

CSC 311: Introduction to Machine Learning

Rahul G. Krishnan & Amanjit Singh Kainth

University of Toronto, Fall 2024

Outline

Tutorial Outline

1. Time Management Advice
2. Breaking down an ML
3. Nearest Neighbors Algorithm Review
4. Probability Review

Time Management

Time Management Advice

As with any rewarding course, this one can be challenging and time-consuming! Some tips to increase efficiency + reduce stress:

- Required readings: Skim before lecture and review after.
- Tutorial: Attempt to go through exercises yourself at least once. If you do it before, you can ask questions during tutorial. If you do it after, you'll cement what you learned.
- Homeworks / Assignments: Start as soon as they are released! Try to determine what parts of lecture + tutorial are relevant.
- Office Hours and Piazza: Ask any questions you didn't get to during lecture or tutorial.

The Machine Learning Problem

Building blocks for a machine learning problem:

- Stages
- Data
- Hyperparameters

- **Learning:** Extract information from data to make predictions.
- **Evaluation:** Check how well the algorithm/model makes predictions.

Data

For supervised learning problems, we often have data of the form:
 n input data samples, each with d features:

$$\begin{bmatrix} x_{11} & x_{12} & x_{13} & \dots & x_{1n} \\ x_{21} & x_{22} & x_{23} & \dots & x_{2n} \\ \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots \\ x_{d1} & x_{d2} & x_{d3} & \dots & x_{dn} \end{bmatrix}$$

n targets, corresponding to each input sample:

$$[y_1 \quad y_2 \quad y_3 \quad \dots \quad y_n]$$

- **Training:** Used at the learning stage to extract information about the data that is relevant to the predictive task and potentially transfer that knowledge from the data such that it can be accessed to make predictions later. (May be used later as well.)
- **Validation:** Used to select one out of a few possible algorithms or learned models by mimicking test time behavior. This serves as a proxy for measuring overfitting. (Hidden during learning.)
- **Test:** Used to evaluate algorithms' performance. (Unavailable to learner before evaluation stage.)

Hyperparameters

We would like to design an algorithm or learn parameters of the model based on the training data.

However, there are some (hyper)parameters that we, the designers of the algorithm, must determine.

- Knobs that we tune to find a right setting for the algorithm.
- Use validation data to choose the setting of a knob.

Hyperparameters

Examples:

- Learning Rate: size of updates made to parameters
- Batch Size: amount of the data used at every step of learning
- k: number of Nearest Neighbors

Classification

Given some data, we want to assign it to meaningful categories by learning the patterns in the training data.

How do we store information about the learned patterns?

- Option 1: We don't! Like the NN algorithm, we can just look at the entire training data.

This is a **Non-Parametric** classifier.

- Option 2: We create a model with some parameters. During the learning stage, we store information in these parameters. During evaluation, we look at the learned model only.

This is a **Parametric** classifier.

Nearest Neighbors

Let's review the stages in the NN algorithm:

- Learning: None! This algorithm holds all the relevant information in the training set.
- Evaluation: For every test point, find the training point “close” to it and assign it the same category.

This needs us to define a notion of “closeness”.

Nearest Neighbors

- “Closeness” is measured as a distance between the input vectors.

For instance, the Euclidean norm:

$$\left\| \begin{array}{ccc} x_{11} & - & x_{12} \\ x_{21} & - & x_{22} \\ \dots & & \dots \\ \dots & & \dots \\ x_{d1} & - & x_{d2} \end{array} \right\|_2$$

- The NN algorithm compares these distances to determine the closest neighbor.
- Curse of Dimensionality: In higher dimensions, common distances are less meaningful.

Probability Review¹

Let's review some probability basics that we will use in future lectures.

¹Following slides adapted from Erdogdu and Zemel

Introduction to Notation

Why do we care about probability?

Uncertainty arises through:

- Noisy measurements
- Variability between samples
- Finite size of data sets

Probability provides a consistent framework for the quantification and manipulation of uncertainty.

Sample Space

Sample space Ω is the set of all possible outcomes of an experiment.

Observations $\omega \in \Omega$ are points in the space also called sample outcomes, realizations, or elements.

Events $E \subset \Omega$ are subsets of the sample space.

