

# Introductory Material and Linear Methods

Robert Schmidt

## Abstract

The following are notes on the key results from the **Elements of Statistical Learning** text. They were primarily derived from course notes and readings in the Stanford STATS 315: *Modern Applied Statistics* series.

## Contents

<b>1</b>	<b>Modeling basics</b>	<b>2</b>
1.1	Linear model bias-variance decomposition . . . . .	2
1.2	KNN bias-variance decomposition . . . . .	3
1.3	Variance of $\hat{\beta}$ . . . . .	3
1.4	Relationship between $t$ and $F$ statistic for dropping coefficients . . . . .	3
<b>2</b>	<b>Ridge regression</b>	<b>4</b>
2.1	Derivation using the SVD . . . . .	4
2.2	Eigenvalue interpretation . . . . .	4
2.3	Relationship between $\lambda$ and the coefficients . . . . .	4
2.4	Bayesian theory: ridge estimate is the mean of the posterior . . . . .	4
<b>3</b>	<b>Degrees of freedom</b>	<b>5</b>
3.1	Ridge df . . . . .	5
3.2	Lasso df . . . . .	5
<b>4</b>	<b>Methods using derived input directions</b>	<b>6</b>
4.1	Principal component regression up to $j$ coefficients . . . . .	6
4.2	PCR connection to OLS/ridge . . . . .	6
4.3	PCR vs. ridge . . . . .	6
<b>5</b>	<b>Zero-residual fits</b>	<b>7</b>
5.1	There are infinitely many OLS solutions with zero residuals . . . . .	7
5.2	Ridge solution has minimum norm among all zero-residual solutions as $\lambda \rightarrow 0$ . . . . .	7

# 1 Modeling basics

## 1.1 Linear model bias-variance decomposition

### 1) Overview

Conceptually, the bias-variance tradeoff is the problem of how closely one should fit a model to presented data. The tradeoff stipulates that there exists some minimum to the total error such that the model has optimal bias (how different the constructed function is from the true function) and variance (how much, on average, the estimated function differs from its expected value – essentially controlling the closeness of the function fit to the training data).

### 2) Mathematical specification

Mathematically, this can be derived from the squared-error loss formula, and similarly for other losses.

We consider squared error for illustrative purposes in this case.

Suppose  $Y = f(X) + \varepsilon$ , for noise  $\varepsilon$  where  $\mathbb{E}[\varepsilon] = 0$  and  $\text{Var}[\varepsilon] = \sigma_\varepsilon^2$ , and errors independent of  $Y$ .

Thus, the *MSE* at a point  $x_0$  is as follows:

$$\begin{aligned} MSE(x_0) &= \mathbb{E}[(Y - \hat{f}(x_0))^2 \mid X = x_0] \\ &= \mathbb{E}[(f(x_0) + \varepsilon - \hat{f}(x_0))^2] \\ &= \mathbb{E}[(\{f(x_0) - \hat{f}(x_0)\} + \varepsilon)^2] \end{aligned}$$

Now, we expand this expression.

$$\begin{aligned} MSE(x_0) &= \mathbb{E}[(f(x_0) - \hat{f}(x_0))^2 + 2\varepsilon(f(x_0) - \hat{f}(x_0)) + \varepsilon^2] \\ &= \mathbb{E}[(f(x_0) - \hat{f}(x_0))^2] + 2\mathbb{E}[\varepsilon(f(x_0) - \hat{f}(x_0))] + \mathbb{E}[\varepsilon^2] \quad \text{by linearity of expectation} \end{aligned}$$

First, consider the middle term. By error independence, we can simplify this term significantly.

$$2\mathbb{E}[\varepsilon(f(x_0) - \hat{f}(x_0))] = 2\mathbb{E}[\varepsilon]\mathbb{E}[f(x_0) - \hat{f}(x_0)] = 0 \quad \text{since } \mathbb{E}[\varepsilon] = 0$$

Since the middle term cancels completely, we see that

$$\begin{aligned} MSE(x_0) &= \mathbb{E}[(f(x_0) - \hat{f}(x_0))^2] + \mathbb{E}[\varepsilon^2] \\ &= \mathbb{E}[(f(x_0) - \hat{f}(x_0))^2] + \sigma_\varepsilon^2 \\ &= \text{Avg}(\text{model error}) + (\text{irreducible error}) \end{aligned}$$

