

Statistics

Jacob Denson

May 6, 2024

Table Of Contents

1	Statistical Models	2
2	Hypothesis Testing	4
2.1	Rejecting the Null Hypothesis	4
3	Point Estimation	9
4	Regression	13
4.1	Parametric Regression	15
4.2	Linear Regression	16
4.3	Additive Models	16
4.4	Tree-Based Models	17
5	Neural Networks	18
5.1	Nets	19
6	Bayesian Networks	20
I	Computational Learning Theory	21

Chapter 1

Statistical Models

Statistics is the theory of inferring features of some probability distribution μ on a set S (normally S is discrete, or equal to \mathbb{R}^d for some d), given a number of independent samples drawn from this probability distribution. In general, it is very difficult to determine properties of μ , so we often work with some assumptions on the distribution μ , a *statistical model*. More precisely, a statistical model \mathcal{F} is a subset of the space of all distributions on a set S . We might write such a set as

$$\mathcal{F} = \{\mu_\theta : \theta \in \Theta\},$$

where Θ is the *parameter space* of the model. If Θ can be seen naturally as a subset of \mathbb{R}^n for some n , we say this is a *parametric model*.

Example. *In many statistical problems, we commonly assume the distribution is normal. That is, we work over the statistical model*

$$\{N(\mu, \Sigma) : \mu \in \mathbb{R}^n, \Sigma \in \mathbb{R}^{n \times n}\}.$$

Thus a normal distribution in \mathbb{R}^n is specifiable by $n^2 + n$ parameters.

Example. *A standard class in nonparametric statistics is the class of all distributions. But some methods might assume some minor regularity about the distribution. For instance, we might work in the model \mathcal{F}_{ABS} , which consists of all probability measures μ which are absolutely continuous (i.e. that have a density function). Or perhaps some additional regularity, i.e. that μ lies in some function space, i.e. a Sobolev space.*

A *statistical functional* is any function of the distribution μ . Examples include the *expected value* of the distribution, the *variance*, or the *median* (only defined for probability measures without atoms). A *statistic* is a function of the data drawn

from this distribution. We can now formally describe the problem of statistics: given a statistical functional, find a statistic which best approximates this functional. This is the problem of *point estimation*. The methods involved depend on the statistical functional employed.

For notational purposes, in statistics we often consider a multitude of distributions μ over the same set S . We therefore introduce the notation \mathbb{P}_μ , \mathbb{E}_μ , and \mathbb{V}_μ , to denote the probabilistic quantities obtained by using the distribution μ . For parametric problems, we might also use the notation \mathbb{P}_θ , \mathbb{E}_θ , and \mathbb{V}_θ .

One of the most important point estimation problem is the problem of *regression*, *prediction*, and *classification*, we are given pairs

$$(X_1, Y_1), \dots, (X_N, Y_N)$$

drawn independently from some product distribution μ . Given μ and some *error*, or *loss function*, we have a given *regression function* $Y = R(X)$, which best approximates the relationship between X and Y given by the distribution μ . Finding a good approximation of R via the data is called *regression*, or *curve estimation*. Given a new input X that has not yet been observed, we can then *predict* from the approximate regression function what the best estimate for the output Y should be.

Related to point estimation is the computation of *confidence intervals*. Given a statistical functional $\theta \in \mathbb{R}$, the problem is to compute an interval C from observed data X_1, \dots, X_N such that for each μ ,

$$\mathbb{P}_\mu(\theta \in C) \geq 1 - \alpha$$

The interval C is then called an α -*confidence interval*. More generally, if θ is a vector we can consider a *confidence set*. It is standard, but not necessary, to take $\alpha = 0.05$, and unless specified we will take this as the choice of α .

Remark. Under the frequentist interpretation of probability, a confidence interval *does not* give a probability associated with the quantity θ , since we assign no probability distribution to the distributions μ we consider. Instead, one can interpret a confidence interval as follows: each time we perform a confidence interval computation, we should be expected to be right 95% of the time.

Another problem is *hypothesis testing*. Here we have a particular fact about a distribution, called the *null hypothesis*. Given some data, the goal of hypothesis testing is to determine whether we should reject the null hypothesis, or if we have grounds to still believe in the null hypothesis.

Chapter 2

Hypothesis Testing

Let us suppose we consider a range of probability distributions on some space \mathcal{X} , parameterized by some set Θ . Our goal is to test whether a given set of data observed in an experiment should lead us to have greater or less belief in a given hypothesis, in which case we say the data is *statistically significant*.

It is important for scientists to follow a methodology in order to determine whether data is statistically significant or not. Otherwise they may be subject to cognitive bias when deciding whether the result of an experiment is statistically significant purely on a whim. There is evidence that banning the use of p-values (a standard method of hypothesis testing) in studies increases the frequency that scientists make erroneous claims (see Frick et. al 2019 for evidence of this as it relates to the null hypothesis significance ban in the journal Basic and Applied Social Psychology).

