

Statistical Machine Learning

Lecture 05: Bayesian Decision Theory

Kristian Kersting

TU Darmstadt

Summer Term 2020

Today's Objectives

- Make you understand how to do an optimal decision!
- Covered Topics:
 - Bayesian Optimal Decisions
 - Classification from a Bayesian point of view
 - Risk-based Classification

Outline

1. Bayesian Decision Theory

2. Risk Minimization

3. Wrap-Up

Outline

1. Bayesian Decision Theory

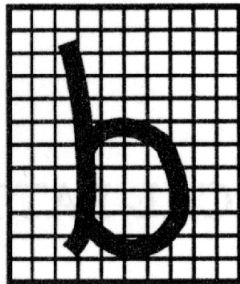
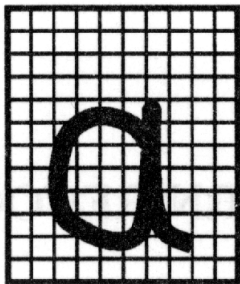
2. Risk Minimization

3. Wrap-Up

Statistical Methods

- **Statistical methods in machine learning** all have in common that they assume that the process that “generates” the data is **governed by the rules of probability**
- The data is understood to be a **set of random samples** from some **underlying probability distribution**
- Today will be all about probabilities. But in future lectures, the use of probability will sometimes be much less explicit
- Nonetheless, **the basic assumption about how the data is generated is always there**, even if you don't see a single probability distribution anywhere

Character Recognition

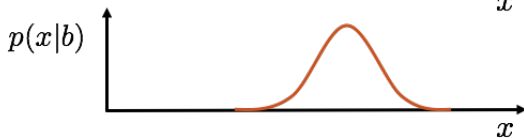
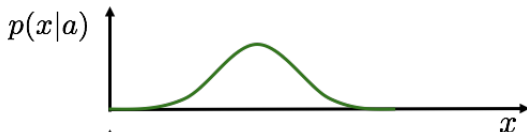


- **Goal:** classify a new letter so that the **probability of a wrong classification is minimized**

Class conditional probabilities

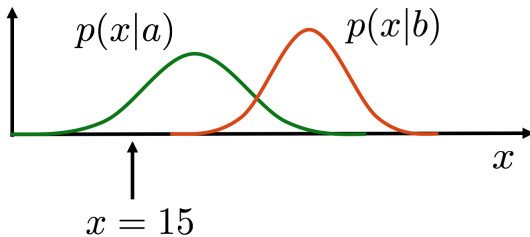
■ Class conditional probabilities

- Probability of making an observation \mathbf{x} knowing that it comes from some class C_k
- Here \mathbf{x} is often a feature vector, which measures/describes properties of the data. E.g.: number of black pixels, height-width ratio, ...



Class conditional probabilities

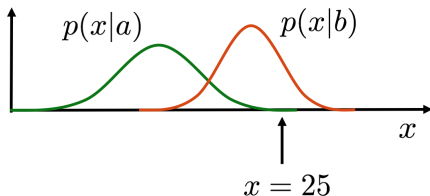
■ Example



- How do we decide which class the data point belongs to?
- Here, we should decide for class **a**

Class conditional probabilities

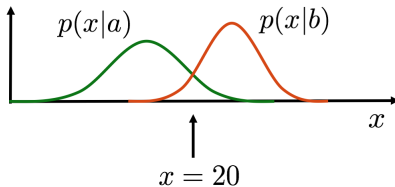
■ Example



- How do we decide which class the data point belongs to?
- Since $p(\mathbf{x}|a)$ is a lot smaller than $p(\mathbf{x}|b)$ we should now decide for class **b**

Class conditional probabilities

■ Example



■ How do we decide which class the data point belongs to?

Class priors

- The *a priori* probability of a data point belonging to a particular class is called the **class prior**
- Example:
 - abaaababaaaabbaaaaa
- What are $p(a)$ and $p(b)$?

