# Expectation

Total expectation theorem:

$$\mathbb{E}[X] = \int_{-inf}^{+inf} f_Y(y) \cdot \mathbb{E}[X|Y=y] \, dy$$

Law of iterated expectation:

$$\mathbb{E}[Y] = \mathbb{E}[\mathbb{E}[Y|X]]$$

Law of total expectation:

$$\mathbb{E}[Y] = \sum_n p_N(n)\mathbb{E}[Y|N=n]$$

Product of **independent** r.vs $X$ and $Y$:

$$\mathbb{E}[X \cdot Y] = \mathbb{E}[X] \cdot \mathbb{E}[Y]$$

Product of **dependent** r.vs $X$ and $Y$:

$$\mathbb{E}[X \cdot Y] \neq \mathbb{E}[X] \cdot \mathbb{E}[Y]$$

$$\mathbb{E}[X \cdot Y] = \mathbb{E}[\mathbb{E}[Y \cdot X|Y]] = \mathbb{E}[Y \cdot \mathbb{E}[X|Y]]$$

Linearity of Expectation where $a$ and $c$ are given scalars:

$$\mathbb{E}[aX + cY] = a\mathbb{E}[X] + c\mathbb{E}[Y]$$

If Variance of $X$ is known:

$$\mathbb{E}[X^2] = var(X) - \mathbb{E}[X]^2$$

## Variance

Variance is the squared distance from the mean.

$$Var(X) = \mathbb{E}[(X - \mathbb{E}[X])^2]$$

Law of total variance
$$Var(X) = \mathbb{E}[Var(X|Y)] + Var(\mathbb{E}[X|Y])$$

Variance of a product with constant $a$:

$$Var(aX) = a^2 Var(X)$$

Variance of sum of two **dependent** r.v.:

$$Var(X+Y) = Var(X) + Var(Y) + 2Cov(X,Y)$$

Variance of sum of n **dependent** r.v.:

$$Var(X_1 + X_2 + ... + X_n) = \sum_{i=1}^{n} Var(X_i) + \sum_{i \neq j} Cov(X_i, X_j)$$

Variance of sum/difference of two **independent** r.v.:

$$Var(X + Y) = Var(X) + Var(Y)$$

$$Var(X - Y) = Var(X) + Var(Y)$$

## Covariance

The Covariance is a measure of how much the values of each of two correlated random variables determine each other

$$Cov(X, Y) = \mathbb{E}[(X - \mu_X)(Y - \mu_Y)]$$

$$Cov(X, Y) = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]$$

Possible notations:

$$Cov(X, Y) = \sigma(X, Y) = \sigma_{(X,Y)}$$

Covariance is commutative:

$$Cov(X, Y) = Cov(Y, X)$$

Covariance with of r.v. with itself is variance:

$$Cov(X, X) = \mathbb{E}[(X - \mu_X)^2] = Var(X)$$

Useful properties:

$$Cov(aX + h, bY + c) = abCov(X, Y)$$

$$Cov(X, X + Y) = Var(X) + cov(X, Y)$$

$$Cov(aX + bY, Z) = aCov(X, Z) + bCov(Y, Z)$$

If $Cov(X, Y) = 0$, we say that X and Y are uncorrelated. If $X$ and $Y$ are independent, their Covariance is zero. The converse is not always true. It is only true if $X$ and $Y$ form a gaussian vector, ie. any linear combination $\alpha X + \beta Y$ is gaussian for all $(\alpha, \beta) \in \mathbb{R}^2$ without $\{0, 0\}$.

## correlation coefficient
$$\rho(X, Y) = \frac{Cov(X,Y)}{\sqrt{Var(X}\sqrt{Var(Y}}$$ Import properties of $\rho$ : $-1 \leq \rho \leq 1$

## Covariance Matrix
Let $X$ be a random vector of dimension $d \times 1$ with expectation $\mu_X$.
Matrix outer products!

$$\Sigma = \mathbb{E}[(X - \mu_X)(X - \mu_X)^T]$$
$$= \mathbb{E}[XX^T] - \mathbb{E}[X]\mathbb{E}[X]^T$$
$$= \mathbb{E}[XX^T] - \mu_X\mu_X^T$$

## LLN and CLT
Let $X_1, ..., X_n \overset{iid}{\sim} P_\mu$, where $E(X_i) = \mu$ and $Var(X_i) = \sigma^2$ for all $i = 1, 2, ..., n$ and $\overline{X}_n = \frac{1}{n}\sum_{i=1}^{n} X_i$.
### Law of large numbers:

$$\overline{X}_n \xrightarrow[n \to \infty]{P, a.s.} \mu$$

$$\frac{1}{n}\sum_{i=1}^{n} g(X_i) \xrightarrow[n \to \infty]{P, a.s.} \mathbb{E}[g(X)]$$

### The Markov Inequality
If $X \geq 0$ and $a > 0$, then

$$P(X \geq a) \leq \frac{\mathbb{E}[X]}{a}$$

"If $X \geq 0$ and $\mathbb{E}[X]$ is small, then X is unlikely to be very large"
### The Chebyshev Inequality
With finite mean $\mu$ and variance $\sigma^2$

$$P(|X \geq \mu| \geq c) \leq \frac{\sigma^2}{c^2}$$

"If variance is small, then X is unlikely to be far from the mean"
### WLLN
$X_1...X_n$ i.i.d. With finite mean $\mu$ and variance $\sigma^2$
For $\epsilon > 0$, $P(|M_n - \mu| \geq \epsilon) = P(|\frac{X_1 + ... + X_n}{n} - \mu| \geq \epsilon) \to 0$, $as\, n \to \infty$

"The sample mean is the empirical frequency of event "
$$M_n \xrightarrow[n \to \infty]{P} \mu \to \text{"convergence in probability"}$$

### Variance of the Mean:

$$Var(\overline{X}_n) = (\frac{\sigma^2}{n})^2 Var(X_1 + X_2, ..., X_n)$$

$$= \frac{\sigma^2}{n}$$

### Expectation of the mean:

$$E[\overline{X}_n] = \frac{1}{n}E[X_1 + X_2, ..., X_n]$$

$$= \mu.$$

## Central Limit theorem
### Central Limit Theorem for Mean:

$$\sqrt{n}\frac{\overline{X}_n - \mu}{\sigma} \xrightarrow[n \to \infty]{d} N(0, 1)$$

$$\sqrt{n}(\overline{X}_n - \mu) \xrightarrow[n \to \infty]{d} N(0, \sigma^2)$$

### Central Limit Theorem for Sums:

$$\sum_{i=1}^{n} X_i \xrightarrow[n \to \infty]{(d)} N(n\mu, \sqrt{n}\sigma)$$

$$\frac{\sum_{i=1}^{n} X_i - n\mu}{\sqrt{n}\sigma} \xrightarrow[n \to \infty]{(d)} N(0, 1)$$

Note: $X_i$ are not necessarily i.i.d

## Bayesian Statistics
Bayesian inference conceptually amounts to weighting the likelihood $L_n(\theta)$ by a prior knowledge we might have on $\theta$. Given a statistical model we technically model our parameter $\theta$ as if it were a random variable. We therefore define the **prior distribution** (PDF):

$$\pi(\theta)$$

Let $X_1, ..., X_n$. We note $f(X_1, ..., X_n|\theta)$ the joint probability distribution of $X_1, ..., X_n$ conditioned on $\theta$ where $\theta \sim \pi$. This is exactly the likelihood from the frequentist approach.

### Bayes' formula
. The posterior distribution verifies:

$$\forall \theta \in \Theta, \pi(\theta|X_1, ..., X_n) \propto$$
$$\pi(\theta)f(X_1, ..., X_n|\theta)$$

The constant is the normalization factor to ensure the result is a proper distribution, and does not depend on $\theta$:

$$\pi(\theta|X_1, ..., X_n) = \frac{\pi(\theta)f(X_1, ..., X_n|\theta)}{\int_\Theta \pi(\theta)f(X_1, ..., X_n|\theta)d\theta}$$

We can often use an **improper prior**, i.e. a prior that is not a proper probability distribution (whose integral diverges), and still get a proper posterior. For example, the improper prior $\pi(\theta) = 1$ on $\Theta$ gives the likelihood as a posterior.

