

ECE 449/590 – OOP and Machine Learning

Lecture 15 Machine Learning Basics II

Professor Jia Wang
Department of Electrical and Computer Engineering
Illinois Institute of Technology

October 17, 2022

Outline

Estimators

Maximum Likelihood Estimation

Supervised Learning Algorithms

Unsupervised Learning Algorithms

Reading Assignment

- ▶ This lecture: Deep Learning 5
- ▶ Next lecture: Deep Learning 6

Outline

Estimators

Maximum Likelihood Estimation

Supervised Learning Algorithms

Unsupervised Learning Algorithms

Point Estimation

- ▶ Input: a set of i.i.d. (independent and identically distributed) data points/examples $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$.
- ▶ Assume the data points are sampled from a parametric distribution with the parameter θ .
 - ▶ E.g. for Gaussian (normal) distribution, $\theta = (\mu, \sigma)$.
 - ▶ The type of the distribution itself, e.g. Gaussian or uniform distribution, is not part of the parameter.
- ▶ A point estimator is any function of the input
 - ▶ $\hat{\theta}_m = g(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)})$
 - ▶ A good estimator provides a close approximation to θ .
- ▶ In practice, the distribution is usually modeled as a function depending on θ and random variables with known distribution.
 - ▶ E.g. $f(x; \theta) = \sigma x + \mu$ generates any Gaussian distribution from x that is sampled from the standard Gaussian distribution.
 - ▶ The point estimator actually solves the function estimation problem in such case.

Bias and Variance

- ▶ Since $\{x^{(1)}, \dots, x^{(m)}\}$ are sampled from a distribution, they should be treated as random variables.
 - ▶ So $\hat{\theta}_m$ is a random variable.
- ▶ For a random variable $(\hat{\theta}_m)$ to approximate a value (θ) well,
 - ▶ The expectation $\mathbb{E}(\hat{\theta}_m)$ should be θ , or the bias $\text{bias}(\hat{\theta}_m) = \mathbb{E}(\hat{\theta}_m) - \theta$ should be 0.
 - ▶ The variance $\text{Var}(\hat{\theta}_m)$ should approach 0.
- ▶ E.g. $\hat{\mu}_m = \frac{1}{m} \sum_{i=1}^m x^{(i)}$ is an estimator of Gaussian mean μ .
 - ▶ $\text{bias}(\hat{\mu}_m) = 0$, $\text{Var}(\hat{\mu}_m) = \frac{\sigma^2}{m}$
 - ▶ A very good estimator for μ but what about estimators for σ ?
- ▶ Machine learning algorithms as estimators need to make trade-off between bias and variance.

Bias and Variance

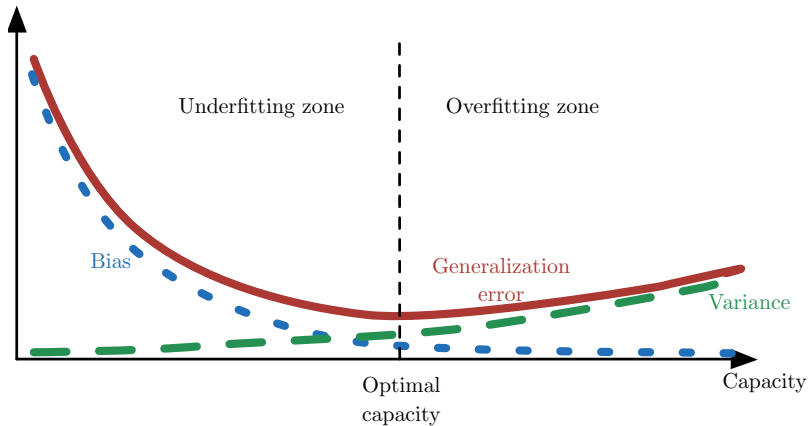


Figure 5.6

Outline

Estimators

Maximum Likelihood Estimation

Supervised Learning Algorithms

Unsupervised Learning Algorithms

Estimator Design via Maximum Likelihood

- ▶ How can we design a good estimator?
 - ▶ Use maximum likelihood as a common principle.
- ▶ Let $\mathbb{X} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ be sampled i.i.d. from p_{data} .
 - ▶ If we further assume no one could know the actual p_{data} , the best we can do is to choose a known parametric distribution $p_{model}(\mathbf{x}; \boldsymbol{\theta})$ and then estimate $\boldsymbol{\theta}$.
- ▶ Maximum likelihood: maximize the probability that you would observe \mathbb{X} assuming they are sampled from $p_{model}(\mathbf{x}; \boldsymbol{\theta})$.

