

LINEAR METHODS FOR REGRESSION: INTRODUCTION

-STATISTICAL MACHINE LEARNING-

Lecturer: Darren Homrighausen, PhD

THE SETUP

Suppose we have data

$$\mathcal{D} = \{(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)\},$$

where

- $X_i \in \mathbb{R}^p$ are the **features**
(or **explanatory variables** or **predictors** or **covariates**. NOT INDEPENDENT VARIABLES!)
- $Y_i \in \mathbb{R}$ are the **response** variables.
(NOT DEPENDENT VARIABLE!)

Our goal for this class is to find a way to explain (at least approximately) the **relationship** between X and Y .

PREDICTION RISK FOR REGRESSION

Given the **training data** \mathcal{D} , we want to predict some independent **test data** $Z = (X, Y) \sim \mathbb{P}$

This means forming a \hat{f} , which is a function of both the range of X and the training data \mathcal{D} , which provides predictions $\hat{Y} = \hat{f}(X)$.

The quality of this prediction is measured via the prediction risk¹

$$R(\hat{f}) = \mathbb{P}_{\mathcal{D}, Z}(Y - \hat{f}(X))^2 = \int l_{\hat{f}} d\mathbb{P}^{\mathcal{D}, Z}$$

We know that the **regression function**, $f_*(X) = \mathbb{P}[Y|X]$, is the best possible predictor.

Note that f_* is *unknown*

$$\mathbb{P}_Z l_{\hat{f}} = \int l_{\hat{f}} d\mathbb{P}$$

"IE"

¹Note: sometimes we integrate with respect to \mathcal{D} only, Z only, neither (loss), or both.

NOTATION RECAP

- X is a vector of **measurements** for each subject
(Example: $X_i = [1, \text{income}_i, \text{education}_i]^\top$)
- x is a vector of **subjects** for each measurement
(Example: $x_j = [\text{income}_1, \text{income}_2, \dots, \text{income}_n]^\top$)
- X_i^j is the j^{th} measurement on the i^{th} subject
(Example: $X_i^j = \text{income}_i$)

Imposing linearity

A LINEAR MODEL: MULTIPLE REGRESSION

If we specify the model: $f_*(X) = X^\top \beta = \sum_{j=1}^p x_j \beta_j$

$$\Rightarrow Y_i = X_i^\top \beta + \epsilon_i$$

Then we recover the usual linear regression formulation

$$\mathbb{X} = \begin{bmatrix} x_1 & \cdots & x_p \end{bmatrix} = \begin{bmatrix} X_1^\top \\ X_2^\top \\ \vdots \\ X_n^\top \end{bmatrix}. \quad Y = \mathbb{X}\beta + \epsilon$$

(When referring to j^{th} entry of any X_i , we write X_i^j)

Commonly, a column $x_0^\top = \underbrace{(1, \dots, 1)}_{n \text{ times}}$ is included

This encodes an intercept term, with intercept parameter β_0

We could (should?) seek to find a β such that $Y \approx \mathbb{X}\beta$

A LINEAR MODEL: POLYNOMIAL EFFECTS

Instead, we may believe

$$f_*(X) = \beta_0 + \sum_{j=1}^p X^j \beta_j + \sum_{j=1}^p \sum_{j'=1}^p X^j X^{j'} \alpha_{j,j'}$$

Then the **feature** matrix is

$$\mathbb{X} = \begin{bmatrix} x_0 & x_1 & \cdots & x_p & x_1^2 & x_1 x_2 & \cdots & x_p^2 \end{bmatrix}$$

(Here, interpret vector multiplication in the entrywise sense, as in **R**: $x * y$)

A LINEAR MODEL: GENERAL FORM

Specify functions $\phi_k : \mathbb{R}^p \rightarrow \mathbb{R}$, $k = 1, \dots, K$

$$\mathbb{X} = [\phi_k(X_i)] = \begin{bmatrix} \Phi(X_1)^\top \\ \Phi(X_2)^\top \\ \vdots \\ \Phi(X_n)^\top \end{bmatrix} \in \mathbb{R}^{n \times K},$$

where $\Phi(\cdot)^\top = (\phi_1(\cdot), \dots, \phi_K(\cdot))$.

EXAMPLE:

$$\phi_k(X) = X^j X^{j'}$$

is an interaction for the j^{th} and j'^{th} covariates

In this case $K = \binom{p}{2} + p = p(p-1)/2 + p = (p^2 + p)/2$

A LINEAR MODEL: GENERAL FORM

We don't know if f_* can actually be expressed as a linear function

Hence, write

$$\Phi = \left\{ f : \exists (\beta_k)_{k=1}^K \text{ such that } f = \sum_{k=1}^K \beta_k \phi_k = \beta^\top \Phi \right\}$$

and

$$f_{*,\Phi} = \operatorname{argmin}_{f \in \Phi} \mathbb{P} \ell_f.$$

The function $f_{*,\Phi}$ is known as the **linear oracle**

This is the object we are estimating when using a linear model

(Alternatively, we are **assuming** $f_* \in \Phi$)

A LINEAR MODEL: MULTIPLE REGRESSION REDUX

Let $K = p$ and define ϕ_k to be the coordinate projection map

That is,

$$\phi_k(X_i) \equiv X_i^k$$

We recover the usual linear regression formulation

$$\mathbb{X} = [\phi_k(X_i)] = \begin{bmatrix} \phi(X_1)^\top \\ \phi(X_2)^\top \\ \vdots \\ \phi(X_n)^\top \end{bmatrix} = \begin{bmatrix} X_1^1 & X_1^2 & \cdots & X_1^p \\ X_2^1 & X_2^2 & \cdots & X_2^p \\ \vdots & \vdots & \ddots & \vdots \\ X_n^1 & X_n^2 & \cdots & X_n^p \end{bmatrix} = \begin{bmatrix} X_1^\top \\ X_2^\top \\ \vdots \\ X_n^\top \end{bmatrix} .$$

