

Statistical Machine Learning: Introduction

Daniel J. McDonald

Indiana University, Bloomington

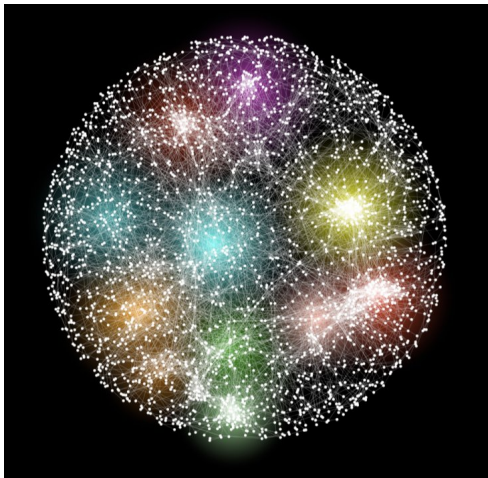
mypage.iu.edu/~dajmcdon

February 24-26, 2015

Introduction and motivation

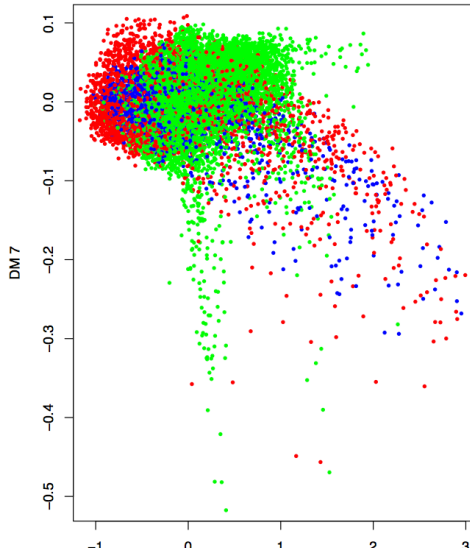
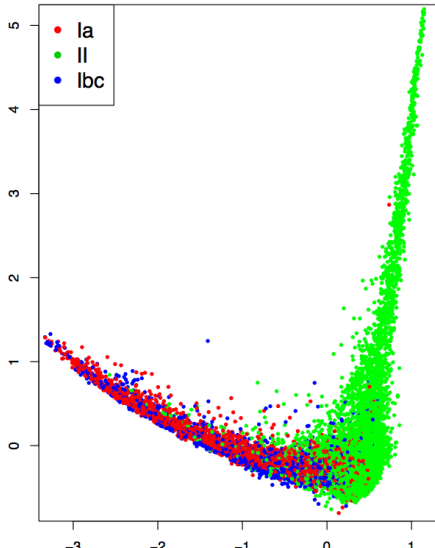
WHAT IS STATISTICS?

DESCRIPTION Collect some data. Give summaries. Make charts, pretty pictures. Also “unsupervised learning”.



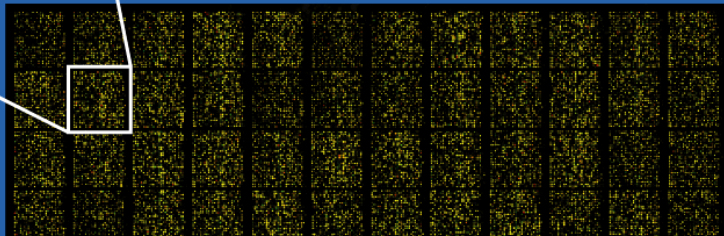
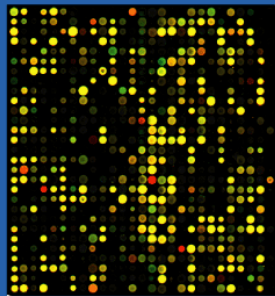
WHAT IS STATISTICS?

ESTIMATION/INFERENCE Try to determine the underlying causal model.



WHAT IS STATISTICS?

PREDICTION Try to predict some of the data using other data.



WHAT IS STATISTICS?

Suppose I come up with some method of description / estimation / prediction. . .

Is it any good?

How can I evaluate my method?

Are there better methods?

