

# Frequentist Statistics

Muchang Bahng

December 2022

## Contents

<b>1</b>	<b>Foundations</b>	<b>3</b>
1.1	Sampling Distributions . . . . .	3
1.2	Concentration of Measure . . . . .	5
1.3	Kullback Leibler Divergence . . . . .	6
1.4	Bounding Maximum of Random Variables . . . . .	7
1.5	Big-O, Little-O Notation . . . . .	7
<b>2</b>	<b>Point Estimation</b>	<b>8</b>
2.1	Sampling from Gaussians . . . . .	11
2.2	Method of moments . . . . .	12
2.3	Maximum likelihood estimation . . . . .	12
2.4	Least squares estimation . . . . .	12
<b>3</b>	<b>Confidence Intervals</b>	<b>12</b>
3.1	CIs for means, proportions, and variances . . . . .	14
3.2	Bootstrap confidence intervals . . . . .	14
<b>4</b>	<b>Hypothesis Testing</b>	<b>14</b>
4.1	One Sample Z and T Tests . . . . .	14
4.2	Power of a test . . . . .	16
4.3	Common tests (t-test, z-test, chi-square test, F-test) . . . . .	16
4.4	Multiple testing problem . . . . .	16
<b>5</b>	<b>Regression Analysis</b>	<b>16</b>
5.1	Ordinary Least Squares . . . . .	16
5.2	Gauss-Markov Theorem . . . . .	18
5.3	Analysis of variance (ANOVA) . . . . .	18
<b>6</b>	<b>Time Series Analysis</b>	<b>18</b>
6.1	Properties of Processes . . . . .	19
6.1.1	Estimation . . . . .	21
6.1.2	Detecting White Noise . . . . .	22
6.2	Autoregressive (AR) Processes . . . . .	22
6.3	Moving Average (MA) Processes . . . . .	24
6.4	Linear Processes . . . . .	25
6.5	ARMA . . . . .	26
6.6	ARIMA . . . . .	26
6.7	Other . . . . .	26
6.8	Components of time series . . . . .	28
6.9	Stationarity and tests for stationarity (including ADF test) . . . . .	28
6.10	Autoregressive (AR) models . . . . .	28

6.11	Moving average (MA) models . . . . .	28
6.12	ARIMA models . . . . .	28
6.13	Forecasting techniques . . . . .	28
<b>7</b>	<b>Advanced Topics</b>	<b>28</b>
7.1	Generalized Linear Models . . . . .	28
7.2	Survival analysis . . . . .	28
7.3	Nonparametric methods . . . . .	28
7.4	Resampling methods (jackknife, bootstrap) . . . . .	28
<b>8</b>	<b>Practical Considerations</b>	<b>28</b>
8.1	Experimental design . . . . .	28
8.2	Sample size and power analysis . . . . .	28
8.3	Dealing with assumptions violations . . . . .	28
8.4	Interpretation and reporting of results . . . . .	28
<b>9</b>	<b>Cross Validation</b>	<b>28</b>
9.1	Leave 1 Out Cross Validation . . . . .	30
9.1.1	Generalized (Approximate) Cross Validation . . . . .	30
9.1.2	Cp Statistic . . . . .	30
9.2	K Fold Cross Validation . . . . .	30
9.3	Data Leakage . . . . .	30
9.4	Information Criterion . . . . .	30
	<b>References</b>	<b>30</b>

We assume the reader is familiar with measure-theoretic probability, and unlike in introductory probability, we throw away the convention that random variables are written with capital Latin letters (so  $x$  can also denote a random variable, which is useful if samples are not fixed). Statistics and probability seem like the same topic, but there are very fundamental differences. In probability, we are given some distribution and must compute certain probabilities. In statistics, we are given the results (the data) and must infer what distribution it came from.

## 1 Foundations

### 1.1 Sampling Distributions

#### Definition 1.1 (Population, Parameters)

When conducting a statistical study, there is a set of items or events which is of interest for some experiment. This can be modeled with some probability space  $(\Omega, \mathcal{F}, \mathbb{P})$ . Usually, we are interested in some numerical property of this population, and so we implicitly define a random variable  $X : \Omega \rightarrow \mathbb{R}$  that induces some distribution  $X \sim P$ , which we call the **population**. A statistical population can be a group of existing objects or a hypothetical and potentially infinite group of objects conceived as a generalization from experience.

With this, we can interpret the population  $X \sim P$  as a random variable, and we often call this the **parent distribution**. We are often interested in its population **parameters**, which can be any measured quantity of a population that summarizes or describes an aspect of it. In generality, the parameter of population  $X$  is denoted  $\theta$ , and it is a fixed value.

#### Example 1.1 (Populations)

Here are some examples of populations:

1. We let  $\Omega$  be the discrete sample space of all hands in poker, and our random variable will assign a numerical ranking to each hand 0 for no hand, 1 for pairs, 2 for two pairs, etc.
2.  $\Omega$  is the sample space of all individuals in the U.S. and we can construct a random variable  $X$  that assigns to each individual their height. Even though  $\Omega$  is finite, we can interpret it as continuous, which leads to a continuous distribution  $X \sim P$ .

#### Example 1.2 (Parameters)

Some population parameters can be:

1. The true mean of the population  $\mu_X = \mathbb{E}[X]$
2. The true variance of the population  $\sigma_X^2 = \text{Var}[X]$
3. In the height example, we can set  $X = (X_1, X_2, X_3)$  and try to construct a linear regression model that predicts  $X_3$  from  $X_1, X_2$ . Theoretically, this is  $\mathbb{E}[X_3 \mid X_1, X_2]$ , and we must find the best function of form  $x_3 = a + b_1x_1 + b_2x_2$  that is closest to the conditional expectation. There does exist a unique one, and so  $a, b_1, b_2$  are all population parameters.

In general, the population is the total set of all relevant things that we are interested in. The specific quantity of the actual population is called the **population parameter**, e.g. the true mean  $\mu$  or the true variance  $\sigma^2$  of  $X$ , usually denoted with  $\theta$ . But usually, these parameters are not known since the population is too big to experimentally measure, so we must try and estimate it with samples. This is the entire point of statistics; otherwise, we would already know everything we want to know.

**Definition 1.2 (Samples)**

From the population  $X \sim P$  (which still has unknown distribution), we can take  $n$  **samples** by considering iid  $x_1, x_2, \dots, x_n \sim P$ .

1. We should note that the samples  $x_i$  are random variables themselves. Not fixing them yet and still considering them in generality as random objects allows us to do more theoretical calculations.
2. Once these samples have been realized (i.e.  $\omega \in \Omega$  is realized, and all  $x_i$ 's are also realized), we can treat them as fixed values.

Sometimes, we may not assume independence, but for most cases we do. A common rule is that if the sample size  $n$  is less than 10% of the population size, then we can assume independence.

**Definition 1.3 (Empirical Distribution)**

Now given that we have these iid samples, we can construct the **empirical distribution**  $\hat{X} \sim \hat{P}$ , defined as the discrete distribution that assigns probability  $1/n$  to each value  $x_i$  for  $i \in [n]$ . In other words, we have

$$\mathbb{P}(\hat{X} = x) = \frac{1}{n} \text{ for } x \in \{x_1, \dots, x_n\} \quad (1)$$

We can write the CDF of the empirical distribution, called the **empirical distribution function**, as the sum of indicators

$$F_X(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}_{[x_i, +\infty)}(x) \quad (2)$$

As expected, we would expect the empirical distribution to converge to the actual distribution.

**Theorem 1.1 (Glivenko–Cantelli theorem)**

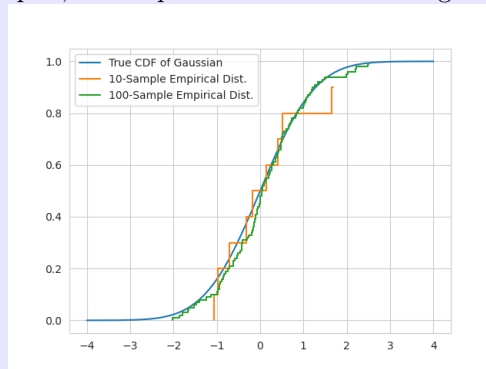
The empirical distribution of iid samples  $x_1, \dots, x_n \sim P_n$  converges almost surely to  $X \sim P$  as  $n \rightarrow \infty$ . More specifically, given that the CDF of  $X$  is  $F$  and the CDF of  $P_n$  is the step function  $F_n$ , we have

$$\|F_n - F\|_\infty = \sup_{x \in \mathbb{R}} |F_n(x) - F(x)| \rightarrow 0 \quad (3)$$

almost surely as  $n \rightarrow \infty$ .

**Example 1.3 (Empirical Distribution of Standard Gaussian)**

We expect the empirical distribution of the standard Gaussian to converge. Indeed, numerical results show that for 10 and 100 samples, the empirical CDF does converge to the true CDF.



## 1.2 Concentration of Measure

Now let's move on to concentration inequalities, which say that the probability that a random variable is greater than something is bounded by something. These probability bounds are extremely useful in of themselves. It allows us to talk about convergence theory, which tells us what happens to a statistic, such as  $\bar{X}$ , as I get more and more data. The first inequality exploits the fact that the tails of a Gaussian RV decay very quickly, and a lot of concentration inequalities attempt to mimic this exponential bound but for non-Gaussian distributions.

### Theorem 1.2 (Gaussian Tail Inequality)

Given  $X \sim \mathcal{N}(0, 1)$ , the inequality says that the probability of  $X$  taking values past a certain  $t$  decays exponentially.

$$\mathbb{P}(|X| > t) \leq \frac{2e^{-t^2/2}}{t}$$

If we have  $x_1, \dots, x_n \sim \mathcal{N}(0, 1)$ , then

$$\mathbb{P}(|\bar{X}| > t) \leq \frac{2}{\sqrt{nt}} e^{-nt^2/2}$$

We can assume that the coefficient is less than 1 if  $n$  is large. The above tells us that this bound exponentially decays with  $t$  but also with the number of samples  $n$ .