Feature \longleftrightarrow Design

A LINEAR MODEL: ORTHOGONAL BASIS EXPANSION

Suppose $f_* \in \mathcal{F}$, where \mathcal{F} is a Hilbert space with norm induced by the inner product $\langle \cdot, \cdot \rangle$.

$$a, b \in \mathbb{R}^Q, \quad \langle a, b \rangle = a^T b$$

Let $(\phi_k)_{k=1}^{\infty}$ be an orthonormal basis for \mathcal{F}

Write

$$e_{\alpha} = (0, \dots, 0, 1, 0, \dots, 0) f_* = \sum_{k=1}^{\infty} \langle f_*, \phi_k \rangle \phi_k = \sum_{k=1}^{\infty} \beta_k \phi_k$$

α^T ENTRY

Then we can estimate $f_{*,\Phi}$ by finding the coefficients of the projection on Φ .

$$a = \sum_{q=1}^Q e_q a_q$$

$\| \cdot \|$
 $\langle e_1, a \rangle$

By Parseval's theorem for Hilbert spaces this induces an **approximation** error of $\sum_{k=K+1}^{\infty} \beta_k^2$.

$$\| f_* - \sum_{k=1}^K \beta_k \phi_k \|_2^2$$

This is small if f_* is smooth

(for instance, if f_* has m derivatives, then $\beta_k \asymp k^{-m}$)

A LINEAR MODEL: NEURAL NETS

Let

$$\phi_k(\mathbf{X}) = \sigma(\alpha_k^\top \mathbf{X} + b_k),$$

where $\sigma(t) = 1/(1 + e^{-t})$ is the **sigmoid** activation function.

Then we can form the **feature** matrix

$$\mathbb{X} = \begin{bmatrix} \phi_1(\mathbf{X}_1) & \phi_2(\mathbf{X}_1) & \cdots \\ \vdots & \vdots & \vdots \\ \phi_1(\mathbf{X}_n) & \phi_2(\mathbf{X}_n) & \cdots \end{bmatrix}$$

For future reference, this is a

“single-layer feed-forward neural network model with linear output”

(It is actually a bit more complicated, as the parameters in the σ map are estimated, and hence this is actually nonlinear)

A LINEAR MODEL: RADIAL BASIS FUNCTIONS

Let

$$\phi_k(\mathbf{X}) = e^{-\|\mu_k - \mathbf{X}\|_2^2 / \lambda_k}.$$

Then $f_{*,\phi}$ is called an²:

“Gaussian radial-basis function estimator”.

This turns out to be a parametric form of a more general technique known as **Gaussian process regression**.

²More on this later

Detour

NOTATION COMMENT

WARNING: It is common to conflate:

- the number of original **covariates** (p)
- the number of created **features** (K)

This means we will always write $\mathbb{X} \in \mathbb{R}^{n \times p}$, regardless of the transformation Φ that generates the matrix \mathbb{X}

The reasons for this are

- multiple regression comes from a particular, degenerate choice of Φ
- the mapping Φ is often not explicitly created (and $K = \infty$)

BOTTOM LINE: Think of X as the vector **after** transformations and $\mathbb{X} \in \mathbb{R}^{n \times p}$ regardless of the choice of Φ

End detour

TURNING THESE IDEAS INTO PROCEDURES

Each of these methods have parameters to choose:

- p could be very large. Do we include all covariates?
- If we include some polynomial (or other function) terms, should we include all of them?
- For neural nets, we need to choose: the activation function σ , the directions α_k , bias terms b_k , as well as the number of units in the hidden layer

Additionally, we need to estimate the associated coefficient vector β , α , or whatever

We would like the data to inform these parameters

TRAINING ERROR AND RISK ESTIMATION

The **linear oracle** is defined to be

$$f_{*,\Phi} = \operatorname{argmin}_{f \in \Phi} \mathbb{P} \ell_f.$$

(**REMINDER:** for regression, $\ell_f(Z) = (f(X) - Y)^2$)

$$R(f) = \mathbb{P} \ell_f$$

Hence, it is intuitive to use $\hat{\mathbb{P}}$ to form the **training error** " $E(f(X)-Y)^2$ "

$$\hat{R}(f) = \hat{\mathbb{P}} \ell_f = \frac{1}{n} \sum_{i=1}^n \ell_f(Z_i) = \frac{1}{n} \sum_{i=1}^n (f(X_i) - Y_i)^2 = \frac{1}{n} \|Y - \mathbb{X}\beta\|_2^2$$

Handwritten notes:
- An arrow points from $\int \ell_f d\mathbb{P}$ to the sum in the equation.
- An arrow points from $f(X) = X^T \beta$ to $\mathbb{X}\beta$ in the equation.

In many statistical applications, this **plug-in** estimator is minimized (Think of how many techniques rely on an unconstrained minimization of squared error, or maximum likelihood, or estimating equations, or ...)