## Jeffreys Prior

$$\pi_J(\theta) \propto \sqrt{det I(\theta)}$$

where $I(\theta)$ is the Fisher information. This prior is **invariant by reparameterization**, which means that if we have $\eta = \phi(\theta)$, then the same prior gives us a probability distribution for $\eta$ verifying:

$$\tilde{\pi}_J(\eta) \propto \sqrt{det\tilde{I}(\eta)}$$

The change of parameter follows the following formula:

$$\tilde{\pi}_J(\eta) = det(\nabla\phi^{-1}(\eta))\pi_J(\phi^{-1}(\eta))$$

### Bayesian confidence region
Let $\alpha \in (0, 1)$. A *Bayesian confidence region with level $\alpha^*$ is a random subset $\mathcal{R} \subset \Theta$ depending on $X_1, ..., X_n$ (and the prior $\pi$) such that:

$$P[\theta \in \mathcal{R}|X_1, ..., X_n] \geq 1 - \alpha$$

Bayesian confidence region and confidence interval are **distinct notions**. The Bayesian framework can be used to estimate the true underlying parameter. In that case, it is used to build a new class of estimators, based on the posterior distribution.
### Bayes estimator (LMS)
### posterior mean:

$$\hat{\theta}_{(\pi)} = \int \theta * \pi(\theta|X_1, ..., X_n)d\theta$$

### Maximum a posteriori estimator (MAP):

$$\hat{\theta}_{(\pi)}^{MAP} = argmax_{\theta \in \Theta}\pi(\theta|X_1, ..., X_n)$$

The MAP is equivalent to the MLE, if the prior is uniform.
### Point estimates of a biased coin
**posterior**:
assuming prior is union distribution:
$$d(n, k)\theta^k(1-\theta)^{n-k}$$
**MAP:**
$$\hat{\theta_{MAP}} = K/n$$ where K is number of heads and n is number of tosses
"find $\hat{\theta}$ to get maximum of posterior probability "
**LMS:**
$$\hat{\theta_{LMS}} = \mathbb{E}[\hat{\theta}|K=k] = \frac{k+1}{n+2}$$, and as $n$ gets large, it approaches MAP estimate
"find $\hat{\theta}$ to get maximum expectation of conditional expectation (posterior mean) of $\hat{\theta}$ based on posterior probability "

## Examples of Parametric Models
### Poisson
Parameter $\lambda$. discrete, approximates the binomial PMF when $n$ is large, $p$ is small, and $\lambda = np$.

$$\mathbf{p_x}(k) = exp(-\lambda)\frac{\lambda^k}{k!}$$ for $k = 0, 1, ...,$
$$\mathbb{E}[X] = \lambda$$
$$Var(X) = \lambda$$
Likelihood:
$$L_n(x_1, ..., x_n, \lambda) = \prod_{i=1}^{n} \frac{\lambda^{\sum_{i=1}^{n} x_i}}{\prod_{i=1}^{n} x_i!} e^{-n\lambda}$$
Loglikelihood:

$$\ell_n(\lambda) =$$
$$= -n\lambda + log(\lambda)(\sum_{i=1}^{n} x_i)) - log(\prod_{i=1}^{n} x_i!)$$
MLE:
$$\hat{\lambda}_{MLE} = \frac{1}{n}\sum_{i=1}^{n}(X_i)$$
Fisher Information:
$$I(\lambda) = \frac{1}{\lambda}$$
Canonical exponential form:
$$f_\theta(y) = \exp\left(y\theta - \underbrace{e^\theta}_{b(\theta)} - \underbrace{\ln y!}_{c(y,\phi)}\right)$$
$$\theta = \ln\lambda$$
$$\phi = 1$$
Poisson process:
k arrivals in t slots $\mathbf{p_x}(k, t) = \mathbb{P}(N_t = k) = e^{-\lambda t}\frac{(\lambda t)^k}{k!}$
$$\mathbb{E}[N_t] = \lambda t$$
$$Var(N_t) = \lambda t$$

### Standard Gaussians
Parameters $\mu$ and $\sigma^2 > 0$, continuous
$$f(x) = \frac{1}{\sqrt{(2\pi\sigma^2)}}exp(-\frac{(x-\mu)^2}{2\sigma^2})$$
$$\mathbb{E}[X] = \mu$$
$$Var(X) = \sigma^2$$
CDF of standard gaussian:
$$\Phi(z) = \int_{-\infty}^{z} \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx$$
Likelihood:
$$L(x_1 ... x_n; \mu, \sigma^2) =$$
$$= \frac{1}{(\sigma\sqrt{2\pi})^n} \exp\left(-\frac{1}{2\sigma^2}\sum_{i=1}^{n}(X_i - \mu)^2\right)$$
Loglikelihood:
$$\ell_n(\mu, \sigma^2) =$$
$$= -nlog(\sigma\sqrt{2\pi}) - \frac{1}{2\sigma^2}\sum_{i=1}^{n}(X_i - \mu)^2$$
MLE:
$$\hat{\mu}_{M}LE = \overline{X}_n$$
$$\hat{\sigma^2}_{M}LE = \frac{1}{n}\sum_{i=1}^{n}(X_i - \overline{X}_n)^2$$ Fisher Information:
$$I(\mu, \sigma^2) = \begin{pmatrix} \frac{1}{\sigma^2} & 0 \\ 0 & \frac{1}{2\sigma^4} \end{pmatrix}$$ Gaussians are invariant under affine transformation:
$$aX + b \sim N(X + b, a^2\sigma^2)$$
Sum of independent gaussians:
Let $X \sim N(\mu_X, \sigma_X^2)$ and $Y \sim N(\mu_Y, \sigma_Y^2)$
If $Y = X + Z$, then $Y \sim N(\mu_X + \mu_Y, \sigma_X^2 + \sigma_Y^2)$
If $U = X - Y$, then $U \sim N(\mu_X - \mu_Y, \sigma_X^2 + \sigma_Y^2)$
Symmetry:
If $X \sim N(0, \sigma^2)$, then $-X \sim N(0, \sigma^2)$
$$\mathbb{P}(|X| > x) = 2\mathbb{P}(X > x)$$
Standardization:
$$Z = \frac{X-\mu}{\sigma} \sim N(0, 1)$$
$$\mathbf{P}(X \leq t) = \mathbf{P}\left(Z \leq \frac{t-\mu}{\sigma}\right)$$ Higher moments:
$$\mathbb{E}[X^2] = \mu^2 + \sigma^2$$
$$\mathbb{E}[X^3] = \mu^3 + 3\mu\sigma^2$$
$$\mathbb{E}[X^4] = \mu^4 + 6\mu^2\sigma^2 + 3\sigma^4$$

### Exponential
Parameter $\lambda$, continuous
$$f_x(x) = \begin{cases} \lambda exp(-\lambda x), & \text{if } x >= 0 \\ 0, & \text{o.w.} \end{cases}$$
$$P(X > a) = exp(-\lambda a)$$
$$F_x(x) = \begin{cases} 1 - exp(-\lambda x), & \text{if } x >= 0 \\ 0, & \text{o.w.} \end{cases}$$
$$\mathbb{E}[X] = \frac{1}{\lambda}$$
$$\mathbb{E}[X^2] = \frac{2}{\lambda^2}$$
$$Var(X) = \frac{1}{\lambda^2}$$

Likelihood:
$$L(X_1 ... X_n; \lambda) = \lambda^n \exp\left(-\lambda\sum_{i=1}^{n} X_i\right)$$
Loglikelihood:
$$\ell_n(\lambda) = nln(\lambda) - \lambda\sum_{i=1}^{n}(X_i)$$
MLE:
$$\hat{\lambda}_{MLE} = \frac{n}{\sum_{i=1}^{n}(X_i)}$$
Fisher Information:
$$I(\lambda) = \frac{1}{\lambda^2}$$
Canonical exponential form:
$$f_\theta(y) = \exp\left(y\theta - \underbrace{(-\ln(-\theta))}_{b(\theta)} + \underbrace{0}_{c(y,\phi)}\right)$$
$$\theta = -\lambda = -\frac{1}{\mu}$$
$$\phi = 1$$

