estimated using  $\mathcal{X}$ ; it is an unsupervised learning approach
- ▶ The estimated model/distribution is then used to make a decision
- ▶ Two approaches to parameter estimation:
  - **Maximum likelihood estimation**:  $\theta$  is a **fixed point** (**point estimation**)
  - **Bayesian estimation**:  $\theta$  is a **random variable** whose prior uncertainty (represented as a **prior distribution**) can be incorporated.

## Parametric Approach to Classification

- Recall Bayes' rule for classification:

$$P(C_i | \mathbf{x}) = \frac{p(\mathbf{x} | C_i)P(C_i)}{\sum_{k=1}^K p(\mathbf{x} | C_k)P(C_k)}$$

Choose  $C_i$  if  $P(C_i | \mathbf{x}) = \max_k P(C_k | \mathbf{x})$

- To compute  $P(C_i | \mathbf{x})$  for classification,  $p(\mathbf{x} | C_i)$  and  $P(C_i)$  have to be estimated from the sample  $\mathcal{X}$ .
- A classifier thus created is often called a generative classifier (as opposed to a discriminative classifier to be studied later in this course), since it specifies how to generate the data based on  $p(\mathbf{x} | C_i)$  and  $P(C_i)$ .

## Maximum Likelihood Estimation - I

- ▶ Maximum likelihood estimation (MLE) seeks to find  $\theta$  that makes sampling  $\mathcal{X}$  from  $p(\mathbf{x} \mid \theta)$  as likely as possible by maximizing the likelihood of  $\theta$  given the sample  $\mathcal{X} = \{\mathbf{x}^t\}_{t=1}^N$ .
- ▶ Likelihood of  $\theta$  given  $\mathcal{X}$  (with the i.i.d. assumption):

$$L(\theta \mid \mathcal{X}) = p(\mathcal{X} \mid \theta) = \prod_{t=1}^N p(\mathbf{x}^t \mid \theta)$$

- ▶ Log likelihood (mainly for computational simplification):

$$\mathcal{L}(\theta \mid \mathcal{X}) = \log L(\theta \mid \mathcal{X}) = \sum_{t=1}^N \log p(\mathbf{x}^t \mid \theta)$$

- ▶ Maximum likelihood (ML) estimate:

$$\hat{\theta} = \arg \max_{\theta} \mathcal{L}(\theta \mid \mathcal{X})$$

## Maximum Likelihood Estimation - II

Local extrema of a log-transformed function is equal to those of the original function.

Proof:

- ▶ Let  $f(\mathbf{x})$  be a positive function.
- ▶ Suppose that  $\mathbf{x}_0$  is the local maximum of  $f(\mathbf{x})$  in a neighborhood of  $\mathbf{x}_0$ , denoted as  $B_\epsilon(\mathbf{x}_0)$ . Then, for any  $\mathbf{y} \in B_\epsilon(\mathbf{x}_0)$ ,  $f(\mathbf{y}) \leq f(\mathbf{x}_0)$ .
- ▶  $\log(\cdot)$  is a monotonically increasing function: if  $z \leq w$ , then  $\log z \leq \log w$ .
- ▶ So for any  $\mathbf{y} \in B_\epsilon(\mathbf{x}_0)$ , since  $f(\mathbf{y}) \leq f(\mathbf{x}_0)$ , we have  $\log f(\mathbf{y}) \leq \log f(\mathbf{x}_0)$ . Hence,  $\mathbf{x}_0$  is also the local maximum for  $\log f(\mathbf{x})$ .
- ▶ The other direction can be proven by noting that the inverse of  $\log(\cdot)$  is also monotonically increasing.



# Outline

## Introduction

## Univariate Data

- Maximum Likelihood Estimation

- Bayesian Estimation

- Parametric Classification

- Regression

- Model Selection

## Multivariate Data

- Parameter Estimation

- Multivariate Normal Distribution

- Parametric Classification

- Discrete Features

- Regression

# Outline

Introduction

Univariate Data

- Maximum Likelihood Estimation

- Bayesian Estimation

- Parametric Classification

- Regression

- Model Selection

Multivariate Data

- Parameter Estimation

- Multivariate Normal Distribution

- Parametric Classification

- Discrete Features

- Regression

## Example: Bernoulli - I

- ▶ Discrete random variable  $x$  with two possible values,  $x \in \{0, 1\}$
- ▶ E.g., in binary classification, use  $P(x = 1)$  to represent  $P(C_1)$  (and hence  $P(x = 0) = 1 - P(x = 1)$  represents  $P(C_2)$ ).
- ▶ Probability mass function (with parameter  $\theta = p = P(x = 1)$ ):

$$P(x \mid p) = p^x(1 - p)^{1-x}$$

- ▶ Log likelihood:

$$\mathcal{L}(p \mid \mathcal{X}) = \sum_{t=1}^N [x^t \log p + (1 - x^t) \log(1 - p)]$$

- ▶ ML estimate:

$$\hat{p} = \frac{1}{N} \sum_{t=1}^N x^t \quad (\text{sample mean})$$

## Example: Bernoulli - II

- ▶ Note that the estimate  $\hat{p}$  is a function of the sample and is another random variable
- ▶ We can discuss the distribution of  $\hat{p}(\mathcal{X}_i)$  given different  $\mathcal{X}_i$ 's sampled from the same  $p(x)$ 
  - the variance of the distribution of  $\hat{p}(\mathcal{X}_i)$  is expected to decrease as  $N$  increases; as the samples get bigger, they (and hence their averages) get more similar.