This sometimes has disastrous results

EXAMPLE

Let's suppose \mathcal{D} is drawn from

$$n = 30$$

$$X = (0:n)/n*2*\pi$$

$$Y = \sin(X) + \text{rnorm}(n, 0, .25)$$

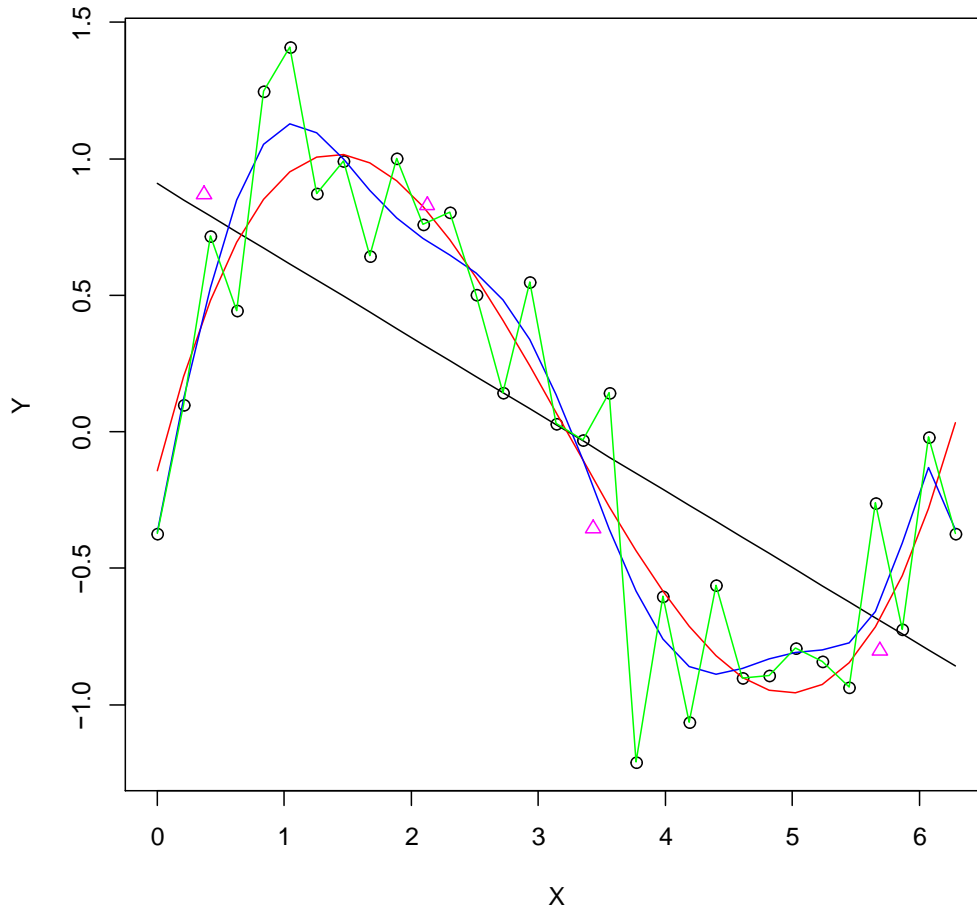
Now, let's fit some polynomials to this data.

We consider the following models:

- Model 1: $f(X_i) = \beta_0 + \beta_1 X_i$
- Model 2: $f(X_i) = \beta_0 + \beta_1 X_i + \beta_2 X_i^2 + \beta_3 X_i^3$
- Model 3: $f(X_i) = \sum_{k=0}^{10} \beta_k X_i^k$
- Model 4: $f(X_i) = \sum_{k=0}^{n-1} \beta_k X_i^k$

Let's look at what happens...

EXAMPLE



The \hat{R} 's are:

$$\hat{R}(\text{Model 1}) = 10.98$$

$$\hat{R}(\text{Model 2}) = 2.86$$

$$\hat{R}(\text{Model 3}) = 2.28$$

$$\hat{R}(\text{Model 4}) = 0$$

What about predicting new observations (\triangle)?

Bias and variance

PREDICTION RISK FOR REGRESSION

Note that $R(\hat{f})$ can be written as

$$R(\hat{f}) = \underbrace{\int \text{bias}^2(x) d\mathbb{P}_X}_{\text{Bias}} + \underbrace{\int \text{var}(x) d\mathbb{P}_X}_{\text{Variance}} + \underbrace{\sigma^2}_{\text{Noise}}$$

where

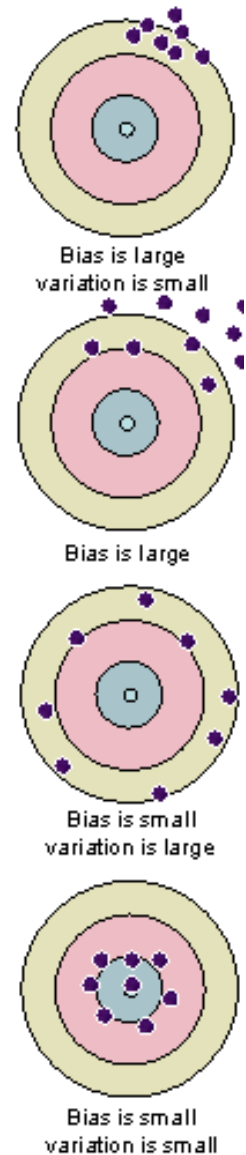
$$\text{bias}(x) = \mathbb{P}\hat{f}(x) - f_*(x)$$

$$\text{var}(x) = \mathbb{V}\hat{f}(x)$$

$$\sigma^2 = \mathbb{P}(Y - f_*(X))^2$$

(As an aside, this decomposition applies to much more general loss functions^a)