### Binomial
Parameters $p$ and $n$, discrete. Describes the number of successes in n independent Bernoulli trials.
$$p_x(k) = \binom{n}{k}p^k(1-p)^{n-k}, k = 0, ..., n$$
$$\mathbb{E}[X] = np$$
$$Var(X) = np(1-p)$$
Likelihood:
$$L_n(X_1, ..., X_n, \theta) =$$
$$= \left(\prod_{i=1}^{n}\binom{K}{X_i}\right)\theta^{\sum_{i=1}^{n} X_i}(1-\theta)^{nK - \sum_{i=1}^{n} X_i}$$
Loglikelihood:
$$\ell_n(\theta) = C + \left(\sum_{i=1}^{n} X_i\right)\log\theta + \left(nK - \sum_{i=1}^{n} X_i\right)\log(1-\theta)$$
MLE:
Fisher Information:
$$I(p) = \frac{n}{p(1-p)}$$
Canonical exponential form:
$$f_p(y) = exp(y\underbrace{(\ln(p) - \ln(1-p))}_{\theta} + \underbrace{n\ln(1-p)}_{-b(\theta)} + \underbrace{\ln(\binom{n}{y})}_{c(y,\phi)})$$

### Geometric
Number of $T$ trials up to (and including) the first success.
$$p_T(t) = (1-p)^{t-1}, t = 1, 2, ...$$
$$\mathbb{E}[T] = \frac{1}{P}$$
$$var(T) = \frac{1-p}{p^2}$$

### Pascal
The negative binomial or Pascal distribution is a generalization of the geometric distribution. It relates to the random experiment of repeated independent trials until observing $m$ successes. I.e. the time of the kth arrival. $Y_k = T_1 + ... T_k$
$T_i \sim iid Geometric(p)$
$$\mathbb{E}[Y_k] = \frac{k}{p}$$
$$Var(Y_k) = \frac{k(1-p)}{p^2}$$ $p_{Y_k}(t) = \binom{t-1}{k-1}p^k(1-p)^{t-k}$
$$t = k, k+1, ...$$

### Multinomial
Parameters $n > 0$ and $p_1, ..., p_r$.
$$p_x(x) = \frac{n!}{x_1!, ..., x_n!}p_1, ..., p_r.$$
$$\mathbb{E}[X_i] = n * p_i$$
$$Var(X_i) = np_i(1-p_i)$$
Likelihood:
$$p_x(x) = \prod_{j=1}^{n} p_j^{T^j}$$, where $T^j = \mathbb{1}(X_i = j)$ is the count how often an outcome is seen in trials.
Loglikelihood:

$\ell_n = \sum_{j=2}^n T_j \ln(p_j)$

## Shifted Exponential

Parameters $\lambda, a \in \mathbb{R}$, continuous

$f_x(x) = \begin{cases} \lambda exp(-\lambda(x-a)), & x >= a \\ 0, & x <= a \end{cases}$

$F_x(x) = \begin{cases} 1 - exp(-\lambda(x-a)), & if x >= a \\ 0, & x <= a \end{cases}$

$\mathbb{E}[X] = a + \frac{1}{\lambda}$

$Var(X) = \frac{1}{\lambda^2}$

Likelihood:

$L(X_1 \ldots X_n; \lambda, \theta) = \lambda^n \exp\left(-\lambda \sum_{i=1}^n (X_i - a)\right) \mathbf{1}_{\min_{i=1,\ldots,n}(X_i) \geq a}$

Loglikelihood:

$\ell(\lambda, a) := n \ln \lambda - \lambda \sum_{i=1}^n X_i + n\lambda a$ MLE:

$\hat{\lambda}_{MLE} = \frac{1}{\overline{X}_n - \hat{a}}$

$\hat{a}_{MLE} = \min_{i=1,\ldots,n}(X_i)$

## Uniform

Parameters $a$ and $b$, continuous.

$\mathbf{f_x}(x) = \begin{cases} \frac{1}{b-a}, & if \ a < x < b \\ 0, & o.w. \end{cases}$

$\mathbf{F_x}(x) = \begin{cases} 0, & for \ x \leq a \\ \frac{x-a}{b-a}, & x \in [a,b) \\ 1, & x \geq b \end{cases}$

$\mathbb{E}[X] = \frac{a+b}{2}$

$Var(X) = \frac{(b-a)^2}{12}$

Likelihood:

$L(x_1 \ldots x_n; b) = \frac{\mathbf{1}(\max_i(x_i \leq b))}{b^n}$

Loglikelihood:

## Cauchy

continuous, parameter $m$,

$f_m(x) = \frac{1}{\pi} \frac{1}{1+(x-m)^2}$

$\mathbb{E}[X] = notdefined!$

$Var(X) = notdefined!$

$med(X) = P(X > M) = P(X < M)$
$= 1/2 = \int_{1/2}^\infty \frac{1}{\pi} \cdot \frac{1}{1+(x-m)^2} dx$

## Chi squared

The $\chi_d^2$ distribution with $d$ degrees of freedom is given by the distribution of $Z_1^2 + Z_2^2 + \cdots + Z_d^2$, where $Z_1, \ldots, Z_d \overset{iid}{\sim} \mathcal{N}(0,1)$ If $V \sim \chi_k^2$:

$\mathbb{E} = \mathbb{E}[Z_1^2] + \mathbb{E}[Z_2^2] + \ldots + \mathbb{E}[Z_d^2] = d$

$Var(V) = Var(Z_1^2) + Var(Z_2^2) + \ldots + Var(Z_d^2) = 2d$

## Student's T Distribution

$T_n := \frac{Z}{\sqrt{V/n}}$ where $Z \sim \mathcal{N}(0,1)$, and $Z$ and $V$ are independent

## Parametric Estimation

An **estimator** $\hat{\theta}_n$ of $\theta$ is any statistic which does not depend on $\theta$. Estimators are random variables if they depend on the data (= realizations of random variables).

---

## bias and variance of estimators
### Jensen's inequality
For any variable $X$ and any convex function $g$

$g(\mathbb{E}[X]) \leq \mathbb{E}[g(X)]$

### Bias of an estimator:

$Bias(\hat{\theta}_n) = \mathbb{E}[\hat{\theta}_n] - \theta$

### Quadratic risk of an estimator

$R(\hat{\theta}_n) = \mathbb{E}[(\hat{\theta}_n - \theta)^2]$
$= Bias^2 + Variance$

"Low quadratic risk means both low bias and low variance"

### Asymptotic normality of an estimator
An estimator $\hat{\theta}_n$ is **weakly consistent**

if: $\lim_{n\to\infty} \hat{\theta}_n = \theta$ or $\hat{\theta}_n \xrightarrow[n\to\infty]{P} \mathbb{E}[g(X)]$. If the convergence is almost surely it is **strongly consistent**.

### Asymptotic normality of an estimator:

$\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow[n\to\infty]{(d)} N(0, \sigma^2)$

$\sigma^2$ is called the **Asymptotic Variance** of the estimator $\hat{\theta}_n$. In the case of the sample mean it is the same variance as as the single $X_i$.
If the estimator is a function of the sample mean the **Delta Method** is needed to compute the asymptotic variance. **Asymptotic Variance** $\neq$ Variance of an estimator.
Note: The asymptotic variance is the limit of a sequence as n goes to infinity. It is a specific real number, not a function of n. It comes from central limit theorem

## Confidence intervals

Confidence Intervals follow the form:

(statistic) $\pm$ (critical value)(estimated standard deviation of statistic)

Let $(E, (\mathbb{P}_\theta)_{\theta \in \Theta})$ be a statistical model based on observations $X_1, \ldots X_n$ and assume $\Theta \subseteq \mathbb{R}$. Let $\alpha \in (0,1)$.
**Non asymptotic** confidence interval of level $1 - \alpha$ for $\theta$:
Any random interval $\mathcal{I}$, depending on the sample $X_1, \ldots X_n$ but not at $\theta$ and such that:
$\mathbb{P}_\theta[\mathcal{I} \ni \theta] \geq 1 - \alpha, \ \forall \theta \in \Theta$
Confidence interval of **asymptotic level** $1 - \alpha$ for $\theta$:
Any random interval $\mathcal{I}$ whose boundaries do not depend on $\theta$ and such that:
$\lim_{n\to\infty} \mathbb{P}_\theta[\mathcal{I} \ni \theta] \geq 1 - \alpha, \ \forall \theta \in \Theta$