2.1 Rejecting the Null Hypothesis

Let us begin with the first approach to hypothesis testing, first advocated by British polymath Ronald Fischer. In this approach, one does not use data to provide evidence *for* a given hypothesis. Perhaps unintuitively, one instead takes data and uses it to provide evidence *against* a given hypothesis (called the *null hypothesis*).

Suppose an experiment will give us some random outcome X valued in a set \mathcal{X} . To perform a hypothesis test, we specify in advance a set $\mathcal{R} \subset \mathcal{X}$, called the *rejection region*. If X takes some value $x \in \mathcal{R}$, we reject the null hypothesis. If X takes some value $x \notin \mathcal{R}$, we do not reject the null hypothesis. If the null hypothesis were true, then this would limit the outcome of X to be drawn from a given family of probability distributions, which we denote by Θ . Usually, we define $\mathcal{R} = \{x \in \mathcal{X} : T(x) > T_0\}$ for a *test statistic* $T : \mathcal{X} \rightarrow \mathbb{R}$ and a *critical value* $T_0 \in \mathbb{R}$.

Note that if a hypothesis does not reject the null hypothesis, then in Fischerian approach we do not view this as evidence that the hypothesis is true. Our goal in this approach is thus only to minimize the number of *Type I errors* we make, where we reject the null hypothesis, despite it being true. If we are able to minimize these errors, then when a given data leads us to reject the null hypothesis, we should treat that as significant evidence that the null hypothesis is false.

This leads us to define the *size* of a test. First, define the *power function* β of the test, given by $\beta(\theta) = \mathbb{P}_\theta(X \in \mathcal{R})$. The *size* of a test is $\alpha = \sup_{\theta \in \Theta_0} \beta(\theta)$, which is precisely the least upper bound on the probability of the test making a *Type I* error. We say a test has *significance level* α if it has size at most α .

Suppose that for a given null hypothesis, we have an increasing family of rejection regions \mathcal{R}_α for each $\alpha \in (0, 1)$, which is increasing in α , and for each α , \mathcal{R}_α corresponds to a hypothesis test of significance level α . Given this setup, we define the *p-value* to be the random quantity $p = \inf \{\alpha : X \in \mathcal{R}_\alpha\}$. For each $\alpha \in [0, 1]$, $\alpha = \sup_{\theta \in \Theta_0} \mathbb{P}_\theta(X \in \mathcal{R}_\alpha)$, so in particular,

$$p = \sup_{\theta \in \Theta_0} \mathbb{P}_\theta(T(X) \geq \mathcal{R}_p).$$

Thus the *p-value* can be thought of as the least upper bound on the probability of observing something equal to or more extreme than the data observed, from distributions in the null hypothesis.

The use of *p-values* as a methodology of checking data for statistical significance was popularized by Ronald Fischer in his 1925 book *Statistical Methods for Research Workers*. Fischer interpreted the *p-value* of a hypothesis test as a continuous measure of how compatible data observed was with a null hypothesis. Fischer's approach interprets *p-values* not as probabilistic statement about the world, but as a continuous measure of how reluctant we should be to believe a null hypothesis - the lower the *p-value*, the more reluctant we should be to believe a null hypothesis.

Neyman and Pearson (1933) expanded upon Fischer's philosophy of hypothesis testing. Rather than simply testing whether to reject the null hypothesis, in the Neyman Pearson approach one instead sets a fixed level at which evidence from an experiment is statistically significant to provide evidence against the null hypothesis: The orthodox approach in 20th century science was to reject the null hypothesis when $p \leq 0.05$, and to accept the null hypothesis when $p > 0.05$. Unlike the Fischerian approach, in the Neyman Pearson approach the particular value of p does not indicate that the null hypothesis is more or less believable.

Example. Let $X = (X_1, \dots, X_n)$ be a vector with components being i.i.d. normally distributed random variables with some unknown mean μ and a known variance σ . Suppose we wish to test the null hypothesis that the mean of these variables is negative. We can parameterize the possible probability distributions that fit this

null hypothesis by the mean of these distributions. Thus we can set $\Theta_0 = (-\infty, 0]$. Consider the test where we reject the null hypothesis if $T(\bar{X}) > T_0$, for some $T_0 \in \mathbb{R}$. Then one can calculate that

$$\beta(\mu) = 1 - \Phi\left(\frac{n^{1/2}(T_0 - \mu)}{\sigma}\right),$$

and thus the size of the test is

$$1 - \Phi\left(\frac{n^{1/2}T_0}{\sigma}\right) \sim \left(\frac{\sigma}{n^{1/2}T_0}\right) e^{-(n^{1/2}T_0/\sigma)^2}.$$

For a given α , we get a test of size α by taking

$$T_0 = \sigma n^{-1/2} \Phi^{-1}(1 - \alpha) \sim \sigma n^{-1/2} \ln(1/\alpha).$$

For a given observation X , the p -value associated with this family of hypothesis tests is thus precisely obtained by the equation

$$X = \sigma n^{-1/2} \Phi^{-1}(1 - \alpha),$$

i.e. we have

$$p = 1 - \Phi(n^{1/2}X/\sigma).$$

The larger X is, the closer the p -value is to 0. For a given mean μ , the CDF of the p -value p is thus

$$t \mapsto 1 - \Phi(\sigma n^{-1/2} \Phi^{-1}(1 - t) - \mu).$$

For $\mu \ll 0$, the distribution of the p -value is concentrated near large values of p , but as $\mu \rightarrow 0$, the distribution of the p -value becomes uniformly distributed on $[0, 1]$.