## Example: Binomial

- ▶ Discrete random variable  $x$  with  $M$  possible values,  $x \in \{0, 1, \dots, M\}$
- ▶ Probability mass function (with parameter  $\theta = p$ ):

$$P(x \mid p) = \binom{M}{x} p^x (1 - p)^{M-x}$$

- ▶ Log likelihood:

$$\mathcal{L}(p \mid \mathcal{X}) = \sum_{t=1}^N \left[ \log \binom{M}{x_t} + x^t \log p + (M - x^t) \log(1 - p) \right]$$

- ▶ ML estimate:

$$\hat{p} = \frac{1}{NM} \sum_{t=1}^N x^t \quad (\text{sample mean})$$

- ▶ can be viewed as  $NM$  iid Bernoulli trials

## Example: Generalized Bernoulli or Multinomial - I

- ▶ Discrete random variable  $x$  with  $K \geq 2$  possible values, e.g., for  $K$  classes.
- ▶ Indicator variables  $x_1, \dots, x_K$ :

$$x_i = \begin{cases} 1 & \text{if outcome is state } i \\ 0 & \text{if outcome is not state } i \end{cases}$$

with  $\sum_{i=1}^K x_i = 1$ .

- ▶ Probability mass function (with parameters  $\theta = (p_1, \dots, p_K)^T$ ):

$$P(x \mid \theta) = P(x_1, \dots, x_K \mid p_1, \dots, p_K) = \prod_{i=1}^K p_i^{x_i}$$

with constraint

$$\sum_{i=1}^K p_i = 1$$

## Example: Generalized Bernoulli or Multinomial - II

- Log likelihood:

$$\mathcal{L}(p_1, \dots, p_K \mid \mathcal{X}) = \sum_{t=1}^N \sum_{i=1}^K x_i^t \log p_i$$

- Constrained optimization problem (with equality constraint  $\sum_{i=1}^K p_i = 1$ )
  - can be solved using the method of Lagrange multipliers (or equivalently KKT optimality conditions).
- ML estimates:

$$\hat{p}_i = \frac{1}{N} \sum_{t=1}^N x_i^t$$

So  $\hat{\theta} = (\hat{p}_1, \dots, \hat{p}_K)^T$  is also the sample mean.

- Another estimation method: view it as  $K$  separate Bernoulli experiments

## Example: Normal (Gaussian)

- ▶ Continuous random variable  $x$  following univariate normal (Gaussian) distribution  $\mathcal{N}(\mu, \sigma^2)$  with mean  $\mu$  and variance  $\sigma^2$  with  $\sigma > 0$ .
- ▶ Probability density function (with parameters  $\theta = (\mu, \sigma)^T$ ):

$$p(x \mid \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[ -\frac{(x - \mu)^2}{2\sigma^2} \right]$$

- ▶ Log likelihood:

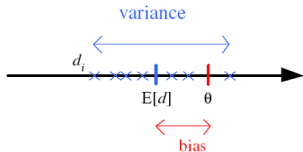
$$\mathcal{L}(\mu, \sigma \mid \mathcal{X}) = -\frac{N}{2} \log(2\pi) - N \log \sigma - \frac{1}{2\sigma^2} \sum_{t=1}^N (x^t - \mu)^2$$

- ▶ ML estimates:

$$m = \frac{1}{N} \sum_{t=1}^N x^t, \quad s^2 = \frac{1}{N} \sum_{t=1}^N (x^t - m)^2$$



## Evaluating an Estimator: Bias and Variance - I



- ▶ **Parameter** to be estimated:  $\theta$ , i.e., the true value
- ▶ **Estimator** based on sample  $\mathcal{X}$ :  $d = d(\mathcal{X})$ ; a random variable depending on  $\mathcal{X}$
- ▶ **Bias**: measures how much the expected value of the estimate  $\mathbb{E}(d) = \mathbb{E}(d(\mathcal{X}))$  deviates from  $\theta$ .

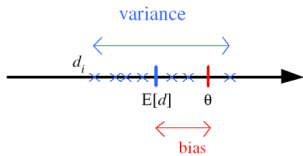
$$b_{\theta}(d) = \mathbb{E}(d) - \theta$$

**Variance**: measures how much, on average, the estimate  $d$  varies from the expected value  $\mathbb{E}[d]$ .

$$\text{Var}(d) = \mathbb{E} \left[ (d - \mathbb{E}(d))^2 \right]$$

- ▶ We would like both to be small.