$$\boldsymbol{\theta}_{ML} = \arg \max_{\boldsymbol{\theta}} p_{model}(\mathbb{X}; \boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta}} \prod_{i=1}^m p_{model}(\mathbf{x}^{(i)}; \boldsymbol{\theta})$$

Log-Likelihood and Cross-Entropy

- ▶ It is more convenient to work with log-likelihood on computers.

$$\begin{aligned}\theta_{ML} &= \arg \max_{\theta} \sum_{i=1}^m \log p_{model}(\mathbf{x}^{(i)}; \theta) \\ &= \arg \max_{\theta} \mathbb{E}_{\mathbf{x} \sim p_{data}} \log p_{model}(\mathbf{x}; \theta)\end{aligned}$$

- ▶ As the negative log-likelihood $-\mathbb{E}_{\mathbf{x} \sim p_{data}} \log p_{model}(\mathbf{x})$ is known as the cross-entropy of p_{model} with respect to p_{data} , to maximize likelihood is equivalent to minimize cross-entropy.

Conditional Log-Likelihood and Supervised Learning

- ▶ So far we assume the examples are $\mathbb{X} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ when applying maximum likelihood estimators.
 - ▶ What about the labels $\{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(m)}\}$ for supervised learning?
- ▶ Estimate $\boldsymbol{\theta}$ in a conditional distribution model $p(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})$.
 - ▶ With the model, for a given \mathbf{x} , the output is predicted as $\arg \max_{\mathbf{y}} p(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})$
- ▶ Conditional maximum likelihood: maximize the probability that you would observe the labels given \mathbb{X} .

$$\boldsymbol{\theta}_{ML} = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^m \log p(\mathbf{y}^{(i)}|\mathbf{x}^{(i)}; \boldsymbol{\theta})$$

Example: Maximum Likelihood for Gaussian

- ▶ Supervised learning: scalar labels $\{y^{(1)}, \dots, y^{(m)}\}$ for the examples $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$.
 - ▶ Assume $p(y|\mathbf{x}; \boldsymbol{\theta})$ is of a Gaussian distribution of known variance σ^2 and unknown mean μ that depends on \mathbf{x} and $\boldsymbol{\theta}$.
 - ▶ The prediction is $\hat{y} = \arg \max_y p(y|\mathbf{x}; \boldsymbol{\theta}) = \mu$.
 - ▶ E.g. if $\mu = \boldsymbol{\theta}^\top \mathbf{x}$ then this is linear regression.
- ▶ The log-likelihood is

$$\sum_{i=1}^m \log p(y^{(i)}|\mathbf{x}^{(i)}; \boldsymbol{\theta}) = -m \log \sigma - \frac{m}{2} \log 2\pi - \sum_{i=1}^m \frac{\|\hat{y}^{(i)} - y^{(i)}\|^2}{2\sigma^2}$$

- ▶ To maximize likelihood is equivalent to minimize MSE.
 - ▶ As long as $p(y|\mathbf{x}; \boldsymbol{\theta})$ is Gaussian with known variance but unknown mean.