The last example should cause us to treat a particular p -value with some skepticism. It is a random quantity, highly dependent on the data we observe, and thus will vary as we replicate a given experiment. The only guarantee we have on the distributions of the p -values is that, if the null hypothesis is true, then a p -value less than α will only be observed at most a fraction α of the time. Thus observing the p -value of one experiment will not necessarily lead to good predictions of p -values of replications. This contrasts with something with more descriptive power, like a confidence interval: Given a 95% confidence interval, 19 out of 20 replications will have confidence intervals overlapping with the confidence interval computed on the first experiment (see Geoff Cumming's [Dance of the \$p\$ -values](#) for more detail).

This behaviour also leads to *Lindlay's Paradox* (Lindsay, 1957). For $\mu \gg 0$ the p -values we will observe from the above experiment will be very low, guaranteed

with very high probability to have p -values smaller than 0.01. For $\mu = 0$, the p -values will be approximately normally distributed. It thus follows that, if we observe a p -value of 0.04, then from a Bayesian maximum likelihood perspective (we should favor hypothesis that are more likely to produce a given data), we are more likely to be in the situation $\mu = 0$ than in the situation $\mu \gg 0$, despite having a p -value less than 0.05 (and thus favoring a rejection of the null hypothesis from the approach of Fischer). Thus different approaches to statistical inference lead to different interpretations of a given data set. One approach which prevents this paradox is decreasing the p -value required for a test to be statistically significant as the sample size increases (see Good, 1992; Leamer, 1978; Maier & Lakens, 2022).

Even if you are using the Neyman-Pearson method for using p -values to refute a null hypothesis, in a paper you should always report the precise p -values obtained from your experiment. This is because these p -values are useful for secondary analysis, and allows researchers to compare p -values if they wish to use another p -value as a baseline.

On the other hand, defined purely in terms of levels, unless we have some information about the power of the tests we are using, the p -value says *nothing* about probabilities associated with the hypothesis. Small p -values indicate that it is unlikely that we could have observed given data if we are observing data from distributions given by the null hypothesis, but *large p -values do not imply that we should reject a given hypothesis*. Large p -values just indicate that, if the null hypothesis is true, it is not unlikely that we have observed data at least as extreme as what was observed. One should perform an *equivalence test* or *minimum effect test* to obtain more information outside of the null hypothesis.

A *hypothesis* is then a statement as to whether a distribution is in some subset Θ_1 of Θ . We call the complement of Θ_1 in Θ the *null hypothesis*, and denote it by Θ_0 . Our goal is to find a method which distinguishes between the two situations. The orthodox philosophy of scientists, in order to obtain greater evidence as to whether a given hypothesis is true, is to instead perform experiments that cause us to believe that the *null hypothesis*, the statement that the hypothesis fails, is highly improbable. The goal of hypothesis testing is to come up with a methodology that allows us to take data, and determine whether this allows us to make a statistically significant statement providing evidence against a given null hypothesis. We want to minimize the number of errors we make: either Type *I* errors, where we reject the null hypothesis, despite it being true, or Type *II* errors, where we retain the null hypothesis when the null hypothesis is false.

For each $\theta \in \Theta$, the hypothesis $\{\theta\}$ is called a *simple hypothesis*. If Θ is a subset of \mathbb{R} , then we call the hypotheses $\{\theta > \theta_0\}$ and $\{\theta < \theta_0\}$ *composite hypotheses*. A *two-sided* test is a test where the null hypothesis is a simple hypothesis, and a *one-sided* test is a test where the null hypothesis is a composite hypothesis.

Mathematically, how does hypothesis testing work? We consider a family of probability distributions over some set \mathcal{X} , indexed by some set Θ , and consider an \mathcal{X} valued random variable X which can be assigned probabilities that fit each distribution in a family parameterized by Θ . A test of the null hypothesis can be given by considering a region $\mathcal{R} \subset \mathcal{X}$, the *rejection region*. If X takes some value $x \in \mathcal{R}$, we reject the null hypothesis, and if X takes some value $x \notin \mathcal{R}$, we do not reject the null hypothesis.

Given a hypothesis test, we define the *power function* of the test is $\beta(\theta) = \mathbb{P}_\theta(X \in \mathcal{R})$. The *size* of the test is $\alpha = \sup_{\theta \in \Theta_0} \beta(\theta)$, and the *power* of the test is $\inf_{\theta \in \Theta_1} \beta(\theta)$. A test has *level* α if its size is less than or equal to α . A test with low size is unlikely to commit Type *I* errors, but the size of a test does not tell us anything about the probability of committing Type *II* errors. Similarly, a test with high power is unlikely to commit Type *II* errors, but the power of a test does not tell us anything about the probability of committing Type *I* errors.