## Evaluating an Estimator: Bias and Variance - II



- Mean squared error (MSE) of estimator  $d$  (measures how much  $d$  is different from  $\theta$ ):

$$\begin{aligned} r(d, \theta) &= \mathbb{E} \left[ (d - \theta)^2 \right] \\ &= \mathbb{E} \left[ (d - \mathbb{E}(d) + \mathbb{E}(d) - \theta)^2 \right] \\ &= \dots \\ &= (\mathbb{E}(d) - \theta)^2 + \mathbb{E} \left[ (d - \mathbb{E}(d))^2 \right] \\ &= (b_{\theta}(d))^2 + \text{Var}(d) = \text{bias}^2 + \text{variance} \end{aligned}$$

## Evaluating an Estimator: Unbiased Estimator

- ▶ If  $b_{\theta}(d) = \mathbb{E}(d) - \theta = 0$  for all  $\theta$ ,  $d$  is an **unbiased estimator** of  $\theta$ .
- ▶ the sample mean  $m = \frac{1}{N} \sum_{t=1}^N x^t$  is an unbiased estimator of the mean  $\mu$ 
  - $x^t$  is from some density with mean  $\mu$

$$\mathbb{E}(m) = \mathbb{E}\left(\frac{1}{N} \sum_{t=1}^N x^t\right) = \frac{1}{N} \sum_{t=1}^N \mathbb{E}(x^t) = \frac{1}{N} N\mu = \mu$$

- $N \rightarrow \infty, m \rightarrow \mu$
- ▶ the normal ML estimator  $s^2 = \frac{1}{N} \sum_{t=1}^N (x^t - m)^2$  is a biased estimator of  $\sigma^2$ 
  - $x^t$  is from some density with mean  $\mu$  and variance  $\sigma^2$

$$\begin{aligned}\mathbb{E}(s^2) &= \mathbb{E}\left(\frac{1}{N} \sum_{t=1}^N (x^t - m)^2\right) = \mathbb{E}\left(\frac{1}{N} \sum_{t=1}^N (x^t)^2\right) - \mathbb{E}(m^2) \\ &= \sigma^2 + \mu^2 - \left(\frac{\sigma^2}{N} + \mu^2\right) = \frac{N-1}{N} \sigma^2\end{aligned}$$

- $\frac{N}{N-1} s^2$  is an unbiased estimator

## Evaluating an Estimator: Consistent Estimator

- ▶ sample mean  $m$  is a **consistent estimator** of  $\mu$ , since  $\text{Var}(m) \rightarrow 0$  as  $N \rightarrow \infty$ 
  - $x^t$  is from some density with variance  $\sigma^2$

$$\text{Var}(m) = \text{Var} \left( \frac{1}{N} \sum_{t=1}^N x^t \right) = \frac{1}{N^2} \sum_{t=1}^N \text{Var}(x^t) = \frac{1}{N^2} N \sigma^2 = \frac{\sigma^2}{N}$$

- as  $N$  gets larger,  $m$  deviates less from  $\mu$

# Outline

Introduction

## Univariate Data

Maximum Likelihood Estimation

Bayesian Estimation

Parametric Classification

Regression

Model Selection

## Multivariate Data

Parameter Estimation

Multivariate Normal Distribution

Parametric Classification

Discrete Features

Regression

## Bayesian Estimation - I

- ▶ Before looking at a sample  $\mathcal{X}$ , we (or experts of the application) may have some prior information on the possible value range that a parameter,  $\theta$ , may take.
- ▶ This information is useful and should be used, especially when the sample is small.
- ▶ The prior information does not tell us exactly what the parameter value is (otherwise we would not need the sample).
- ▶ Unlike MLE which treats  $\theta$  as a fixed (but unknown) point, the Bayesian estimation approach treats it as a **random variable** with **prior density**  $p(\theta)$  (i.e., our a prior **uncertainty** about  $\theta$ , e.g. normal distribution).
- ▶ Combining the likelihood density  $p(\mathcal{X} | \theta)$  (i.e., information from the sample  $\mathcal{X}$ ), using Bayes' rule, we obtain **posterior density** of  $\theta$  (i.e., uncertainty about  $\theta$  after observing the sample  $\mathcal{X}$ ):

$$p(\theta | \mathcal{X}) = \frac{p(\mathcal{X} | \theta)p(\theta)}{p(\mathcal{X})} = \frac{p(\mathcal{X} | \theta)p(\theta)}{\int p(\mathcal{X} | \theta')p(\theta')d\theta'}$$