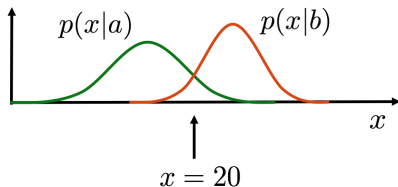
$$C_1 = a \quad p(C_1) = 0.75$$

$$C_2 = b \quad p(C_2) = 0.25$$

$$\sum_k p(C_k) = 1$$

Back to our problem...

■ Example



- How do we decide which class the data point belongs to?
- If $p(a) = 0.75$ and $p(b) = 0.25$, we should decide for class **a**

Bayesian Decision Theory

- **Bayes Theorem** lets us formalize the previous intuitive decision
- We want to find the **a-posteriori probability** (posterior) of the class C_k given the observation (feature) \mathbf{x}

$$p(C_k|\mathbf{x}) = \frac{p(\mathbf{x}|C_k)p(C_k)}{p(\mathbf{x})} = \frac{p(\mathbf{x}|C_k)p(C_k)}{\sum_j p(\mathbf{x}|C_j)p(C_j)}$$

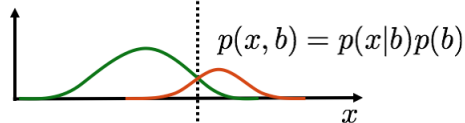
- class prior: $p(C_k)$
- class-conditional probability (likelihood): $p(\mathbf{x}|C_k)$
- class posterior: $p(C_k|\mathbf{x})$
- normalization term: $p(\mathbf{x})$

Bayesian Decision Theory



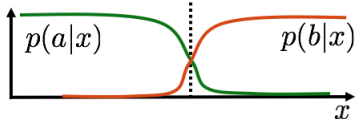
Likelihood

$$p(x, a) = p(x|a)p(a)$$



Likelihood \times Prior

Decision boundary



$$\text{Posterior} = \frac{\text{Likelihood} \times \text{Prior}}{\text{Normalization factor}}$$

Bayesian Decision Theory

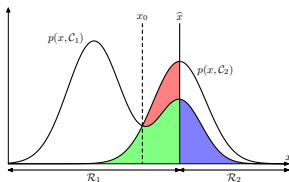
- Why is it called this way?
 - To some extent, because it involves applying Bayes' rule
 - But this is not the whole story...
 - The real reason is that it is **built on so-called Bayesian probabilities**

Bayesian Probabilities

- Probability is not just interpreted as a frequency of a certain event happening
- Rather, it is seen as a **degree of belief** in an outcome
- Only this allows us to assert a **prior belief** in a data point coming from a certain class
- Even though this might seem easy to accept to you now, this interpretation was quite contentious in statistics for a long time

Bayesian Decision Theory

- Goal: **Minimize the misclassification rate** (the probability of classifying wrongly)



$$\begin{aligned}
 p(\text{error}) &= p(x \in R_1, C_2) + p(x \in R_2, C_1) \\
 &= \int_{R_1} p(x, C_2) dx + \int_{R_2} p(x, C_1) dx \\
 &= \int_{R_1} p(x|C_2)p(C_2) dx + \int_{R_2} p(x|C_1)p(C_1) dx
 \end{aligned}$$

Bayesian Decision Theory

- **Decision rule:** decide C_1 if $p(C_1|x) > p(C_2|x)$
- Equivalent to

$$\frac{p(x|C_1)p(C_1)}{p(x)} > \frac{p(x|C_2)p(C_2)}{p(x)}$$
$$p(x|C_1)p(C_1) > p(x|C_2)p(C_2)$$
$$\frac{p(x|C_1)}{p(x|C_2)} > \frac{p(C_2)}{p(C_1)}$$

- A classifier obeying this rule is called a **Bayes Optimal Classifier**

Bayesian Decision Theory

$$\frac{p(x|C_1)}{p(x|C_2)} > \frac{p(C_2)}{p(C_1)}$$

■ Special cases

- If $p(x|C_1) = p(x|C_2)$, then use $p(C_1) > p(C_2)$
- If $p(C_1) = p(C_2)$, then use $p(x|C_1) > p(x|C_2)$

More than two Classes

■ Generalization to more than 2 classes:

- Decide for class k iff it has the highest a-posteriori probability

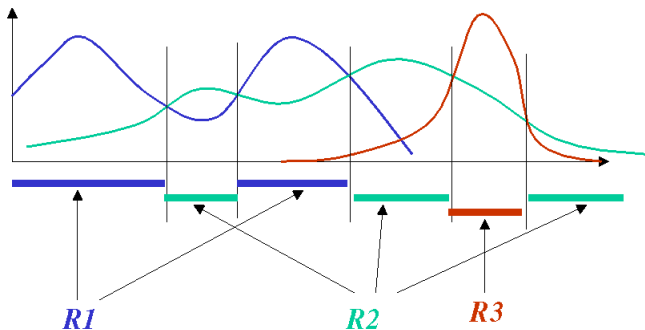
$$p(C_k|x) > p(C_j|x) \quad \forall j \neq k$$

- Equivalent to

$$p(x|C_k)p(C_k) > p(x|C_j)p(C_j) \quad \forall j \neq k$$
$$\frac{p(x|C_k)}{p(x|C_j)} > \frac{p(C_j)}{p(C_k)} \quad \forall j \neq k$$

More than two Classes

- Decision regions: R_1, R_2, R_3, \dots



High Dimensional Features

- So far we have only considered one-dimensional features, i.e., $x \in \mathbb{R}$
- We can use more features and generalize to an arbitrary D -dimensional feature space, i.e., $\mathbf{x} \in \mathbb{R}^D$

- For instance, in the salmon vs. sea-bass classification task

$$\mathbf{x} = [x_1 \quad x_2]^T \in \mathbb{R}^2$$

- Where x_1 is the width, and x_2 is the lightness
- The decision boundary we devised still applies to $\mathbf{x} \in \mathbb{R}^D$. We just need to use **multivariate class-conditional densities** $p(\mathbf{x}|C_k)$

Dummy Classes

- There are also applications, where it may be advantageous to have a **dummy class** denoted “don't know” or “don't care”
 - Also called a **reject option**
- Not a common case though and we will not cover this in this class

Outline

1. Bayesian Decision Theory

2. Risk Minimization

3. Wrap-Up

2. Risk Minimization

- So far, we have tried to **minimize the misclassification rate**
- There are many cases when **not every misclassification is equally bad**
- Smoke detector
 - If there is a fire, we need to be very sure that we classify it as such
 - If there is no fire, it is ok to occasionally have a false alarm
- Medical diagnosis
 - If the patient is sick, we need to be very sure that we report them as sick
 - If they are healthy, it is ok to classify them as sick and order further testing that may help clarifying this up

Loss Functions

- **Key idea:** we have to construct a loss function in a way that expresses what we want to achieve

$$\begin{aligned} \text{loss}(\text{decision} = \text{healthy} | \text{patient} = \text{sick}) &>> \\ \text{loss}(\text{decision} = \text{sick} | \text{patient} = \text{healthy}) \end{aligned}$$

- Possible decisions: α_j
- True classes: C_j
- Loss function: $\lambda(\alpha_j | C_j)$
- Expected loss of making a decision α_j

$$R(\alpha_j | X) = \mathbb{E}_{C_k \sim p(C_k | X)} [\lambda(\alpha_j | C_k)] = \sum_j \lambda(\alpha_j | C_j) p(C_j | X)$$

Risk Minimization

- The expected loss of a decision is also called the **risk of making a decision**
- **Instead of minimizing the Misclassification rate**

$$\begin{aligned} p(\text{error}) &= p(x \in R_1, C_2) + p(x \in R_2, C_1) \\ &= \int_{R_1} p(x, C_2) dx + \int_{R_2} p(x, C_1) dx \\ &= \int_{R_1} p(x|C_2) p(C_2) dx + \int_{R_2} p(x|C_1) p(C_1) dx \end{aligned}$$

- **We minimize the Overall Risk**

$$R(\alpha_i|x) = \mathbb{E}_{C_k \sim p(C_k|x)} [\lambda(\alpha_i|C_k)] = \sum_j \lambda(\alpha_i|C_j) p(C_j|x)$$