In 1933, Neyman and Pearson, inspired by insights into p-values due to Gosset and Fischer, developed an approach called *statistical hypothesis testing*. In this framework, the goal of statistical tests is to guide researchers with respect to two different hypothesis, a null hypothesis, and an alternate hypothesis. Given a particular value α , classically equal to 0.05, one computes the p-value of a null hypothesis. If the value is smaller than α , one should be more prone to believing in the alternate hypothesis of the null hypothesis. Thus the exact quantitative p-value is not viewed as a measure of how reluctant we should be to accept the null hypothesis.

Chapter 3

Point Estimation

Let us consider point estimation. We consider a statistical functional θ of some distribution μ , from which we observe data X_1, \dots, X_N . We must now find a statistic $\hat{\theta}_N$, called the *estimator*, which is our best guess of θ . We say the estimator $\hat{\theta}$ is *unbiased* if $\mathbb{E}_\mu[\hat{\theta} - \theta_N] = 0$ for any distribution μ in our statistical model. It is *consistent* if $\hat{\theta}_N$ converges to θ in probability as $N \rightarrow \infty$, for any μ .

As a function of random data, the estimator $\hat{\theta}$ is a random variable. Its distribution is called the *sampling distribution*. The standard deviation of the estimator is called the *standard error* SE_N . Often, this error will depend on μ , and will thus be unknown. But we can often find an *estimator* for the standard error, which will be denoted by \widehat{SE}_N .

To determine how good an estimator is, we produce a *loss function* L , and then to determine the *expected loss*

$$\mathbb{E}_\mu[L(\theta, \hat{\theta}_N)].$$

The function L is normally selected so that if $\hat{\theta}_N = \theta$ then no loss is incurred, and more loss is incurred the ‘further away’ $\hat{\theta}_N$ is from being correct. There are several natural choices for the loss function:

- The most common is the loss function

$$L(\theta_1, \theta_2) = (\theta_1 - \theta_2)^2.$$

The resulting expected loss is the *mean square error*, denoted MSE. A nice feature of this formula is the *bias variance* decomposition, i.e. for any estimator, if we define

$$\text{BIAS} = \mathbb{E}[\hat{\theta}_N - \theta_N],$$

then

$$\text{MSE} = \text{BIAS}^2 + SE^2.$$

- Another common choice is the L^1 loss function

$$L(\theta_1, \theta_2) = |\theta_1 - \theta_2|.$$

- For discrete problems, another choice is the 0 – 1 loss function

$$L(\theta_1, \theta_2) = \mathbf{I}(\theta_1 \neq \theta_2).$$

Given a loss function, our goal is often therefore to find an estimator which minimizes the expected loss over all other possible estimators.

Expanding out the mean square error, we conclude that

$$\text{MSE} = \text{BIAS}^2 + \text{SE}^2,$$

where $\text{BIAS} = \mathbb{E}_\mu[\hat{\theta}_N] - \theta$. In particular, we conclude that if $\text{BIAS} \rightarrow 0$ and $\text{SE} \rightarrow 0$ as $N \rightarrow \infty$, then $\hat{\theta}_N$ is a consistent estimator, since then $\hat{\theta}_N$ converges in the L^2 norm to θ for any distribution μ , and convergence in L^2 implies convergence in probability.

Example. Consider $X_1, \dots, X_N \sim \text{Bernoulli}(p)$, where $p \in [0, 1]$ is unknown. A natural estimator to use is

$$\hat{p}_N = \frac{X_1 + \dots + X_N}{N}.$$

It is simple to check via linearity of expectation that \hat{p}_N is unbiased. It's standard error is

$$\text{SE} = \sqrt{\frac{p(1-p)}{N}}.$$

A natural estimator for this standard error is therefore

$$\widehat{\text{SE}}_N = \sqrt{\frac{\hat{p}_N(1 - \hat{p}_N)}{N}}.$$

For any value p , the standard error converges to zero, so this is a consistent estimator.

A common method for point estimation, especially in parametric problems where the point we are estimating determines the distribution under study, is the *maximum likelihood estimator*. We begin by assuming the statistical model we are working in is contained in the space of all distributions which are either continuous or are discrete. Given data X_1, \dots, X_N , we define the *likelihood function*

$$L(\mu) = f(X_1) \cdots f(X_N)$$

where f is the density function of μ . To find the maximum likelihood estimator, we find

$$\mu_* = \arg \max_{\mu \in \mathcal{F}} (L(\mu)),$$

the ‘most likely’ distribution to generate a given dataset. We then set $\hat{\theta}_N = \theta^*$, where θ^* is the value of the statistical function θ we are estimating which is associated with the distribution μ^* . To compute μ^* , it is often easier to work with the *log likelihood*

$$l(\mu) = \log f(X_1) + \log f(X_2) + \cdots + \log f(X_N),$$

since the distribution maximizing likelihood also maximizes the log-likelihood.