Let us now expand the model error via the bias-variance decomposition.

$$\begin{aligned} \text{Model error} &= \mathbb{E}[(f(x_0) - \hat{f}(x_0))^2] \\ &= \mathbb{E}[(f(x_0) - \mathbb{E}\hat{f}(x_0) + \mathbb{E}\hat{f}(x_0) - \hat{f}(x_0))^2] \\ &= \mathbb{E}[(f(x_0) - \mathbb{E}\hat{f}(x_0))^2 + 2(f(x_0) - \mathbb{E}\hat{f}(x_0))(\mathbb{E}\hat{f}(x_0) - \hat{f}(x_0)) + (\mathbb{E}\hat{f}(x_0) - \hat{f}(x_0))^2] \end{aligned}$$

Again, we will focus on the middle term. If we note that  $f(x_0)$  and  $\mathbb{E}\hat{f}(x_0)$  are not random in regards to the expectation over the data, we can cancel this term:

$$\begin{aligned} 2\mathbb{E}[(f(x_0) - \mathbb{E}\hat{f}(x_0))(\mathbb{E}\hat{f}(x_0) - \hat{f}(x_0))] &= 2(f(x_0) - \mathbb{E}\hat{f}(x_0))\mathbb{E}[\mathbb{E}\hat{f}(x_0) - \hat{f}(x_0)] \\ &= 2(f(x_0) - \mathbb{E}\hat{f}(x_0))[\mathbb{E}\hat{f}(x_0) - \mathbb{E}\hat{f}(x_0)] \\ &= 0 \end{aligned}$$

So, we finally get:

$$\begin{aligned} \text{Model error} &= (\mathbb{E}\hat{f}(x_0) - f(x_0))^2 + \mathbb{E}[(\hat{f}(x_0) - \mathbb{E}\hat{f}(x_0))^2] \\ &= \text{Avg}(\text{bias}^2) + \text{Avg}(\text{variance}) \end{aligned}$$

And, returning to the overall average error, we see that:

$$\begin{aligned} MSE(x_0) &= (\mathbb{E}\hat{f}(x_0) - f(x_0))^2 + \mathbb{E}[(\hat{f}(x_0) - \mathbb{E}\hat{f}(x_0))^2] + \sigma_\varepsilon^2 \\ &= (\text{bias}^2) + (\text{variance}) + (\text{irreducible error}) \\ &= \text{bias-variance decomposition} \end{aligned}$$

## 1.2 KNN bias-variance decomposition

## 1.3 Variance of $\hat{\beta}$

### 1) Assumptions

1.  $y_i$  uncorrelated with constant variance  $\sigma^2$
2.  $x_i$  fixed/not random

### 2) Derivation

Consider the linear model  $y = \mathbf{X}\beta + \varepsilon$ .

- Let  $\mathbf{X}$  be fixed. This implies  $\mathbb{E}(y) = \mathbf{X}\beta$ .
- $\mathbb{E}(\hat{\beta}) = \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T y] = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \beta = \beta$ .
- So,

$$\begin{aligned}\hat{\beta} - \mathbb{E} \hat{\beta} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T y - (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \beta \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (y - \mathbf{X} \beta) \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \varepsilon\end{aligned}$$

$$\begin{aligned}\implies \text{Var}(\hat{\beta}) &= \mathbb{E}(\hat{\beta} - \mathbb{E} \hat{\beta})(\hat{\beta} - \mathbb{E} \hat{\beta})^T \\ &= \mathbb{E}((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \varepsilon)(\varepsilon^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}) \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{E}(\varepsilon \varepsilon^T) \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \text{Var}(\varepsilon) \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \quad \text{since } \mathbb{E}(\varepsilon) = 0\end{aligned}$$

- Recall the assumption that all entries of  $y$  are uncorrelated with variance  $\sigma^2$ . Hence,

$$\implies \text{Var}(\varepsilon) = \sigma^2 \mathbf{I}_N$$

$$\begin{aligned}\implies \text{Var}(\hat{\beta}) &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \sigma^2 \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \sigma^2\end{aligned}$$

This matches equation (3.8) on page 47 of *ESL*.