Modeling your data

STATISTICAL MODELS

We observe data Z_1, Z_2, \dots, Z_n generated by some probability distribution P . We want to use the data to learn about P .

A **statistical model** is a set of distributions \mathbb{P} .

Some examples:

1 $\mathbb{P} = \{P(z = 1) = p, P(z = 0) = 1 - p, 0 < p < 1\}.$

2 $\mathbb{P} = \{Y \sim N(X^\top \beta, \sigma^2), \beta \in \mathbb{R}^p, \sigma > 0, X \text{ fixed}\}.$

3 $\mathbb{P} = \{\text{all CDF's } F\}.$

4 $\mathbb{P} = \{\text{all smooth functions } f : \mathbb{R}^p \rightarrow \mathbb{R}\}$

STATISTICAL MODELS

We observe data Z_1, Z_2, \dots, Z_n generated by some probability distribution P . We want to use the data to learn about P .

$$\mathbb{P} = \{P(z = 1) = p, P(z = 0) = 1 - p, 0 < p < 1\}$$

To completely characterize P , I just need to estimate p .

Need to assume that $P \in \mathbb{P}$.

This assumption is mostly empty: need independent, can't see $z = 12$.

STATISTICAL MODELS

We observe data $Z_i = (Y_i, X_i)$ generated by some probability distribution P . We want to use the data to learn about P .

$$\mathbb{P} = \{Y \sim N(X^\top \beta, \sigma^2), \beta \in \mathbb{R}^p, \sigma > 0, X \text{ fixed}\}$$

To completely characterize P , I just need to estimate β .

Need to assume that $P \in \mathbb{P}$.

This time, I have to assume a lot more: **Linearity, independence, Gaussian noise, no ignored variables, no collinearity, etc**

Things you may have seen: Notation

NECESSARY BACKGROUND: NOTATION

- We will write **vectors** as

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

We write this as $x \in \mathbb{R}^n$, which is “x is a member of ar-en.”

- We commonly will need to “turn” the vector, which we write as

$$x^\top = [x_1 \quad x_2 \quad \dots \quad x_n]$$

Here, the superscript “T” takes a vector and flips it on its side.

NECESSARY BACKGROUND: NOTATION

If we have two vectors, we will double subscript them

Suppose $x_1, x_2 \in \mathbb{R}^n$, then

$$x_1 = \begin{bmatrix} x_{11} \\ x_{12} \\ \vdots \\ x_{1n} \end{bmatrix} \quad \text{and} \quad x_2 = \begin{bmatrix} x_{21} \\ x_{22} \\ \vdots \\ x_{2n} \end{bmatrix}$$

NECESSARY BACKGROUND: NOTATION

Often, we will combine many vectors into a **matrix**

$$\mathbb{X} = \begin{bmatrix} x_1 & x_2 & \cdots & x_p \end{bmatrix} = \begin{bmatrix} X_1^\top \\ X_2^\top \\ \vdots \\ X_n^\top \end{bmatrix}$$

As much as possible

- **Lower case** Roman letters will be columns
- **Upper case** Roman letters will be rows

NECESSARY BACKGROUND: LENGTHS

We will need to measure the **size** of both vectors and matrices.

The most common is the one we use every day **Euclidean distance**

(Think: the Pythagorean theorem)

$$||x||_2 = \sqrt{\sum_{k=1}^p x_k^2}$$

We call this a **norm** and refer to this as the “ell two norm”

Additionally, we will need the **Manhattan distance**

$$||x||_1 = \sum_{k=1}^p |x_k|$$

We call this the “ell one norm”

NECESSARY BACKGROUND: LENGTHS

For matrices, we will just define something very related to ‘length’
(but it doesn’t technically qualify)

Many times, we are interested in the size of the diagonal of a matrix

This is known as the **trace** and is defined to be

$$\text{trace}(\mathbb{X}) = \sum_{j=1}^p \mathbb{X}_{jj}$$

That is, the trace is the sum of the diagonal entries.

Things you may have seen: Probability

NECESSARY BACKGROUND: PROBABILITY

For this class, we need to recall (some) probability.

Again, I would be satisfied with you accepting that certain manipulations are reasonable, rather than ‘understanding’ everything.