## Bayesian Estimation - II

- Full Bayesian estimation approach (through marginalization over  $\theta$ ):
  - Estimation of the density at any  $x$  (i.e., the probability that any sample occurs):

$$\begin{aligned}p(x | \mathcal{X}) &= \int p(x, \theta | \mathcal{X}) d\theta \\&= \int p(x | \theta, \mathcal{X}) p(\theta | \mathcal{X}) d\theta \\&= \int p(x | \theta) p(\theta | \mathcal{X}) d\theta\end{aligned}$$

taking an average over predictions using all  $\theta$ , weighted by their probabilities ( $p(x | \theta, \mathcal{X}) = p(x | \theta)$  because once we know  $\theta$ , the sufficient statistics, we know everything about the distribution)

- Prediction, as in regression, based on a function  $g(x | \theta)$ :

$$y = \int g(x | \theta) p(\theta | \mathcal{X}) d\theta$$

## Computational Considerations

- ▶ Evaluating the integrals (for marginalization over  $\theta$ ) may be computationally difficult, except in cases where the posterior  $p(\theta | \mathcal{X})$  has a nice form.
- ▶ So the full Bayesian estimation approach is sometimes replaced by other methods for computational considerations.
- ▶ **Maximum a posteriori (MAP)** estimation - mode of the posterior density:

$$\theta_{\text{MAP}} = \arg \max_{\theta} p(\theta | \mathcal{X}) = \arg \max_{\theta} p(\mathcal{X} | \theta)p(\theta)$$

$$p(x | \mathcal{X}) = p(x | \theta_{\text{MAP}}) \quad y = g(x | \theta_{\text{MAP}})$$

- **Maximum likelihood (ML) estimation** – MAP with uniform/flat prior:

$$\theta_{\text{ML}} = \arg \max_{\theta} p(\mathcal{X} | \theta)$$

- ▶ **Bayes' estimator** – expectation w.r.t. posterior density:

$$\theta_{\text{Bayes'}} = \mathbb{E}(\theta | \mathcal{X}) = \int \theta p(\theta | \mathcal{X}) d\theta$$

$$p(x | \mathcal{X}) = p(x | \theta_{\text{Bayes'}}) \quad y = g(x | \theta_{\text{Bayes'}})$$



## Bayes' Estimator

- Bayes' estimator – expectation w.r.t. posterior density:

$$\theta_{\text{Bayes'}} = \mathbb{E}(\theta \mid \mathcal{X}) = \int \theta p(\theta \mid \mathcal{X}) d\theta$$

- The reason for taking the expected value is that the best estimate of a random variable is its mean.
  - Let us say  $\theta$  is the variable we want to predict with  $\mathbb{E}(\theta) = \mu$ . For any estimate  $c$  (a constant value), we have the mean squared error (MSE)

$$\begin{aligned}\mathbb{E}[(\theta - c)^2] &= \mathbb{E}[(\theta - \mu + \mu - c)^2] \\ &= \mathbb{E}[(\theta - \mu)^2] + (\mu - c)^2\end{aligned}$$

which is minimum if  $c$  is taken as  $\mu$ .

- Strictly speaking,  $\mathbb{E}(\theta \mid \mathcal{X})$  is the Bayes' estimator under the **squared error loss function**, i.e., a minimum mean square error (MMSE) estimator.
- In the case of a normal density, the mode is the mean.
  - If  $p(\theta \mid \mathcal{X})$  is normal, then  $\theta_{\text{MAP}} = \theta_{\text{Bayes'}}$ .

## Example - I

- Bayesian estimation with known  $\mu_0, \sigma_0$ , and  $\sigma$ :

$$x^t \sim \mathcal{N}(\theta, \sigma^2), \theta \sim \mathcal{N}(\mu_0, \sigma_0^2)$$

- ML estimator:

$$\theta_{\text{ML}} = \frac{1}{N} \sum_{t=1}^N x^t = m \quad (\text{sample mean})$$

- The  $p(\theta | \mathcal{X})$  is normal.

$$p(\theta | \mathcal{X}) \propto p(\mathcal{X} | \theta)p(\theta)$$

$$\propto \frac{1}{(2\pi)^{\frac{N}{2}} \sigma^N} \exp \left[ -\frac{\sum_{t=1}^N (x^t - \theta)^2}{2\sigma^2} \right] \times \frac{1}{(2\pi)^{\frac{1}{2}} \sigma_0} \exp \left[ -\frac{(\theta - \mu_0)^2}{2\sigma_0^2} \right]$$

$$\propto \exp \left[ -\frac{\sum_{t=1}^N (\theta - x^t)^2}{2\sigma^2} \right] \exp \left[ -\frac{(\theta - \mu_0)^2}{2\sigma_0^2} \right]$$

$$\text{Univariate Data} \propto \exp \left[ -\frac{\sigma_0^2 \sum_{t=1}^N (\theta - x^t)^2 + \sigma^2 (\theta - \mu_0)^2}{2\sigma_0^2 \sigma^2} \right] \propto \exp \left[ -\frac{\left( \theta - \frac{\sigma_0^2 \sum_{i=1}^N x_i + \sigma^2 \mu_0}{N\sigma_0^2 + \sigma^2} \right)^2}{2C} \right]$$