Risk Minimization

- 2 classes: C_1, C_2
- 2 decisions: α_1, α_2
- Loss function: $\lambda(\alpha_i|C_j) = \lambda_{ij}$
- Risk of both decisions

$$R(\alpha_1|x) = \lambda_{11}p(C_1|x) + \lambda_{12}p(C_2|x)$$

$$R(\alpha_2|x) = \lambda_{21}p(C_1|x) + \lambda_{22}p(C_2|x)$$

- Goal: Create a decision rule so that overall risk is minimized
 - Decide α_1 if $R(\alpha_2|x) > R(\alpha_1|x)$

Risk Minimization

$$\begin{aligned}
 R(\alpha_2|x) &> R(\alpha_1|x) \\
 \lambda_{21}p(C_1|x) + \lambda_{22}p(C_2|x) &> \lambda_{11}p(C_1|x) + \lambda_{12}p(C_2|x) \\
 (\lambda_{21} - \lambda_{11})p(C_1|x) &> (\lambda_{12} - \lambda_{22})p(C_2|x)
 \end{aligned}$$

$$\begin{aligned}
 \frac{\lambda_{21} - \lambda_{11}}{\lambda_{12} - \lambda_{22}} &> \frac{p(C_2|x)}{p(C_1|x)} = \frac{p(x|C_2)p(C_2)}{p(x|C_1)p(C_1)} \\
 \frac{p(x|C_1)}{p(x|C_2)} &> \frac{\lambda_{12} - \lambda_{22}}{\lambda_{21} - \lambda_{11}} \frac{p(C_2)}{p(C_1)}
 \end{aligned}$$

- It is reasonable to assume that the loss of a correct decision is smaller than that of a wrong decision: $\lambda_{ij} > \lambda_{ji} \quad \forall j \neq i$

Risk Minimization 0-1 Loss

$$\frac{p(x|C_1)}{p(x|C_2)} > \frac{\lambda_{12} - \lambda_{22} p(C_2)}{\lambda_{21} - \lambda_{11} p(C_1)}$$

- Decide α_1 if

$$\lambda(\alpha_i|C_j) = \begin{cases} 0 & i = j \\ 1 & i \neq j \end{cases}$$

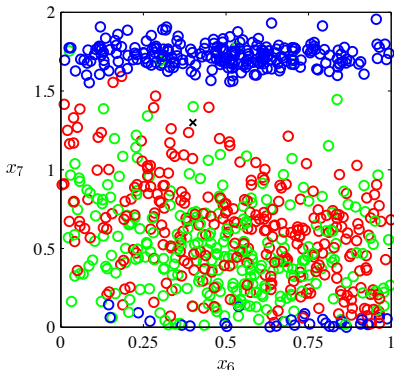
$$\frac{p(x|C_1)}{p(x|C_2)} > \frac{p(C_2)}{p(C_1)}$$

- The **0-1 loss** leads to the **same decision rule** that minimized the **misclassification rate**

Bayesian Decision Theory

- Are we done with classification?
 - We have decision rules for simple and general loss functions
 - Even “Bayes optimal”
 - We can deal with 2 or more classes
 - We can deal with high dimensional feature vectors
 - We can incorporate prior knowledge on the class distribution
- What are we going to do the rest of the semester? Where is the catch?
- Where do we get the probability distributions from?

Training Data



- How do we get the probability distributions from this so that we can classify with them?

Outline

1. Bayesian Decision Theory

2. Risk Minimization

3. Wrap-Up

3. Wrap-Up

You know now:

- What class conditional probabilities, class priors and class posteriors are
- What Bayesian Decision Theory is
- How to use Bayes Theorem for classification
- What misclassification rate is
- What a Bayes Optimal Classifier is
- How to generalize decision to more than 2 classes
- What risk is, and how it relates to misclassification

Self-Test Questions

- How can we decide on classifying a query based on simple and general loss functions?
- What does “Bayes optimal” mean?
- How to deal with 2 or more classes?
- How to deal with high dimensional feature vectors?
- How to incorporate prior knowledge on the class distribution?
- What are the equations for misclassification rate and risk

Homework

- Reading Assignment for next lecture
 - Bishop ch. 2 (Probability Distributions), 9 (Mixture Models and EM)