### Two-sided asymptotic CI

Let $X_1, \ldots, X_n = \bar{X}$ and $\bar{X} \overset{iid}{\sim} P_\theta$. A two-sided CI is a function depending on $\bar{X}$ giving an upper and lower bound in which the estimated parameter lies $\mathcal{I} = [l(\bar{X}, u(\bar{X})]$ with a certain probability $\mathbb{P}(\theta \in \mathcal{I}) \geq 1 - q_\alpha$ and conversely $\mathbb{P}(\theta \notin \mathcal{I}) \leq \alpha$
Since the estimator is a r.v. depending on

---

$\bar{X}$ it has a variance $Var(\hat{\theta}_n)$ and a mean $\mathbb{E}[\hat{\theta}_n]$. Since the CLT is valid for every distribution standardizing the distributions and massaging the expression yields an an asymptotic CI:

$\mathcal{I} = [\hat{\theta}_n - \frac{q_{\alpha/2}\sqrt{Var(X_i)}}{\sqrt{n}},$
$\hat{\theta}_n + \frac{q_{\alpha/2}\sqrt{Var(X_i)}}{\sqrt{n}}]$

This expression depends on the real variance $Var(X_i)$ of the r.vs, the variance has to be estimated.
Three possible methods: plugin (use sample mean or empirical variance), solve (solve quadratic inequality), conservative (use the theoretical maximum of the variance.

### Sample Mean and Sample Variance

Let $X_1, \ldots, X_n \overset{iid}{\sim} P_\mu$, where $E(X_i) = \mu$ and $Var(X_i) = \sigma^2$ for all $i = 1, 2, \ldots, n$
**Sample Mean:**

$\overline{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$

**Sample Variance:**

$S_n = \frac{1}{n} \sum_{i=1}^n (X_i - \overline{X}_n)^2$
$= \frac{1}{n}(\sum_{i=1}^n X_i^2) - \overline{X}_n^2$

**Unbiased estimator of sample variance:**

$\tilde{S}_n = \frac{1}{n-1} \sum_{i=1}^n \left(X_i - \overline{X}_n\right)^2$
$= \frac{n}{n-1} S_n$

## Confidence intervals Approaches
### Solution 1 - Conservative Bound
1. use WLLN to get boundaries dependent on $\theta$
2. Find the most conservative bounds $argmax(f(\theta))$

### Solution 2 - Solve

transfer the 2 sided inequality equations to quadratic equation on $\theta$
$I_{solve} = [\hat{\lambda}(1 + \frac{q_{\alpha/2}}{\sqrt{n}})^{-1}, \hat{\lambda}(1 - \frac{q_{\alpha/2}}{\sqrt{n}})^{-1}]$

### Solution 3 - plug-in

replace $\theta$ with $\hat{\theta}$ (slutzky theorem)
$I_{plug-in} = [\hat{\lambda}(1 - \frac{q_{\alpha/2}}{\sqrt{n}}), \hat{\lambda}(1 + \frac{q_{\alpha/2}}{\sqrt{n}})]$

## Delta Method

To find the asymptotic CI if the estimator is a function of the mean. Goal is to find an expression that converges a function of the mean using the CLT. Let $Z_n$ be a sequence of r.v. $\sqrt{n}(Z_n - \theta) \xrightarrow[n\to\infty]{(d)} N(0, \sigma^2)$ and let $g: R \longrightarrow R$ be continuously differentiable at $\theta$, then:

$\sqrt{n}(g(Z_n) - g(\theta)) \xrightarrow[n\to\infty]{(d)}$
$\mathcal{N}(0, g'(\theta)^2 \sigma^2)$

**Example:** let $X_1, \ldots, X_n \ exp(\lambda)$ where $\lambda > 0$. Let $\overline{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ denote the sample mean. By the CLT, we know that

---

$\sqrt{n}\left(\overline{X}_n - \frac{1}{\lambda}\right) \xrightarrow[n\to\infty]{(d)} N(0, \sigma^2)$ for some value of $\sigma^2$ that depends on $\lambda$.
If we set $g: \mathbb{R} \to \mathbb{R}$ and $x \mapsto 1/x$, then by the Delta method:

$\sqrt{n}\left(g(\overline{X}_n) - g\left(\frac{1}{\lambda}\right)\right)$
$\xrightarrow[n\to\infty]{(d)} N(0, g'(\mathbb{E}[X])^2 VarX)$
$\xrightarrow[n\to\infty]{(d)} N(0, g'\left(\frac{1}{\lambda}\right)^2 \frac{1}{\lambda^2})$
$\xrightarrow[n\to\infty]{(d)} N(0, \lambda^2)$

## Asymptotic Hypothesis tests (Theory)

Two hypotheses ($\Theta_0$ disjoint set from $\Theta_1$):
$\begin{cases} H_0 : \theta \epsilon \Theta_0 \\ H_1 : \theta \epsilon \Theta_1 \end{cases}$ . Goal is to reject $H_0$ using a test statistic.

### Statistic
A (statistical) test is an statistic whose output is always either 0 or 1, and like an estimator, does not depend explicitly on the value of true unknown parameter.
A test $\psi$ has **level** $\alpha$ if $\alpha_\psi(\theta) \leq \alpha, \forall \theta \in \Theta_0$. and **asymptotic level** $\alpha$ if $\lim_{n\to\infty} P_\theta(\psi = 1) \leq \alpha$.

**A hypothesis-test** has the form

$\psi = \mathbf{1}\{T_n \geq c\}$

for some test statistic $T_n$ and threshold $c \in \mathbb{R}$. Threshold $c$ is usually $q_{\alpha/2}$
### Rejection region:

$R_\psi = \{T_n > c\}$

### Symmetric about zero and acceptance Region interval:

$\psi = \mathbf{1}\{|T_n| - c > 0\}.$

### Power of the test:
*worst possible result of Type II*:

$\pi_\psi = \inf_{\theta \in \Theta_1} (1 - \beta_\psi(\theta))$

Where $\beta_\psi$ is the probability of making a Type2 Error and $inf$ is the minimum.
### Two-sided test:

$H_1 : \theta \neq \Theta_0$
$\mathbf{1}(|T_n| > q_{\alpha/2})$

### One-sided tests:

$H_1 : \theta > \Theta_0$
$\mathbf{1}(T_n < -q_\alpha)H_1 \qquad : \theta < \Theta_0$
$\mathbf{1}(T_n > q_\alpha)$

### Type1 Error:
Test rejects null hypothesis $\psi = 1$ but it is actually true $H_0 = TRUE$ also known as the level of a test.
### Type2 Error:
Test does not reject null hypothesis $\psi = 0$ but alternative hypothesis is true

---

$H_1 = TRUE$

**Example:** Let $X_1, \ldots, X_n \overset{i.i.d.}{\sim} Ber(p^*)$.
Question: is $p^* = 1/2$.
$H_0 : p^* = 1/2; H_1 : p^* \neq 1/2$
If asymptotic level $\alpha$ then we need to standardize the estimated parameter $\hat{p} = \overline{X}_n$ first.

$T_n = \sqrt{n} \frac{|\overline{X}_n - 0.5|}{\sqrt{0.5(1-0.5)}}$
$\psi_n = \mathbf{1}(T_n > q_{\alpha/2})$

where $q_{\alpha/2}$ denotes the $q_{\alpha/2}$ quantile of a standard Gaussian, and $\alpha$ is determined by the required level of $\psi$. Note the absolute value in $T_n$ for this two sided test.
### Pivot:
Let $T_n$ be a function of the random samples $X_1, \ldots, X_n, \theta$. Let $g(T_n)$ be a random variable whose distribution is the same for all $\theta$. Then, $g$ is called a pivotal quantity or a pivot.
**Example:** let $X$ be a random variable with mean $\mu$ and variance $\sigma^2$. Let $X_1, \ldots, X_n$ be iid samples of $X$. Then,

$g_n \triangleq \frac{\overline{X}_n - \mu}{\sigma}$

is a pivot with $\theta = \begin{bmatrix} \mu & \sigma^2 \end{bmatrix}^T$ being the parameter vector (not the same set of paramaters that we use to define a statistical model).

