# Bayesian Decision Theory

Prof. Ziping Zhao

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

CS182: Introduction to Machine Learning (Fall 2021)
http://cs182.sist.shanghaitech.edu.cn

# Outline

# Outline

# Coin Tossing Example

▶ Outcome of tossing a coin $\in$ {head, tail}
▶ Random variable $X$:
$$X = \begin{cases} 1 & \text{if outcome is head} \\ 0 & \text{if outcome is tail} \end{cases}$$

▶ $X$ is Bernoulli-distributed:

$$P(X) = p_0^X (1 - p_0)^{1-X}$$

where the parameter $p_0$ is the probability that the outcome is head, i.e., $p_0 = P(X = 1)$.

# Estimation and Prediction

▶ Estimation of parameter $p_0$ from sample $\mathcal{X} = \{x^{(\ell)}\}_{\ell=1}^{N}$:

$$\hat{p}_0 = \frac{\#\text{heads}}{\#\text{tosses}}$$
$$= \frac{\sum_{\ell=1}^{N} x^{(\ell)}}{N}$$

▶ Prediction of outcome of next toss:

$$\text{Predicted outcome} = \begin{cases} \text{head} & \text{if } p_0 > 1/2 \\ \text{tail} & \text{otherwise} \end{cases}$$

by choosing the more probable outcome, which minimizes the probability of error ($= 1 -$ probability of our choice for the predicted outcome).

# Outline

# Classification as Bayesian Decision

► Credit scoring example:
  – Inputs: income and savings, or $\mathbf{x} = (x_1, x_2)^T$
  – Output: risk $\in$ {low,high}, or $C \in \{0, 1\}$ (a Bernoulli random variable)

► Prediction:

$$\text{Choose} \begin{cases} C = 1 & \text{if } P(C = 1 \mid \mathbf{x}) > 0.5 \\ C = 0 & \text{otherwise} \end{cases}$$

or equivalently

$$\text{Choose} \begin{cases} C = 1 & \text{if } P(C = 1 \mid \mathbf{x}) > P(C = 0 \mid \mathbf{x}) \\ C = 0 & \text{otherwise} \end{cases}$$

► Probability of error:

$$1 - \max(P(C = 1 \mid \mathbf{x}), P(C = 0 \mid \mathbf{x})) = \min(P(C = 1 \mid \mathbf{x}), P(C = 0 \mid \mathbf{x}))$$

► Similar to coin tossing except that $C$ is conditioned on two observable variables $\mathbf{x}$

# Bayes' Rule

▶ Bayes' rule:

$$\text{Posterior } P(C \mid \mathbf{x}) = \frac{\text{likehihood} \times \text{prior}}{\text{evidence}} = \frac{p(\mathbf{x} \mid C)P(C)}{p(\mathbf{x})}$$

▶ prior probability: knowledge we have as to $C$ before looking at the observables $\mathbf{x}$
▶ class likelihood: derived from data
▶ evidence: the marginal probability that an observation $\mathbf{x}$ is seen
▶ Some useful properties to note:
  – $P(C = 0) + P(C = 1) = 1$
  – $p(\mathbf{x}) = p(\mathbf{x} \mid C = 1)P(C = 1) + p(\mathbf{x} \mid C = 0)P(C = 0)$
  – $P(C = 0 \mid \mathbf{x}) + P(C = 1 \mid \mathbf{x}) = 1$
▶ we will discuss the estimation of $p(C)$ and $p(\mathbf{x}|C)$ from training samples in later lectures

## Bayes' Rule for $K > 2$ Classes

▶ Bayes' rule for general case ($K$ mutually exclusive and exhaustive classes):

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

▶ Optimal decision rule for Bayes' classifier:

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

# Outline

# Losses and Risks

▶ In general different decisions or actions may not be equally good or costly.

