# Review of Basic Concepts

Yan Liu

CSCI 567 Machine Learning

January 10, 2023

# Outline

- **Probability and Statistics**: basic concepts
- **Information theory**: entropy, information gain
- **Convex Optimization**: convex, concave, basic algorithms

**Probability and Statistics**

# Probability

**Sample Space**: set of all possible outcomes or realizations.
*Example*: Toss a coin twice; the sample space is
$\Omega = \{HH, HT, TH, TT\}$.

**Event**: A subset of sample space
*Example*: the event that at least one toss is a head is
$A = \{HH, HT, TH\}$.

**Probability**: We assign a real number P(A) to each event A, called the probability of A.

**Probability Axioms**: The probability $P$ must satisfy three axioms:

1. $P(A) \geq 0$ for every A;
2. $P(\Omega) = 1$;
3. If $A_1, A_2, \ldots$ are disjoints, then $P(\cup_{i=1}^{\infty} A_i) = \sum_{i=1}^{\infty} P(A_i)$

# Random Variable

**Definition**: A random variable is a measurable function that maps a probability space into a measurable space, i.e. $X : \Omega \to R$, that assigns a real number $X(\omega)$ to each outcome $\omega$.

*Example*: if $\Omega = \{(x, y) : x^2 + y^2 \leq 1\}$ and our outcomes are samples $(x, y)$ from the unit disk, then these are some examples of random variables: $X(\omega) = x$, $Y(\omega) = y$, $Z(\omega) = x + y$.

**Data and Statistics** The data are specific realizations of random variables; A statistics is just any function of the data or random variables.

# Distribution Function

**Definition**: Suppose $X$ is a random variable, $x$ is a specific value of it, *Cumulative distribution function (CDF)* is the function $F : R \to [0, 1]$, where $F(x) = P(X \le x)$.

If $X$ is discrete $\Rightarrow$ *probability mass function*: $f(x) = P(X = x)$.
If $X$ is continuous $\Rightarrow$ *probability density function* for $X$ if there exists a function $f$ such that $f(x) \ge 0$ for all x, $\int_{-\infty}^{\infty} f(x)dx = 1$ and for every $a \le b$,

$$P(a \le X \le b) = \int_a^b f(x)dx.$$

If $F(x)$ is differentiable everywhere, $f(x) = F'(x)$.

# Expectation

**Expected Values**

- Discrete random variable X, $E[g(X)] = \sum_{x \in \mathcal{X}} g(x)f(x)$;
- Continuous random variable X, $E[g(X)] = \int_{-\infty}^{\infty} g(x)f(x)$

**Mean and Variance** $\mu = E[X]$ is the mean; $var[X] = E[(X - \mu)^2]$ is the variance.

We also have $var[X] = E[X^2] - \mu^2$.

# Common Distributions

| Discrete variable | Probability function | Mean | Variance |
|---|---|---|---|
| **Uniform** $X \sim U[1, \ldots, N]$ | $1/N$ | $\frac{N+1}{2}$ | |
| **Binomial** $X \sim Bin(n, p)$ | $\binom{x}{n} p^x (1-p)^{(n-x)}$ | np | |
| **Geometric** $X \sim Geom(p)$ | $(1-p)^{x-1} p$ | $1/p$ | |
| **Poisson** $X \sim Poisson(\lambda)$ | $\frac{e^{-\lambda} \lambda^x}{x!}$ | $\lambda$ | |
| Continuous variable | Probability density function | Mean | Variance |
| **Uniform** $X \sim U(a, b)$ | 1/ (b-a) | (a + b)/2 | |
| **Gaussian** $X \sim N(\mu, \sigma^2)$ | $\frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{1}{2\sigma^2}(x-\mu)^2)$ | $\mu$ | |
| **Gamma** $X \sim \Gamma(\alpha, \beta)$ $(x \geq 0)$ | $\frac{1}{\Gamma(\alpha)\beta^a} x^{a-1} e^{-x/\beta}$ | $\frac{\alpha}{\beta}$ | |
| **Exponential** $X \sim exponen(\beta)$ | $\frac{1}{\beta} e^{-\frac{x}{\beta}}$ | $\beta$ | |