### P-Value

The (asymptotic) p-value of a test $\psi_\alpha$ is the smallest (asymptotic) level $\alpha$ at which $\psi_\alpha$ rejects $H_0$. It is random since it depends on the sample. It can also interpreted as the probability that the test-statistic $T_n$ is realized given the null hypothesis.

If $pvalue \leq \alpha$ , $H_0$ is rejected by $\psi_\alpha$ at the (asymptotic) level $\alpha$

The smaller the p-value, the more confidently one can reject $H_0$.
### Left-tailed p-values:

$pvalue = \mathbb{P}(X \leq x | H_0)$
$= \mathbf{P}(Z < T_{n,\theta_0}(\overline{X}_n)))$
$= \Phi(T_{n,\theta_0}(\overline{X}_n))$
$Z \sim \mathcal{N}(0,1)$

### Right-tailed p-values:

$pvalue = \mathbb{P}(X \geq x | H_0)$

**Two-sided p-values:** If asymptotic, create normalized $T_n$ using parameters from $H_0$. Then use $T_n$ to get to probabilities.

$pvalue = 2min\{\mathbb{P}(X \leq x | H_0), \mathbb{P}(X \geq x | H_0)\}$
$\mathbb{P}(|Z| > |T_{n,\theta_0}(\overline{X}_n)|) = 2(1 - \Phi(T_n))$
$Z \sim N(0,1)$

---

## Comparisons of two proportions

Let $X_1, \ldots, X_n \overset{iid}{\sim} Bern(p_x)$ and $Y_1, \ldots, Y_n \overset{iid}{\sim} Bern(p_y)$ and be $X$ independent of $Y$. $\hat{p}_x = 1/n \sum_{i=1}^n X_i$ and $\hat{p}_x = 1/n \sum_{i=1}^n Y_i$

$H_0 : p_x = p_y; H_1 : p_x \neq p_y$
To get the asymptotic Variance use multi-variate Delta-method. Consider $\hat{p}_x - \hat{p}_y = g(\hat{p}_x, \hat{p}_y); g(x, y) = x - y$, then

$\sqrt{(n)}(g(\hat{p}_x, \hat{p}_y) - g(p_x - p_y)) \xrightarrow{(d)}$
$N(0, \nabla g(p_x - p_y)^T \Sigma \nabla g(p_x - p_y))$
$\Rightarrow N(0, p_x(1 - p_x) + p_y(1 - py))$

## Non-asymptotic Hypothesis tests (Tests)

### Chi squared

The $\chi_d^2$ distribution with $d$ degrees of freedom is given by the distribution of $Z_1^2 + Z_2^2 + \cdots + Z_d^2$, where $Z_1, \ldots, Z_d \overset{iid}{\sim} \mathcal{N}(0,1)$
If $V \sim \chi_d^2$:

$\mathbb{E} = \mathbb{E}[Z_1^2] + \mathbb{E}[Z_2^2] + \ldots + \mathbb{E}[Z_d^2] = d$

$Var(V) = Var(Z_1^2) + Var(Z_2^2) + \ldots + Var(Z_d^2) = 2d$

### Cochranes Theorem:
If $X_1, \ldots, X_n \overset{iid}{\sim} N(\mu, \sigma^2)$, then sample mean $\overline{X}_n$ and the sample variance $S_n$ are independent. The sum of squares of $n$ variables follows a chi squared distribution with (n-1) degrees of freedom:

$\frac{nS_n}{\sigma^2} \sim \chi_{n-1}^2$

If formula for unbiased sample variance is used:

$\frac{(n-1)S_n}{\sigma^2} \sim \chi_{n-1}^2$

### Student's T Test

Non-asymptotic hypothesis test for small samples (works on large samples too), data must be gaussian.

**Student's T distribution** with $d$ degrees of freedom: $t_d := \frac{Z}{\sqrt{V/n}}$ where $Z \sim \mathcal{N}(0,1)$ and $V \sim \chi_k^2$ are independent.

**Student's T test (one sample + two-sided):**

Let $X_1, \ldots, X_n \overset{iid}{\sim} N(\mu, \sigma^2)$ and suppose we want to test $H_0 : \mu = \mu_0 = 0$ vs. $H_1 : \mu \neq 0$.

Test statistic follows Student's T distribution:

$$T_n = \frac{Z}{S}$$
$$= \frac{\overline{X} - \mu}{\frac{\hat{\sigma}}{\sqrt{n}}}$$
$$= \frac{\sqrt{n}\frac{\overline{X}_n - \mu_0}{\sigma}}{\sqrt{\frac{\hat{S}_n}{\sigma^2}}}$$
$$\sim \frac{N(0,1)}{\sqrt{\frac{\chi^2_{n-1}}{n-1}}}$$
$$\sim t_{n-1}$$

Works bc. under $H_0$ the numerator $N(0,1)$ and the denominator $\frac{\hat{S}_n}{\sigma^2} \sim \frac{1}{n-1}\chi^2_{n-1}$ are independent by Cochran's Theorem.

Student's T test at level $\alpha$:
$$\psi_\alpha = \mathbf{1}\{|T_n| > q_{\alpha/2}(t_{n-1})\}$$

**Student's T test (one sample, one-sided):**
$$\psi_\alpha = \mathbf{1}\{T_n > q_\alpha(t_{n-1})\}$$

**Student's T test (two samples, two-sided):**
Let $X_1,...,X_n \overset{iid}{\sim} N(\mu_X, \sigma_X^2)$ and $Y_1,...,Y_n \overset{iid}{\sim} N(\mu_Y, \sigma_Y^2)$, suppose we want to test $H_0 : \mu_X = \mu_Y$ vs $H_1 : \mu_X \neq \mu_Y$.

$$T_{n,m} = \frac{\overline{X}_n - \overline{Y}_m}{\sqrt{\frac{\hat{\sigma}_X^2}{n} + \frac{\hat{\sigma}_Y^2}{m}}}$$

**Welch-Satterthwaite formula:**

When samples are different sizes we need to finde the Student's T distribution of: $T_{n,m} \sim t_N$

Calculate the degrees of freedom for $t_N$ with:

$$N = \frac{\left(\frac{\hat{\sigma}_X^2}{n} + \frac{\hat{\sigma}_Y^2}{m}\right)^2}{\frac{\hat{\sigma}_X^4}{n^2(n-1)} + \frac{\hat{\sigma}_Y^4}{m^2(m-1)}} \geq \min(n,m)$$

$N$ should be rounded down.

**Walds Test**

Squared distance of $\widehat{\theta}_n^{MLE}$ to true $\theta_0$ using the fisher information $I(\widehat{\theta}_n^{MLE})$ as metric.
Let $X_1,...,X_n \overset{iid}{\sim} \mathbf{P}_{\theta^*}$ for some true parameter $\theta^* \in \mathbb{R}^d$ and the maximum likelihood estimator $\widehat{\theta}_n^{MLE}$ for $\theta^*$.