▶ Action $\alpha_i$: decision to assign the input **x** to class $C_i$

▶ Loss $\lambda_{ik}$: loss incurred for taking action $\alpha_i$ when the actual state is $C_k$

▶ Expected risk/loss or conditional risk for taking action $\alpha_i$ given input **x**:

$$R(\alpha_i \mid \mathbf{x}) = \sum_{k=1}^{K} \lambda_{ik} P(C_k \mid \mathbf{x})$$

▶ Optimal decision rule with minimum expected risk:

$$\text{Choose } \alpha_i \text{ if } R(\alpha_i \mid \mathbf{x}) = \min_k R(\alpha_k \mid \mathbf{x})$$

## 0-1 Loss Function

▶ All correct decisions have zero loss and all errors have unit cost (i.e., are equally costly):

$$\lambda_{ik} = \begin{cases} 0 & \text{if } i = k \\ 1 & \text{if } i \neq k \end{cases}$$

▶ Expected risk:

$$\begin{aligned} R(\alpha_i \mid \mathbf{x}) &= \sum_{k=1}^{K} \lambda_{ik} P(C_k \mid \mathbf{x}) \\ &= \sum_{k \neq i} P(C_k \mid \mathbf{x}) \\ &= 1 - P(C_i \mid \mathbf{x}) \end{aligned}$$

▶ Optimal decision rule with minimum expected risk (or, equivalently, highest posterior probability):

$$\text{Choose } \alpha_i \text{ if } P(C_i \mid \mathbf{x}) = \max_k P(C_k \mid \mathbf{x})$$

# Reject Option

▶ If the certainty of a decision is low but misclassification has very high cost, the action of reject or doubt ($\alpha_{K+1}$) may be more desirable.

▶ A possible loss function:

$$\lambda_{ik} = \begin{cases} 0 & \text{if } i = k \\ \lambda & \text{if } i = K + 1 \\ 1 & \text{otherwise} \end{cases}$$

where $0 < \lambda < 1$ is the loss incurred for choosing the action of reject.

▶ Expected risk:

$$R(\alpha_i \mid \mathbf{x}) = \begin{cases} \sum_{k=1}^{K} \lambda P(C_k \mid \mathbf{x}) = \lambda & \text{if } i = K + 1 \\ \sum_{k \neq i} P(C_k \mid \mathbf{x}) = 1 - P(C_i \mid \mathbf{x}) & \text{if } i \in \{1, \dots, K\} \end{cases}$$

## Reject Option (2)

▶ Optimal decision rule:

$$\begin{cases} \text{Choose } C_i & \text{if } R(\alpha_i \mid \mathbf{x}) = \min_{1 \leq k \leq K} R(\alpha_k \mid \mathbf{x}) < R(\alpha_{K+1} \mid \mathbf{x}) \\ \text{Reject} & \text{otherwise} \end{cases}$$

▶ Equivalent form of optimal decision rule:

$$\begin{cases} \text{Choose } C_i & \text{if } P(C_i \mid \mathbf{x}) = \max_{1 \leq k \leq K} P(C_k \mid \mathbf{x}) > 1 - \lambda \\ \text{Reject} & \text{otherwise} \end{cases}$$

▶ This approach is meaningful only if $0 < \lambda < 1$:
  – If $\lambda = 0$, we always reject (a reject is as good as a correct classification).
  – If $\lambda \geq 1$, we never reject (a reject is at least as costly as, or costlier than, a misclassification error).

# Outline

## Discriminant Functions

▶ One way of specifying a classifier for classification is through a set of discriminant functions, $g_i(\mathbf{x})$, $i = 1, \ldots, K$.

▶ Classification rule:

$$\text{Choose } C_i \text{ if } g_i(\mathbf{x}) = \max_k g_k(\mathbf{x})$$

▶ Some ways of defining the discriminant functions:
- $g_i(\mathbf{x}) = -R(\alpha_i \mid \mathbf{x})$ (generally for Bayes' classifier)
- $g_i(\mathbf{x}) = P(C_i \mid \mathbf{x})$
- $g_i(\mathbf{x}) = p(\mathbf{x} \mid C_i)P(C_i)$

▶ For the two-class case, it suffices to use just one discriminant function:

$$g(\mathbf{x}) = g_1(\mathbf{x}) - g_2(\mathbf{x})$$

with the following classification rule:

$$\text{Choose } \begin{cases} C_1 & \text{if } g(\mathbf{x}) > 0 \\ C_2 & \text{otherwise} \end{cases}$$

# Decision Regions

▶ The feature space is divided into $K$ decision regions $\mathcal{R}_1, \ldots, \mathcal{R}_K$, where

$$\mathcal{R}_i = \left\{ \mathbf{x} \mid g_i(\mathbf{x}) = \max_k g_k(\mathbf{x}) \right\}$$

▶ The decision region corresponding to a class may consist of noncontiguous subregions.

▶ The decision regions are separated by decision boundaries (a.k.a. decision surfaces) where ties occur among the discriminant functions with the largest values.