Example. Consider $X_1, \dots, X_N \sim \text{Bernoulli}(p)$, where p is unknown, and suppose we want to estimate p . We will find that the estimator we saw above is the maximum likelihood estimator for this example. Indeed, we have

$$\begin{aligned} l(p) &= \sum_{i=1}^N \log \mathbb{P}_p(X_i) \\ &= \sum_{i=1}^N \log (p^{X_i} (1-p)^{1-X_i}) \\ &= \sum_{i=1}^N X_i \log(p) + (1-X_i) \log(1-p). \end{aligned}$$

The derivative of l as a function of p is equal to

$$\sum_{i=1}^N \frac{X_i}{p} - \frac{1-X_i}{1-p} = \frac{1}{1-p} \left(\frac{S}{p} - N \right),$$

where $S = \sum X_i$. Setting this quantity equal to zero gives that we should set $\hat{p}_N = S/N$, which was the estimator we considered before.

Example. Consider $X_1, \dots, X_N \sim N(\mu, \sigma^2)$, where μ and σ are unknown. The problem is to estimate $\theta = (\mu, \sigma)$. The likelihood function is

$$L(\theta) = \frac{1}{\sigma^N (2\pi)^{N/2}} \exp \left(-\frac{1}{2} \sum_{i=1}^N \left(\frac{X_i - \mu}{\sigma} \right)^2 \right),$$

and thus the log likelihood is proportional to

$$l(\theta) = -N \log(\sigma) - \frac{1}{2} \sum_{i=1}^N \left(\frac{X_i - \mu}{\sigma} \right)^2.$$

Taking partial derivatives in σ and μ gives that the maximum likelihood parameters μ^* and σ^* satisfy

$$\mu^* = \frac{X_1 + \cdots + X_N}{N}$$

and

$$\sigma^* = \left(\frac{(X_1 - \mu)^2 + \cdots + (X_N - \mu)^2}{N} \right)^{1/2},$$

which are the natural estimators we might choose given the observed data.

Example. Consider $X_1, \dots, X_N \sim \text{Uniform}([0, \theta])$, for some $\theta \in [0, \infty)$. What is the maximum likelihood estimator of θ ? The likelihood function is

$$L(\theta) = \theta^{-N} \mathbf{I}(\max(X_1, \dots, X_N) \leq \theta).$$

This function is equal to zero for $\theta \in [0, \max(X_1, \dots, X_N))$, and decreases for values bigger than $\max(X_1, \dots, X_N)$. Thus we conclude that the maximum likelihood estimator for θ is

$$\hat{\theta}_N = \max(X_1, \dots, X_N).$$

Example. Consider a practical example. Let us suppose that there are N fish in a lake, a quantity unknown to us, and which we would like to estimate. We capture A fish, paint them with red spots, and then release them back into the wild. We then capture B fish, and observe that X of them are red. Let us use this information to compute the maximum likelihood estimate of N . The only random quantity here is X , the other quantities A and B determining the statistical model we use. Given N , the distribution of X is a hypergeometric distribution, and so the likelihood function is

$$L(N) = \frac{\binom{A}{X} \binom{N-A}{n-X}}{\binom{N}{B}},$$

To determine the maximum quantity here, we compute there for $N > A$,

$$\frac{L(N)}{L(N-1)} = \frac{N-A}{N-A-B+X} \frac{N-B}{N}.$$

This quantity is greater than one if $N < AB/X$, and smaller than one if $N > AB/X$. The maximum likelihood estimator for N is therefore the greatest integer less than or equal to AB/X , i.e.

$$\hat{N} = \left\lfloor \frac{AB}{X} \right\rfloor.$$

As an example, if we capture 1000 fish and mark them with red spots, capture another set of 1000 fish, and observe that 100 have red spots, we should guess that there are 10000 fish in the lake.

Chapter 4

Regression

The most basic task of statistics is regression. Given a sequence of independent experiments $(X_1, Y_1), \dots, (X_N, Y_N)$ drawn from some distribution (X, Y) , we want to determine the a function f such that $f(X)$ approximates Y . To measure the approximations success, for each value (x, y) , we consider a *loss function* $L(y, f(x)) \geq 0$. The *estimated prediction error* is then defined to be $EPE(f) = \mathbb{E}(L(Y, f(X)))$. The goal of *regression* is to find a function f which minimizes the estimated prediction error. In parametric problems, the goal is to find the function f lying in a given finite dimensional class of candidate functions minimized estimated prediction error. In non parametric problems, we must find the function f from an infinite dimensional class of candidate functions.

Example. *In most cases, the most analytically convenient loss functions is the square error loss*

$$L(y, f(x)) = (y - f(x))^2.$$

We then have

$$EPE(f) = \mathbb{E}(L(Y, f(X))) = \mathbb{E}((Y - f(X))^2) = \mathbb{E}(\mathbb{E}(Y - f(X))^2 | X)).$$

In this case, the expected prediction error is known as the mean square error, denoted $MSE(f)$. Thus regression reduces to minizing $(Y - f(X))^2$ pointwise, given X . The minimizer of this quantity is the conditional expectation

$$f(X) = \mathbb{E}(Y|X).$$

This is because if f is any function, then $\mathbb{E}(f(X)|X) = f(X)$, and so

$$\mathbb{E}((Y - \mathbb{E}(Y|X))f(X)|X) = \mathbb{E}(Y|X)f(X) - \mathbb{E}(Y|X)f(X) = 0.$$