### Theorem 1.3 (Markov's Inequality)

Given a nonnegative random variable  $X > 0$ , we have

$$\mathbb{P}(X > t) \leq \frac{\mathbb{E}[X]}{t}$$

#### Proof.

We have

$$\begin{aligned} \mathbb{E}[X] &= \int_0^\infty x p_X(x) dx \\ &\geq \int_t^\infty x p_X(x) dx \\ &= t \int_t^\infty p_X(x) dx \\ &= t \mathbb{P}(X > t) \end{aligned}$$

### Theorem 1.4 (Chebyshev's Inequality)

Given a random variable  $X$  with mean  $\mu = \mathbb{E}[X]$ , we have

$$\mathbb{P}(|x - \mu| > t) \leq \frac{\text{Var}(X)}{t^2}$$

**Theorem 1.5 (Hoeffding's Inequality)**

Let  $x_1, x_2, \dots, x_n$  be independent (not necessarily identical) random variables s.t.  $a_i \leq X_i \leq b_i$  almost surely. Consider the random variable  $\bar{X} = \frac{1}{n}(x_1 + \dots + x_n)$ . Then, for all  $t > 0$ ,

$$\mathbb{P}(|\bar{X} - \mathbb{E}[\bar{X}]| \geq t) \leq 2 \exp\left(-\frac{2n^2 t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right)$$

Now in addition to bounding probabilities, we would like to bound expectations.

**Theorem 1.6 (Cauchy-Schwartz)**

Given random variables  $X, Y$ , it is often hard to compute the expectation of  $XY$  since it is hard to compute the distribution of it (sums are easy). But we can bound it as

$$|\mathbb{E}[XY]| \leq \mathbb{E}[|XY|] \leq \sqrt{\mathbb{E}[X] \mathbb{E}[Y]}$$

**Theorem 1.7 (Jensen's Inequality)**

Given  $g$  a convex function and  $X$  a random variable, we have

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

**1.3 Kullback Leibler Divergence**

Now a popular metric between PMFs/PDFs is the KL divergence.

**Definition 1.4 (Kullback-Leibler Divergence)**

Given random variable  $X$  and  $Y$ ,

1. If they are discrete with PMFs  $P$  and  $Q$ , the KL-divergence is defined

$$D_{KL}(P \parallel Q) := \sum_x P(x) \log \left( \frac{P(x)}{Q(x)} \right) = \mathbb{E} \left[ \log \frac{P(x)}{Q(x)} \right]$$

where we can interpret the expectation as  $X \sim p$ .

2. If they are continuous with PDFs  $p$  and  $q$ , the KL-divergence is defined

$$D_{KL}(p \parallel q) := \int p(x) \log \left( \frac{p(x)}{q(x)} \right) dx = \mathbb{E} \left[ \log \frac{P(x)}{Q(x)} \right]$$

where we can interpret the expectation as  $X \sim p$ .

We should prove that this is indeed a metric.

1. The fact that  $D_{KL}(p \parallel p) = 0$  is obvious.
2. To prove that  $D_{KL}(p \parallel q) \geq 0$ , we use Jensen's inequality

$$-D_{KL}(p \parallel q) = \mathbb{E} \log \frac{q(X)}{p(X)} \leq \log \mathbb{E} \frac{q(X)}{p(X)} = \log \int \frac{q(x)}{p(x)} p(x) dx = \log(1) = 0$$

It is a common trick to switch the log and the expectation using Jensen's.

## 1.4 Bounding Maximum of Random Variables

Given that we have  $n$  samples  $x_1, \dots, x_n \sim P$ , it is conventional to index them with open brackets to denote order

$$X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$$

Our goal is now to find the distribution of  $X_{(n)} = \max_i X_i$ . If we know the distribution of  $P$ ,  $\mathbb{P}(X_{(n)} \leq x)$  is just the probability that all the  $X_i$ 's are less than  $x$ , and by independence we can product out the CDFs, differentiate to get the PDF, and compute. So it's not too hard to do this theoretically, but in practice this is hard to do since we don't exactly know  $P$ .

So if you didn't know the  $X_i$ 's, the best you can assume is that

$$\mathbb{E} \max\{x_1, \dots, x_n\}$$

is going to grow like  $n$ . But if we can bound the MGF with  $\mathbb{E}e^{tX} \leq e^{t^2\sigma^2/2}$ , then we can show that  $\mathbb{E} \max X_i$  doesn't grow like  $n$ , but rather like  $\log n$ .

$$\mathbb{E} \max X_i \leq \sigma \sqrt{2 \log n}$$

## 1.5 Big-O, Little-O Notation

Going back to calculus, if we have a function  $f : X \rightarrow \mathbb{R}$ , we can say that

1.  $f(x) = O(g(x))$  if  $f$  is of the same order as  $g$ . That is, they grow at the same rate

$$\frac{f(x)}{g(x)} \rightarrow c \text{ as } x \rightarrow \infty$$

for some constant  $c$ .

2.  $f(x) = o(g(x))$  if  $f$  is negligible w.r.t.  $g$ . That is,  $f$  is infinitesimal w.r.t.  $g$ .

$$\frac{f(x)}{g(x)} \rightarrow 0 \text{ as } x \rightarrow \infty$$

for some constant  $c$ .

Now there is a probabilistic notation as well. The concept of boundedness translates to being able to capture most of the mass of the random variable within some interval, and infinitesimality translates to the probability mass concentrating around 0.

### Definition 1.5 ( $O_p, o_p$ Notation)

Let  $x_1, x_2, \dots$  be a sequence of random variables.

1.  $x_n = o_p(1)$  if

$$\mathbb{P}(|Y_n| > \epsilon) \rightarrow 0 \text{ as } n \rightarrow \infty$$

for all  $\epsilon > 0$ . This means that  $x_n$  gets more and more concentrated around 0.

2.  $x_n = O_p(1)$  if for all  $\epsilon > 0$ , then there exists a  $C$  s.t.

$$\mathbb{P}(|x_n| > C) \leq \epsilon$$

for all large  $n$ . That is, we can always trap the majority of the probability mass of  $x_n$  within the interval  $[-C, C]$ . This must hold for all  $x_n$  with  $n > N$ , so the mass can't "escape" to infinity. We can think of it as the distribution is "settling down" and not shooting off to somewhere.

3.  $x_n = o_p(a_n)$  means that

$$\frac{x_n}{a_n} = o_p(1)$$

4.  $x_n = O_p(a_n)$  means that

$$\frac{x_n}{a_n} = O_p(1)$$

### Theorem 1.8 ()

Given  $Y_1, \dots, Y_n \sim \text{Bernoulli}(p)$ , let  $\hat{p} = \frac{1}{n} \sum Y_i$ . Then

$$\hat{p}_n - p = o_p(1)$$

which is also written  $\hat{p}_n = p + o_p(1)$ , which means that the random variable  $\hat{p}_n$  is some constant  $p$  plus a random variable that is going to 0.

### Proof.

By Hoeffding's inequality,

$$\mathbb{P}(|\hat{p}_n - p| > \epsilon) \leq 2e^{-2n\epsilon^2}$$

which goes to 0 as  $n \rightarrow \infty$ .

### Example 1.4 ()

Given  $Y_1, \dots, Y_n \sim \text{Bernoulli}(p)$ , let  $\hat{p} = \frac{1}{n} \sum Y_i$ . Then,

$$\hat{p} - p = O_p\left(\frac{1}{\sqrt{n}}\right)$$

## 2 Point Estimation

### Definition 2.1 (Sample Statistic, Estimators and Estimates)

Now given a population  $X$ , we would like to use the  $n$  iid samples  $x_1, \dots, x_n$  to estimate a parameter  $\theta$  of interest with our own random variable/value  $\hat{\theta}_n$ , called a **sample statistic**. We must note the dual nature of the sample statistic as a random variable and a value is similar to that of samples.

1. The statistic  $\hat{\theta}_n$  is a random variable itself, referred to as the **estimator**. More specifically, it is a function  $\hat{\theta}_n : \mathbb{R}^n \rightarrow \mathbb{R}$  of the  $n$  samples, i.e. a transformation of random variables

$$\hat{\theta}_n = \hat{\theta}_n(x_1, x_2, \dots, x_n) \quad (4)$$

This makes  $\hat{\theta}_n$  also a random variable, which attempts to estimate the true  $\theta$ , which is some unknown fixed value. Since  $\hat{\theta}_n$  is a random variable, it has its own distribution, called the **sampling distribution** of  $\hat{\theta}_n$ .

2. Once these samples  $x_i$  have been realized, the estimator realizes and the value realized is now called the **estimate**.

This sampling distribution is a distribution of the statistic  $\hat{\theta}_n$ , and this forms a separate distribution with its own mean and variance.

1. the mean of the sampling distribution is denoted  $\mu_{\hat{\theta}_n}$
2. the standard deviation of the sampling distribution is denoted  $\sigma_{\hat{\theta}_n}^2$ , also called the **standard error**.

We would want these estimators to have three properties:

1. unbiasedness



2. consistency
3. efficiency

We would like the sampling distribution of our statistic to give us good estimate in two ways.  $\hat{\theta}_n$  should not be too far off from the actual parameter  $\theta$  (bias is small), and  $\hat{\theta}_n$  should not fluctuate too widely (variance of  $\hat{\theta}_n$  should be small).

**Definition 2.2 (Bias, Variance of Estimator)**

Given an estimator  $\hat{\theta}$  of a sample  $x_1, \dots, x_n$  estimating population parameter  $\theta$ , the **sampling bias** refers to

$$\text{Bias}(\hat{\theta}) = |\mathbb{E}[\hat{\theta}] - \theta| \quad (5)$$

and the **sampling variance** refers to

$$\text{Var}(\hat{\theta}) = \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2] \quad (6)$$

A good rule of thumb to remember is that statistics is about replacing expectations with averages.

$$\mathbb{E} \mapsto \frac{1}{n} \sum_i \quad (7)$$

This is really the fundamental quality of statistics. Then after that we can do some fancy things, like minimizing something or manipulating another, but every single time we see an expectation just replace it with an average.