## Example - II

- MAP and Bayes' estimator:

$$\begin{aligned}\theta_{\text{MAP}} = \theta_{\text{Bayes'}} &= \mathbb{E}(\theta \mid \mathcal{X}) = \frac{N/\sigma^2}{N/\sigma^2 + 1/\sigma_0^2} m + \frac{1/\sigma_0^2}{N/\sigma^2 + 1/\sigma_0^2} \mu_0 \\ &= \frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2} m + \frac{\sigma^2}{N\sigma_0^2 + \sigma^2} \mu_0\end{aligned}$$

weighted average of the sample mean  $m$  and the prior mean  $\mu_0$ , with weights being inversely proportional to their variances

- As  $N$  increases, the Bayes' estimator gets closer to the sample average, using more the information provided by the sample.
- When  $\sigma_0^2$  is small (we have little prior uncertainty regarding the correct value of  $\theta$ ), or when  $N$  is small, our prior guess  $\mu_0$  has a higher effect.

## Bayesian Estimation

- ▶ Both MAP and Bayes' estimators reduce the whole posterior density to a **single point** and lose information unless the posterior is unimodal and makes a narrow peak around these points.
- ▶ With computation getting cheaper, we can use a Monte Carlo approach that generates samples from the posterior density.
- ▶ There also are many approximation methods one can use to evaluate the full integral in the full Bayesian estimation approach.

# Outline

Introduction

## Univariate Data

Maximum Likelihood Estimation

Bayesian Estimation

**Parametric Classification**

Regression

Model Selection

## Multivariate Data

Parameter Estimation

Multivariate Normal Distribution

Parametric Classification

Discrete Features

Regression

## Classification with Discriminant Functions

- ▶ In Bayes' decision rule for classification, the discriminant function for of class  $C_i$  is

$$p(x | C_i)P(C_i) \text{ or } \log [p(x | C_i)P(C_i)]$$

- ▶ Assume Gaussian density for each class:

$$p(x | C_i) = \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp \left[ -\frac{(x - \mu_i)^2}{2\sigma_i^2} \right]$$

- ▶ Discriminant functions:

$$\begin{aligned} g_i(x) &= \log [p(x | C_i)P(C_i)] \\ &= \log p(x | C_i) + \log P(C_i) \\ &= -\frac{1}{2} \log 2\pi - \log \sigma_i - \frac{(x - \mu_i)^2}{2\sigma_i^2} + \log P(C_i) \end{aligned}$$

## Discriminant Functions based on ML Estimates

- ▶ Samples  $\mathcal{X} = \{(x^t, \mathbf{r}^t)\}_{t=1}^N$  where

$$x^t \in \mathbb{R} \quad \mathbf{r} \in \{0, 1\}^K \text{ with } r_i^t = \begin{cases} 1 & \text{if } x^t \text{ belongs to } C_i \\ 0 & \text{if } x^t \text{ belongs to } C_k, k \neq i \end{cases}$$

- ▶ ML estimates:

$$\hat{P}(C_i) = \frac{1}{N} \sum_{t=1}^N r_i^t \text{ (generalized Bernoulli density)}$$

$$m_i = \frac{\sum_{t=1}^N x^t r_i^t}{\sum_{t=1}^N r_i^t} \quad s_i^2 = \frac{\sum_{t=1}^N (x^t - m_i)^2 r_i^t}{\sum_{t=1}^N r_i^t} \text{ (normal density)}$$

- ▶ Discriminant functions (with constant term  $-\frac{1}{2} \log 2\pi$  dropped):

$$g_i(x) = -\log s_i - \frac{(x - m_i)^2}{2s_i^2} + \log \hat{P}(C_i)$$

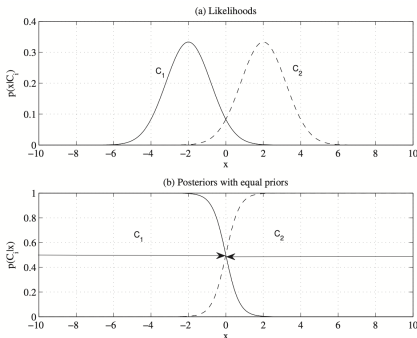
- ▶ Decision threshold between two adjacent classes  $C_i$  and  $C_j$  ( $i \neq j$ ):  $g_i(x) = g_j(x)$

## Case I: Equal Priors and Variances

- ▶ Simplified discriminant functions:  $g_i(x) = -(x - m_i)^2$ 
  - the quadratic term  $x^2$  can be reduced since it is common in all discriminants, leading to a linear discriminant
- ▶ Classification rule (a.k.a. **nearest centroid (mean) classifier**):