^aVariance and Bias for General Loss Functions; , Machine Learning 2003



BIAS-VARIANCE TRADEOFF

This can be heuristically thought of as

$$\text{Prediction risk} = \text{Bias}^2 + \text{Variance}.$$

There is a natural conservation between these quantities

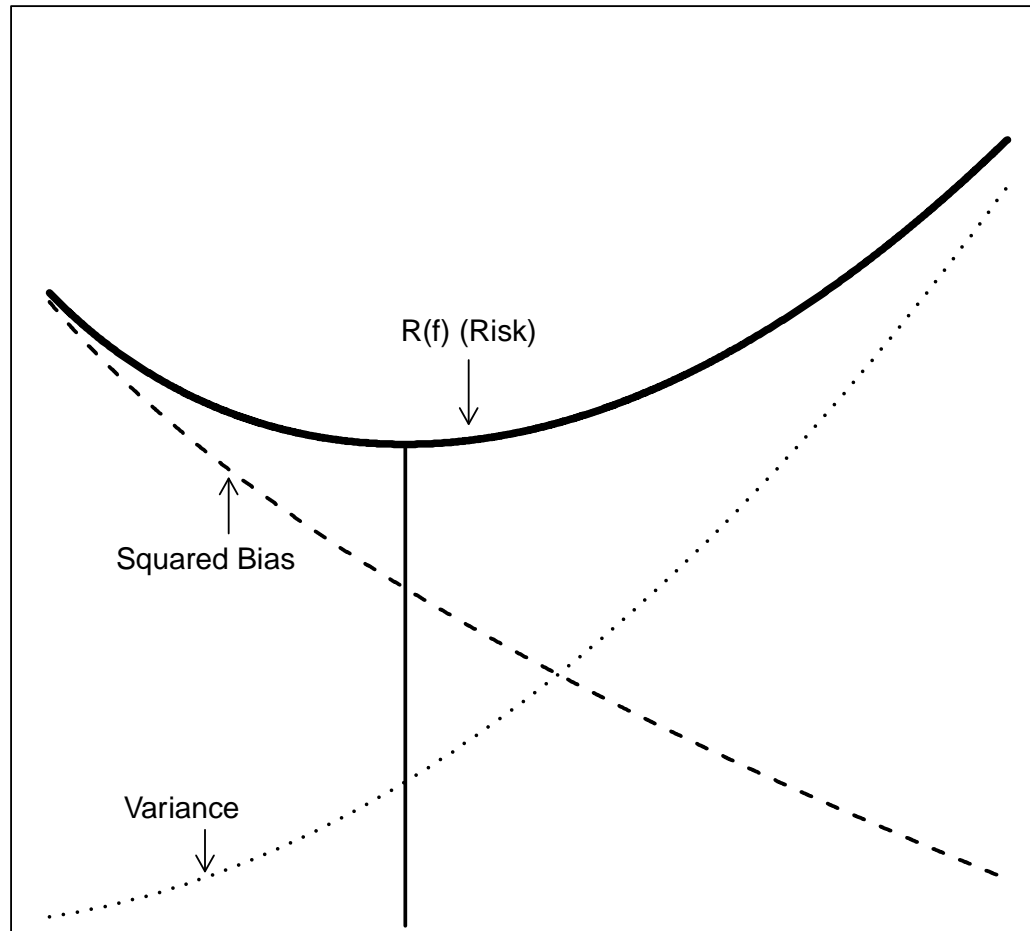
Low bias \rightarrow complex model \rightarrow many parameters \rightarrow high variance

The opposite also holds

(Think: $\hat{f} \equiv 0$.)

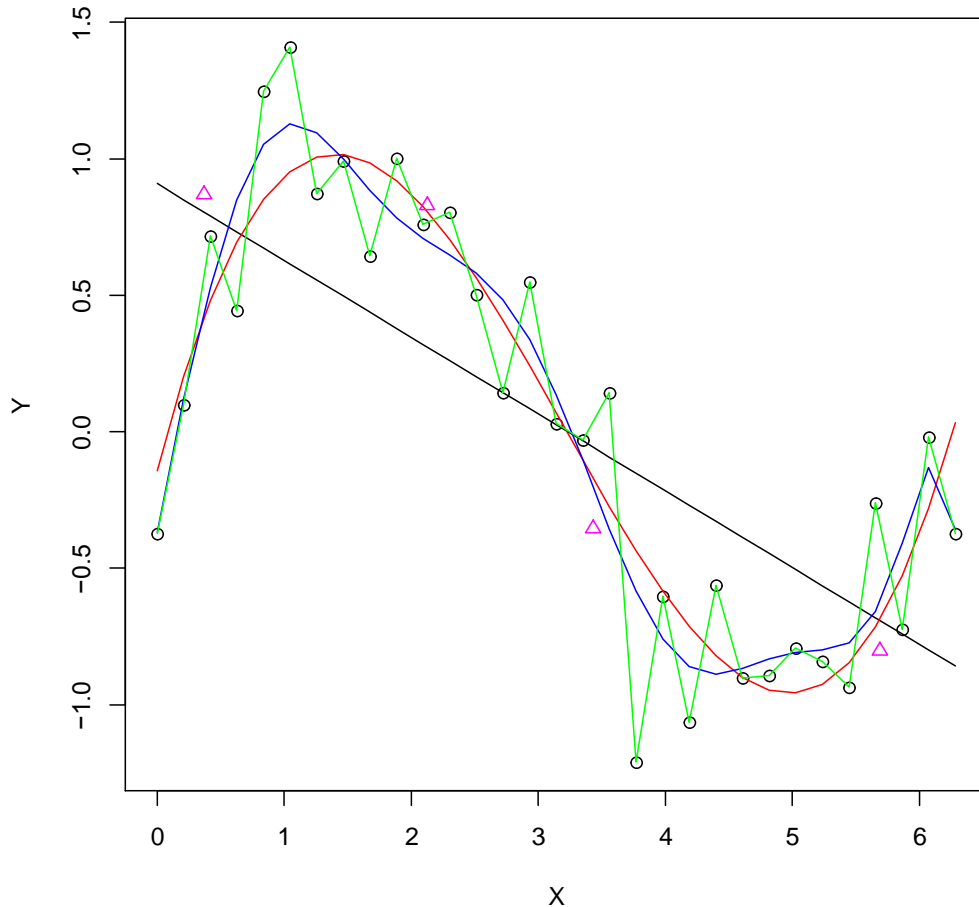
We'd like to 'balance' these quantities to get the best possible predictions