**Definition 2.3 (Sample Mean)**

Given a population  $X$  with  $\mu = \mathbb{E}[X]$  and  $\sigma^2 = \text{Var}(X)$ , our estimator for  $\mu$  is simply the average of the  $n$  samples  $x_1, \dots, x_n$ , called the **sample mean** or the **sampling distribution of the sample mean**.

$$\bar{x}_n = \hat{\mu}_n = \frac{1}{n}(x_1 + \dots + x_n) \quad (8)$$

This gives us the sampling distribution of the sample means. The mean and standard deviation (i.e. standard error) of  $\bar{x}_n$  is denoted  $\mu_{\bar{x}_n}$  and  $\sigma_{\bar{x}_n}$ .

1. The mean of  $\bar{x}_n$  is  $\mu$ .

$$\mu_{\bar{x}_n} = \mu \quad (9)$$

because

$$\mathbb{E}[\bar{x}_n] = \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n x_i\right] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[x_i] = \mathbb{E}[x] = \mu \quad (10)$$

2. The variance of  $\bar{x}_n$  is  $\sigma^2/n$ , i.e. the standard error of  $\bar{x}_n$  is  $\sigma_{\bar{x}_n} = \sigma/\sqrt{n}$ .

$$\sigma_{\bar{x}_n} = \frac{\sigma}{\sqrt{n}} \quad (11)$$

because

$$\sigma_{\bar{x}_n}^2 = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(x_i) = \frac{1}{n} \text{Var}(x) = \frac{\sigma^2}{n} \quad (12)$$

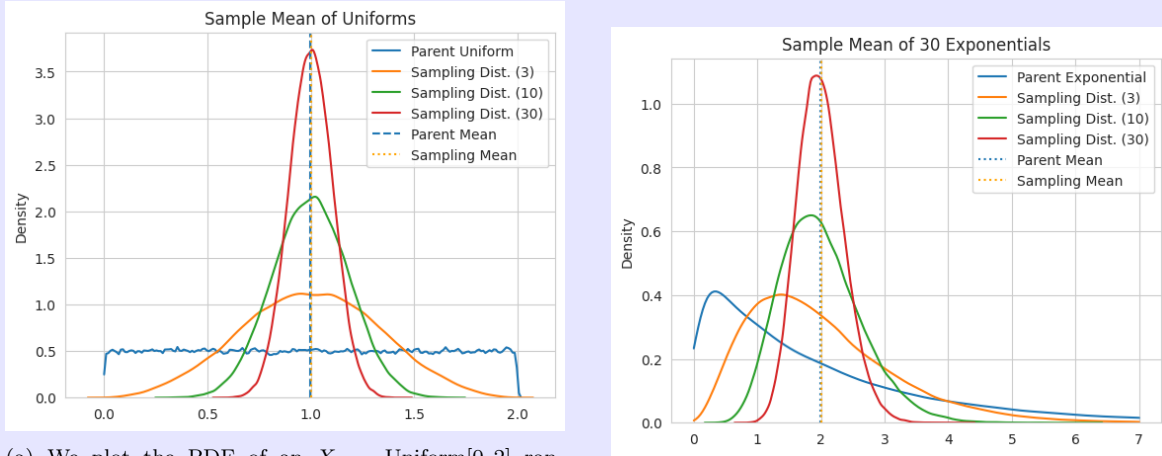
Practically, this tells us that when trying to estimate the value of a population mean, due to the factor of  $1/\sqrt{n}$ , reducing the error on the estimate by a factor of 2 requires acquiring 4 times as many observations in the sample. But realistically, the true standard deviation  $\sigma$  is unknown, and so the standard error of the mean is usually estimated by replacing  $\sigma$  with the sample standard deviation  $S$  instead.

$$\sigma_{\bar{x}_n} \approx \frac{S}{\sqrt{n}} \quad (13)$$

3. By CLT,  $\bar{x}_n$  converges to  $\mathcal{N}(\mu, \sigma^2/n)$  in distribution as  $n \rightarrow +\infty$  (but in practicality, we assume this for  $n \geq 30$ ). The fact that its mean and variance is  $\mu$  and  $\sigma^2/n$  isn't that impressive. What is really impressive is that no matter what the distribution of  $x$  is, the sampling distribution of the mean will be Gaussian.

### Example 2.1 (Sample Means)

Here are some figures of sample means. Note that with a uniform parent distribution, the sampling distribution of its mean looks like a Gaussian even without a large  $n$ . However, this is not necessarily true for different parent distributions, such as the exponential.



(a) We plot the PDF of an  $X \sim \text{Uniform}[0, 2]$  random variable by taking 100k samples. We also take 100k samples from the sampling distribution of the mean  $\bar{X}_3, \bar{X}_{10}, \bar{X}_{30}$ . We can see that the standard deviation decreases by a factor of  $\sqrt{n}$ .

(b) We plot the PDF of an  $X \sim \text{Exponential}(1.5)$  random variable by taking 100k samples. We also take 100k samples from the sampling distribution of the mean  $\bar{X}_3, \bar{X}_{10}, \bar{X}_{30}$ .

Figure 1

If the parent distribution is normal, then we don't even need CLT to claim that the sampling distribution of the sample mean is normal, since sums of normals are normal.

Now the variance of the population is defined to be  $\sigma^2 = \mathbb{E}[(X - \mathbb{E}[X])^2]$ , and by our rule of thumb, we can replace the expectations with sample means, by first setting  $\mathbb{E}[X] = \hat{\mu}$  and averaging out the values  $(X - \hat{\mu})^2$ .

### Definition 2.4 (Sample Variance)

Given a population  $X$ , our estimator for  $\sigma^2 = \mathbb{E}[(X - \mathbb{E}[X])^2]$  is simply the average of the squared distances of the  $n$  samples  $\{(x_i - \hat{\mu})^2\}_{i=1}^n$ .

$$S_n^2 = \hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)^2 \quad (14)$$

The mean and standard deviation of  $S_n^2$  is denoted  $\mu_{S_n^2}$  and  $\sigma_{S_n^2}$ . Note that there is a small difference that the sum for variance is divided by  $n - 1$  rather than  $n$ , since we want it to be unbiased, but we will correct this later.

While the CLT states that the sampling distribution of the sample mean will look approximately Gaussian, we do not have this luxury when looking at the sampling distribution of sample variance.

**Example 2.2 (Sample Variance)**

Take a look at the following sampling distributions of the sample variance. There does not seem to be strong signs of convergence to a Gaussian. Their means do not align either.

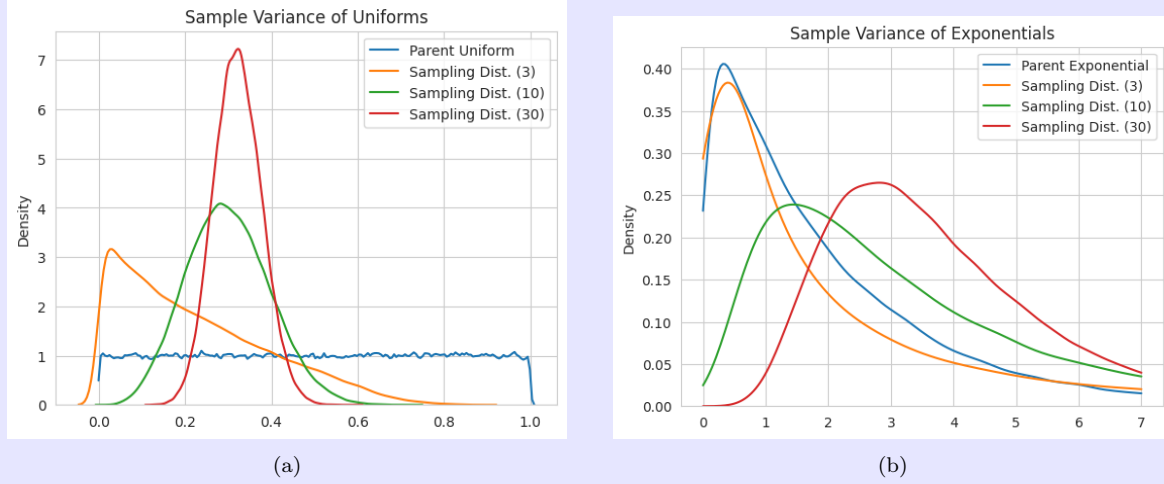


Figure 2

**2.1 Sampling from Gaussians**

Now if we assume that the parent distribution is Gaussian, then we can conclude some extra things and more kinds of distributions arise. Let  $x_1, \dots, x_n \sim \mathcal{N}(\mu, \sigma^2)$ , with  $\bar{x}_n$  the sample mean and  $S_n^2$  the sample variance. Say that we want to find the distribution of  $\bar{x}_n$ .

1. In the unrealistic case where we know the true  $\sigma^2$ , we don't even need to consider the sample variance. From the basic property of Gaussians, we know that  $\bar{x}_n \sim \mathcal{N}(\mu, \sigma^2/n)$ , or after standardizing,

$$\frac{\bar{x}_n - \mu}{\sigma/\sqrt{n}} \sim \mathcal{N}(0, 1) \quad (15)$$

2. In the realistic case where we don't know the true  $\sigma^2$ , we should replace it with our sample variance  $S^2$ , and it turns out that because of this extra uncertainty in the variance, our sampling distribution follows the student-t distribution, which can be interpreted as a mixture of Gaussians with differing variances.

$$\frac{\bar{x}_n - \mu}{S/\sqrt{n}} \sim \text{StudentT}(n - 1) \quad (16)$$

Now if we are interested in finding the distribution of  $S_n^2$ :

1. In the unrealistic case where the know the true  $\mu$ , we don't need to consider the sampling distribution of  $\bar{x}_n$ . We have

$$S_n^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 \sim \text{Gamma}\left(\frac{n}{2}, \frac{n}{2\sigma^2}\right) \quad (17)$$

2. In the realistic case where we don't know  $\mu$ , we have

$$\frac{n-1}{\sigma^2} S_n^2 = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \bar{x}_n)^2 \sim \chi^2(n-1) \quad (18)$$

## 2.2 Method of moments

## 2.3 Maximum likelihood estimation

## 2.4 Least squares estimation

# 3 Confidence Intervals

Recall that the central limit theorem says that given a sequence of iid random variables  $x_1, \dots, x_n$  coming from a random variable with true mean  $\mu$  and variance  $\sigma^2$ , the sample mean is similar to a  $\mathcal{N}(\mu, \sigma^2/n)$  random variable. That is, the sample mean converges in distribution

$$\bar{X}_n \xrightarrow{\text{dist}} \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right) \quad (19)$$

as  $n \rightarrow \infty$ . Another way to state it is that the normalized sample mean is similar to a standard Gaussian.

$$\frac{\bar{x}_n - \mu}{\sigma_{\bar{x}_n}} = \frac{\bar{x}_n - \mu}{\sigma/\sqrt{n}} \xrightarrow{\text{dist}} \mathcal{N}(0, 1) \quad (20)$$

So, given that we have enough samples, I will perfectly understand its fluctuations. Now let's introduce some definitions that will allow us to unify some ideas into simpler notation: the realized value  $x$ , the number of standard deviations it is away from the mean, and the probability that it takes that value (or more extreme).