Thus $Y - \mathbb{E}(Y|X)$ is orthogonal to the subspace of random variables measurable with respect to the sigma algebra generated by X . Thus we can apply the Pythagorean theorem to conclude that

$$\begin{aligned} \text{MSE}(f) &= \mathbb{E}((Y - f(X))^2) = \mathbb{E}((Y - \mathbb{E}(Y|X))^2) + \mathbb{E}((\mathbb{E}(Y|X) - f(X))^2) \\ &\geq \mathbb{E}((Y - \mathbb{E}(Y|X))^2) = \text{MSE}(\mathbb{E}(Y|X)). \end{aligned}$$

Thus regression with respect to squared loss is equivalent to estimating the conditional expectation of one variable with respect to one another. The decomposition above using Pythagoras' theorem is very useful, which is really what makes the square error loss most useful in a given situation.

Example. Another standard loss function is the L^1 loss function, given by

$$L(y, f(x)) = |y - f(x)|.$$

As with the squared loss it suffices to choose $f(X)$ which pointwise minimizes

$$\mathbb{E}(|Y - f(X)||X).$$

Fix y . We note that if $\varepsilon > 0$, then

$$|Y - y + \varepsilon| - |Y - y| = \begin{cases} \varepsilon & : Y \geq y \\ -\varepsilon & : Y \leq y - \varepsilon \\ 2(Y - y) + \varepsilon & : y - \varepsilon < Y < y \end{cases}$$

Thus

$$\begin{aligned} &|\mathbb{E}(|Y - y + \varepsilon| - |Y - y||X) - \varepsilon[\mathbb{P}(Y \geq y|X) - \mathbb{P}(Y \leq y - \varepsilon|X)]| \\ &\leq \varepsilon \mathbb{P}(y - \varepsilon < Y < y|X). \end{aligned}$$

Provided we are working with a regular probability measure, this means that for each ω , the function $y \mapsto \mathbb{E}(|Y - y||X = X(\omega))$ is right differentiable, with derivative $\mathbb{P}(Y \geq y|X = X(\omega)) - \mathbb{P}(Y < y|X = X(\omega))$. In particular, a choice of y which minimizes $\mathbb{E}(|Y - y||X = X(\omega))$ must satisfy

$$\mathbb{P}(Y \geq y|X = X(\omega)) = \mathbb{P}(Y < y|X = X(\omega)) = 0.5.$$

If $f(X)$ is a function such that almost surely,

$$\mathbb{P}(Y \geq f(X)|X) = \mathbb{P}(Y < f(X)|X) = 0.5$$

then we say it is a conditional median, and it will be a pointwise minimizer of the expected loss. We normally denote a conditional median by $\mathbb{M}(Y|X)$. Unlike the conditional expectation, the conditional median need not be unique if the underlying distribution of X is not continuous.

Example. Suppose that the values of Y lie in some discrete set of values. The problem of regression in this setting is normally called classification. A natural loss function to use here is the 0-1 loss function

$$L(y, f(x)) = \mathbf{I}(y \neq f(x)).$$

As with the previous examples, the regression function in this setting can be easily proved to be

$$\mathbf{B}(Y|X) = \arg \min \mathbb{P}(Y = y|X).$$

This is known as the Bayes classifier. The value $EPE(\mathbf{B}(Y|X))$ is known as the Bayes rate.

Even in these examples, where we can calculate an explicit formula for the regression function, in practice we cannot compute the regression function from sample data. Thus we must come up with an approximation \hat{f} of the regression function $f(x)$, where for each x , $\hat{f}(x)$ is a random variable determined by the data $(X_1, Y_1), \dots, (X_N, Y_N)$. Such a random variable is known as a *statistic*.

There is a very useful decomposition result for the expected prediction error of \hat{f} , where $f(x) = \mathbb{E}(Y|X = x)$. We can write $Y = f(X) + \varepsilon$, where ε has mean zero and variance σ^2 . If we define the *bias*

$$\text{Bias}(x) = \mathbb{E}(\hat{f}(x)) - f(x).$$

then

$$\begin{aligned} \text{MSE}(f) &= \mathbb{E}[(Y - f(X))^2] + \mathbb{E}[(f(X) - \hat{f}(X))^2] \\ &= \sigma^2 + \mathbf{V}[\hat{f}] + \mathbb{E}[\text{Bias}(X)^2]. \end{aligned}$$

This is referred to as the *Bias-Variance decomposition*. The error σ^2 is unavoidable for any function \hat{f} , whereas to minimize the mean square error, we must make a tradeoff between making the variance of our estimator small, and the bias small.