That being said, let's take it away...

WHAT'S A RANDOM VARIABLE?

Let X be a **random variable**. That is, X ...

- Has a **probability density function** p_X such that the **probability** (denote this by \mathbb{P}) that X takes on a set of values A is given by¹

$$\mathbb{P}(A) = \int_A p_X(x) dx$$

- And p_X has certain properties such as $p_X \geq 0$ and $\int p_X = 1$.

¹Anyone who has studied probability would have serious problems with this statement. If this is you, don't quibble; we're trying to avoid unnecessary complications.

WHAT ARE THE PROPERTIES OF A RANDOM VARIABLE?

In this class, we really only care about X 's

- **mean** (alternatively known as its **expectation**)

(This is all about finding its **center**)

- and **variance**.

(This is all about finding its **spread**)

WHAT'S EXPECTATION?

Imagine taking a metal rod of a certain mass.

However, its mass isn't necessarily even along its length.

Attempt to balance the rod on your finger. The balancing point is the **center of mass** of the rod.



FIGURE: A family calculates expectations

WHAT'S EXPECTATION?

Crucial connection: If we think about the density of the random variable determining where the rod's mass is **distributed**, then the “center of mass” is the **expectation**.

$$\mathbb{E}[X] = \int x p_X(x) dx$$

WHAT'S VARIANCE?

For variance, I'll just give you the definition

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

In words:

“variance is the average squared deviation from the average”

Note: $\text{Var}[X] = \mathbb{E}[X^2] - \mathbb{E}[X]^2$

CAN YOU TAKE ME HIGHER?

Implicitly, we were assuming that $X \in \mathbb{R}$.

What happens if $X \in \mathbb{R}^p$?

The expectation is going to look the same, but be a vector

$$\mathbb{E}[X] \in \mathbb{R}^p$$

For variance, we need to use some matrix notation:

$$\text{Var}[X] = \mathbb{E}[(X - \mathbb{E}[X])(X - \mathbb{E}[X])^\top] \in \mathbb{R}^{p \times p}$$

If you write this out, you'll see that this a matrix with

- The variances of the components on the diagonal
- The covariances of any two components in the off-diagonal entries.

COMBINE MATRICES AND PROBABILITY

We will commonly combine matrix multiplication with probability statements

Suppose that $Y \in \mathbb{R}^n$ is a random variable such that $\mathbb{E}[Y] = \mu$ and $\text{Var}[Y] = \Sigma$.

What is the distribution of $\mathbb{X}Y$ (if \mathbb{X} is a fixed matrix, not random)?

It turns out expectation is **linear** and hence we can rearrange ‘ \mathbb{E} ’ and ‘ \mathbb{X} ’

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

Variance is little more complicated, but not much

$$\text{Var}[\mathbb{X}Y] = \mathbb{X} \text{Var}[Y] \mathbb{X}^\top = \mathbb{X}\Sigma\mathbb{X}^\top \quad (\text{check this!})$$

MORE ON RANDOM VARIABLES

- If X_1, \dots, X_n are independent then $\text{Var}[\sum_{i=1}^n a_i X_i] = \sum_i a_i^2 \text{Var}[X_i]$
- The **covariance** is

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

- The **correlation** is

$$\rho(X, Y) = \text{Cov}[X, Y] / \sigma_X \sigma_Y.$$

The **conditional expectation of Y given X** is the **random variable** $\mathbb{E}[Y|X]$.

$$\mathbb{E}[Y|X = x] = \int y p(y|x) dy$$

IMPORTANT DISTRIBUTIONS

NORMAL $X \sim N(\mu, \sigma^2)$ if

$$p(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2\sigma^2} (x - \mu)^2 \right\}$$

BERNOULLI $X \sim \text{Bernoulli}(\theta)$ if $P(X = 1) = \theta$ and $P(X = 0) = 1 - \theta$. Thus the pdf is

$$p(x; \theta) = \theta^x (1 - \theta)^{1-x} I_{\{0,1\}}(x).$$

CONVERGENCE

Let X_1, X_2, \dots be a sequence of random variables, and let X be another random variable with distribution P . Let F_n be the cdf of X_n and let F be the cdf of X .

