## Probability

- The goal of probability is to quantify our belief about the potential occurrence of an event.

- $\Omega$ = sample space = set of all possible outcomes

- event = subset of $\Omega$

- $A \subseteq \Omega$ occurs if we do an experiment and get $x \in A$.

- We quantify our belief in the occurrence of $A$ via a probability measure.

$\Omega = \{rain, no\ rain\}$

$A = \{rain\} =$ event that it rains

$A = \Omega =$ event that it either rains or doesn't.

**Definition ③** A probability measure on a countable sample space $\Omega$ is a map

$$P: \mathcal{P}(\Omega) \longrightarrow [0,1] \quad ; \quad P: A \longmapsto P(A)$$

(power set of $\Omega$)

such that $P(\Omega) = 1$ and

$$P\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} P(A_i) \quad \text{if} \quad A_i \cap A_j = \emptyset \quad \forall i \neq j.$$

- If $\Omega$ uncountable, replace $\mathcal{P}(\Omega)$ with $\mathcal{U} \subseteq \mathcal{P}(\Omega)$ called a $\sigma$-algebra. Each $A \in \mathcal{U}$ is called a measurable set.

- $\mathcal{U}$ usually contain sets that play well with Lebesgue integration. This is because we define many of our favorite probability measures via integration of a density function:

**Definition ④** A density function is an integrable function $f: \mathbb{R}^m \longrightarrow [0, \infty)$ with

$$\int_{\mathbb{R}^m} f = 1.$$

Given density $f$, we define a probability measure $P$ by

$$P(A) := \int_A f \quad , \quad A \in \mathcal{U}.$$

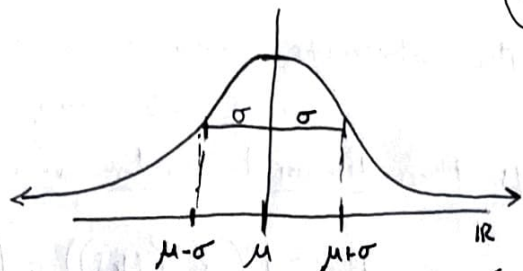The triple $(\Omega, \mathcal{U}, P)$ is a probability space. $\mathbb{R}^m \sim \mathcal{B}_m$.

**Example** (Normal Distribution) $\Omega = \mathbb{R}$

$\mathcal{U} = \mathcal{B} := \left\{ \begin{array}{l} \text{subsets of } \mathbb{R} \text{ given by countable union, intersection and relative complement} \\ \text{of open sets (intervals)} \end{array} \right\}$

- $\mu \in \mathbb{R}$ is the mean

- $\delta > 0$ is the variance

Define the density

$$\phi_{\mu,\sigma}(x) := \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right), \quad x \in \mathbb{R}$$

• Denote distribution as $N(\mu,\sigma)$.



• Given a probability space we can also ask. if one event is more or less likely given knowledge about another event. This is the notion of dependence / independence...

**Definition ⑤** $A, B \in \mathcal{U}$ such that $P(B) > 0$. The <u>conditional probability</u> of

● $A$ given $B$ is $P(A \mid B)$ where

$$P(A \cap B) = P(A \mid B) P(B).$$

● We say $A$ and $B$ are <u>independent</u> if either $P(B) > 0$ and $P(A \mid B) = P(A)$ or $P(B) = 0$.

Exercise: $A, B$ independent $\iff P(A \cap B) = P(A) P(B)$.

## Random Variables

- Most experiments we do yield numerical quantities that we may want to modify. Random variables give us the language for doing that.

**Definition ⑥** $(\Omega, \mathcal{U}, P)$ a probability space. A random variable is a

● function $X: \Omega \to \mathbb{R}$ such that for every measurable set $B \in \mathcal{B}$ $X^{-1}(B) \in \mathcal{U}$. The distribution $P^X$ associated to $X$ is

$$P^X(B) := P(X^{-1}(B)), \quad B \in \mathcal{B}.$$

●

**Example**. ① <u>Continuous random variable</u> $\Rightarrow \Omega$ uncountable. For example,

$X: \mathbb{R} \to \mathbb{R}$ the identity and $P = \int_B \phi_{\mu,\sigma}$ for $B \in \mathcal{B}$ is

called a <u>Gaussian</u> random variable; i.e., $X \sim N(\mu,\sigma)$

② <u>Discrete random variable</u> $\Rightarrow \Omega$ countable. For example:

• A coin has chance $p \in (0,1)$ of landing Heads up $(H)$.