4.1 Parametric Regression

In parametric statistics, we are given independent and identically distributed data $\{(X_1, Y_1), \dots, (X_N, Y_N)\}$. We do not know the distribution they are drawn from, but we *do* know the distribution lies in some

4.2 Linear Regression

In some case, we assume our regression functions take the form $\beta^* X$, where $\beta^* \in (\mathbf{R}^n)^*$. Given the data $(X_1, Y_1), \dots, (X_k, Y_k)$, we determine the best estimate $\hat{\beta}$ of β^* by evaluating it against the loss function $\mathcal{L}(\beta) = \mathbf{E}[(Y - \beta X)^2]$. Of course, we cannot calculate \mathcal{L} directly, but we may estimate it with our samples. Because the loss function is a differentiable function of β , we may take derivatives to determine β :

$$\nabla \mathcal{L}(\beta) = 2\mathbf{E}[(Y - \beta X)X^T] = 2\mathbf{E}(YX^T) - 2\beta\mathbf{E}(XX^T)$$

This is optimized when the derivative of this function is zero. i.e., when

$$\beta\mathbf{E}(XX^T) = \mathbf{E}(YX^T)$$

Assuming $\mathbf{E}(XX^T)$ is invertible, we may invert, and determine that the optimal value β^* can be calculated as

$$\beta^* = \mathbf{E}(YX^T)\mathbf{E}(XX^T)^{-1}$$

Now if we only have the samples (X_i, Y_i) , we may approximate this value by forming the conglomerate matrices $\mathbf{X} = (X_1|X_2|\dots|X_n)$ and $\mathbf{Y} = (Y_1, \dots, Y_n)$, and calculating $\hat{\beta} = \mathbf{YX}^T(\mathbf{XX}^T)^{-1}$. This minimizes the error over the training data set $\sum (Y_i - \beta X_i)^2 = \|\mathbf{y} - \beta\mathbf{X}\|^2$.

How do we estimate how accurate our prediction is. First, assume that each Y_i is independant, with the same variance σ^2 . Then

$$\mathbf{V}(\hat{\beta}) = \mathbf{V}(\mathbf{YX}^T(\mathbf{XX}^T)^{-1})$$

INSERT THEORETICAL ESTIMATES, Gauss Markov theorem, etc.

4.3 Additive Models

A generalized additive model has a regression function of the form

$$\mathbf{E}(Y|X) = \alpha + f_1(X^1) + \dots + f_n(X^n)$$

where the f_i are unspecified smooth (C^∞) functions, and $X = (X^1, \dots, X^n)$ is a random vector. To fit an additive model, given a sample $(X_1, Y_1), \dots, (X_m, Y_m)$, we take as a cost function the penalized sum of squares to find the constant α and functions f_i ,

$$\sum_{i=1}^m \left(Y_i - \alpha - \sum_{j=1}^n f_j(X_i^j) \right)^2 + \sum_{j=1}^m \lambda_j \int (f_j'')^2$$

Where the $\lambda_j \geq 0$ are arbitrary parameters. The minimizer of this cost function is not unique – it is standard convention to require that $\sum_{i,j} f_i(X_i^j) = 0$. One can apply an iterative cubic smoothing spline solution to find this minimum.

4.4 Tree-Based Models

Tree based methods partition the feature space, and then fit a simple model (normally a constant) into each one. If \mathcal{S} is such a space, and we partition it into S_1, \dots, S_n , each with a model f_1, \dots, f_n , then the model is

$$\mathbf{E}(Y|X) = \sum_{i=1}^n f_i(X) [X \in S_i]$$

If S_i has already been decided, and we are using constants for the f_i , then the best choice of constants (according to the least squares cost function) given a sample $(X_1, Y_1), \dots, (X_n, Y_n)$ are just the mean values of the Y_j with $X_i \in S_i$. The question remains, however, on how to choose our partitions.

To find optimal partitions, we assume our feature space is \mathbf{R}^n , and our regions formed by ‘binary splits’. We start with the whole space, pick a ‘splitting coordinate’ i and ‘splitting point’ $t \in \mathbf{R}$, and partition our region into two sets $A = \{x \in \mathbf{R}^n : x_i < t\}$ and $B = \{x \in \mathbf{R}^n : x_i \geq t\}$. We then recursively partition A and B up in this manner, until we are satisfied with our splits.

Finding the best choice of partition using the method above is generally computationally infeasible. We shall proceed with a greedy approximation. Given a region A containing features X_1, \dots, X_n , we seek to find a splitting variable i and split point t which minimize the cost function

$$\arg \min_{i,t} \min_a \sum_{X_j^i \leq t} (Y_j - a)^2 + \min_b \sum_{X_j^i > t} (Y_j - b)^2$$

Given i and t , the minimum values of a and b are just obtained by taking the mean of the results Y_j . By doing a linear scan on each coordinate, it is fairly simple to find i and t . Then we recursively perform this greedy approach on the subpartitions.

Now when do we stop splitting? If we split far enough, then we will only have very few examples in each subregion, and we will have overfitted our training data! Furthermore, it will be very difficult to interpret the model we have created.

Chapter 5

Neural Networks

Neural Networks arise from the solution of a certain model, known as the Projection Pursuit Regression model. Assume we have an input vector $X \in \mathbf{R}^n$, with target Y . The projection pursuit regression model has the form

$$f(X) = \sum_{i=1}^M g_i(\beta_i X)$$

Where the g_i are unspecified, and $\beta_i \in (\mathbf{R}^n)^*$. This is an additive model, but in the features $V_i = \beta_i X$. Each $g_i(\beta_i X)$ is called a ridge function in \mathbf{R}^n .