BIAS-VARIANCE TRADEOFF



Model Complexity ↗

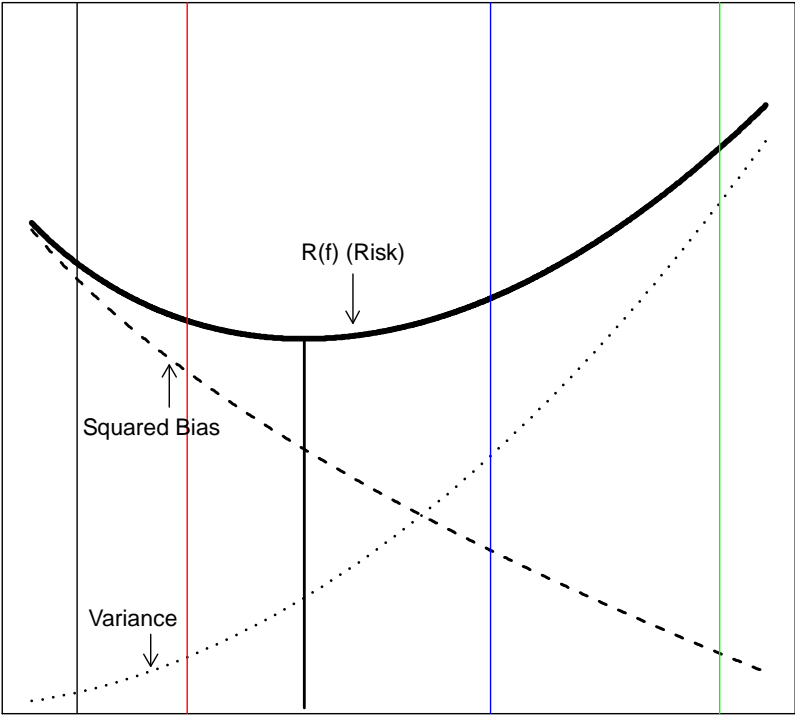
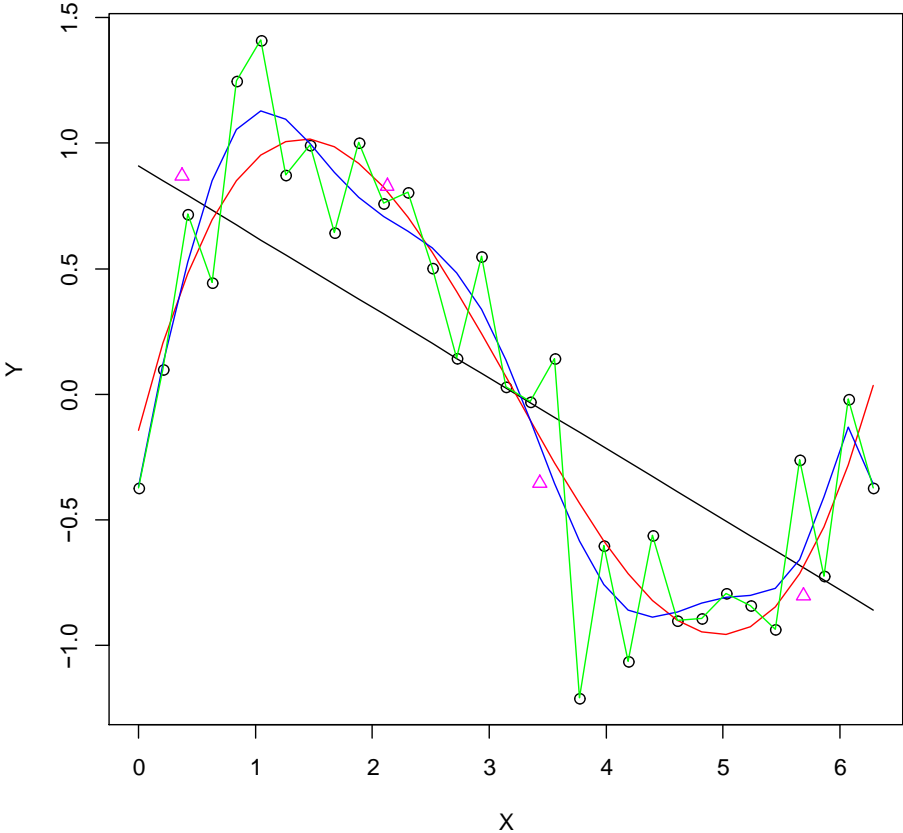
EXAMPLE



- Black model has low variance, high bias
- Green model has low bias, but high variance
- Red model and Blue model have intermediate bias and variance.

We want to balance these two quantities.