• Flip coin $n$ times get outcome $\underline{\omega} = (\omega_1, \omega_2, \ldots, \omega_n) \in \Omega := \{H, T\}^n$.

• Since outcomes $\omega_i$ independent

$$P(\underline{\omega}) = P(\omega_1) P(\omega_2) \cdots P(\omega_n) = p^k (1-p)^{n-k}; \quad k = \# \text{ of } H\text{'s in } \underline{\omega}$$

- The discrete random variable

$$X : \Omega \longrightarrow \mathbb{R} \quad ; \quad \omega \longmapsto \sum_{i=1}^{n} X(\omega_i = H)$$

is the **binomial** random variable.

$$P_k := P(X^{-1}(\{k\})) = \binom{n}{k} p^k (1-p)^{n-k} \quad, \quad k \in \{0, \dots, n\}$$

is $p^X$, called the binomial distribution $B(n, p)$.

**Definition ⑦** A **random vector** in $\mathbb{R}^m$ is a vector of random variables $(X_1, \dots, X_m)$ all defined on the same $(\Omega, \mathcal{U}, P)$. The joint distribution of $(X_1, \dots, X_m)$, denoted $p^{(X_1, \dots, X_m)}$, defined like $p^X$ but for $\mathcal{B}_m$.

- Given jointly distributed random variables, we will often be interested in the dependencies amongst the variables encoded by the distribution, generalizing how we each ask about the dependencies of one event upon another.

**Definition ⑧** Given random variables (or vectors) $X_1, \dots, X_m$, on the probability space $(\Omega, \mathcal{U}, P)$, we say $X_1, \dots, X_m$ are **mutually independent** if

$$P(X_1 \in B_1, \dots, X_m \in B_m) = \prod_{i=1}^{m} P(X_i \in B_i) \quad \forall \text{ measurable } B_1, \dots, B_m.$$

**Lemma ⑥** Let $X_1, \dots, X_m$ random vectors on $(\Omega, \mathcal{U}, P)$ with densities $f_1, \dots, f_m$. If $X_1, \dots, X_m$ mutually independent then $(X_1, \dots, X_m)$ has density

$$f(x) = \prod_{i=1}^{m} f_{X_i}(x_i). \quad \circledast$$

If $(X_1, \dots, X_m)$ has density $f$ then $X_1, \dots, X_n$ independent if $\circledast$ holds almost surely.

**Example** (Multivariate Normal (Gaussian) Distributions) The random vector $(X_1, \dots, X_m)$ is distributed according to the **standard multivariate normal distribution** if

① $X_i \sim N(0, 1) \quad \forall i \in [m]$, and

② $X_1, \dots, X_m$ mutually independent.

Lemma ① $\Longrightarrow f_{(X_1, \dots, X_m)}(x) = \prod_{i \in [m]} \phi_{0,1}(x_i) = \prod_{i \in [m]} \frac{1}{\sqrt{2\pi}} \exp\left(-\tfrac{1}{2} x_i^2\right) = \frac{1}{(\sqrt{2\pi})^m} \exp\left(\tfrac{1}{2} x^T x\right)$

- More generally, take $\underline{\mu} \in \mathbb{R}^m$ and
  $$\Sigma \in PD_m := \{ m \times m \text{ positive definite matrices} \}.$$

- $(x_1, ..., x_m)$ distributed according to the **multivariate normal distribution** $N(\underline{\mu}, \Sigma)$ if it has density

$$f_{(x_1, ..., x_m)}(\underline{x}) = \phi_{\mu, \Sigma}(\underline{x}) := \frac{1}{(2\pi)^{m/2} |\Sigma|^{1/2}} \exp\left( -\frac{1}{2}(y - \underline{\mu})^T \Sigma^{-1}(y - \underline{\mu}) \right).$$

## Expectation, Variance, and Covariance

- We we start working with data (i.e. statistics) we will want numerical features of the data-generating distribution we can try and estimate.

- The first such feature is the "center" or expectation of the distribution...