1 X_n converges in probability to X , $X_n \xrightarrow{P} X$, if for every $\epsilon > 0$,

$$\lim_{n \rightarrow \infty} P(|X_n - X| > \epsilon) = 0.$$

2 X_n converges in distribution to X , $X_n \rightsquigarrow X$, if

$$\lim_{n \rightarrow \infty} F_n(t) = F(t)$$

Heuristics: the sequence is somehow getting “closer” to some limit. But things are random...

CONVERGENCE RULES

Suppose X_1, X_2, \dots are independent random variables, each with mean μ and variance σ^2 .
Let $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$.

WEAK LAW OF LARGE NUMBERS

$$\bar{X}_n \xrightarrow{P} \mu$$

The law of large numbers tell us that the probability mass of an average of random variables “piles up” near its expectation.

CENTRAL LIMIT THEOREM

$$\frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \rightsquigarrow N(0, 1).$$

The CLT tells us about the shape of the “piling”, when appropriately normalized.

Who cares?

STATISTICAL MODELS

We observe data $Z_i = (Y_i, X_i)$ generated by some probability distribution P . We want to use the data to learn about P .

$$\mathbb{P} = \{Y \sim N(X^\top \beta, \sigma^2), \beta \in \mathbb{R}^p, \sigma > 0, X \text{ fixed}\}$$

To completely characterize P , I just need to estimate β .

Need to assume that $P \in \mathbb{P}$.

This time, I have to assume a lot more: **Linearity, independence, Gaussian noise, no ignored variables, no collinearity, etc**

EVALUATION

Once I choose some way to “learn” a statistical model, I need to decide if I’m doing a good job.

How do I decide if I’m doing anything good?

PROPERTIES

Lots of ways to evaluate estimators, $\hat{\mu}$ of parameters μ .

- Consistency: $\hat{\mu} \xrightarrow{P} \mu$.
- Asymptotic Normality: $\hat{\mu} \xrightarrow{D} N(\mu, \Sigma)$
- Efficiency: how large is Σ
- Unbiased: $\mathbb{E}[\hat{\mu}] \stackrel{?}{=} \mu$
- etc.

None of these things make sense unless **your model is correct**.

Your model is wrong!

[unless you are flipping coins, gambling in a casino, or running randomized, controlled trials on cereal grains]

MIS-SPECIFIED MODELS

What happens when your model is wrong? And it **IS** wrong.

None of those evaluation criteria make any sense. The parameters no longer have any meaning.

[The criteria still hold in some sense: I can demand that I get close to the projection of the truth onto \mathbb{P}]

PREDICTION

Prediction is easier: your model may not actually represent the true state of nature, but it may still predict well.

*Over an 13-year period, [David Leinweber] found, [that annual **butter production** in Bangladesh] “explained” 75% of the variation in the annual returns of the Standard & Poor’s 500-stock index.*

*By tossing in **U.S. cheese production** and the **total population of sheep** in both Bangladesh and the U.S., Mr. Leinweber was able to “predict” past U.S. stock returns with 99% accuracy.*

via Carl Richards, NYT 3/26/2012

The predictive viewpoint

THE SETUP

What do we mean by good predictions?

We make observations and then attempt to “predict” new, unobserved data.

Sometimes this is the same as estimating the mean.

Mostly, we observe $(y_1, x_1), \dots, (y_n, x_n)$, and we want some way to predict Y from X .

EVALUATING PREDICTIONS

Of course, both Y and \hat{Y} are **random**

I want to know how well I can predict **on average**

Let \hat{f} be some way of making predictions \hat{Y} of Y using covariates X

In fact, suppose I observe a dataset $\mathcal{D}_n = \{(Y_1, X_1), \dots, (Y_n, X_n)\}$. Then I want to **choose** some \hat{f} using \mathcal{D}_n .

Is \hat{f} good on average?