Choose  $C_i$  if  $|x - m_i| = \min_k |x - m_k|$

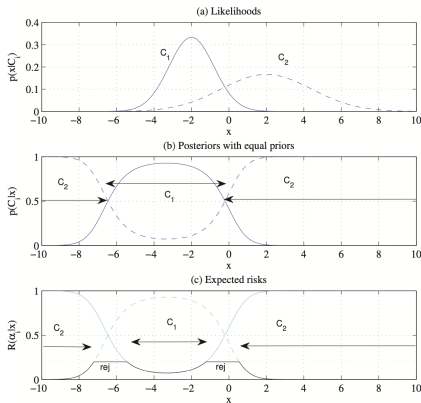
- ▶ Likelihood functions and posterior densities:





## Case II: Equal Priors but Different Variances

- ▶ Simplified discriminant functions:  $g_i(x) = -\log s_i - \frac{(x-m_i)^2}{2s_i^2}$
- ▶ Likelihood functions, posterior densities, and expected risks (for reject with  $\lambda = 0.2$ ):



## Case III: Different Priors

- ▶ Discriminant functions:

$$g_i(x) = -\log s_i - \frac{(x - m_i)^2}{2s_i^2} + \log \hat{P}(C_i)$$

- a quadratic concave function in  $x$
- ▶ If the priors are different, this has the effect of moving the threshold of decision toward the mean of the less likely class.

## Discriminant Functions Based on Bayesian Estimates

- ▶ We have used the ML estimators for the parameters.
- ▶ If we have some prior information about them, for example, for the means, we can use a Bayesian estimate of  $p(x \mid C_i)$  with prior on  $\mu_i$ .

## Remarks

### Gaussianity

- ▶ When  $x$  is continuous,  $p(x | C_i)$  cannot be immediately assumed as Gaussian densities. The classification algorithm—that is, the threshold points—will be wrong if the densities are not Gaussian.

### Generative vs. Discriminative Classifiers

- ▶ This likelihood-based approach is a **generative** approach to classification where we use data to estimate the densities separately, calculate posterior densities using Bayes' rule, and then get the discriminant.
- ▶ We will discuss the **discriminative** approach (discriminant-based approach) in later lectures where we bypass the estimation of densities and directly estimate the discriminants.

# Outline

Introduction

## Univariate Data

Maximum Likelihood Estimation

Bayesian Estimation

Parametric Classification

Regression

Model Selection

## Multivariate Data

Parameter Estimation

Multivariate Normal Distribution

Parametric Classification

Discrete Features

Regression

## Additive Parametric Model

- Functional relationship between numeric output (**dependent variable**) and input (**independent variable**) expressed in an **additive form**:

$$r = f(x) + \epsilon$$

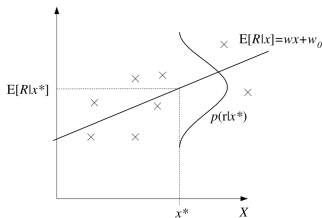
where  $f(x)$  is an unknown but **deterministic** regression function and  $\epsilon$  is **random** noise independent of the input.

- **Parametric modeling**:
  - $f(x) \approx$  estimator  $g(x | \theta)$  (defined up to a set of parameters  $\theta$ )
  - $\epsilon \sim \mathcal{N}(0, \sigma^2)$

## Maximum Likelihood Estimation - I

- Conditional probability of output given input:

$$p(r | x, \theta) \sim \mathcal{N}(g(x | \theta), \sigma^2)$$



- We again use maximum likelihood to learn the parameters  $\theta$
- Joint density:

$$p(r, x | \theta) = p(r | x, \theta)p(x | \theta)$$

where  $p(x)$  is the input density.

## Maximum Likelihood Estimation - II

- Log likelihood of  $\theta$  given an i.i.d. sample  $\mathcal{X} = \{(x^t, r^t)\}_{t=1}^N$ :

$$\begin{aligned}\mathcal{L}(\theta \mid \mathcal{X}) &= \log \prod_{t=1}^N p(x^t, r^t \mid \theta) \\ &= \sum_{t=1}^N \log p(r^t \mid x^t, \theta) + \sum_{t=1}^N \log p(x^t \mid \theta) \\ &= \sum_{t=1}^N \log p(r^t \mid x^t, \theta) + \sum_{t=1}^N \log p(x^t) \\ &= \sum_{t=1}^N \log p(r^t \mid x^t, \theta) + \text{const.}\end{aligned}$$

- We can ignore the second term since it does not depend on our estimator



## Maximum Likelihood Estimation - III

- We have

$$\begin{aligned}\mathcal{L}(\theta \mid \mathcal{X}) &= \sum_{t=1}^N \log \frac{1}{\sqrt{2\pi}\sigma} \exp \left[ -\frac{[r^t - g(x^t \mid \theta)]^2}{2\sigma^2} \right] + \text{const.} \\ &= -\frac{1}{2\sigma^2} \sum_{t=1}^N [r^t - g(x^t \mid \theta)]^2 + \text{const.}\end{aligned}$$

- Maximizing  $\mathcal{L}(\theta \mid \mathcal{X})$  is equivalent to minimizing the following error function:

$$E(\theta \mid \mathcal{X}) = \frac{1}{2} \sum_{t=1}^N [r^t - g(x^t \mid \theta)]^2$$

which is the most frequently used squared error function

- So the ML estimate of  $\theta$  under Gaussian error assumption is also the least squares estimate.