### Definition 3.1 (z-score)

Given a  $\mathcal{N}(\mu, \sigma^2)$  distribution, the **z-score** of a number  $x \in \mathbb{R}$  is defined to be the number of standard deviations away from the mean.

$$z = \frac{x - \mu}{\sigma} \quad (21)$$

### Definition 3.2 (Percentile)

Given  $X \sim \mathcal{N}(0, 1)$  and significance level  $\alpha \in [0, 1]$ , let us define  $q_\alpha \in \mathbb{R}$  as the point where

$$\mathbb{P}(X \geq q_\alpha) = \alpha \quad (22)$$

i.e. the  $100\alpha$ th percentile of the standard normal. Note that given  $X \sim \mathcal{N}(0, 1)$ , we have

$$\mathbb{P}(|X| > q_{\alpha/2}) = \alpha \quad (23)$$

Now given  $x_1, \dots, x_n$  from a population  $X$  with mean  $\mu$  and standard deviation  $\sigma$ , let  $\bar{x}_n$  be the sampling distribution of the mean. By virtue of the central limit theorem, we can write

$$\mathbb{P}\left(\left|\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}}\right| \geq q_{\alpha/2}\right) \approx \alpha \iff \mathbb{P}\left(\left|\frac{\bar{X}_n - \mu}{\sigma\sqrt{n}}\right| \leq q_{\alpha/2}\right) \approx 1 - \alpha \quad (24)$$

which implies that with probability  $1 - \alpha$ , we have

$$\bar{X}_n \in \left[\mu - q_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \mu + q_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right] \iff \mu \in \left[\bar{X}_n - q_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{X}_n + q_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right] \quad (25)$$

This is how we construct a confidence interval. In other words, as  $n$  becomes large (ideally at least 30), the probability that an interval around our sample mean contains the actual mean  $\mu$  can be approximated by a Gaussian. But note that CI requires to know the actual standard deviation  $\sigma$ . There are three ways to deal with this:

1. This may actually be known from the start, especially if we are working with calibrated devices with standard devices that have been experimentally verified.

2. We can simply bound  $\sigma$ , depending on what kind of random variable we are working with. For example, given  $X \sim \text{Bernoulli}(p)$ , its standard deviation is bounded by  $\sigma = \sqrt{p(1-p)} \leq \frac{1}{2}$ , so we can create a confidence interval that is larger than any other confidence interval we can make if we had known the true  $\sigma$ .

$$p \in \left[ \bar{X}_n - q_{\alpha/2} \frac{1}{2\sqrt{n}}, \bar{X}_n + q_{\alpha/2} \frac{1}{2\sqrt{n}} \right] \quad (26)$$

3. We can approximate  $\sigma$  with the sample standard deviation  $S$ , which turns out to be an unbiased estimator.

### Example 3.1 (Proportion of Right-Side Kissers)

We have observed 80 out of 124 right-side kisses, resulting in a sample estimate of  $\hat{p} = 0.645$ . Given that we want a confidence interval of 95%, we want an  $\alpha = 0.05$ , implying a the value  $q_{\alpha/2} = q_{0.025} = 1.96$ . So, with probability 0.95, we have

$$p \in \left[ 0.645 - \frac{1.96}{2\sqrt{124}}, 0.645 + \frac{1.96}{2\sqrt{124}} \right] = [0.56, 0.73] \quad (27)$$

If we had, say 3 observations, rather than 124, we would have a 95% confidence interval of  $p \in [0.10, 1.23]$ , which is terrible, but in this case even CLT is not valid.

### Example 3.2 (Proportion of Voters)

Given that we sample  $n = 100$  people from a city's population to ask whether they support candidate A or B, we have 54 people who support candidate A, so  $\hat{p} = 0.54$ . Say that we want a 95% confidence interval, which leads to  $q_{\alpha/2} = q_{0.025} = 1.96$ . So, with probability 0.95, we have

$$p \in \left[ 0.54 - 1.96 \frac{\sigma}{\sqrt{100}}, 0.54 + 1.96 \frac{\sigma}{\sqrt{100}} \right] \quad (28)$$

and by substituting  $\sigma$  for  $S = \sqrt{0.54(1 - 0.54)} \approx 0.5$ , we get

$$p \in \left[ 0.54 - 1.96 \frac{0.284}{\sqrt{100}}, 0.54 + 1.96 \frac{0.284}{\sqrt{100}} \right] = [0.44, 0.64] \quad (29)$$

An interpretation of confidence intervals is that if you keep on sampling  $\bar{x}$  or  $\hat{p}$  and construct 95% CIs, then 95% of the time these intervals will contain the true mean  $\mu$  or proportion  $p$  (or more if we had bounded the CI with a bigger interval).

### Example 3.3 ()

We survey 6250 teachers to ask whether they think computers are essential for teaching. 250 were randomly selected and 142 felt that they were essential. Let's construct a 99% confidence interval for the proportion of teachers who felt that computers were essential. We would like to construct a CI for the true  $\mu = p$ , and we have  $\bar{x} = 142/250 = 0.568$ .

1. 99% confidence corresponds to  $\alpha = 0.01$ , which corresponds to a z-score of  $q_{\alpha/2} = 2.576$ .
2. The parent distribution is  $\text{Bernoulli}(p)$ , with  $\mu = p$  and  $\sigma = \sqrt{p(1-p)}$ . The sampling distribution of  $\bar{x}$  has  $\mu_{\bar{x}} = p$  also and  $\sigma_{\bar{x}} = \sigma/\sqrt{n}$ .
3. We need to know the details of the sampling distribution, but we don't know  $\sigma$ , which is needed to calculate  $\sigma_{\bar{x}}$ . However, we can estimate it using the sample standard deviation  $S = \sqrt{0.568(1 - 0.568)} = 0.5$ .
4. Our sampling distribution has standard deviation  $\sigma_{\bar{x}} \approx S/\sqrt{n} = 0.5/\sqrt{250} = 0.031$ , and our z-score was 2.576, so our 99% confidence interval is 2.576 standard deviations from our mean.

That is, with probability 0.99,

$$p \in [0.568 - 2.576 \cdot 0.031, 0.568 + 2.576 \cdot 0.031] = [0.488144, 0.647856] \quad (30)$$

### 3.1 CIs for means, proportions, and variances

### 3.2 Bootstrap confidence intervals

## 4 Hypothesis Testing

A significance test is a method used to decide whether the data at hand sufficiently supports a particular hypothesis. The hypothesis to be tested is called the **alternative hypothesis**, denoted  $H_1$  or  $H_a$ , and the status quo is called the **null hypothesis**, denoted  $H_0$ . Assuming that  $H_0$  is true, we compute the likelihood of the data happening. If the sample is not too unlikely (past some significance level), we fail to reject  $H_0$ , and if there is strong evidence, we reject  $H_0$ .  $H_0$  and  $H_a$  can be devised in countless ways.

#### Example 4.1 ()

There are countless test statistics we can build, but here are some common examples,

1. Proportion: Company A produces circuit boards, but 10% of them are defective. Company B claims that they produce fewer defective circuit boards.

$$H_0 : p = 0.10 \text{ versus } H_a : p < 0.10 \quad (31)$$

2. Means: It is known that the average height of boys in KIS is 176cm. Ben claims that the average height is lower than this.

$$H_0 : \mu = 176 \text{ versus } H_a : \mu < 176 \quad (32)$$

3. Difference of Means: If  $\mu_1$  and  $\mu_2$  denote the true average breaking strengths of the same type of twine produced by two different companies. Jenny claims that the  $\mu_1 - \mu_2 > 5$ .

$$H_0 : \mu_1 - \mu_2 = 0 \text{ versus } H_a : \mu_1 - \mu_2 > 5 \quad (33)$$

### 4.1 One Sample Z and T Tests

Let us have some population  $X \sim P$  and a null hypothesis that claims  $H_0 : \mu = \theta_0$ . Since we are interested in the mean, we would like to use CLT or some other theorem to determine what the distribution of the mean of  $n$  samples  $\bar{x}_n$  looks like (either Normal or Student T centered around  $\theta_0$  and scaled down by factor of  $\sqrt{n}$ ). When we actually sample, the value  $\bar{x}_n = \hat{\theta}$  is realized, and we would like to see if sampling  $\hat{\theta}$  from the distribution centered around  $\theta_0$  is likely, usually after normalizing. If it isn't, then we reject  $H_0$ .

How do we decide whether to use the z-test or the t-test? It is known that  $\text{StudentT}(n-1)$  converges to  $\mathcal{N}(0, 1)$  in distribution as  $n \rightarrow +\infty$ . Therefore, depending on the context of the problem, at a certain point  $N$  (usually  $N = 30$  or perhaps higher for skewed distributions), the difference between these two are negligible.

1. Z-test: if we know the population variance  $\sigma^2$ , but it is rarely the case that we actually know  $\sigma^2$ .
2. T-test: if we do not know the population variance  $\sigma^2$ , which we then substitute for the sample variance  $S^2$ .
3. Z-test: if we do not know the population variance (which we substitute for  $S^2$ ), but our sample size is greater than  $N$ , then we can approximate the  $t$ -distribution with our normal, allowing us to use the Z-test again.