EVALUATING PREDICTIONS

Choose some **loss function** that measures prediction quality: $\ell : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$. We predict Y with \hat{Y}

Examples:

SQUARED-ERROR: $\ell(y, \hat{y}) = (y - \hat{y})^2$

ABSOLUTE-ERROR: $\ell(y, \hat{y}) = |y - \hat{y}|$

ZERO-ONE: $\ell(y, \hat{y}) = I(y \neq \hat{y}) = \begin{cases} 0 & y = \hat{y} \\ 1 & \text{else} \end{cases}$

Can be generalized to Y in arbitrary spaces.

PREDICTION RISK

PREDICTION RISK

$$R_n(\hat{f}) = \mathbb{E}[\ell(Y, \hat{f}(X))]$$

where the expectation is taken over the new data point (Y, X) and \mathcal{D}_n (everything that is random).

For **regression** applications, we will use squared-error loss:

$$R_n(\hat{f}) = \mathbb{E}[(Y - \hat{f}(X))^2]$$

For **classification** applications, we will use zero-one loss:

$$R_n(\hat{f}) = \mathbb{E}[I(Y \neq \hat{f}(X))]$$

Example 1: The mean

ESTIMATING THE MEAN

Suppose we know that we want to predict a quantity Y , where $\mathbb{E}[Y] = \mu \in \mathbb{R}$ and $\text{Var}[Y] = 1$. That is, $Y \sim P \in \mathbb{P}$, where

$$\mathbb{P} = \{P : \mathbb{E}[Y] = \mu \text{ and } \text{Var}[Y] = 1\}.$$

Our data is $\mathcal{D}_n = \{Y_1, \dots, Y_n\}$ such that $Y_i \stackrel{i.i.d.}{\sim} P$, and we want to estimate μ (and hence P).

ESTIMATING THE MEAN

Let $\hat{Y} = \bar{Y}_n$ be the sample mean.

We can ask about the **estimation risk** (since we're estimating μ):

$$\begin{aligned} R_n(\bar{Y}_n; \mu) &= \mathbb{E}[(\bar{Y}_n - \mu)^2] \\ &= \mathbb{E}[\bar{Y}_n^2] - 2\mu\mathbb{E}[\bar{Y}_n] + \mu^2 \\ &= \mu^2 + \frac{1}{n} - 2\mu^2 + \mu^2 \\ &= \frac{1}{n} \end{aligned}$$

PREDICTING NEW Y 'S

Let $\hat{Y} = \bar{Y}_n$ be the sample mean. This is a better idea.

What is the **prediction risk** of \bar{Y} ?

$$\begin{aligned} R_n(\bar{Y}_n) &= \mathbb{E}[(\bar{Y}_n - Y)^2] \\ &= \mathbb{E}[\bar{Y}_n^2] - 2\mathbb{E}[\bar{Y}_n Y] + \mathbb{E}[Y^2] \\ &= \mu^2 + \frac{1}{n} - 2\mu^2 + \mu^2 + 1 \\ &= 1 + \frac{1}{n} \end{aligned}$$

PREDICTING NEW Y 'S

What is the prediction risk of guessing $Y = 0$?

You can probably guess that this is a stupid idea.

Let's show why it's stupid.

$$\begin{aligned}R_n(0) &= \mathbb{E}[(0 - Y)^2] \\&= \mathbb{E}[(Y - \mu + \mu)^2] \\&= \mathbb{E}[(Y - \mu)^2] + 2\mu\mathbb{E}[(Y - \mu)] + \mu^2 \\&= 1 + \mu^2\end{aligned}$$

PREDICTING NEW Y 'S

What is the prediction risk of guessing $Y = \mu$?

This is a great idea, but we don't know μ .

Let's see what happens anyway.

$$\begin{aligned} R_n(\mu) &= \mathbb{E}[(Y - \mu)^2] \\ &= 1 \end{aligned}$$

ESTIMATING THE MEAN

Prediction risk: $R(\bar{Y}_n) = 1 + \frac{1}{n} = \sigma^2 + \frac{\sigma^2}{n}$

Estimation risk: $R(\bar{Y}_n; \mu) = \frac{1}{n}$

There is actually a nice interpretation here:

The common $1/n$ term is $\text{Var}[\bar{Y}_n]$

The extra factor of 1 in the prediction risk is **irreducible error** — Y is a random variable, and hence noisy. We can never eliminate its intrinsic variance. In other words, even if we knew μ , we could never get closer than 1, on average.