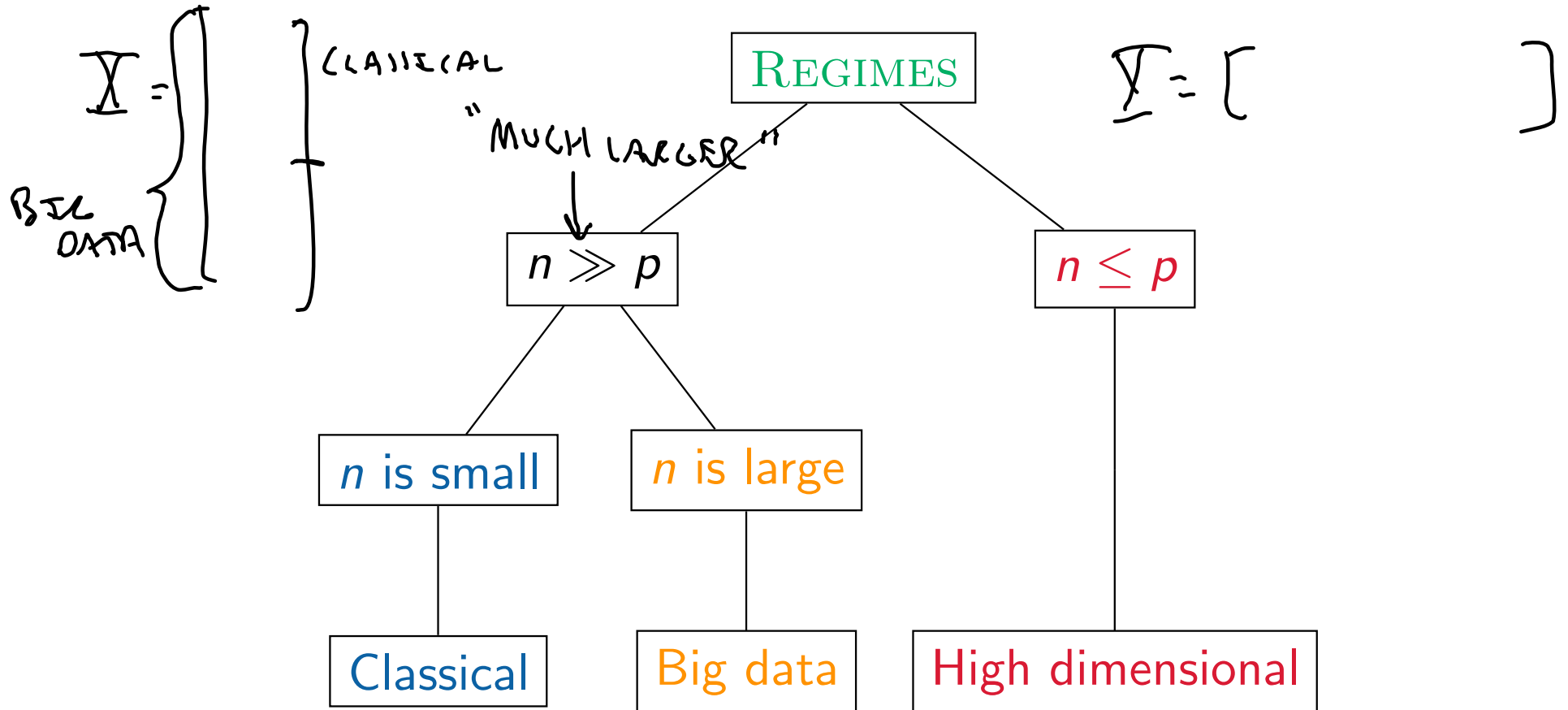
BIAS VS. VARIANCE



Model Complexity ↗

TURNING THESE IDEAS INTO PROCEDURES

There are roughly **three** regimes of interest, assuming $\mathbb{X} \in \mathbb{R}^{n \times p}$



CLASSICAL REGIME

Suppose we have the matrix \mathbb{X} with the features we're considering

Now, we want to estimate a parameter vector β in the model

$$Y = \mathbb{X}\beta + \epsilon$$

(E.g. we are modeling the **regression function** as (globally) linear in these features)

Minimize the training error $\hat{R}(f)$ over all functions $f_\beta(X) = X^\top \beta$

$$\hat{\beta}_{LS} = \operatorname{argmin}_{\beta} \hat{R}(f_\beta) = \operatorname{argmin}_{\beta} \underbrace{\|Y - \mathbb{X}\beta\|_2^2}_{= \Phi(X)^\top \beta}$$

(Though we write this as equality, there is only a unique solution if $\operatorname{rank}(\mathbb{X}) = p$)

CLASSICAL REGIME

In this case,

$$\hat{f}(X) = X^T \hat{\beta}_{LS} = X^T \mathbb{X}^\dagger Y \stackrel{\text{LEAST SQUARES SOLUTION}}{=} X^T (\mathbb{X}^T \mathbb{X})^{-1} \mathbb{X}^T Y$$

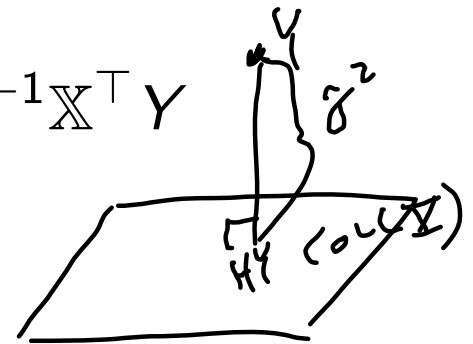
$\int_{\beta}(X) = X^T \beta$ PREDICTED $\hat{\beta}_{LS}$ $\text{rank}(\mathbb{X})=p$

(\mathbb{X}^\dagger is the Moore-Penrose pseudo inverse)

↳ SVD

The fitted values are $\mathbb{X} \hat{\beta}_{LS} = HY$, where H is the orthogonal projection onto the column space of \mathbb{X}

(Contrary to $\hat{\beta}_{LS}$, the fitted values are always unique)



CLASSICAL REGIME

We can examine the first and second moment properties of $\hat{\beta}_{LS}$

$$\mathbb{E}\hat{\beta}_{LS} = \beta \quad (\text{unbiased}) \quad (1)$$

$$\mathbb{V}\hat{\beta}_{LS} = \mathbb{X}^\dagger (\mathbb{V}Y) (\mathbb{X}^\dagger)^\top \quad \underbrace{=} \quad \mathbb{V}[Y_i] (\mathbb{X}^\top \mathbb{X})^{-1} \quad (2)$$

$\text{rank}(\mathbb{X})=p, \mathbb{V}Y \propto I_n$

NOTE: Here is where we need to be more careful:

The ‘true’ parameter β we are estimating is a coefficient vector of the linear oracle with respect to

$$\{f : \text{There exists } \beta \text{ where } f(X) = \beta^\top X\}$$

There is no reason to believe this approximation error is zero, hence ‘bias’ really references the linear oracle

CLASSICAL REGIME

The Gauss-Markov theorem assures us that this is the best linear **unbiased** estimator of β

(Effectively, equation (2) is minimized subject to equation (1))

Also, it is the maximum likelihood estimator under a homoskedastic, independent Gaussian model

(Hence, it is asymptotically efficient)

Does that necessarily mean it is any good?

CLASSICAL REGIME

Write $\mathbb{X} = UDV^T$ for the SVD of \mathbb{X}

Then $\mathbb{V}\hat{\beta}_{LS} \propto (\mathbb{X}^T \mathbb{X})^{-1} = VD^{-1} \underbrace{U^T U}_{=I} D^{-1} V^T = VD^{-2} V^T$

Handwritten annotations: "RIGHT SVEC" with an arrow pointing to V^T , "SVAL" with an arrow pointing to D , and "LEFT SVEC" with an arrow pointing to U .