## Linear Regression - I

- ▶ Linear regression function:

$$g(x^t \mid w_0, w_1) = w_1 x^t + w_0$$

- ▶ In a “linear model”,  $g$  is a linear function of the parameter  $\theta$  but not necessarily of the data  $x^t$ .
- ▶ Error function:

$$E(w_0, w_1 \mid \mathcal{X}) = \frac{1}{2} \sum_{t=1}^N (r^t - w_1 x^t - w_0)^2$$

- ▶ Setting the derivatives of  $E(w_0, w_1 \mid \mathcal{X})$  w.r.t.  $w_0, w_1$  to 0 gives two equations

$$\begin{aligned} \sum_{t=1}^N r^t &= Nw_0 + \sum_{t=1}^N x^t w_1 \\ \sum_{t=1}^N r^t x^t &= \sum_{t=1}^N x^t w_0 + \sum_{t=1}^N (x^t)^2 w_1 \end{aligned}$$

## Linear Regression - II

- ▶ Linear system of equations in vector-matrix form:

$$\mathbf{A}\mathbf{w} = \mathbf{y}$$

where

$$\mathbf{A} = \begin{bmatrix} N & \sum_t x^t \\ \sum_t x^t & \sum_t (x^t)^2 \end{bmatrix}$$

$$\mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$$

$$\mathbf{y} = \begin{bmatrix} \sum_t r^t \\ \sum_t r^t x^t \end{bmatrix}$$

- ▶ Least squares estimate (assuming that  $\mathbf{A}$  is invertible):

$$\hat{\mathbf{w}} = \mathbf{A}^{-1}\mathbf{y}$$

- ▶ Linear models can be solved algebraically in closed form, while many non-linear models need to be solved numerically.

## Polynomial Regression

- Polynomial regression function is a polynomial in  $x$  of order  $k$ :

$$g(x^t \mid w_0, w_1, \dots, w_k) = w_k(x^t)^k + \dots + w_2(x^t)^2 + w_1x^t + w_0$$

- The model is still linear with respect to the parameters.
  - As before, in the least squares estimate, taking the derivatives, we get  $k + 1$  equations in  $k + 1$  unknowns.
- An alternative approach for the least squares estimate (assuming that  $\mathbf{D}^T \mathbf{D}$  is invertible):

$$\hat{\mathbf{w}} = (\underbrace{\mathbf{D}^T \mathbf{D}}_{\mathbf{A}})^{-1} \underbrace{\mathbf{D}^T \mathbf{r}}_{\mathbf{y}}$$

where

$$\mathbf{D} = \begin{bmatrix} 1 & x^1 & (x^1)^2 & \dots & (x^1)^k \\ 1 & x^2 & (x^2)^2 & \dots & (x^2)^k \\ \vdots & & & & \vdots \\ 1 & x^N & (x^N)^2 & \dots & (x^N)^k \end{bmatrix} \quad \mathbf{r} = \begin{bmatrix} r^1 \\ r^2 \\ \vdots \\ r^N \end{bmatrix}$$

## Other Error Measures

- Squared error:

$$E(\theta | \mathcal{X}) = \frac{1}{2} \sum_t [r^t - g(x^t | \theta)]^2$$

- Relative squared error (RSE):

$$E(\theta | \mathcal{X}) = \frac{\sum_t [r^t - g(x^t | \theta)]^2}{\sum_t (r^t - \bar{r})^2}$$

- Absolute error:

$$E(\theta | \mathcal{X}) = \sum_t |r^t - g(x^t | \theta)|$$

- $\epsilon$ -insensitive error (for support vector machines to be covered later):

$$E(\theta | \mathcal{X}) = \sum_t \mathbf{1}(|r^t - g(x^t | \theta)| > \epsilon)(|r^t - g(x^t | \theta)| - \epsilon)$$

# Outline

Introduction

## Univariate Data

Maximum Likelihood Estimation

Bayesian Estimation

Parametric Classification

Regression

## Model Selection

## Multivariate Data

Parameter Estimation

Multivariate Normal Distribution

Parametric Classification

Discrete Features

Regression

## Tuning Model Complexity

- ▶ Remember that for best generalization, we should adjust the complexity of our learner model to the complexity of the data.
- ▶ In polynomial regression, the complexity parameter is the order of the fitted polynomial.
- ▶ Therefore we need to find a way to choose the best order that minimizes the generalization error, that is, tune the complexity of the model to best fit the complexity of the function inherent in the data.