Intuitively, \bar{Y}_n is the obvious thing to do.

PREDICTING NEW Y 'S

Let's try one more: $\hat{Y}_a = a\bar{Y}_n$ for some $a \in (0, 1]$.

$$R_n(\hat{Y}_a) = \mathbb{E}[(\hat{Y}_a - Y)^2] = (1 - a)^2\mu^2 + \frac{a^2}{n} + 1$$

We can minimize this in a to get the best possible prediction risk for an estimator of the form \hat{Y}_a :

$$\operatorname{argmin}_a R_n(\hat{Y}_a) = \left(\frac{\mu^2}{\mu^2 + 1/n} \right) \bar{Y}_n$$

What happens if $\mu \ll 1$?

Wait a minute! You're saying there is a **better** estimator than \bar{Y}_n ?

PREDICTING NEW Y 'S

Let's try one more: $\hat{Y}_a = a\bar{Y}_n$ for some $a \in (0, 1]$.

$$R_n(\hat{Y}_a) = \mathbb{E}[(\hat{Y}_a - Y)^2] = (1 - a)^2\mu^2 + \frac{a^2}{n} + 1$$

We can minimize this in a to get the best possible prediction risk for an estimator of the form \hat{Y}_a :

$$\operatorname{argmin}_a R_n(\hat{Y}_a) = \left(\frac{\mu^2}{\mu^2 + 1/n} \right) \bar{Y}_n$$

What happens if $\mu \ll 1$?

Wait a minute! You're saying there is a **better** estimator than \bar{Y}_n ?

Why deal with prediction risk?

PREDICTION RISK

$$R_n(f) = \mathbb{E}[\ell(Y, f(X))]$$

Why care about $R_n(f)$?

Measures predictive accuracy on average.

How much confidence should you have in f 's predictions.

Compare with other models.

This is hard:

Don't know P (if I knew the truth, this would be easy)

WHAT IF YOU REALLY WANT TO MAKE INFERENCES?

- 1 You don't really care about predicting what will happen next year / quarter / millisecond
- 2 But you do want to offer an explanation / evaluate counterfactuals / describe the world
- 3 So you need the structure of your model to be at least approximately right
- 4 If you cannot predict well, then your model cannot be correct at all
- 5 Therefore, a necessary condition to believe your counterfactuals is good predictive accuracy

- Step 4 is about prediction error. (Sum of squared residuals??)
- Best not be fooling yourself in step 5.

RISK FOR GENERAL MODELS

We just saw that when you know the true model, and you have a nice estimator, the prediction risk has a nice decomposition

(this generalizes to much more complicated situations)

- Suppose we have a class of prediction functions \mathcal{F} ,

$$\text{e.g. } \mathcal{F} = \left\{ \beta : f(x) = \beta^\top x \right\}$$

- We use the data to choose some $\hat{f} \in \mathcal{F}$ and set $\hat{Y} = \hat{f}(X)$
- The true model is g (not necessarily in \mathcal{F}). Then:

$$R_n(\hat{f}) = \int \left[\text{bias}^2(\hat{f}(x)) + \text{var}(\hat{f}(x)) \right] p(x) dx + \sigma^2$$

where $X \sim p$ and

$$\text{bias}(\hat{f}(x)) = \mathbb{E}[\hat{f}(x)] - g(x)$$

$$\text{var}(\hat{f}(x)) = \mathbb{E}[(\hat{f}(x) - \mathbb{E}\hat{f}(x))^2]$$

$$\sigma^2 = \mathbb{E}[(Y - g(X))^2]$$

BIAS-VARIANCE DECOMPOSITION

So,

$$\begin{aligned}\text{prediction risk} &= \text{bias}^2 + \text{variance} + \text{irreducible error} \\ \text{estimation risk} &= \text{bias}^2 + \text{variance}\end{aligned}$$

What is $R(a)$ for our estimator $\hat{Y}_a = a\bar{Y}_n$?