## 1.4 Relationship between $t$ and $F$ statistic for dropping coefficients

## 2 Ridge regression

### 2.1 Derivation using the SVD

As derived in ESL,  $\hat{\beta}(\lambda) = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$ .

We will re-write this estimate in terms of the SVD of  $\mathbf{X} = \mathbf{U} \mathbf{D} \mathbf{V}^T$ .

$$\begin{aligned}
 \hat{\beta}(\lambda) &= (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \\
 &= (\mathbf{V} \mathbf{D} \mathbf{U}^T \mathbf{U} \mathbf{D} \mathbf{V}^T + \lambda \mathbf{I})^{-1} \mathbf{V} \mathbf{D} \mathbf{U}^T \mathbf{y} \\
 &= (\mathbf{V} \mathbf{D}^2 \mathbf{V}^T + \lambda \mathbf{V} \mathbf{V}^T)^{-1} \mathbf{V} \mathbf{D} \mathbf{U}^T \mathbf{y} \quad \text{since } \mathbf{U}^T \mathbf{U} = \mathbf{I}_{p \times p} = \mathbf{V} \mathbf{V}^T \\
 &= (\mathbf{V} (\mathbf{D}^2 + \lambda \mathbf{I}) \mathbf{V}^T)^{-1} \mathbf{V} \mathbf{D} \mathbf{U}^T \mathbf{y} \\
 &= \mathbf{V} (\mathbf{D}^2 + \lambda \mathbf{I})^{-1} \mathbf{V}^{-1} \mathbf{V} \mathbf{D} \mathbf{U}^T \mathbf{y} \quad \text{since } (\mathbf{V}^T)^{-1} = \mathbf{V} \text{ given that } \mathbf{V} \text{ is orthogonal} \\
 &= \mathbf{V} (\mathbf{D}^2 + \lambda \mathbf{I})^{-1} \mathbf{D} \mathbf{U}^T \mathbf{y}
 \end{aligned}$$

Hence, in terms of the SVD,  $\boxed{\hat{\beta}(\lambda) = \mathbf{V} (\mathbf{D}^2 + \lambda \mathbf{I})^{-1} \mathbf{D} \mathbf{U}^T \mathbf{y}}$ .

### 2.2 Eigenvalue interpretation

### 2.3 Relationship between $\lambda$ and the coefficients

### 2.4 Bayesian theory: ridge estimate is the mean of the posterior

Assume that  $\mathbf{Y} | \beta \sim N(\mathbf{X}\beta, \sigma^2 \mathbf{I})$  and  $\beta \sim N(0, \tau^2 \mathbf{I})$ . Thus:

$$f(\mathbf{Y} | \beta) \propto \exp \left( -\frac{\|\mathbf{Y} - \mathbf{X}\beta\|^2}{2\sigma^2} \right) \quad \text{and} \quad f(\beta) \propto \exp \left( -\frac{\|\beta\|^2}{2\tau^2} \right)$$

Note that  $f(\beta | \mathbf{Y}) = \frac{f(\beta, \mathbf{Y})}{f(\mathbf{Y})} = \frac{f(\beta)f(\mathbf{Y}|\beta)}{f(\mathbf{Y})} \propto f(\beta)f(\mathbf{Y} | \beta)$  since  $f(\mathbf{Y})$  is a constant in terms of the function.

So, plugging in the expressions we found for these functions,

$$f(\beta | \mathbf{Y}) \propto \exp \left( -\frac{\|\mathbf{Y} - \mathbf{X}\beta\|^2}{2\sigma^2} - \frac{\|\beta\|^2}{2\tau^2} \right) = \exp \left( -\left( \frac{(\mathbf{Y} - \mathbf{X}\beta)^T (\mathbf{Y} - \mathbf{X}\beta)}{2\sigma^2} + \frac{\beta^T \beta}{2\tau^2} \right) \right)$$

We will focus our simplification efforts on the function inside the exponential.