Outline

Estimators

Maximum Likelihood Estimation

Supervised Learning Algorithms

Unsupervised Learning Algorithms

Probabilistic Supervised Learning

- ▶ Choose a parametric family $p_{model}(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})$ of conditional distributions.
- ▶ Apply maximum likelihood estimation to obtain $\boldsymbol{\theta}$.
- ▶ We have seen the example for \mathbf{y} being continuous and p_{model} being Gaussian.
 - ▶ But for the more general case, it is difficult to specify the conditional pdf (probability density function).
- ▶ For \mathbf{y} being discrete, it is possible to specify the conditional probability for individual cases.
 - ▶ E.g. logistic regression when \mathbf{y} can be either 0 or 1:

$$p_{model}(y = 1|\mathbf{x}; \boldsymbol{\theta}) = \text{sigmoid}(\boldsymbol{\theta}^\top \mathbf{x})$$

where $\text{sigmoid}(a) = \frac{e^a}{e^a + 1}$.

- ▶ When \mathbf{y} may take more values, one may use a vector-valued function of $\boldsymbol{\theta}$ and \mathbf{x} for such purpose.
 - ▶ The challenge is to find the proper form of the function to be able to approximate the unknown p_{data} .

Other Supervised Learning Algorithms

- ▶ Support vector machine: mainly for binary classification
- ▶ k -nearest neighbor: memorizing training set
- ▶ Decision tree: branching on features one by one
- ▶ Has their limitations but may work for specific problems and settings.

Outline

Estimators

Maximum Likelihood Estimation

Supervised Learning Algorithms

Unsupervised Learning Algorithms

Representation and Unsupervised Learning

- ▶ Find the “best” representation for the data set \mathbb{X} .
 - ▶ Simpler or more accessible than \mathbb{X} .
 - ▶ While preserving as much information about \mathbb{X} as possible.
- ▶ Common approaches
 - ▶ Lower-dimensional representations where information is compressed.
 - ▶ Sparse representations where only non-zero elements need to be stored.
 - ▶ Independent representations where a joint pdf is decomposed into products of marginal pdfs.

Principal Components Analysis (PCA)

- ▶ Consider \mathbb{X} as the design matrix $\mathbf{X} \in \mathbb{R}^{m \times n}$.
 - ▶ Assume $\mathbb{E}[\mathbf{x}] = 0$ for \mathbf{x} be an example (a column of \mathbf{X}^\top).
- ▶ Unbiased sample covariance: $\text{Var}[\mathbf{x}] = \frac{1}{m-1} \mathbf{X}^\top \mathbf{X}$
- ▶ Eigendecomposition: $\mathbf{X}^\top \mathbf{X} = \mathbf{W} \mathbf{\Lambda} \mathbf{W}^\top$.
 - ▶ where \mathbf{W} is orthogonal and $\mathbf{\Lambda}$ is diagonal.
- ▶ Let $\mathbf{Z} = \mathbf{X} \mathbf{W}$, then $\text{Var}[\mathbf{z}] = \frac{1}{m-1} \mathbf{\Lambda}$.
- ▶ Now \mathbf{Z} works as a representation of \mathbf{X} while all \mathbf{z} are mutually uncorrelated.

Principal Components Analysis

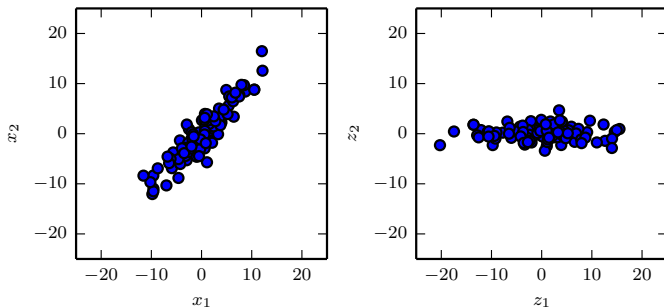


Figure 5.8

k -means Clustering

- ▶ Divide \mathbb{X} into k different clusters.
 - ▶ Each example can be represented as a length- k one hot code.
- ▶ k -means algorithm: each cluster has a centroid and examples in a cluster are closer to its own centroid than to other centroids.
- ▶ When k is not given and thus is a hyperparameter, how to decide if it is good?

Summary

- ▶ Bias and variance matter for estimators.
- ▶ Maximum Likelihood is widely used to design estimators.
- ▶ Simple supervised and unsupervised learning algorithms.