REMINDER: Elements of D , d_j , are the axes lengths of the ellipse induced by \mathbb{X}



Also, suppose we are interested in estimating β ,

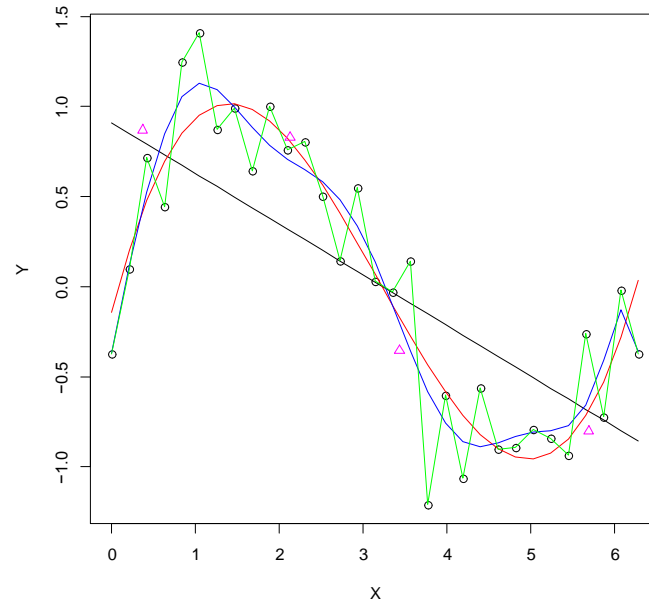
$$\mathbb{E} \|\hat{\beta}_{LS} - \beta\|_2^2 = \text{trace}(\mathbb{V}\hat{\beta}) \propto \sum_{j=1}^p \frac{1}{d_j^2}$$

Handwritten annotation: "trace(V^T V D^{-2})" with an arrow pointing to the trace term in the equation.

(Can you show this? Hint: add and subtract $\mathbb{E}\hat{\beta}_{LS}$)

IMPORTANT: Even in the classical regime, we can do arbitrarily badly if $d_p \approx 0$!

RETURNING TO POLYNOMIAL EXAMPLE: BIAS



Using a Taylor's series, for all X

$$\sin(X) = \sum_{q=0}^{\infty} \frac{(-1)^q X^{2q+1}}{(2q+1)!} = \Phi(X)^T \beta$$

Higher order polynomial models will **reduce** the bias part

RETURNING TO POLYNOMIAL EXAMPLE: VARIANCE

The least squares solution is given by solving $\min \|\mathbb{X}\beta - Y\|_2^2$

$$\mathbb{X} = \begin{bmatrix} 1 & X_1 & \dots & X_1^{p-1} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & X_n & \dots & X_n^{p-1} \end{bmatrix},$$

is the associated Vandermonde[#] matrix.

This matrix is well known for being numerically unstable

(Letting $\mathbb{X} = UDV^\top$, this means that $d_1/d_p \rightarrow \infty$)

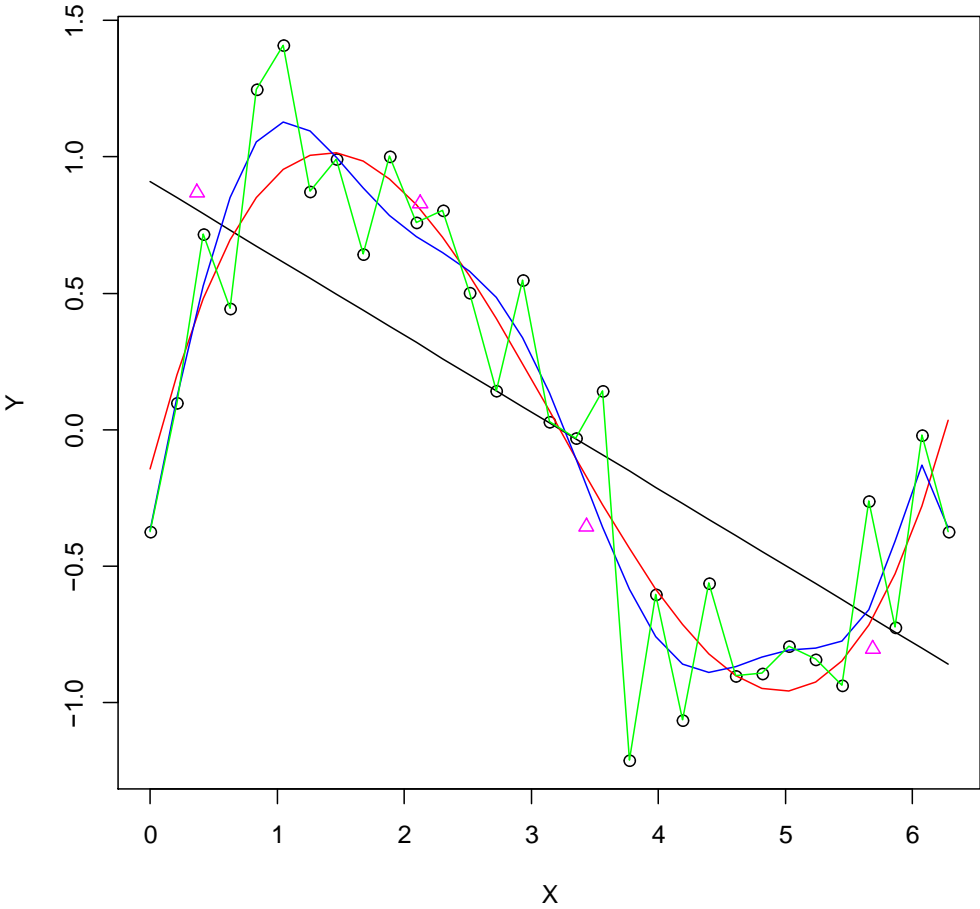
Hence³

$$\|(\mathbb{X}^\top \mathbb{X})^{-1}\|_2 = \frac{1}{d_p^2}$$

grows larger, where here $\|\cdot\|_2$ is the spectral (operator) norm[#]

³This should be compared with the variance computation in equation (2)

RETURNING TO THE POLYNOMIAL EXAMPLE



CONCLUSION

CONCLUSION: Fitting the full least squares model, even in the classical regime, can lead to poor prediction/estimation performance

In the other regimes, we encounter even for sinister problems

BIG DATA REGIME

Big data: The computational complexity scales extremely quickly. This means that procedures that are feasible classically are not for large data sets

EXAMPLE: Fit $\hat{\beta}_{LS}$ with $\mathbb{X} \in \mathbb{R}^{n \times p}$. Next fit $\hat{\beta}_{LS}$ with $\mathbb{X} \in \mathbb{R}^{3n \times 4p}$

The second case will take $\approx (3 * 4^2) = 48$ times longer to compute, as well as ≈ 12 times as much memory!