In general, the alternative to the null hypothesis  $H_0 : \theta = \theta_0$  will look like one of the following three assertions:

1. Two-Sided Test:  $H_a : \theta \neq \theta_0$
2. One-Sided Test:  $H_a : \theta > \theta_0$  (in which case the null hypothesis is  $\theta \leq \theta_0$ )
3. One-Sided Test:  $H_a : \theta < \theta_0$  (in which case the null hypothesis is  $\theta \geq \theta_0$ )

Now we must still quantify *how* unlikely our sample mean  $\theta$  must be compared to  $\theta_0$  in order to reject the null hypothesis. This is where we specify our **significance level**, denoted by  $\alpha$  (common values 0.10, 0.05, 0.01). This specifies the tail-regions in which  $\theta$  will land in with probability  $\alpha$ . Usually, working with general normal/t distributions is tedious, so we can rescale them and use their z/t-scores.

#### Definition 4.1 (Z-score)

Given a value  $x$  sampled from distribution  $X \sim \mathcal{N}(\mu, \sigma^2)$ , its **z-score** is defined to be the number of standard deviations away from the mean.

$$z := \frac{x - \mu}{\sigma} \quad (34)$$

Now given a significance level  $\alpha \in [0, 1]$ , let  $z_\alpha$  be the value such that the measure of a standard normal distribution past  $z_\alpha$  is  $1 - \alpha$  (i.e. the  $100\alpha$  percentile).  $z_\alpha$  is called the **critical z-value**.

#### Definition 4.2 (T-score)

Given a value  $x$  sampled from distribution  $X \sim \text{StudentT}(n)$ , its **t-score** is defined to be the number of standard deviations away from the mean.

#### Example 4.2 ()

A factory has a machine that dispenses 80mL of fluid in a bottle. An employee believes the average amount of fluid is not 80mL. Using 40 samples, he measures the average amount dispensed by the machine to be 78mL with a sample standard deviation of 2.5.

1. Let the true mean be  $\mu$  and true standard deviation be  $\sigma$ . The null hypothesis is  $H_0 : \mu = 80$  and the alternative is  $H_1 : \mu \neq 80$ , making this a two-sided test.
2. We don't know the true standard deviation  $\sigma$ , so we must use the sample standard deviation  $S$ . This requires us to use the  $t$ -test, but since  $n > 30$ , we can invoke CLT and state that  $\bar{x}_{40}$  is (approximately) Gaussian with mean  $\mu$  and standard deviation  $S/\sqrt{n}$ . So, we use the  $z$ -test.
3. At a 95% confidence level, we have  $\alpha = 0.05$ , and our rejection region is  $(-\infty, z_{0.025}] \cup [z_{0.975}, +\infty)$ . Since we are looking at a standard Gaussian, we have by symmetry  $z_{0.025} = -1.96$  and  $z_{0.975} = 1.96$ , and our critical  $z$ -value is  $z^* = 1.96$ .
4. So the  $z$ -score for 78 is

$$z = \frac{\bar{x} - \mu_0}{s/\sqrt{n}} = \frac{78 - 80}{2.5/\sqrt{40}} = -5.06 \quad (35)$$

which is definitely in the reject region. So this tells us that we can reject the null hypothesis with a 95% level of confidence.

#### Example 4.3 ()

A company manufactures car batteries with an average life span of 2 or more years. An engineer believes this value to be less. Using 10 samples, he measures the average life span to be 1.8 years with a standard deviation of 0.15.

1. Let the true mean be  $\mu$  and true standard deviation be  $\sigma$ . The null hypothesis is  $H_0 : \mu \geq 2$

- and the alternative is  $H_1 : \mu < 2$ , making this a one-sided test.
2. We don't know the true standard deviation  $\sigma$ , so we must use the sample standard deviation  $S$ . This requires us to use the  $t$ -test, especially since  $n = 10$  is not large enough for us to invoke CLT.
  3. At a 99% confidence level, we have  $\alpha = 0.01$ , and our rejection region is  $(-\infty, t_{0.01}] = (-\infty, -2.82]$ .
  4. The  $t$ -score for the observed mean value is

$$t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}} = \frac{1.8 - 2}{0.15/\sqrt{10}} = -4.22 \quad (36)$$

which is definitely in the reject region. So this tells us that we can reject the null hypothesis with a 99% level of confidence.

We may have to account for errors. There is always a chance that our evidence leads us to an incorrect conclusion, and we have names for this.

#### Definition 4.3 (Errors)

Given a hypothesis test where we look for evidence supporting our alternative claim,

1. A **type 1 error** is when the null hypothesis is rejected, but it is true (false positive).
2. A **type 2 error** is when we fail to reject the null hypothesis, when it is false (false negative).

## 4.2 Power of a test

## 4.3 Common tests (t-test, z-test, chi-square test, F-test)

## 4.4 Multiple testing problem

# 5 Regression Analysis

Now we will talk about regression analysis from a statistical point of view. Regression can be used to approximate the relationship between two random variables (through a smooth function) and can be used for casual inference. Essentially, linear regression attempts to model the conditional distribution  $Y | X$ .

## 5.1 Ordinary Least Squares

If we use a squared loss function, this is called **ordinary least squares**. It is a well known fact that the true regressor that minimizes this loss is

$$f^*(x) = \mathbb{E}[Y | X = x] \quad (37)$$

which is the conditional expectation of  $Y$  given  $X$ . This is the true regressor function, which is the best approximation of  $Y$  over the  $\sigma$ -algebra generated by  $X$ . This may or may not be linear.

### Theorem 5.1 (Least Squares Solution For Linear Regression)

Given the design matrix  $\mathbf{X}$ , we can present the linear model in vectorized form:

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}) \quad (38)$$

The solution that minimizes the squared loss is

$$\begin{aligned} \boldsymbol{\beta} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \in \mathbb{R}^d \\ \text{Var}(\hat{\boldsymbol{\beta}}) &= \hat{\sigma}^2 (\mathbf{X}^T \mathbf{X})^{-1} \in \mathbb{R}^{d \times d} \end{aligned}$$



**Proof.**

The errors can be written as  $\epsilon = \mathbf{Y} - \mathbf{X}\beta$ , and you have the following total sum of squared errors:

$$S(\beta) = \epsilon^T \epsilon = (\mathbf{Y} - \mathbf{X}\beta)^T (\mathbf{Y} - \mathbf{X}\beta)$$

We want to find the value of  $\beta$  that minimizes the sum of squared errors. In order to do this, remember the following matrix derivative rules when differentiating with respect to vector  $\mathbf{x}$ .

1.  $\mathbf{x}^T \mathbf{A} \mapsto \mathbf{A}$
2.  $\mathbf{x}^T \mathbf{A} \mathbf{x} \mapsto 2\mathbf{A}\mathbf{x}$

Now this should be easy.

$$\begin{aligned} S(\beta) &= \mathbf{Y}^T \mathbf{Y} - \beta^T \mathbf{X}^T \mathbf{Y} - \mathbf{Y}^T \mathbf{X} \beta + \beta^T \mathbf{X}^T \mathbf{X} \beta \\ &= \mathbf{Y}^T \mathbf{Y} - 2\mathbf{Y}^T \mathbf{X} \beta + \beta^T \mathbf{X}^T \mathbf{X} \beta \\ \frac{\partial}{\partial \beta} S(\beta) &= -2\mathbf{X}^T \mathbf{Y} + 2\mathbf{X}^T \mathbf{X} \beta \end{aligned}$$

and setting it to  $\mathbf{0}$  gives

$$2\mathbf{X}^T \mathbf{X} \beta - 2\mathbf{X}^T \mathbf{Y} = 0 \implies \mathbf{X}^T \mathbf{X} \beta = \mathbf{X}^T \mathbf{Y}$$

and the variance of  $\beta$ , by using the fact that  $\text{Var}[\mathbf{A}\mathbf{X}] = \mathbf{A} \text{Var}[\mathbf{X}] \mathbf{A}^T$ , is

$$\text{Var}(\hat{\beta}) = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \sigma^2 \mathbf{I} \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} = \sigma^2 (\mathbf{X}'\mathbf{X})^{-1} (\mathbf{X}'\mathbf{X}) (\mathbf{X}'\mathbf{X})^{-1} = \sigma^2 (\mathbf{X}'\mathbf{X})^{-1}$$

But we don't know the true  $\sigma^2$ , so we estimate it with  $\hat{\sigma}^2$  by taking the variance of the residuals. Therefore, we have

$$\begin{aligned} \beta &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \in \mathbb{R}^d \\ \text{Var}(\hat{\beta}) &= \hat{\sigma}^2 (\mathbf{X}^T \mathbf{X})^{-1} \in \mathbb{R}^{d \times d} \end{aligned}$$

Note that we have assumed that  $\mathbf{X}^T \mathbf{X}$  was invertible in order for such a solution to be unique, i.e.  $\mathbf{X}$  must be full rank. This process breaks down when it isn't invertible, e.g. if there are repetitions in the features (one feature is a linear combination of the others and hence not full column rank). We will talk more about this soon.

**Definition 5.1 (Hat Matrix)**

For convenience of notation, let's call

$$\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \tag{39}$$

the  $n \times n$  **hat matrix**, which is essentially a projection of the observed  $y_i$ 's to the predictions.

$$\hat{\mathbf{y}} = \mathbf{H}\mathbf{y} \tag{40}$$

**Lemma 5.1 (Properties)**

The hat matrix is an orthogonal projection matrix that projects to the column space of  $\mathbf{X}$ .

Note that this parallels the orthogonal projection of conditional expectation to the true function onto the subspace of  $X$  measurable functions. Except that we are not doing this in function space, but rather the sample space  $\mathbb{R}^n$ .

We can also see that the residuals  $\hat{\epsilon}_i = y_i - \hat{y}_i$  has the property that

$$\hat{\epsilon} = \mathbf{y} - \hat{\mathbf{y}} = (\mathbf{I}_n - \mathbf{H})\mathbf{y} \quad (41)$$

Now if we look back to the derivative of the loss  $S$ , we really want to set

$$\mathbf{X}^T(\mathbf{y} - \mathbf{X}\hat{\beta}) = \mathbf{X}^T\hat{\epsilon} = \mathbf{0} \quad (42)$$

## 5.2 Gauss-Markov Theorem

At this point, we have only talked about the mathematical properties of the least squares regression, but now let's talk about some statistical properties. In machine learning, we talk about some assumptions (homoscedacity, uncorrelated residuals, etc.), and we now formalize their need.