Sample Space

Example: Flip a coin twice:

Sample space includes all possible outcomes

$$\Omega = \{HH, HT, TH, TT\}$$

Observation is any single element of the sample space

$$\omega = HT \in \Omega$$

Event is a subset of the sample space (eg. the event where both flips have the same outcome)

$$E = \{HH, TT\} \subset \Omega$$

Probability

The probability of an event E , $P(E)$, satisfies three axioms:

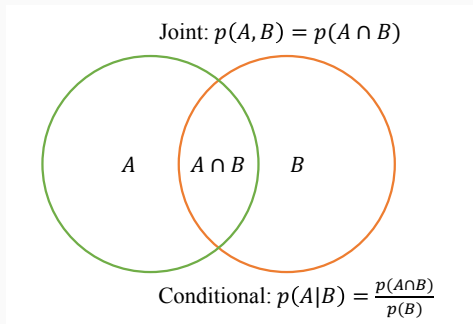
- 1: $P(E) \geq 0$ for every E
- 2: $P(\Omega) = 1$
- 3: If E_1, E_2, \dots are disjoint then

$$P\left(\bigcup_{i=1}^{\infty} E_i\right) = \sum_{i=1}^{\infty} P(E_i)$$

Joint and Conditional Probabilities

Joint Probability of A and B is denoted $P(A, B)$.

Conditional Probability of A given B is denoted $P(A|B)$.



$$p(A, B) = p(A|B)p(B) = p(B|A)p(A)$$

Conditional Example

Probability of passing the midterm is 60% and probability of passing both the final and the midterm is 45%.

What is the probability of passing the final given the student passed the midterm?

$$\begin{aligned}P(F|M) &= P(M, F)/P(M) \\&= 0.45/0.60 \\&= 0.75\end{aligned}$$

Independence

Events A and B are independent if $P(A, B) = P(A)P(B)$.

Independence

Suppose you have 2 coins. Coin 1 always comes up Heads and Coin 2 always comes up Tails. You close your eyes, pick a coin and toss it. Then you replace it, pick again and toss again.

- Independent: Before seeing the result of any toss, you wonder about 2 events; A : first toss is Head, B : second toss is Head.

$$P(A, B) = 0.5 \times 0.5 = P(A)P(B)$$

- Not Independent: Now you wonder about the same events A and B but you toss the same coin twice.

$$P(A, B) = 0.5 \neq P(A)P(B)$$

Conditional Independence

Events A and B are **conditionally independent** given C if

$$P(A, B|C) = P(B|C)P(A|C)$$

Consider two coins²: A regular coin and a coin which always outputs heads.

A = The first toss is heads;

B = The second toss is heads;

C = The regular coin is used

D = The biased coin is used

Then A and B are conditionally independent given C and given D .

²www.probabilitycourse.com/chapter1/1_4_4_conditional_independence.php

Conditional Dependence

Events A and B are **conditionally independent** given C if

$$P(A, B|C) = P(A|C)P(B|C)$$

Consider a coin which outputs heads if the first toss was heads, and tails otherwise.

A = The first toss is heads;

B = The second toss is heads;

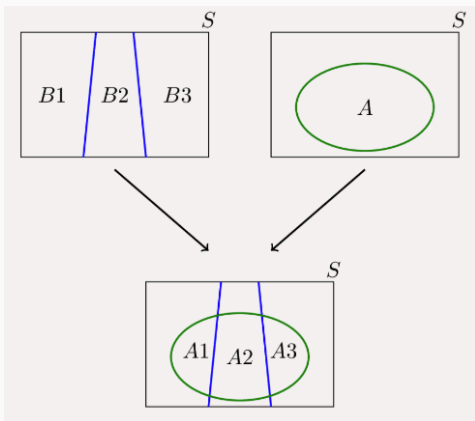
E = The eventually biased coin is used

Then A and B are conditionally dependent given E .

Marginalization and Law of Total Probability

Law of Total Probability³

$$P(A) = \sum_B P(A, B) = \sum_B P(A|B)P(B)$$



³www.probabilitycourse.com/chapter1/1_4_2_total_probability.php

Bayes' Rule

Bayes' Rule

Bayes' Rule:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

$$P(\theta|x) = \frac{P(x|\theta)P(\theta)}{P(x)}$$

$$\text{Posterior} = \frac{\text{Likelihood} * \text{Prior}}{\text{Evidence}}$$

$$\text{Posterior} \propto \text{Likelihood} \times \text{Prior}$$

Bayes' Example

Suppose you have tested positive for a disease. What is the probability you actually have the disease?