## Bias and Variance - I

- ▶ A sample  $\mathcal{X} = \{(x^t, r^t)\}_{t=1}^N$  is drawn from some unknown joint pdf  $p(x, r)$ .
- ▶ Using this sample  $\mathcal{X}$ , we construct our estimate  $g(\cdot)$ .
- ▶ Expected squared error of estimator  $g(\cdot)$  at  $x$  over  $p(x, r)$ :

$$\mathbb{E}[(r - g(x))^2 | x] = \underbrace{\mathbb{E}[(r - \mathbb{E}[r | x])^2 | x]}_{\text{noise}} + \underbrace{(\mathbb{E}[r | x] - g(x))^2}_{\text{squared error}}$$

Only the second term depends on the estimator  $g(\cdot)$ .

- The first term is the variance of  $r$  given  $x$ ; it does not depend on  $g(\cdot)$  or  $\mathcal{X}$ . It is the variance of noise added,  $\sigma^2$ . This is the part of error that can never be removed, no matter what estimator we use.
- The second term quantifies how much  $g(x)$  deviates from the regression function,  $f(x) = \mathbb{E}[r | x]$ . It depends on the estimator  $g(\cdot)$  and the training set  $\mathcal{X}$ .



## Bias and Variance - II

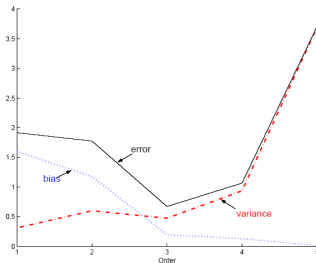
- ▶ It may be the case that for one sample,  $g(x)$  may be a very good fit; and for some other sample, it may make a bad fit.
- ▶ We average over different samples  $\mathcal{X}$  to quantify how well an estimator  $g(\cdot)$  is, giving the expected value (average over samples  $\mathcal{X}$ , all of size  $N$  and drawn from the same joint density  $p(x, r)$ ):

$$\begin{aligned} & \mathbb{E}_{\mathcal{X}} \left[ (\mathbb{E}[r | x] - g(x))^2 | x \right] \\ &= \underbrace{\left( \mathbb{E}[r | x] - \mathbb{E}_{\mathcal{X}}[g(x)] \right)^2}_{\text{bias}^2} + \underbrace{\mathbb{E}_{\mathcal{X}} \left[ (g(x) - \mathbb{E}_{\mathcal{X}}[g(x)])^2 \right]}_{\text{variance}} \end{aligned}$$

- bias measures how much  $g(x)$  is wrong disregarding the effect of varying samples
- variance measures how much  $g(x)$  fluctuate around the expected value,  $\mathbb{E}_{\mathcal{X}}[g(x)]$ , as the sample varies
- we want both to be small

## Bias/Variance Dilemma

- ▶ As model complexity increases (e.g., order of polynomial increases):
  - bias decreases (better fit to data)
  - variance increases (fit varies more with data)
  - called the **bias/variance dilemma** which is true for any machine learning system
- ▶ If there is bias, this indicates that our model class does not contain the solution; this is underfitting.
- ▶ If there is variance, the model class is too general and also learns the noise; this is overfitting.
- ▶ **Optimal model**: best **tradeoff** between bias and variance.



- ▶ **Model selection**: search for optimal or suboptimal model.

## Common Model Selection Methods - I

- ▶ **Cross-validation**: Measure generalization accuracy by testing on data unused during model training.
- ▶ **Regularization**: Penalize complex models by minimizing an **augmented error function**:

$$E' = \text{error on data} + \lambda \cdot \text{model complexity}$$

The **regularization parameter**  $\lambda$ , which can be determined by cross-validation, controls the weight of penalty.

## Common Model Selection Methods - II

- ▶ **Bayesian model selection** is used when there exists prior knowledge about the appropriate class of approximating functions, represented as a **prior distribution** over models,  $p(\text{model})$ .
- ▶ **Posterior distribution** over models given the data:

$$p(\text{model} \mid \text{data}) = \frac{p(\text{data} \mid \text{model})p(\text{model})}{p(\text{data})}$$

- ▶ From the posterior distribution, we may:
  - choose the model with the **highest** posterior probability, or
  - choose multiple models with **high** posterior probabilities, or
  - use **all** models weighted by their posterior probabilities.
- ▶ Taking the log of the above equation

$$\log p(\text{model} \mid \text{data}) = \log p(\text{data} \mid \text{model}) + \log p(\text{model}) - \text{const.}$$

- ▶ **Regularization** may be seen as a Bayesian approach in which the “prior” favors simpler models.

## Common Model Selection Methods - III

- ▶ Cross-validation is different from the other methods for model selection in that it makes no prior assumption about the model or parameters.
- ▶ If there is a large enough validation dataset, cross-validation is the best approach.
- ▶ When the data sample is small, the other models become useful.