$$\text{bias}(\hat{Y}_a) = \mathbb{E}[a\bar{Y}_n] - \mu = (a - 1)\mu$$

$$\text{var}(\hat{f}(x)) = \mathbb{E}[(a\bar{Y}_n - \mathbb{E}[a\bar{Y}_n])^2] = a^2\mathbb{E}[(\bar{Y}_n - \mu)^2] = \frac{a^2}{n}$$

$$\sigma^2 = \mathbb{E}[(Y - \mu)^2] = 1$$

$$\left(\text{That is: } R_n(\hat{Y}_a) = (a - 1)^2\mu^2 + \frac{a^2}{n} + 1 \right)$$

BIAS-VARIANCE DECOMPOSITION

So,

$$\begin{aligned}\text{prediction risk} &= \text{bias}^2 + \text{variance} + \text{irreducible error} \\ \text{estimation risk} &= \text{bias}^2 + \text{variance}\end{aligned}$$

What is $R(a)$ for our estimator $\hat{Y}_a = a\bar{Y}_n$?

$$\text{bias}(\hat{Y}_a) = \mathbb{E}[a\bar{Y}_n] - \mu = (a - 1)\mu$$

$$\text{var}(\hat{f}(x)) = \mathbb{E}[(a\bar{Y}_n - \mathbb{E}[a\bar{Y}_n])^2] = a^2\mathbb{E}[(\bar{Y}_n - \mu)^2] = \frac{a^2}{n}$$

$$\sigma^2 = \mathbb{E}[(Y - \mu)^2] = 1$$

$$\left(\text{That is: } R_n(\hat{Y}_a) = (a - 1)^2\mu^2 + \frac{a^2}{n} + 1 \right)$$

BIAS-VARIANCE DECOMPOSITION

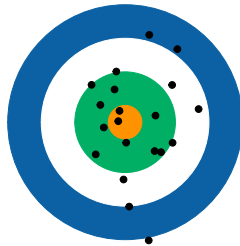
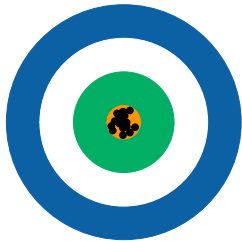
Important implication: prediction risk is proportional to estimation risk. However, defining estimation risk requires stronger assumptions.

In order to make good predictions, we want our prediction risk to be small. This means that we want to ‘balance’ the bias and variance.

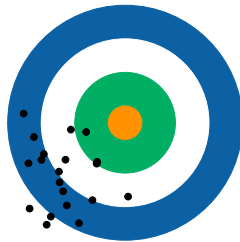
low variance

high variance

low bias

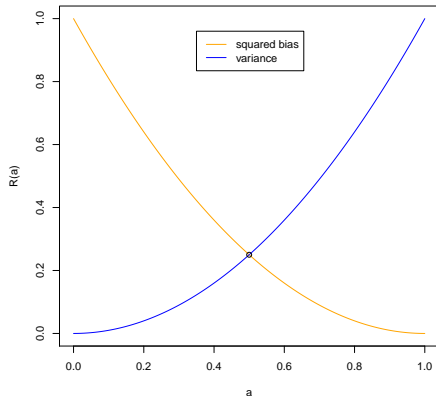


high bias



BIAS-VARIANCE TRADEOFF: ESTIMATING THE MEAN

$$R(a) = R_n(\hat{Y}_a) = (a - 1)^2 \mu^2 + \frac{a^2}{n} + 1$$



Here, $\mu = 1 (= \sigma^2)$ and $n = 1$, and hence optimal $a = .5$. This corresponds to the intersection of these curves.

WHAT?

Just to restate:

If $\mu = 1$ and $n = 1$ then it is better to predict with $.5Y$ than with Y itself.

In this case

$$R(a) = R_1(aY) = 1.5$$

$$R(Y) = 1 + 1 = 2$$

BIAS-VARIANCE TRADEOFF: OVERVIEW

- bias: how well does \hat{f} approximate the truth g
- more complicated \mathcal{F} , lower bias. Flexibility \Rightarrow Parsimony
- more flexibility \Rightarrow larger variance
- complicated models are hard to estimate precisely for fixed n
- irreducible error

Example 2: Normal means

NORMAL MEANS

Suppose we observe the following data:

$$Y_i = \beta_i + \epsilon_i, \quad i = 1, \dots, n$$

where $\epsilon_i \stackrel{iid}{\sim} \mathcal{N}(0, 1)$.