**Definition 9** $X$ a random variable taking values in $\mathcal{X} \subseteq \mathbb{R}$ such that $\int_{\mathcal{X}} |x| f(x) < \infty$. The _expected value_, _expectation_ or _mean_ of $X$ is

$$E[X] := \int_{\mathbb{R}} x f(x) \qquad \left( E[X] := \sum_{x \in \mathcal{X}} x \, P(X=x) \quad \text{if } X \text{ discrete} \right)$$

The vector $X = (x_1, ..., x_m)$ has expectation $E[X] = (E[x_1], ..., E[x_m])$.

**Lemma 2** Expectation is linear: $X = (x_1, ..., x_m)$, $a_0 \in \mathbb{R}$, $\underline{a} \in \mathbb{R}^m$ then if $Y = a_0 + a^t X$,

$$E[Y] = a_0 + \underline{a}^t E[X] = a_0 + \sum_{i=1}^m a_i E[x_i].$$

**Examples** $E[X] = \mu$ if $X \sim N(\mu, \sigma)$. (check).

- We also measure the spread of a distribution around its mean:

**Definition 10** The _variance_ of a random variable $X$ is
$$Var[X] := E[(X - E[X])^2] \longleftarrow \begin{array}{l} \text{the expectation of the squared} \\ \text{deviation from the mean...} \end{array}$$

- Lemma 2 $\Rightarrow$ $Var[X] = E[X^2] - E[X]^2$. (check!)

- By generalizing variance, we get a measure of dependency between variables

**Definition 11** The _covariance_ of two random variables $X$ and $Y$ is
$$Cov[X, Y] := E[(X - E[X])(Y - E[Y])].$$

The covariance matrix of $X = (x_1, ..., x_m)$ is $cov[X] := [Cov[x_i, x_j]]_{i,j=1}^m$.

**Example** $Cov[X] = \Sigma$ for $X \sim N(\underline{\mu}, \Sigma)$.

**Exercise** Show that $Cov[X]$ is positive semidefinite for any random vector $X$.

# Statistics

- While probability gives a language for describing our belief in the occurrence of an event when the probability space is known, statistics focuses on how to infer the potential occurrence of an event when it is unknown...

- In statistics, start with only data and try to infer the probability of an event.

- Two Perspectives:

① Frequentist: There is a "true" unknown probability distribution an we want to use our data to learn this distribution precisely as # samples $\to \infty$.

② Bayesian: Use data to update our prior beliefs about an unknown parameter.

- We will mainly consider the frequentist perspectives...

## Statistical Models

**Definition ⑫/** A statistical model $M$ is a family of probability distributions. A parametric statistical model $M_\Theta$ is a mapping from a finite dimensional parameter space $\Theta \subseteq \mathbb{R}^d$ to a space of distributions or densities:

$$P_\bullet : \Theta \to M_\Theta \quad ; \quad \Theta \longmapsto P_\Theta.$$

$M_\Theta := P_\bullet(\Theta) = \{P_\Theta : \Theta \in \Theta\}$. $M_\Theta$ is _identifiable_ if $P_\bullet$ is injective.

**Examples** ① The _Binomial_ Random Variable Model: $\Theta := [0,1] \subseteq \mathbb{R}^1$.

Fix $r \in \mathbb{Z}_{\geq 0}$. Then

$$M_\Theta := \left\{ \left( \binom{r}{i} \theta^i (1-\theta)^{r-i} \mid i \in [0,\dots,r] \right) : \Theta \in [0,1] \right\}.$$

② _Multivariate Gaussian (Normal) Model_: $X$ an $m$-dim'l random vector.

$\Theta := \mathbb{R}^m \times PD_m$. $M_\Theta$ is the image of

$$P_\bullet : \Theta \ni (\mu, \Sigma) \longmapsto \emptyset_{\mu, \Sigma}(\underline{x}).$$

**Definition ⑬/** A statistical model $M$ is _implicit_ if it is defined by a set of constraints.

**Example** conditional independence models. (Markov chain)

**Note** A model can be ~~both~~ implicit and parametric...

- When trying to estimate how our data-generating distribution fits a model's parameters or defining constraints we will rely on features of our data that we call statistics...