$$\begin{aligned}
 \frac{(\mathbf{Y} - \mathbf{X}\beta)^T (\mathbf{Y} - \mathbf{X}\beta)}{2\sigma^2} + \frac{\beta^T \beta}{2\tau^2} &= \frac{\mathbf{Y}^T \mathbf{Y} - 2\mathbf{Y}^T \mathbf{X}\beta + \beta^T \mathbf{X}^T \mathbf{X}\beta}{2\sigma^2} + \frac{\beta^T \beta}{2\tau^2} \\
 &= \beta^T \left( \frac{\mathbf{X}^T \mathbf{X}}{2\sigma^2} + \frac{\mathbf{I}}{2\tau^2} \right) \beta - \left( \frac{\mathbf{Y}^T \mathbf{X}}{\sigma^2} \right) \beta + c
 \end{aligned}$$

This expression is quadratic in  $\beta$ ; we will factor it into the form  $(\beta - \mu)^T \left( \frac{\mathbf{X}^T \mathbf{X}}{2\sigma^2} + \frac{\mathbf{I}}{2\tau^2} \right) (\beta - \mu) + c$ .

Such a factorization would show that  $f(\beta | \mathbf{Y}) \propto \exp \left( (\beta - \mu)^T \left( \frac{\mathbf{X}^T \mathbf{X}}{2\sigma^2} + \frac{\mathbf{I}}{2\tau^2} \right) (\beta - \mu) \right)$ , which would imply that  $E[f(\beta | \mathbf{Y})] = \mu$  since  $f(\beta | \mathbf{Y})$  would be distributed multivariate Gaussian.

Since  $(\beta - \mu)^T \left( \frac{\mathbf{X}^T \mathbf{X}}{2\sigma^2} + \frac{\mathbf{I}}{2\tau^2} \right) (\beta - \mu) + c = \beta^T \left( \frac{\mathbf{X}^T \mathbf{X}}{2\sigma^2} + \frac{\mathbf{I}}{2\tau^2} \right) \beta - \left( \frac{\mathbf{Y}^T \mathbf{X}}{\sigma^2} \right) \beta + c'$ ,

we need  $-2\mu^T \left( \frac{\mathbf{X}^T \mathbf{X}}{2\sigma^2} + \frac{\mathbf{I}}{2\tau^2} \right) \beta = -\left( \frac{\mathbf{Y}^T \mathbf{X}}{\sigma^2} \right) \beta$ .

Focusing on this expression to solve for  $\mu$ :

$$\begin{aligned}
 -2\mu^T \left( \frac{\mathbf{X}^T \mathbf{X}}{2\sigma^2} + \frac{\mathbf{I}}{2\tau^2} \right) \beta &= -\left( \frac{\mathbf{Y}^T \mathbf{X}}{\sigma^2} \right) \beta \\
 \implies \mu &= \left( \frac{\mathbf{X}^T \mathbf{X}}{2\sigma^2} + \frac{\mathbf{I}}{2\tau^2} \right)^{-1} \left( \frac{\mathbf{Y}^T \mathbf{X}}{2\sigma^2} \right)^T \\
 \implies \mu &= \left( \frac{\mathbf{X}^T \mathbf{X}}{2\sigma^2} + \frac{\mathbf{I}}{2\tau^2} \right)^{-1} \left( \frac{\mathbf{X}^T \mathbf{Y}}{2\sigma^2} \right) = \left( \mathbf{X}^T \mathbf{X} + \frac{\sigma^2}{\tau^2} \mathbf{I} \right)^{-1} \mathbf{X}^T \mathbf{Y}
 \end{aligned}$$

Thus,  $\mu = (\mathbf{X}^T \mathbf{X} + \frac{\sigma^2}{\tau^2} \mathbf{I})^{-1} \mathbf{X}^T \mathbf{Y}$  is the mean of the posterior distribution,  $f(\beta | \mathbf{Y})$ .

Recall that the ridge regression solution is  $\hat{\beta}(\lambda) = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{Y}$ . This solution is identical to our result with the stipulation that  $\boxed{\lambda = \sigma^2 / \tau^2}$ . Hence, the ridge estimate is indeed the mean of the posterior distribution.

### 3 Degrees of freedom

#### 3.1 Ridge df

#### 3.2 Lasso df

## 4 Methods using derived input directions

### 4.1 Principal component regression up to $j$ coefficients

### 4.2 PCR connection to OLS/ridge

### 4.3 PCR vs. ridge

## 5 Zero-residual fits

5.1 There are infinitely many OLS solutions with zero residuals

5.2 Ridge solution has minimum norm among all zero-residual solutions as  $\lambda \rightarrow 0$