### Theorem 5.2 (Gauss-Markov Theorem)

Given a dataset with

1. mean zero residuals  $\mathbb{E}[\epsilon_i] = 0$ , i.e.  $\mathbb{E}[\mathbf{Y} | \mathbf{X}] = \mathbf{X}\beta$ .
2. homoscedacity  $\text{Var}[\epsilon_i] = \sigma^2 < \infty$  for all  $i$ ,
3. uncorrelated residuals  $\text{Cov}(\epsilon_i, \epsilon_j) = 0$  for all  $i \neq j$ . This and the previous assumption can be combined into  $\text{Cov}[\mathbf{Y} | \mathbf{X}] = \sigma^2 \mathbf{I}_n$ .

We were concerned with estimating the parameters  $\beta_1, \dots, \beta_d$ . Now let's generalize this and consider the problem of estimating, for some known constants  $c_1, \dots, c_{d+1}$ , the point estimator

$$\theta = c_1\beta_1 + c_2\beta_2 + \dots + c_d\beta_d + c_{d+1} \quad (43)$$

Then the estimator

$$\hat{\theta} = c_1\hat{\beta}_1 + c_2\hat{\beta}_2 + \dots + c_d\hat{\beta}_d + c_{d+1} \quad (44)$$

where  $\hat{\beta}_i$  is clearly an unbiased estimator of  $\theta$  and it is a linear estimator of  $\theta$ , i.e.

$$\hat{\theta} = \sum_{i=1}^n b_i y_i \quad (45)$$

for some known (given  $\mathbf{X}$ ) constants  $b_i$ . Then, the Gauss-Markov theorem states that the estimator  $\hat{\theta}$  has the smallest (best) variance among *all* linear unbiased estimators of  $\theta$ , i.e.  $\hat{\theta}$  is BLUE.

## 5.3 Analysis of variance (ANOVA)

## 6 Time Series Analysis

If we try sticking to linear algebra, we hope to model time series of the form

$$X_t = f(t) + w_t \quad (46)$$

so that we can decompose to a deterministic process followed by some white noise. There are several ways to approach this, including kernel smoothing, moving average smoothing, or cubic spline smoothing. However, this falls short when you look the residuals. They will follow some pattern that must be removed due to autocorrelation.

In linear regression, one of the fundamental assumptions was independence of errors. Ideally, we would also like independence of features, but this is usually not true (in fact, in extreme cases, multicollinearity can screw us up). The relaxation of these assumptions helps us transition from linear regression to time series analysis. Let's go over some basic things with new terms.

**Definition 6.1 (Time Series)**

A stochastic process

$$\{X_1, \dots, X_t, \dots\} \quad (47)$$

of random variables indexed by time  $t$  is a **time series**. The stochastic behavior of  $\{X_t\}$  is determined by specifying the PDF/PMF

$$p(x_{t_1}, \dots, x_{t_m}) \quad (48)$$

for all finite collections of time indices

$$\{(t_1, \dots, t_m), m < \infty\} \quad (49)$$

i.e. all finite-dimensional distributions of  $X_t$ .

**Definition 6.2 (White Noise)**

**White noise**  $w_t$  is a random variable indexed by time  $t$  satisfying

1.  $\mathbb{E}[w_t] = 0$
2.  $\text{Var}[w_t] = \sigma^2$
3.  $\text{Cov}[w_t, w_s] = 0$  for  $s \neq t$ . That is, they are uncorrelated but not necessarily independent.

Note that this third condition can be strengthened to independence or uncorrelated Gaussians, which automatically imply independence.

**6.1 Properties of Processes**

Now let's define some properties. We will start with the time series analogue of covariance and correlation.

**Definition 6.3 (Autocovariance)**

The **autocovariance** between two time steps  $t, s$  of process  $\{X_t\}$  is defined

$$K_X(s, t) = \text{Cov}(X_t, X_s) \quad (50)$$

**Definition 6.4 (Autocorrelation)**

The **autocorrelation** is

$$\rho_X(s, t) = \frac{K_X(s, t)}{\sqrt{K_X(s, s) K_X(t, t)}} \quad (51)$$

**Definition 6.5 (Cross Covariance)**

Given two stochastic processes  $\{X_t\}, \{Y_t\}$ , the **cross covariance** is

$$K_{XY}(t, s) = \text{Cov}(X_t, Y_s) \quad (52)$$

and the **cross correlation** is

$$\rho_{XY}(t, s) = \frac{K_{XY}(t, s)}{K_X(t, s) K_Y(s, s)} \quad (53)$$

It is used to model the correlations between two related products with a certain time lag perhaps.

**Definition 6.6 (Stationarity)**

There are two types of stationarity.

1. A **weakly stationary** or **covariance stationary** process means that its mean and autocovariance are invariant to time shifts. That is, for all  $r$ ,

$$\mathbb{E}[X_t] = \mathbb{E}[X_{t+r}] = \mu \quad (54)$$

$$\text{Var}[X_t] = \text{Var}[X_{t+r}] = \sigma_X^2 \quad (55)$$

$$K_X(t, s) = K_X(t + r, s + r) \quad (56)$$

$$(57)$$

2. A **strongly stationary** process means that any joint distribution function of a finite set of time steps is invariant to time shifts. That is, for any  $r > 0$  and finite collection of time points  $t_1, \dots, t_k$ ,

$$F(X_{t_1}, \dots, X_{t_k}) = F(X_{t_1+r}, \dots, X_{t_k+r}) \quad (58)$$

where  $F$  is the joint pdf and equality means almost everywhere equality.

Clearly, weakly stationary implies strongly stationary, and the difference is that weakly stationary has invariance in the first two moments while strongly stationary holds for all moments.

**Theorem 6.1 ()**

It immediately follows that for a stationary process  $X_t$ , the autocovariance function can be defined

$$K_X(s, t) = K_X(s - t, 0) = K_X(\tau) \quad (59)$$

for some difference between the time points, called the lag. From this, we can see that  $\text{Var}[X_t] = K_X(0)$ , so the autocorrelation can be defined as

$$\rho_X(\tau) = \frac{K_X(\tau)}{K_X(0)} \quad (60)$$

Stationary time series are very desirable, since if we do parameter estimation, we don't want to estimate parameters that are always changing. For example, in stationary processes, we know that the mean never changes, so we have a bunch of sample points to choose from, and if every wasn't stationary, then every  $X_t$  would have its own mean and we won't be able to estimate it. Similarly, we also know that for some fixed  $\tau$ , the autocorrelation does not change, so we can estimate  $K_X(\tau)$  with a bunch of fixed intervals of length  $\tau$ . Therefore, if we want to test for stationary of a fixed time process, we want to conduct a test where we want to find whether the autocovariance is relatively invariant. This gives us a bit of intuition.

**Theorem 6.2 ()**

Note the following properties.

1.  $K_X(\tau) = K_X(-\tau)$
2. By Cauchy-Schwartz,  $K_X(0)^2 = \text{Var}[X_t] \text{Var}[X_{t+r}] \geq \text{Cov}(X_t, X_{t+r}) = K_X(r)^2$ , so  $|K_X(\tau)| \leq K_X(0)$ .

Therefore, we would like to decompose a general time series to a stationary component and a nonstationary simple component, and do some statistics on the stationary one.

**Definition 6.7 (Joint Stationarity)**

Two processes  $X_t, Y_t$ , are said to be jointly stationary if both are individually stationary and also if the cross covariance function is also stationary. That is, for all  $r$ ,

$$K_{XY}(t, s) = K_{XY}(t + r, s + r) \quad (61)$$

**Definition 6.8 (Backshift Operator)**

The backshift operator  $B$  acts on time series by

$$BX_t = X_{t-1} \quad (62)$$

It can be iterated to get  $B^k X_t = X_{t-k}$  and can also be inverted to get a forward shift  $B^{-k} X_t = X_{t+k}$ . We can just think of this as (not necessarily linear?) operators between the function space of  $X$ -measurable functions.

**6.1.1 Estimation**

We should now try to estimate some parameters of a weakly stationary process.

**Theorem 6.3 (Sampling Distribution of Mean)**

We can already estimate the mean. We should get the mean of the mean and the variance of the mean.

1. The mean is trivial, since by linearity of expectation we can get

$$\hat{\mu} = \bar{X} = \frac{1}{T} \sum_{t=1}^T X_t \quad (63)$$

2. The variance is a bit more involved since there are covariance terms, so

$$\text{Var}[\bar{X}] = \text{Var}\left(\frac{1}{T} \sum_{t=1}^T X_t\right) \quad (64)$$

$$= \frac{1}{T^2} \sum_{t=1}^T \sum_{s=1}^T \text{Cov}(X_t, X_s) \quad (65)$$

$$= \frac{1}{T^2} \sum_{t=1}^T \sum_{s=1}^T K_X(|t-s|) \quad (66)$$

$$= \frac{1}{T} K_X(0) + \frac{2}{T} \sum_{z=1}^{T-1} \left(1 - \frac{z}{T}\right) K_X(z) \quad (67)$$

In the unrealistic situation where the  $X_t$ 's are uncorrelated, we have  $K_X(0) = \sigma^2$  and  $K_X(z) = 0$  for all  $z > 0$ , leaving us with  $\sigma^2/T$ .

**Theorem 6.4 (Sampling Distribution of Autocovariance)**

To estimate the autocovariance of a weakly stationary process, we can define the sample autocovariance function to be

$$\hat{K}_X(h) = \frac{1}{T} \sum_{t=1}^{T-h} (X_{t+h} - \bar{X})(X_t - \bar{X}) \quad (68)$$

Note that we divide by  $T$  rather than  $T-h$  so that this covariance is positive semidefinite. Note that as  $h$  gets bigger, the number of terms in the sum decreases giving less accurate estimation. Similarly, the sample autocorrelation function is

$$\hat{\rho}(h) = \frac{\hat{K}_X(h)}{\hat{K}_X(0)} \quad (69)$$

The sample cross covariance and cross correlation are

$$\hat{K}_{XY}(h) = \frac{1}{T} \sum_{t=1}^{T-h} (X_{t+h} - \bar{X})(Y_t - \bar{Y}) \text{ and } \hat{\rho}_{XY}(h) = \frac{\hat{K}_{XY}(h)}{\sqrt{\hat{K}_X(0) \hat{K}_Y(0)}} \quad (70)$$

Note that even though we can just plug these formulas and get the sample estimators for any time series, these don't mean anything if they are not stationary.