We want to estimate $\beta = (\beta_1, \dots, \beta_n)$.

The maximum likelihood estimator is $\hat{\beta}^{MLE} = (Y_1, \dots, Y_n)$.

This estimator has lots of nice properties: consistent, unbiased, UMVUE, (asymptotic) normality

NORMAL MEANS

But the MLE **STINKS!** It's a bad estimator.

It has no bias, but big variance.

$$R_n(\hat{\beta}^{MLE}) = \text{bias}^2 + \text{var} = 0 + n \cdot 1 = n$$

What if we use a biased estimator?

Consider the following estimator instead:

$$\hat{\beta}_i^S = \begin{cases} Y_i & i \in S \\ 0 & \text{else.} \end{cases}$$

Here $S \subseteq \{1, \dots, n\}$.

NORMAL MEANS

$$R_n(\hat{\beta}^S) = \sum_{i \notin S} \beta_i^2 + |S|.$$

In other words, if some $\beta_i < 1$, then don't bother estimating them!

In general, introduced parameters like S will be called **tuning parameters**.

Of course we don't know which $\beta_i < 1$.

But we could try to estimate $R_n(\hat{\beta}^S)$, and choose S to minimize our estimate.

ESTIMATING R_n

By definition, for any estimator $\hat{\beta}$,

$$R_n(\hat{\beta}) = \mathbb{E} \left[\sum_{i=1}^n (\hat{\beta}_i - \beta_i)^2 \right]$$

An intuitive estimator of R_n is

$$\hat{R}_n(\hat{\beta}) = \sum_{i=1}^n (\hat{\beta}_i - Y_i)^2.$$

This is known as the **training error** and it can be shown that

$$\hat{R}_n(\hat{\beta}) \approx R_n(\hat{\beta}).$$

Also,

$$\hat{\beta}^{MLE} = \underset{\beta}{\operatorname{argmin}} \hat{R}_n(\hat{\beta}^{MLE}).$$

What could possibly go wrong?

DANGERS OF USING THE TRAINING ERROR

Although

$$\hat{R}_n(\hat{\beta}) \approx R_n(\hat{\beta}),$$

this approximation can be very bad. In fact:

$$\begin{array}{llll} \text{TRAINING ERROR:} & \hat{R}_n(\hat{\beta}^{MLE}) & = & 0 \\ \text{RISK:} & R_n(\hat{\beta}^{MLE}) & = & n \end{array}$$

In this case, the **optimism** of the training error is n .

NORMAL MEANS

What about $\hat{\beta}^S$?

$$\hat{R}_n(\hat{\beta}^S) = \sum_{i=1}^n (\hat{\beta}_i - Y_i)^2 = \sum_{i \notin S} Y_i^2$$

Well

$$\mathbb{E} \left[\hat{R}_n(\hat{\beta}^S) \right] = R_n(\hat{\beta}^S) - 2|S| + n.$$

So I can choose S by minimizing $\hat{R}_n(\hat{\beta}^S) + 2|S|$.

Estimate of Risk = training error + penalty.

The penalty term corrects for the optimism.

Where we're going

THEMES OF COURSE

BIAS IS GOOD

- 1 Very often, we can trade some bias for (much) lower variance
- 2 Bias is controlled by setting **tuning parameters** (e.g. S)
- 3 Choosing tuning parameters carefully gives good risk properties
- 4 To know how to choose the tuning parameters, we need an estimate of the risk
- 5 Training error (which a risk estimator) is a bad choice (optimistic)
- 6 Unbiased estimators of parameters in correct models **may** have nice properties
- 7 Unbiased estimators of parameters in mis-specified models rarely have nice properties
- 8 All models are mis-specified

UP NEXT...

- 1 Regression and regularization
- 2 Economic forecasting and time series
- 3 Classification
- 4 More fun stuff. . .