Test $H_0 : \theta^* = \mathbf{0}$ vs $H_1 : \theta^* \neq \mathbf{0}$

Under $H_0$, the asymptotic normality of the MLE $\widehat{\theta}_n^{MLE}$ implies that:

$$\left\|\sqrt{n}\mathcal{I}(\mathbf{0})^{1/2}\left(\widehat{\theta}_n^{MLE} - \mathbf{0}\right)\right\|^2 \xrightarrow[n\to\infty]{(d)} \chi_d^2$$

**Test statistic:**
$$T_n = n(\widehat{\theta}_n^{MLE} - \theta_0)^\top I(\widehat{\theta}_n^{MLE})(\widehat{\theta}_n^{MLE} - \theta_0)$$
$$\xrightarrow[n\to\infty]{(d)} \chi_d^2$$

**Wald test** of level $\alpha$:
$$\psi_\alpha = \mathbf{1}\{T_n > q_\alpha(\chi_d^2)\}$$

**Likelihood Ratio Test**

Parameter space $\Theta \subseteq \mathbb{R}^d$ and $H_0$ is that parameters $\theta_{r+1}$ through $\theta_d$ have values $\theta_c^{r+1}$ through $\theta_d^c$ leaving the other $r$ unspecified. That is:
$H_0 : (\theta_{r+1},...,\theta_d)^T = \theta_{r+1...d} = \theta_0$

**Construct two estimators:**
$$\widehat{\theta}_n^{MLE} = argmax_{\theta \in \Theta}(\ell_n(\theta))$$
$$\widehat{\theta}_n^c = argmax_{\theta \in \Theta_0}(\ell_n(\theta))$$

**Test statistic:**
$$T_n = 2(\ell(X_1,..X_n|\widehat{\theta}_n^{MLE}) - \ell(X_1,..X_n|\widehat{\theta}_n^c)))$$

**Wilk's Theorem:** under $H_0$, if the MLE conditions are satisfied:
$$T_n \xrightarrow[n\to\infty]{(d)} \chi_{d-r}^2$$

**Likelihood ratio test** at level $\alpha$:
$$\psi_\alpha = \mathbf{1}\{T_n > q_\alpha(\chi_{d-r}^2)\}$$

**Implicit Testing**
Todo

**Goodness of Fit Discrete Distributions**

Let $X_1,...,X_n$ be iid samples from a categorical distribution. Test $H_0 : p = p^0$ against $H_1 : p \neq p^0$. Example: against the uniform distribution $p^0 = (1/K,...,1/K)^\top$.

**Test statistic** under $H_0$:
$$T_n = n\sum_{k=1}^K \frac{(\hat{p}_k - p_k^0)^2}{p_k^0} \xrightarrow[n\to\infty]{(d)} \chi_{K-1}^2$$

**Test at level alpha:**
$$\psi_\alpha = \mathbb{1}\{T_n > q_\alpha(\chi_{K-1}^2)\}$$

**Kolmogorov-Smirnov test**

**Kolmogorov-Lilliefors test**

**QQ plots**

**Heavier tails**: below > above the diagonal.
**Lighter tails**: above > below the diagonal.
**Right-skewed**: above > below > above the diagonal.
**Left-skewed**: below > above > below the diagonal.

# Distances between distributions
## Total variation distance
The total variation distance TV between the propability measures $P$ and $Q$ with a sample space $E$ is defined as:
$TV(\mathbf{P},\mathbf{Q}) = \max_{A \subseteq E}|\mathbf{P}(A) - \mathbf{Q}(A)|$,
Calculation with $f$ and $g$:
$$TV(\mathbf{P},\mathbf{Q}) =$$
$$\begin{cases} \frac{1}{2}\sum_{x \in E}|f(x) - g(x)|, \text{discr} \\ \frac{1}{2}\int_{x \in E}|f(x) - g(x)|dx, \text{cont} \end{cases}$$

Symmetry: $TV(\mathbf{P},\mathbf{Q}) = TV(\mathbf{Q},\mathbf{P})$
Positive: $TV(\mathbf{P},\mathbf{Q}) \geq 0$
Definite: $TV(\mathbf{P},\mathbf{Q}) = 0 \iff \mathbf{P} = \mathbf{Q}$
Triangle inequality: $TV(\mathbf{P},\mathbf{V}) \leq TV(\mathbf{P},\mathbf{Q}) + TV(\mathbf{Q},\mathbf{V})$

If the support of $\mathbf{P}$ and $\mathbf{Q}$ is disjoint:
$$TV(\mathbf{P},\mathbf{V}) = 1$$

TV between continuous and discrete r.v:
$$TV(\mathbf{P},\mathbf{V}) = 1$$

## KL divergence
The KL divergence (aka relative entropy) KL between probability measures $P$ and $Q$ with the common space $E$ and pmf/pdf functions $f$ and $g$ is defined as:
$$KL(\mathbf{P},\mathbf{Q}) =$$
$$\begin{cases} \sum_{x \in E} p(x)\ln\left(\frac{p(x)}{q(x)}\right), & \text{discr} \\ \int_{x \in E} p(x)\ln\left(\frac{p(x)}{q(x)}\right)dx, & \text{cont} \end{cases}$$

The KL divergence is not a distance measure! Always sum over the support of $P$!
Asymetric in general: $KL(\mathbf{P},\mathbf{Q}) \neq KL(\mathbf{Q},\mathbf{P})$
Nonnegative: $KL(\mathbf{P},\mathbf{Q}) \geq 0$
Definite: if $\mathbf{P} = \mathbf{Q}$ then $KL(\mathbf{P},\mathbf{Q}) = 0$
Does not satisfy triangle inequality in general: $KL(\mathbf{P},\mathbf{V}) \not\leq KL(\mathbf{P},\mathbf{Q}) + KL(\mathbf{Q},\mathbf{V})$

**Estimator of KL divergence:**
$$KL(\mathbf{P}_{\theta^*},\mathbf{P}_\theta) = \mathbb{E}_{\theta^*}\left[\ln\left(\frac{p_{\theta^*}(X)}{p_\theta(X)}\right)\right]$$
$$\widehat{KL}(\mathbf{P}_{\theta_*},\mathbf{P}_\theta) = const - \frac{1}{n}\sum_{i=1}^n log(p_\theta(X_i))$$

*maximize the above equation can derive the maximum likelihood*

# Maximum likelihood estimation
Let $\left\{E, (\mathbf{P}_\theta)_{\theta \in \Theta}\right\}$ be a statistical model associated with a sample of i.i.d. random variables $X_1, X_2,...,X_n$. Assume that there exists $\theta^* \in \Theta$ such that $X_i \sim \mathbf{P}_{\theta^*}$.
The **likelihood** of the model is the product of the $n$ samples of the pdf/pmf:

$$L_n(X_1, X_2,...,X_n,\theta) =$$
$$\begin{cases} \prod_{i=1}^n p_\theta(x_i) & \text{if } E \text{ is discrete} \\ \prod_{i=1}^n f_\theta(x_i) & \text{if } E \text{ is continous} \end{cases}$$

The maximum likelihood estimator is the (unique) $\theta$ that minimizes

$\widehat{KL}(\mathbf{P}_{\theta^*},\mathbf{P}_\theta)$ over the parameter space. (The minimizer of the KL divergence is unique due to it being strictly convex in the space of distributions once is fixed.)

$$\widehat{\theta}_n^{MLE} = argmin_{\theta \in \Theta}\widehat{KL}_n(\mathbf{P}_{\theta^*},\mathbf{P}_\theta)$$
$$= argmax_{\theta \in \Theta}\sum_{i=1}^n \ln p_\theta(X_i)$$
$$= argmax_{\theta \in \Theta}\ln\left(\prod_{i=1}^n p_\theta(X_i)\right)$$

Since taking derivatives of products is hard but easy for sums and $exp()$ is very common in pdfs we usually take the log of the likelihood function before maximizing it.

$$\ell((X_1, X_2,...,X_n,\theta) = ln(L_n(X_1, X_2,...,X_n,\theta))$$
$$= \sum_{i=1}^n ln(L_i(X_i,\theta))$$

Cookbook: set up the likelihood function, take log of likelihood function. Take the partial derivative of the loglikelihood function wrt. the parameter(s). Set the partial derivative(s) to zero and solve for the parameter.
If an indicator function on the pdf/pmf does not depend on the parameter, it can be ignored. If it depends on the parameter it can't be ignored because there is an discontinuity in the loglikelihood function. The maximum/minimum of the $X_i$ is then the maximum likelihood estimator.

**Fisher Information**
The Fisher information is the covariance matrix of the gradient of the loglikelihood function. It is equal to the negative expectation of the Hessian of the loglikelihood function and captures the negative of the expected curvature of the loglikelihood function.

Let $\theta \in \Theta \subset \mathbb{R}^d$ and let $\left(E, \{\mathbf{P}_\theta\}_{\theta \in \Theta}\right)$ be a statistical model. Let $f_\theta(\mathbf{x})$ be the pdf of the distribution $\mathbf{P}_\theta$. Then, the Fisher information of the statistical model is.

$$\mathcal{I}(\theta) = Cov(\nabla\ell(\theta)) =$$
$$= \mathbb{E}[\nabla\ell(\theta))\nabla\ell(\theta)^T] - \mathbb{E}[\nabla\ell(\theta)]\mathbb{E}[\nabla\ell(\theta)] =$$
$$= -\mathbb{E}[\mathbf{H}\ell(\theta)]$$

Where $\ell(\theta) = \ln f_\theta(\mathbf{X})$. If $\nabla\ell(\theta) \in \mathbb{R}^d$ it is a $d \times d$ matrix. The definition when the distribution has a pmf $p_\theta(\mathbf{x})$ is also the same, with the expectation taken with respect to the pmf.