This model is really general. For instance, the product of the coordinates can be written

$$X_1 X_2 = \frac{(X_1 + X_2)^2 - (X_1 - X_2)^2}{4}$$

In fact, if we let M be large enough, for appropriate choices of g_i can approximate arbitrary continuous functions on \mathbf{R}^n (this model is a universal approximator). Unfortunately, this means this model will be hard to fit exactly, and thus the model is better for estimating data rather than obtaining an understandable model.

Given some data $(X_1, Y_1), \dots, (X_n, Y_n)$ from the regression model, we thus seek to minimize the error

$$\sum_{i=1}^n \left(y_i - \sum_{j=1}^M g_j(\beta_j x_i) \right)$$

as a choice of g_j and β_j . We need to impose constraints on g_j to prevent overfitting.

Suppose we have $M = 1$, and that $g = g_1$ is differentiable, and $\beta = \beta_1$ is the linear functional. Then, taking the initial terms around the Taylor series,

$$g(\beta x_i) \approx g(\alpha x_i) + g'(\alpha x_i)(\alpha - \beta)x_i$$

$$\begin{aligned}
\sum_{i=1}^n [y_i - g(\beta x_i)]^2 &\approx \sum_{i=1}^n [y_i - g(\alpha x_i) - g'(\alpha x_i)(\alpha - \beta)x_i]^2 \\
&= \sum_{i=1}^n g'(\alpha x_i)^2 \left[\alpha x_i + \frac{y_i - g(\alpha x_i)}{g'(\alpha x_i)} - \beta x_i \right]^2
\end{aligned}$$

We can minimize the right-hand side by carrying out a least squares regression with target

$$\alpha x_i + \frac{y_i - g(\alpha x_i)}{g'(\alpha x_i)}$$

We can then iterate this regression until convergence. With more than one term in the model, we just perform forward stage-wise regression.

The projection pursuit regression model has not been widely used in the field of statistics, possibly due to the lack of computational resources when it was created. Nonetheless, it leads to the field of neural networks, which are much more useful.

5.1 Nets

There is a lot of mysticism surrounding neural networks (perhaps for the same reason ‘the god particle’ is so controversial) but they are really just non-linear statistical models. Here we will discuss the most basic kind of neural nets, the single hidden layer back-propagation network.

Suppose we are given a set of inputs $X = (X_1, \dots, X_n)$. A neural net creates layers of derived features $Z = (Z_1, \dots, Z_m)$ as affine transformations of the X_i , ‘flattened’ by some activation function σ . In the single layer approach, we have one layer of these derived features, and then these derived features are used to generate the target $Y = (Y_1, \dots, Y_k)$ as a function of the Z_i , again modified by an output function. In terms of formulas, our mathematical model is

$$Z = \sigma(\Lambda X + v) \quad W = \Delta Z + w \quad f(X) = g(W)$$

where Λ and Δ are linear transformations, and our regression function is f .

For regression, we normally choose not to modify our outputs via an output function (that is, we let $g = \mathbf{1}$). For classification, we need to choose an output function which results in reasonable results. The sigmoid function is often chosen as the activation, $\sigma(t) = (1 + e^{-t})^{-1}$. Sometime Gaussian radial basis functions are used, producing a radial basis function network. Note that if we let σ and the output regularization function be the identity, we obtain a linear model. Thus in this way, a neural network is a generalization of the linear model.

Chapter 6

Bayesian Networks

Let X, Y, Z be random variables. We say that X and Y are *conditionally independent* given Z , if for any measurable $A, B \subset \mathbf{R}$,

$$\mathbf{P}(X \in A, Y \in B | Y, Z) = \mathbf{P}(X \in A | Z) \mathbf{P}(Y \in B | Z)$$

This just means that once you know Z , you can gain no information about X from information about Y . Bayesian networks are a model of information which allow us to measure the conditional independence of random variables.

Given a set X_1, \dots, X_n of random variables, suppose we form a directed, acyclic graph whose certices are the random variables. We say the variables are *Markov* with respect to the graph if for any random variable X_i ,

$$\mathbf{P}(X_1 \in A_1, \dots, X_n \in A_n) = \prod_{i=1}^n \mathbf{P}(X_i \in A_i | \text{parents}(X_i))$$

In other words, this means exactly that X_i is independent of all variables once we condition on the parents of X_i .

Part I

Computational Learning Theory

Computational Learning Theory is the study of Machine Learning, from the perspective of theoretical computing science. In the study of computational learning theory, we specify learning models, which limit the *complexity* of the kinds of algorithms one can use to study data in machine learning. Since we can often obtain huge inputs in statistics, we often limit our study to data with low *sample complexity*, i.e. the minimum number of samples required to solve an algorithm, up to a certain degree of error. We might also naturally limit algorithms in *time complexity*, since, despite a problem having low sample complexity, there might not exist an efficient algorithm to analyze a given set of data. We might also limit the *hypothesis complexity*.