This depends on the prior probability of the disease:

- $P(T = 1|D = 1) = 0.95$ (likelihood)
- $P(T = 1|D = 0) = 0.10$ (likelihood)
- $P(D = 1) = 0.1$ (prior)

So $P(D = 1|T = 1) = ?$

Bayes' Example

$$P(D = 1|T = 1) = ?$$

Random Variables and Statistics

Random Variable

How do we connect sample spaces and events to data?

A **random variable** is a mapping which assigns a real number $X(\omega)$ to each observed outcome $\omega \in \Omega$

For example, let's flip a coin 10 times. $X(\omega)$ counts the number of Heads we observe in our sequence. If $\omega = HHTHTHHTHT$ then $X(\omega) = 6$. We often shorten this and refer to the random variable X .

Probability Distribution Statistics

Expectation: First Moment, μ

$$E[x] = \sum_{i=1}^{\infty} x_i p(x_i) \quad (\text{univariate discrete r.v.})$$

$$E[x] = \int_{-\infty}^{\infty} x p(x) dx \quad (\text{univariate continuous r.v.})$$

Variance: Second (central) Moment, σ^2

$$\begin{aligned} \text{Var}[x] &= \int_{-\infty}^{\infty} (x - \mu)^2 p(x) dx \\ &= E[(x - \mu)^2] \\ &= E[x^2] - E[x]^2 \end{aligned}$$

Expectations

From our example, we see that X does not have a fixed value, but rather a distribution of values it can take. It is natural to ask questions about this distribution, such as “What is the average number of heads in 10 coin tosses?”

This average value is called the expectation and denoted as $E[X]$. It is defined as

$$E[x] = \sum_{a \in \mathcal{A}} P[X = a] \times a$$

where \mathcal{A} represents the set of all possible values $X(w)$ can take.

Expectation Practice

- What is the expected value of a fair die?

•

Linearity of Expectations

There are two powerful properties regarding expectations.

1. $E[X + Y] = E[X] + E[Y]$.

This holds even if the random variables are dependent.

2. $E[cX] = cE[X]$, where c is a constant.

Note we cannot say anything in general about $E[XY]$.

Linearity of Expectation Practice

What is the expected value of the sum of two dice?

X_1 = value of roll 1

X_2 = value of roll 2

Linearity of Expectation Practice 2

Suppose there are n students in class, and they each complete an assignment. We hand back assignments randomly. What is the expected number of students that receive the correct assignment?

When $n = 3$? In general?

X = Number of students that get their assignment back

X_i = Student i gets their assignment back

Variances

Knowing the expectation can only tell us so much. We have another quantity used to describe how far off we are from the expected value. It is defined as follows for a random variable X with $E[X] = \mu$:

$$\text{Var}[x] = E[(X - \mu)^2]$$

The variance can be simplified as:

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

Variance Properties

Constants get squared:

$$\text{Var}[cX] = c^2 \text{Var}[X]$$

For independent random variables X and Y , we have

$$E[XY] = E[X]E[Y]$$

and

$$\text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y]$$

Variance Practice

Consider a particle that starts at position 0. At each time step, the particle moves one step to the left or one step to the right with equal probability. What is the variance of the particle at time step n ?

$$X = X_1 + X_2 + \dots + X_n$$

Discrete and Continuous Random Variables

Discrete Random Variables

- Takes countably many values, e.g., number of heads
- Distribution defined by probability mass function (PMF)
- Marginalization: $p(x) = \sum_y p(x, y)$

Continuous Random Variables

- Takes uncountably many values, e.g., time to complete task
- Distribution defined by probability density function (PDF)
- Marginalization: $p(x) = \int_y p(x, y) dy$

Random variables are said to be **independent and identically distributed** (i.i.d.) if they are sampled from the same probability distribution and are mutually independent.

This is a common assumption for observations. For example, coin flips are assumed to be i.i.d.

Problem: Eigendecomposition of a 2×2 matrix

Given the symmetric positive-definite matrix

$$\Sigma = \begin{bmatrix} a & b \\ b & c \end{bmatrix}, \quad a, c > 0,$$

Compute its (two) eigenvalues λ_1, λ_2 , and the corresponding *orthonormal* (pairwise-orthogonal and unit-normalized) eigenvectors \mathbf{v}_1 and \mathbf{v}_2 . Orthogonal means that $\langle \mathbf{v}_1, \mathbf{v}_2 \rangle = 0$, and unit-normalized implies $\|\mathbf{v}_i\|_2 = 1$.