| Definition (14) | A *statistic* is a function from the state space of a random variable (i.e., its image) to another set $T: \mathcal{X} \to \mathbb{T}$. For a parametric model $M_\Theta$, $T$ is *sufficient* if its densities $f_\Theta(x)$ satisfies

$$f_\Theta(x) = h(x) g(T(x), \Theta) \text{ for some } g, h.$$

[Example] $\mathcal{X}_F$ is the binomial random variable where $\mathcal{X} = 0$ if an event happens and $\mathcal{X} = 1$ if it doesn't. Do $n$ i.i.d experiments

① $\mathcal{X}^{(1)}, ..., \mathcal{X}^{(n)}$. Then if $\Theta \in (0,1)$ is $P(\mathcal{X} = 0)$,

② $P(\mathcal{X}^{(1)} = x_1, ..., \mathcal{X}^{(n)} = x_n) = \prod_{i \in [n]} P(\mathcal{X}_i = x_i) = \Theta^{\Sigma x_i} (1-\Theta)^{n - \Sigma x_i}.$

So $h(x) \equiv 1$, $g(t, \Theta) = \Theta^t (1-\Theta)^{n-t}$ and so $T(\mathcal{X}) = \sum_{i=1}^{n} \mathcal{X}^{(i)}$ is sufficient.

- Here, our data was the outcome of an event and the # of samples was $h$. If the event was "a flipped coin landed heads up" then each sample would have been drawn from the same distribution and independent.

⊛ Our data will usually be iid = independent and identically distributed.

• To estimate the probability distribution assuming it comes from our parametric model, we will often use sufficient statistics to estimate the model parameters.

| Definition (15) | A *parameter* of a statistical model $M_\Theta$ is a function

① $s: \Theta \to \mathbb{R}$. An *estimator* of $s$ is a function $\hat{s}: \mathbb{D} \to \mathbb{R}$, where $\mathbb{D}$ is the data. $\hat{s}$ is *consistent* if $\hat{s} \to s$ as $|\mathbb{D}| \to \infty$.

[Example] The maximum likelihood estimator (MLE)

$$\hat{\Theta} = \arg\max_{\Theta \in \Theta} L(\Theta | \mathbb{D}) = \arg\max_{\Theta \in \Theta} P(\mathbb{D} | \Theta).$$

(see last lecture!)

• Parameter estimation only gives an approximation of the model parameters assuming the data-generating distribution is actually in the model. This may be a bad assumption.

# Hypothesis Testing

- Hypothesis tests are the statistical framework for rejecting an assumption (such as "the data-generating distribution is in this model") base on data.

- $H_0$ = the null hypothesis

[Q] Given the data can we reject the null hypothesis $H_0$?

[Example]. $X, Y$ are two random variables.

- $H_0 : (X = Y)$; i.e. $X$ and $Y$ are equidistributed (same distribution).
- $H_a : (X \neq Y)$ is the alternative hypothesis ($\neg H_0$).

- Can reject $H_0$ if $E[X] \neq E[Y]$.

- Collect samples $X^{(1)}, \ldots, X^{(n_1)}$ from $X$ and $Y^{(1)}, \ldots, Y^{(n_2)}$ from $Y$.

- Compute the sample means (estimators of the means):

$$\hat{\mu}_1 = \frac{1}{n_1} \sum_{i=1}^{n_1} X^{(i)} \quad \text{and} \quad \hat{\mu}_2 = \frac{1}{n_2} \sum_{i=1}^{n_2} Y^{(i)}$$

- Permutation test (exact test):

  - set $d = |\hat{\mu}_1 - \hat{\mu}_2|$ and $Z := \{ z_1, \ldots, z_{n_1 + n_2} \} = \{ X^{(1)}, \ldots, X^{(n_1)}, Y^{(1)}, \ldots, Y^{(n_2)} \}$

  - for each partition $A \cup B$ of $Z$ where $|A| = n_1$ compute

$$d_{A \cup B} := \left| \frac{1}{n_1} \sum_{a \in A} a - \frac{1}{n_2} \sum_{b \in B} b \right|.$$

  - If $p := \dfrac{|\{ A \cup B | d_{A \cup B} \geq d \}|}{\binom{n_1 + n_2}{n_1}} < \alpha = 0.05$ then reject $H_0$.

  - We can do this because if few $d_{A \cup B}$ are larger than $d$, it suggests that $\hat{\mu}_1$ and $\hat{\mu}_2$ are further away that would be possible due to sampling error...