Let $(\mathbb{R}, \{\mathbf{P}_\theta\}_{\theta \in \mathbb{R}})$ denote a continuous statistical model. Let $f_\theta(x)$ denote the pdf (probability density function) of the continuous distribution $\mathbf{P}_\theta$. Assume that $f_\theta(x)$ is twice-differentiable as a function of the parameter $\theta$.

Formula for the calculation of Fisher Information of $X$:

$$\mathcal{I}(\theta) = \int_{-\infty}^\infty \frac{\left(\frac{\partial f_\theta(x)}{\partial\theta}\right)^2}{f_\theta(x)} dx$$

Models with one parameter (ie. Bernoulli):

$$\mathcal{I}(\theta) = Var(\ell'(\theta))$$
$$\mathcal{I}(\theta) = -\mathbf{E}(\ell''(\theta))$$

Models with multiple parameters (ie. Gaussians):
$$\mathcal{I}(\theta) = -\mathbb{E}[\mathbf{H}\ell(\theta)]$$

Cookbook:

Better to use 2nd derivative.

- Find loglikelihood
- Take second derivative (=Hessian if multivariate)
- Massage second derivative or Hessian (isolate functions of $X_i$ to use with $-\mathbf{E}(\ell''(\theta))$ or $-\mathbb{E}[\mathbf{H}\ell(\theta)]$.
- Find the expectation of the functions of $X_i$ and subsitute them back into the Hessian or the second derivative. Be extra careful to subsitute the right power back. $\mathbb{E}[X_i] \neq \mathbb{E}[X_i^2]$.
- Don't forget the minus sign!

**Asymptotic normality of the maximum likelihood estimator**
Under certain conditions the MLE is asymptotically normal and consistent. This applies even if the MLE is not the sample average.
Let the true parameter $\theta^* \in \Theta$. Necessary assumptions:

- The parameter is identifiable
- For all $\theta \in \Theta$, the support $\mathbb{P}_\theta$ does not depend on $\theta$ (e.g. like in $Unif(0,\theta)$);
- $\theta^*$ is not on the boundary of $\Theta$;
- Fisher information $\mathcal{I}(\theta)$ is invertible in the neighborhood of $\theta^*$
- A few more technical conditions

The asymptotic variance of the MLE is the inverse of the fisher information.
$$\sqrt{n}(\widehat{\theta}_n^{MLE} - \theta^*) \xrightarrow[n\to\infty]{(d)} N_d(0,\mathcal{I}(\theta^*)^{-1})$$

# Method of Moments
Let $X_1,...,X_n \overset{iid}{\sim} \mathbf{P}_{\theta^*}$ associated with model $(\mathbb{E}, \{\mathbf{P}_\theta\}_{\theta \in \Theta})$, with $\mathbb{E} \subseteq \mathbb{R}$ and $\Theta \subseteq \mathbb{R}$, for some $d \geq 1$
Population moments:
$$m_k(\theta) = \mathbb{E}_\theta[X_1^k], 1 \leq k \leq d$$

Empirical moments:
$$\widehat{m_k}(\theta) = \overline{X_n^k} = \frac{1}{n}\sum_{i=1}^n X_i^k$$
Convergence of empirical moments:
$$\widehat{m_k} \xrightarrow[n\to\infty]{P,a.s.} m_k$$

$$(\widehat{m_1},...,\widehat{m_d}) \xrightarrow[n\to\infty]{P,a.s.} (m_1,...,m_d)$$
MOM Estimator $M$ is a map from the parameters of a model to the moments of its distribution. This map is invertible, (ie. it results into a system of equations that can be solved for the true parameter vector $\theta^*$). Find the moments (as many as parameters), set up system of equations, solve for parameters, use empirical moments to estimate.
$$\psi : \Theta \to \mathbb{R}^d$$
$$\theta \mapsto (m_1(\theta), m_2(\theta),...,m_d(\theta))$$
$$M^{-1}(m_1(\theta^*), m_2(\theta^*),...,m_d(\theta^*))$$

The MOM estimator uses the empirical moments:
$$M^{-1}\left(\frac{1}{n}\sum_{i=1}^n X_i, \frac{1}{n}\sum_{i=1}^n X_i^2,..., \frac{1}{n}\sum_{i=1}^n X_i^d\right)$$

Assuming $M^{-1}$ is continuously differentiable at $M(0)$, the asymptotical variance of the MOM estimator is:

$$\sqrt{(n)}(\widehat{\theta_n^{MM}} - \theta) \xrightarrow[n\to\infty]{(d)} N(0,\Gamma)$$

where,
$$\Gamma(\theta) = \left[\frac{\partial M^{-1}}{\partial\theta}(M(\theta))\right]^T \Sigma(\theta)\left[\frac{\partial M^{-1}}{\partial\theta}(M(\theta))\right]$$
$$\Gamma(\theta) = \nabla_\theta(M^{-1})^T \Sigma \nabla_\theta(M^{-1})$$

$\Sigma_\theta$ is the covariance matrix of the random vector of the moments $(X_1^1, X_1^2..., X_1^d)$.

# Algebra
Absolute Value Inequalities:
$|f(x)| < a \Rightarrow -a < f(x) < a$
$|f(x)| > a \Rightarrow f(x) > a$ or $f(x) < -a$

# Matrixalgebra
$\|\mathbf{A}\mathbf{x}\|^2 = (\mathbf{A}\mathbf{x})^T(\mathbf{A}\mathbf{x}) = \mathbf{x}^T\mathbf{A}^T\mathbf{A}\mathbf{x}$
$\mathbf{x}^T\mathbf{A}^T\mathbf{A}\mathbf{x}$

# Calculus
Differentiation under the integral sign
$$\frac{d}{dx}\left(\int_{a(x)}^{b(x)} f(x,t)dt\right) = f(x,b(x))b'(x) - f(x,a(x))a'(x) + \int_{a(x)}^{b(x)} f_x(x,t)dt.$$

**Concavity in 1 dimension**

If $g : I \to \mathbb{R}$ is twice differentiable in the interval $I$:
concave:
if and only if $g''(x) \leq 0$ for all $x \in I$

strictly concave:
if $g''(x) < 0$ for all $x \in I$

convex:
if and only if $g''(x) \geq 0$ for all $x \in I$

strictly convex if:
$g''(x) > 0$ for all $x \in I$

## Multivariate Calculus

The Gradient $\nabla$ of a twice differntiable function $f$ is defined as:

$$\nabla f : \mathbb{R}^d \to \mathbb{R}^d$$

$$\theta = \begin{pmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_d \end{pmatrix} \mapsto \left.\begin{pmatrix} \frac{\partial f}{\partial \theta_1} \\ \frac{\partial f}{\partial \theta_2} \\ \vdots \\ \frac{\partial f}{\partial \theta_d} \end{pmatrix}\right|_\theta$$

### Hessian

The Hessian of $f$ is a symmetric matrix of second partial derivatives of $f$

$$\mathbf{H}h(\theta) = \nabla^2 h(\theta) =$$
$$\begin{pmatrix} \frac{\partial^2 h}{\partial \theta_1 \partial \theta_1}(\theta) & \cdots & \frac{\partial^2 h}{\partial \theta_1 \partial \theta_d}(\theta) \\ & \vdots & \\ \frac{\partial^2 h}{\partial \theta_d \partial \theta_1}(\theta) & \cdots & \frac{\partial^2 h}{\partial \theta_d \partial \theta_d}(\theta) \end{pmatrix} \in$$
$$\mathbb{R}^{d \times d}$$

A symmetric (real-valued) $d \times d$ matrix $\mathbf{A}$ is:

Positive semi-definite:
$\mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0 \quad$ for all $\mathbf{x} \in \mathbb{R}^d$.