**6.1.2 Detecting White Noise**

Ultimately, the main goal of time series analysis is to transform the data into a white noise process. We want to first identify trends and patterns in the process, remove them, and hopefully get white noise. To actually detect if we have white noise, one way to do this is to look at the estimated autocorrelation function across  $h$ . Note that for white noise, we have a spike at  $h = 0$  to be 1 (since it is just the correlation of a variable with itself), and then it drops to 0 immediately (since by definition,  $w_s, w_t$  are uncorrelated). We would like to see this behavior within a certain confidence interval.

**6.2 Autoregressive (AR) Processes**

The assumptions are:

1. the data must be stationary (though it is not always stationary as it may contain a unit root)
2. the relationship between the variables and their lagged values must be linear (nonlinear gives large language models like LSTMs)
3. the error term should be white noise

**Definition 6.9 (Autoregressive Process)**

An **AR(p)** process encodes causality<sup>a</sup> into the white noise process. It is a stochastic process with mean 0 and of the form

$$X_t = w_t + \sum_{i=1}^p \phi_i X_{t-i} \quad (71)$$

where  $p$  is the hyperparameter of steps to look back,  $w_t$  is white noise with variance  $\sigma^2$ , and  $\phi_i$  are constants  $\neq 0$ . Using the backshift operator  $B$ , we can write the AR(p) process as

$$\Phi(B)X_t = w_t \quad (72)$$

where

$$\Phi(B) = \left(1 - \sum_{i=1}^p \phi_i B^i\right) \quad (73)$$

In fact, we have already seen this process many times.

### Example 6.1 (AR(p) Processes)

Consider the following.

1. AR(O) is simply a white noise process

$$X_t = w_t \quad (74)$$

2. AR(1) with  $\theta = 1$  gives us the formula

$$X_t = X_{t-1} + w_t \quad (75)$$

which is a random walk. It is also a Markov process and a martingale.

3. AR(1) of the form

$$X_t = a + X_{t-1} + w_t \quad (76)$$

is a random walk with drift.

4. AR(2) can be of form

$$X_t = X_{t-1} - 0.2X_{t-2} + w_t \quad (77)$$

5. AR(3) can be of form

$$X_t = X_{t-1} - 0.2X_{t-2} + 0.13X_{t-3} + w_t \quad (78)$$

Occasionally, it may be hard to determine the difference between the difference of AR(p) processes.

### Example 6.2 (AR(1) Processes)

Let's focus on the AR(1) process. Later on in linear processes, we see that the AR(1) process has a causal representation as a linear process.

$$X_t = \phi_1 X_{t-1} + w_t = \sum_{i=0}^{\infty} \phi_1^i w_{t-i} \quad (79)$$

This is stationary under certain conditions.

1. If  $\phi < 1$ , then the series is stationary.
2. If  $\phi = 1$ , this is a random walk which is not stationary.
3. If  $\phi > 1$ , then this process grows exponentially fast.

Now to determine weak stationarity, let's go back to the equation. Talk about unit root test.

### Definition 6.10 (Augmented Dickey-Fuller Test)

The Augmented Dickey-Fuller (ADF) test is a statistical test used to determine whether a time series is stationary or not. Here's a step-by-step explanation of how the ADF test is typically implemented:

1. **Model Specification.** The ADF test is based on an autoregressive model. The general form is:

$$\Delta Y_t = \alpha + \beta t + \gamma Y_{t-1} + \delta_1 \Delta Y_{t-1} + \cdots + \delta_{p-1} \Delta Y_{t-p+1} + \varepsilon_t \quad (80)$$

Where:

- $\Delta Y_t$  is the first difference of the series at time  $t$
- $\alpha$  is the constant term
- $\beta t$  is the time trend term
- $\gamma$  and  $\delta$  are coefficients
- $\varepsilon_t$  is the error term

<sup>a</sup>on how a random variable  $Y$  is *caused* by another RV  $X$ .

- $p$  is the lag order
- 2. **Determine the lag order ( $p$ ):**
  - This can be done using information criteria like AIC or BIC
  - Or by starting with a maximum lag and testing down
- 3. **Estimate the model:**

$$\Delta X_t = \alpha + \beta t + \gamma X_{t-1} + \sum_{i=1}^{p-1} \delta_i \Delta X_{t-i} + \varepsilon_t \quad (81)$$

Where  $\Delta X_t = X_t - X_{t-1}$  is the first difference of the series. To apply OLS, we rewrite this in matrix form:

$$Y = X\beta + \varepsilon \quad (82)$$

Where:

- $Y$  is an  $(n - p) \times 1$  vector of  $\Delta X_t$  values
- $X$  is an  $(n - p) \times (p + 2)$  matrix of explanatory variables
- $\beta$  is a  $(p + 2) \times 1$  vector of coefficients  $(\alpha, \beta, \gamma, \delta_1, \dots, \delta_{p-1})$
- $\varepsilon$  is an  $(n - p) \times 1$  vector of error terms
- $n$  is the number of observations
- $p$  is the lag order

The OLS estimator for  $\beta$  is given by:

$$\hat{\beta} = (X'X)^{-1}X'Y \quad (83)$$

This estimator minimizes the sum of squared residuals:

$$\sum_{t=p+1}^n \varepsilon_t^2 = (Y - X\hat{\beta})'(Y - X\hat{\beta}) \quad (84)$$

- Use Ordinary Least Squares (OLS) to estimate the coefficients of the model
- 4. **Calculate the test statistic:**
  - The test statistic is the t-statistic for  $\gamma$ :

$$t = \frac{\hat{\gamma} - 0}{SE(\hat{\gamma})} \quad (85)$$

Where  $\hat{\gamma}$  is the estimated coefficient and  $SE(\hat{\gamma})$  is its standard error

- 5. **Determine the critical values:**
  - These depend on the sample size and the model specification (whether it includes a constant and/or trend)
  - They're typically obtained from statistical tables or through simulation
- 6. **Compare the test statistic to the critical values:**
  - If the test statistic is less than (more negative than) the critical value, reject the null hypothesis
  - The null hypothesis is that the series has a unit root (is non-stationary)
- 7. **Interpret the results:**
  - If we reject the null, we conclude the series is stationary
  - If we fail to reject the null, we cannot conclude the series is stationary

Once this is settled, our job is now to estimate the parameters. We can use MLE.

### 6.3 Moving Average (MA) Processes

The key assumptions are:



1. The random shocks are white noise, mutually independent and coming from the same distribution with mean 0 and constant variance.

#### Definition 6.11 (Moving Average Process)

The **MA(q)** process is a smoother type of noise than the white noise process. It is expressed by the formula

$$X_t = \sum_{j=1}^q \phi_j w_{t-j} + w_t \quad (86)$$

for  $\phi_j \in \mathbb{R}$ . Compared to the AR formula, the MA formula averages over the noise terms  $w_t$ . It focuses on the ripples of the process; if there is a shock to the process  $w_{t-1}$ , then that shock is still felt at time  $t$  by the term  $\phi_1 w_{t-1}$ .

Alternatively, the MA model can be written as an overall average of both the past and future white noise.

$$X_t = \sum_{j=-q/2}^{q/2} \phi_j w_{t+j} \quad (87)$$

#### Theorem 6.5 ()

A nice property of MA(q) is that autocovariance vanishes beyond a certain point. More specifically, it decays *linearly* and vanishes after  $q$  steps behind.

## 6.4 Linear Processes

Many time series fall under the category of linear processes.

#### Definition 6.12 (Linear Processes)

A **linear process** is defined as

$$X_t = \mu + \sum_{j=-\infty}^{+\infty} \theta_j w_{t-j} \quad (88)$$

which means that every  $X_t$  is a linear combination of the terms in the white noise process with some mean  $\mu$  added on. To ensure that this series doesn't blow up, we add the constraint that

$$\sum_j \theta_j^2 < \infty \quad (89)$$

However, since we are more interested in causal inference, to use the past to predict the future, we use the form

$$X_t = \mu + \sum_{j=0}^{\infty} \theta_j w_{t-j} \quad (90)$$

In fact, some AR processes are linear processes.

#### Example 6.3 (AR(1) as a Linear Process)

Note that AR(1) has a causal representation as a linear process. We can use the formula  $X_t =$

$\theta X_{t-1} + w_t$  and recursively define

$$X_t = \theta(\theta X_{t-2} + w_{t-1}) + w_t = \dots = \sum_{j=0}^{\infty} \theta^j w_{t-j} \quad (91)$$

Going back to analysis, infinite series are just limits.

$$\lim_{N \rightarrow \infty} \sum_{j=0}^N \theta^j w_{t-j} \quad (92)$$

So this sum may not converge. Letting  $S_N(\theta)$  be defined as above, we can compute that

$$\mathbb{E}[S_N(\theta)] = 0 \text{ and } \text{Var}[S_N] = \sigma^2 \sum_{j=0}^N \theta^{2j} = \sigma^2 \left( \frac{1 - \theta^{2N+2}}{1 - \theta^2} \right) \quad (93)$$

Thus, if  $|\theta| < 1$ , then  $\text{Var}[S_N(\theta)] \rightarrow \sigma^2/(1 - \theta^2)$ , and if  $w_t$  is Gaussian noise, then

$$S_N(\theta) \xrightarrow{d} \mathcal{N}(0, \sigma^2/(1 - \theta^2)) \quad (94)$$

If  $|\theta| = 1$ , the series does not converge and is not stationary, and if  $|\theta| > 1$ , then the random walk will grow exponentially fast.

## 6.5 ARMA

We can combine both the AR and MA processes to make a more sophisticated model.