# Multivariate Distributions

**Definition**:

$$F_{X,Y}(x,y) := P(X \le x, Y \le y),$$

and

$$f_{X,Y}(x,y) := \frac{\partial^2 F_{X,Y}(x,y)}{\partial x \partial y},$$

**Marginal Distribution** of $X$ (Discrete case):

$$f_X(x) = P(X = x) = \sum_y P(X = x, Y = y) = \sum_y f_{X,Y}(x,y)$$

or $f_X(x) = \int_y f_{X,Y}(x,y)dy$ for continuous variable.

**Conditional probability** of $X$ given $Y = y$ is

$$f_{X|Y}(x|y) = P(X = x|Y = y) = \frac{P(X = x, Y = y)}{P(Y = y)} = \frac{f_{X,Y}(x,y)}{f_Y(y)}$$

# Bayes Rule

**Law of total Probability**: $X$ takes values $x_1, \ldots, x_n$ and $y$ is a value of $Y$, we have

$$f_Y(y) = \sum_j f_{Y|X}(y|x_j) f_X(x_j)$$

**Bayes Rule**:
(Simple Form)

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

(Discrete Random Variables)

$$f_{X|Y}(x_i|y) = \frac{f_{Y|X}(y|x_i) f_X(x_i)}{\sum_j f_{Y|X}(y|x_j) f_X(x_j)}$$

(Continuous Random Variables)

$$f_{X|Y}(x|y) = \frac{f_{Y|X}(y|x) f_X(x)}{\int_x f_{Y|X}(y|x) f_X(x) dx}$$

# Independence

**Independent Variables** $X$ and $Y$ are *independent* if and only if:

$$P(X = x, Y = y) = P(X = x)P(Y = y)$$

or $f_{X,Y}(x, y) = f_X(x)f_Y(y)$ for all values $x$ and $y$.

**IID variables**: *Independent and identically distributed* (IID) random variables are drawn from the same distribution and are all mutually independent.

If $X_1, \ldots, X_n$ are independent, we have

$$E[\prod_{i=1}^{n} X_i] = \prod_{i=1}^{n} E[X_i], \ \ var[\sum_{i=1}^{n} a_i X_i] = \sum_{i=1}^{n} a_i^2 var[X_i]$$

**Linearity of Expectation**: Even if $X_1, \ldots, X_n$ are not independent,

$$E[\sum_{i=1}^{n} X_i] = \sum_{i=1}^{n} E[X_i].$$

# Correlation

**Covariance**

$$cov(X, Y) = E[(X - \mu_x)(Y - \mu_y)],$$

**Correlation coefficients**

$$corr(X, Y) = Cov(X, Y)/\sigma_x \sigma_y$$

- Independence $\Rightarrow$ Uncorrelated ($corr(X, Y) = 0$).

However, the reverse is generally not true.
The important special case: multi-variate Gaussian distribution.

# Exponential family

**Definition** A family of pdf or pmfs is called an exponential family if

$$f(x|\theta) = h(x)c(\theta)\exp(\sum_{i=1}^{k} w_i(\theta)t_i(x))$$

*Natural parameterization Form*: For $k = 1$, we have

$$f(x|\eta) = h(x)\exp(\eta t(x) - A(\eta)),$$

where $A(\eta) = log \int h(x)\exp(\eta t(x))dx$ and:

- $t(x)$ is a *sufficient statistics* of the distribution,

- $\eta$ is called the *natural parameter*,

- $A(\eta)$ is a *normalization factor*, or *log-partition function*.