Positive definite:
$\mathbf{x}^T \mathbf{A} \mathbf{x} > 0$ for all non-zero vectors $\mathbf{x} \in \mathbb{R}^d$

Negative semi-definite (resp. negative definite):

$\mathbf{x}^T \mathbf{A} \mathbf{x}$ is negative for all $\mathbf{x} \in \mathbb{R}^d - \{\mathbf{0}\}$.

Positive (or negative) definiteness implies positive (or negative) semi-definiteness.

If the Hessian is positive definite then $f$ attains a local minimum at $a$ (convex).

If the Hessian is negative definite at $a$, then f attains a local maximum at $a$ (concave).

If the Hessian has both positive and negative eigenvalues then $a$ is a saddle point for $f$.

## Multivariate Related

A random vector $\mathbf{X} = \left(X^{(1)},\ldots,X^{(d)}\right)^T$ of dimension $d \times 1$ is a vector-valued function from a probability space $\omega$ to $\mathbb{R}^d$:

$$\mathbf{X} : \Omega \longrightarrow \mathbb{R}^d$$

$$\omega \longrightarrow \begin{pmatrix} X^{(1)}(\omega) \\ X^{(2)}(\omega) \\ \vdots \\ X^{(d)}(\omega) \end{pmatrix}$$

where each $X^{(k)}$, is a (scalar) random variable on $\Omega$.

PDF of $\mathbf{X}$: joint distribution of its components $X^{(1)},\ldots,X^{(d)}$.

CDF of $\mathbf{X}$:

$$\mathbb{R}^d \to [0,1]$$

$$\mathbf{x} \mapsto \mathbf{P}(X^{(1)} \leq x^{(1)},\ldots, X^{(d)} \leq x^{(d)}).$$

The sequence $\mathbf{X}_1, \mathbf{X}_2,\ldots$ converges in probability to $\mathbf{X}$ if and only if each component of the sequence $X_1^{(k)}, X_2^{(k)},\ldots$ converges in probability to $X^{(k)}$.

### Expectation of a random vector

The expectation of a random vector is the elementwise expectation. Let $\mathbf{X}$ be a random vector of dimension $d \times 1$.

$$\mathbb{E}[\mathbf{X}] = \begin{pmatrix} \mathbb{E}[X^{(1)}] \\ \vdots \\ \mathbb{E}[X^{(d)}] \end{pmatrix}.$$

The expectation of a random matrix is the expected value of each of its elements. Let $X = \{X_{ij}\}$ be an $n \times p$ random matrix. Then $\mathbb{E}[X]$, is the $n \times p$ matrix of numbers (if they exist):

$$\mathbb{E}[X] =$$
$$\begin{bmatrix} \mathbb{E}[X_{11}] & \mathbb{E}[X_{12}] & \cdots & \mathbb{E}[X_{1p}] \\ \mathbb{E}[X_{21}] & \mathbb{E}[X_{22}] & \cdots & \mathbb{E}[X_{2p}] \\ \vdots & \vdots & \ddots & \vdots \\ \mathbb{E}[X_{n1}] & \mathbb{E}[X_{n2}] & \cdots & \mathbb{E}[X_{np}] \end{bmatrix}$$

Let $X$ and $Y$ be random matrices of the same dimension, and let $A$ and $B$ be conformable matrices of constants.

$$\mathbb{E}[X + Y] = \mathbb{E}[X] + \mathbb{E}[Y]$$
$$\mathbb{E}[AXB] = A\mathbb{E}[X]B$$

### Covariance Matrix

Let $X$ be a random vector of dimension $d \times 1$ with expectation $\mu_X$.
Matrix outer products!

$$\Sigma = \mathbb{E}[(X - \mu_X)(X - \mu_X)^T] =$$

$$\mathbb{E}\left[\begin{pmatrix} X_1 - \mu_1 \\ X_2 - \mu_2 \\ \cdots \\ X_d - \mu_d \end{pmatrix}[X_1 - \mu_1, X_2 - \mu_2,\ldots, X_d - \mu_d]\right]$$

$$\Sigma = Cov(X) = \begin{bmatrix} \sigma_{11} & \sigma_{12} & \cdots & \sigma_{1d} \\ \sigma_{21} & \sigma_{22} & \cdots & \sigma_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{d1} & \sigma_{d2} & \cdots & \sigma_{dd} \end{bmatrix}$$

The covariance matrix $\Sigma$ is a $d \times d$ matrix. It is a table of the pairwise covariances of the elements of the random vector. Its diagonal elements are the variances of the elements of the random vector, the off-diagonal elements are its covariances. Note that the covariance is commutative e.g. $\sigma_{12} = \sigma_{21}$

Alternative forms:

$$\Sigma = \mathbb{E}[XX^T] - \mathbb{E}[X]\mathbb{E}[X]^T =$$
$$= \mathbb{E}[XX^T] - \mu_X \mu_X^T$$

Let the random vector $X \in \mathbb{R}^d$ and $A$ and $B$ be conformable matrices of constants.

$$Cov(AX + B) = Cov(AX) = ACov(X)A^T = A\Sigma A^T$$

Every Covariance matrix is positive definite.

$$\Sigma < 0$$

### Gaussian Random Vectors

A random vector $\mathbf{X} = (X^{(1)},\ldots,X^{(d)})^T$ is a Gaussian vector, or multivariate Gaussian or normal variable, if any linear combination of its components is a (univariate) Gaussian variable or a constant (a "Gaussian"variable with zero variance), i.e., if $\alpha^T \mathbf{X}$ is (univariate) Gaussian or constant for any constant non-zero vector $\alpha \in \mathbb{R}^d$.

### Multivariate Gaussians

The distribution of, $X$ the $d$-dimensional Gaussian or normal distribution, is completely specified by the vector mean $\mu = \mathbb{E}[\mathbf{X}] = (\mathbb{E}[X^{(1)}],\ldots,\mathbb{E}[X^{(d)}])^T$ and the $d \times d$ covariance matrix $\Sigma$. If $\Sigma$ is invertible, then the pdf of $X$ is:

$$f_{\mathbf{X}}(\mathbf{x}) =$$
$$\frac{1}{\sqrt{(2\pi)^d \det(\Sigma)}} e^{-\frac{1}{2}(\mathbf{x}-\mu)^T \Sigma^{-1}(\mathbf{x}-\mu)},$$
$$\mathbf{x} \in \mathbb{R}^d$$

Where $\det(\Sigma)$ is the determinant of $\Sigma$, which is positive when $\Sigma$ is invertible. If $\mu = 0$ and $\Sigma$ is the identity matrix, then $X$ is called a standard normal random vector .

If the covariant matrix $\Sigma$ is diagonal, the pdf factors into pdfs of univariate Gaussians, and hence the components are independent.

The linear transform of a gaussian $X \sim N_d(\mu, \Sigma)$ with conformable matrices $A$ and $B$ is a gaussian:

$$AX + B = N_d(A\mu + b, A\Sigma A^T)$$

### Multivariate CLT

Let $X_1,\ldots,X_d \in \mathbb{R}^d$ be independent copies of a random vector $X$ such that $\mathbb{E}[x] = \mu$ ($d \times 1$ vector of expectations) and $Cov(X) = \Sigma$

$$\sqrt{(n)}(\overline{X_n} - \mu) \xrightarrow[n\to\infty]{(d)} N(0, \Sigma)$$

$$\sqrt{(n)}\Sigma^{-1/2}\overline{X_n} - \mu \xrightarrow[n\to\infty]{(d)} N(0, I_d)$$

Where $\Sigma^{-1/2}$ is the $d \times d$ matrix such that $\Sigma^{-1/2}\Sigma^{-1/2} = \Sigma^1$ and $I_d$ is the identity matrix.

### Multivariate Delta Method

Given a sequence of r.v $(\mathbb{T}_n)_{n\geq 1}$ satisfying $\sqrt{n}(\mathbb{T}_n - \Theta) \xrightarrow[(d)]{n\to\infty} \mathbb{T}$ and $g$ is continuously differentiable at $\Theta$ then

$$\sqrt{n}(g(\mathbb{T}_n) - g(\Theta)) \xrightarrow[(d)]{n\to\infty} \Delta g(\Theta^T)\mathbf{T}$$

where $\Delta g = [\Delta g_1 \Delta g_2 \Delta g_3 \ldots \Delta g_k]$