### Definition 6.13 (ARMA)

The time series  $X_t$  is an ARMA( $p, q$ ) process if  $X_t$  has 0-mean and if we can write it as

$$X_t = w_t + \sum_{i=1}^p \phi_i X_{t-i} + \sum_{j=1}^q \theta_j w_{t-j} \quad (95)$$

where  $w_t$  is white noise with variance  $\sigma^2$  and  $\phi, \theta$  do not have any zero elements. Using the backshift operator, we can write it as

$$\Phi(B)X_t = \Theta(B)w_t \quad (96)$$

where

$$\Phi(B) = \left( 1 + \sum_{i=1}^p \phi_i B^i \right) \text{ and } \Theta(B) = \left( 1 + \sum_{j=1}^q \theta_j B^j \right) \quad (97)$$

## 6.6 ARIMA

## 6.7 Other

### Theorem 6.6 (Wold Representation Theorem)

Any 0-mean covariance stationary time series  $\{X_t\}$  can be decomposed into two time series

$$X_t = V_t + S_t \quad (98)$$

where

1.  $V_t$  is a linear combination of past variables of  $V_t$  with constant coefficients.
2.  $S_t = \sum_{i=0}^{\infty} \psi_i \eta_{t-i}$  is an infinite moving average process of error terms, where
  - (a)  $\psi_0 = 1, \sum_{i=0}^{\infty} \psi_i^2 < \infty$ .
  - (b)  $\{\eta_t\}$  is linearly unpredictable white noise, i.e.

$$\mathbb{E}[\eta_t] = 0 \quad (99)$$

$$\mathbb{E}[\eta_t^2] = \sigma^2 \quad (100)$$

$$\mathbb{E}[\eta_t \eta_s] = 0 \text{ for } s \neq t \quad (101)$$

and  $\eta_t$  is uncorrelated with  $\{V_t\}$ .

$$\mathbb{E}[\eta_t V_s] = 0 \text{ for all } t, s \quad (102)$$

#### Example 6.4 (Construction on Dataset)

Say that we have data  $\{X_t\}_{t=1}^T$  that we want to model and we have evidence that it is covariance stationary. We can do the following.

1. Initialize a parameter  $p$ , the number of parameters in the linearly deterministic term of the Wold decomposition of  $\{X_t\}$ .
2. By assumption we would like to estimate the linear projection of  $X_t$  on  $(X_{t-1}, X_{t-2}, \dots, X_{t-p})$ . Therefore, let us index the  $n$  subseries of length  $p+1$  by  $y$  and we can write the OLS equation

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \quad \mathbf{Z} = \begin{bmatrix} 1 & y_0 & y_{-1} & \cdots & y_{-(p-1)} \\ 1 & y_1 & y_0 & \cdots & y_{-(p-2)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & y_{n-1} & y_{n-2} & \cdots & y_{n-p} \end{bmatrix} \quad (103)$$

and we apply OLS to the problem  $\mathbf{y} = \mathbf{Z}\beta$  to give

$$\hat{\mathbf{y}} = \mathbf{Z}(\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z} \mathbf{y} \quad (104)$$

$$= \hat{P}(Y_t | Y_{t-1}, \dots, Y_{t-p}) \quad (105)$$

$$= \hat{\mathbf{y}}^{(p)} \quad (106)$$

We can compute the projection residuals

$$\epsilon^{(p)} = \mathbf{y} - \hat{\mathbf{y}}^{(p)} \quad (107)$$

and apply time series analysis to the sequence  $\epsilon^{(p)} = \{\epsilon_t^{(p)}\}$  to specify a moving average model.

$$\epsilon_t^{(p)} = \sum_{i=0}^{\infty} \psi_i \eta_{t-i} \quad (108)$$

yielding  $\{\hat{\psi}_j\}$  and  $\{\hat{\eta}_t\}$  estimates of parameters and innovations. We then check these estimates and see if they are consistent with the model assumptions. If not, we can add additional legs or modify  $p$ .

Theoretically, as we increase  $p$ , the projection of  $Y_t$  over the past  $p$ th history should approach the true linear projection  $Y_t$  over the whole history.

$$\lim_{p \rightarrow \infty} \hat{\mathbf{p}}^{(p)} = \hat{\mathbf{y}} \quad (109)$$

But if  $p$  is too large compared to  $n$ , you run out of freedom to estimate your models. You generally want to have more data than the number of parameters.

**Definition 6.14 (Lag Operator)**

The **lag operator**  $L^k$  simply maps

$$L^k(X_t) = X_{t-k} \quad (110)$$

Inverses also exist, so  $L^{-k}(X_t) = X_{t+k}$ .

Therefore, the Wold representation for a covariance stationary time series  $\{X_t\}$  can be expressed as

$$X_t = \sum_{i=0}^{\infty} \psi_i \eta_{t-i} + V_t \quad (111)$$

$$= \sum_{i=0}^{\infty} \psi_i L^i(\eta_t) + V_t \quad (112)$$

$$= \psi(L)\eta_t + V_t \quad (113)$$

where  $\psi(L) = \sum_{i=0}^{\infty} \psi_i L^i$ .

**6.8 Components of time series****6.9 Stationarity and tests for stationarity (including ADF test)****6.10 Autoregressive (AR) models****6.11 Moving average (MA) models****6.12 ARIMA models****6.13 Forecasting techniques****7 Advanced Topics****7.1 Generalized Linear Models****7.2 Survival analysis****7.3 Nonparametric methods****7.4 Resampling methods (jackknife, bootstrap)****8 Practical Considerations****8.1 Experimental design****8.2 Sample size and power analysis****8.3 Dealing with assumptions violations****8.4 Interpretation and reporting of results****9 Cross Validation**

We have understood the theoretical foundations of overfitting and underfitting with the bias variance decomposition. But in practice, we don't have an ensemble of datasets; we just have one. Therefore, we don't actually know what the bias, the variance, or the noise is at all. Therefore, how do we actually *know* in practice when we are underfitting or overfitting? Easy. We just split our dataset into 2 different parts: the training set and testing sets.

$$\mathcal{D} = \mathcal{D}_{train} \sqcup \mathcal{D}_{test} \quad (114)$$

What we usually have is a **training set** that allows us to train the model, and then to check its performance we have a **test set**. We would train the model on the training set, where we will always minimize the loss, and then we would look at the loss on the test set. Though we haven't made a testing set, since we know the true model let us just generate more data and use that as our testing set. For each model, we can calculate the optimal  $\theta$ , which we will denote  $\theta^*$ , according to the **root mean squared loss**

$$h_{\theta^*} = \operatorname{argmin}_{h_{\theta}} \sqrt{\frac{1}{N} \sum_{i=1}^N (y^{(i)} - h_{\theta}(\mathbf{x}^{(i)}))^2} \quad (115)$$

where division of  $N$  allows us to compare different sizes of datasets on equal footing, and the square root ensures that this is scaled correctly. Let us see how well these different order models perform on a separate set of data generated by the same function with Gaussian noise.

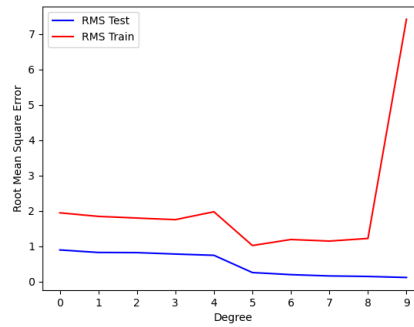


Figure 3: We can see that the RMS decreases monotonically on the training error as more complex functions become more fine-tuned to the data. However, when we have a 9th degree polynomial the RMS for the testing set dramatically increases, meaning that this model does not predict the testing set well, and performance drops.

Now we know that a more complex model (i.e. that captures a greater set of functions) is not necessarily the best due to overfitting. Therefore, researchers perform **cross-validation** by taking the training set  $(\mathcal{X}, \mathcal{Y})$ . We divide it into  $S$  equal pieces

$$\bigcup_{s=1}^S D_s = (\mathcal{X}, \mathcal{Y}) \quad (116)$$

Then, we train the model  $\mathcal{M}$  on  $S - 1$  pieces of the data and then test it across the final piece, and do this  $S$  times for every test piece, averaging its performance across all  $S$  test runs. Therefore, for every model  $\mathcal{M}_k$ , we must train it  $S$  times, for all  $K$  models, requiring  $KS$  training runs. If data is particularly scarce, we set  $S = N$ , called the **leave-one-out** technique. Then we just choose the model with the best average test performance.

The following result shows that cross-validation (data splitting) leads to an estimator with risk nearly as good as the best model in the class.

**Theorem 9.1 (Gyorfi, Kohler, Krzyak, Walk (2002))**

Let  $\mathcal{M} = \{m_h\}$  be a finite class of regression estimators indexed by a parameter  $h$ , with  $m$  being the true risk minimizer,  $m_{\hat{h}}$  being the empirical risk minimizer over the whole dataset  $\mathcal{D}$ , and  $m_H$  being

the empirical risk minimizer over the test set  $\mathcal{D}_{\text{test}}$  for ordinary least squares loss.

$$m_H = \operatorname{argmin}_{m_h} \frac{1}{N} \sum_{i \in \mathcal{D}_{\text{test}}} (y_i - m_h(x_i))^2 \quad (117)$$

$$m_{\hat{h}} = \operatorname{argmin}_{m_h} \frac{1}{N} \sum_{i \in \mathcal{D}} (y_i - m_h(x_i))^2 \quad (118)$$

If the data  $Y_i$  and estimators are bounded by  $L$ , then for any  $\delta > 0$ , we have

$$\mathbb{E} \int |m_H(x) - m(x)|^2 d\mathbb{P}(x) \leq (1 + \delta) \mathbb{E} \int |m_{\hat{h}}(x) - m(x)|^2 d\mathbb{P}(x) + \frac{C(1 + \log |M|)}{n} \quad (119)$$

where  $c = L^2(16/\delta + 35 + 19\delta)$ .

### Code 9.1 (Minimal Example of Train Test Split in scikit-learn)

To implement this in scikit-learn, we want to use the `train_test_split` class. We can also set a random state parameter to reproduce results.

```
1 from sklearn.model_selection import train_test_split
2
3 # Split into training (80\%) and test (20\%) data
4 X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2,
    random_state=66)
```

However, this process requires a lot of training runs and therefore may be computationally infeasible. Therefore, various **information criterion** has been proposed to efficiently select a model.

## 9.1 Leave 1 Out Cross Validation

### 9.1.1 Generalized (Approximate) Cross Validation

#### 9.1.2 Cp Statistic

## 9.2 K Fold Cross Validation

## 9.3 Data Leakage

## 9.4 Information Criterion

## References