**Properties**:

$$E[t(x)] = A'(\eta), \ \ Var[t(x)] = A''(\eta).$$

**Examples**: Gaussian, Exponential, Poisson, Bionomial distributions. Note that uniform distribution is NOT.

# Statistics

Suppose $X_1, \ldots, X_n$ are random variables:

**Sample Mean**:

$$\bar{X} = \frac{1}{N} \sum_{i=1}^{N} X_i$$

**Sample Variance**:

$$S^2 = \frac{1}{N-1} \sum_{i=1}^{N} (X_i - \bar{X})^2.$$

If $X_i$ are iid:

$$E[\bar{X}] = E[X_i] = \mu,$$
$$Var(\bar{X}) = \sigma^2/N,$$
$$E[S^2] = \sigma^2$$

# Point Estimation

**Definition** The *point estimator* $\hat{\theta}_N$ is a function of samples $X_1, \ldots, X_N$ that approximates a parameter $\theta$ of the distribution of $X_i$.

**Sample Bias**: The bias of an estimator is

$$bias(\hat{\theta}_N) = E_\theta[\hat{\theta}_N] - \theta$$

An estimator is *unbiased estimator* if $E_\theta[\hat{\theta}_N] = \theta$

**Standard error** The standard deviation (i.e. the square-root of variance) of $\hat{\theta}_N$ is called the *standard error*

$$se(\hat{\theta}_N) = \sqrt{Var(\hat{\theta}_N)}.$$

**Information Theory**

# Information Theory

Suppose $X$ can have on of the m values: $x_1, \ldots, x_m$. The probability distribution is $P(X = x_i) = p_i$.

**Entropy** is the smallest possible number of bits, on average, per symbol, needed to transmit a steam of symbols drawn from distribution of $X$.

$$H(X) = -\sum_{j=1}^{m} p_i \log p_i$$

- "High entropy" means X is from a uniform (boring) distribution;
- "Low entropy" means X is from varied (peaks and valleys) distribution.

# Information Theory

**Conditional Entropy** is the remaining entropy of a random variable $Y$ given that the value of another random variable $X$ is known.

$$H(Y|X) = \sum_{i=1}^{m} p_i H(Y|X = x_i) = -\sum_{i=i}^{m} \sum_{j=1}^{n} p(x_i, y_j) \log p(y_j|x_i)$$

**Mutual Information**: if $Y$ must be transmitted, how many bits on average would be saved if both ends of the line knew $X$?

$$I(Y; X) = H(Y) - H(Y|X).$$

Notice that $I(Y; X) = I(X; Y) = H(X) + H(Y) - H(X, Y)$

**Kullback-Leibler divergence** is a measure of distance between two distributions: a "true" distribution $p(X)$, and an arbitrary distribution $q(X)$.

$$\mathsf{KL}(p||q) = \sum_{x} p(x) \log \frac{p(x)}{q(x)}$$

We can write $I(X; Y) = KL(p(x, y)||p(x)p(y))$.

**Optimization**

# Optimization

**Definition**: Optimization refers to choosing the best element from some set of available alternatives. A general form is as follows:

$$
\begin{aligned}
\text{minimize} \quad & f_0(x) \\
\text{subject to} \quad & f_i(x) \leq 0, i = 1, \ldots, m \\
& h_i(x) = 0, i = 1, \ldots, p.
\end{aligned}
\tag{1}
$$

**Difficulties**:

1. Local or global optimimum?
2. Difficulty to find a feasible point,
3. Stopping criteria,
4. Poor convergence rate,
5. numerical issues

# Convex Optimization

**Convex Functions**: if for any two points $x_1$ and $x_2$ in its domain $X$ and any $t \in [0, 1]$,

$$f(tx_1 + (1-t)x_2) \leq tf(x_1) + (1-t)f(x_2).$$

A function $f$ is said to be *concave* if $-f$ is convex.

**Convex Set** a set $S$ is convex if and only if for any $x_1, x_2 \in S$,
$tx_1 + (1-t)x_2 \in S$ for any $t \in [0, 1]$,

**Convex Optimization** is minimization (maximization) of a convex (concave) function over a convex set.

*Examples*: Linear Programming (LP), Quadratic Programming (QP), and Semi-Definite Programming (SDP).

**Popular convex optimization algorithms**:

- Gradient descent
- Conjugate gradient
- Newton's method
- Quasi-Newton method
- Subgradient method