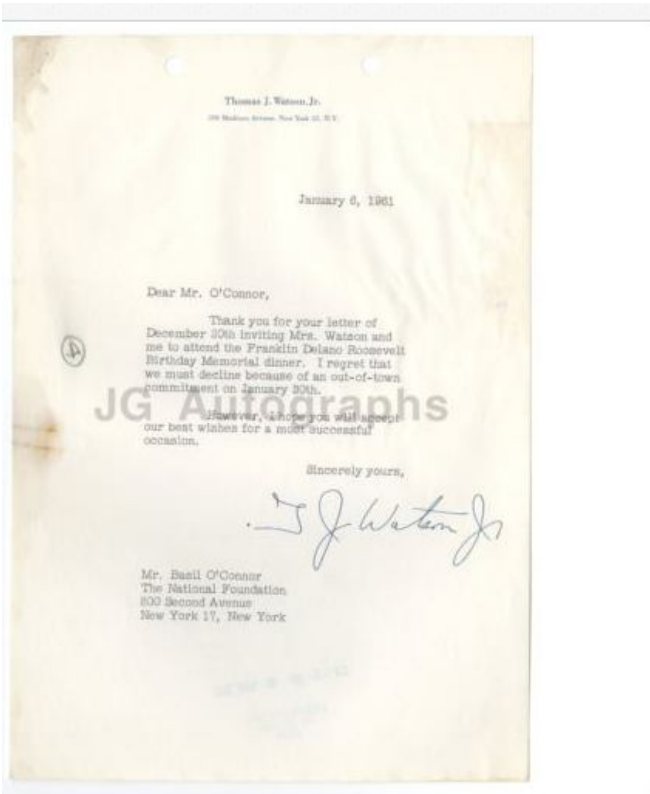
(Actually, for software such as **R** it might take 36 times as much memory, though there are data structures specifically engineered for this purpose that update objects 'in place')

CONCLUSION

```
p = 300; n = 10000
Y = rnorm(n); X = matrix(rnorm(n*p),nrow=n,ncol=p)
start = proc.time()[3]
out    = lm(Y~.,data=data.frame(X))
end    = proc.time()[3]
smallTime = end - start
```

```
n = nMultiple*n; nMultiple = 3
p = pMultiple*p; pMultiple = 4
Y = rnorm(n); X = matrix(rnorm(n*p),nrow=n,ncol=p)
start = proc.time()[3]
out    = lm(Y~.,data=data.frame(X))
end    = proc.time()[3]
bigTime = end - start
> print(bigTime/smallTime)
  elapsed
38.61458
> print(nMultiple*pMultiple**2)
[1] 48
```

EXAMPLE BIG DATA PROBLEM



Thomas Watson, Jr. - IBM Chairman - Authentic Autographed Letter (TLS)

Item condition: --

Ended: May 27, 2014 16:59:11 PDT

Winning bid: **US \$11.61** [6 bids]

Shipping: **\$3.99** Standard Shipping | [See details](#)

Item location: **United States**

Ships to: **Worldwide**

Delivery: Estimated within 3-6 business days 📅

Payments: **PayPal** | [See details](#)

Returns: 14 days money back, buyer pays return shipping | [See details](#)

Guarantee: **ebay** MONEY BACK GUARANTEE | [See details](#)

Get the item you ordered or get your money back.
Covers your purchase price and original shipping.

Seller information

jgautographs (64927 ★)
100% Positive feedback

[+ Follow this seller](#)

[See other items](#)

Visit store: [JG Autograph](#)

EXAMPLE BIG DATA PROBLEM

Buyer:

	Always a pleasure! Smooth & pleasant transaction!	f**a (3618 ★)	Jun-10-14 13:52
	Thomas Watson, Jr. - IBM Chairman - Authentic Autographed Letter (TLS) (#390846670600)	US \$11.61	View Item

Seller:

	Great communication. A pleasure to do business with.	Buyer: f**a (3618 ★)	Jun-05-14 18:59
	Thomas Watson, Jr. - IBM Chairman - Authentic Autographed Letter (TLS) (#390846670600)	–	View Item

The data (~750 Gb, millions of rows, thousands of columns):

User	ID	Rating	Comment	Role	WinBid	SellerID
dorkyporky	134	1	fast delivery.....very good seller...AAA++	B	15.51	princesskitten2001

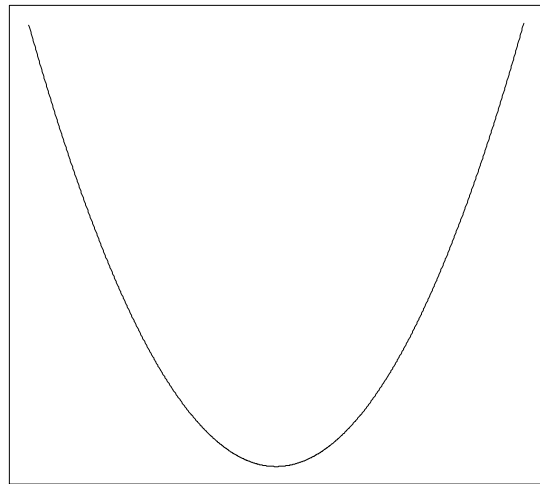
HIGH DIMENSIONAL REGIME

High dimensional: These problems tend to have many of the computational problems as **Big data**, as well as a **rank problem**:

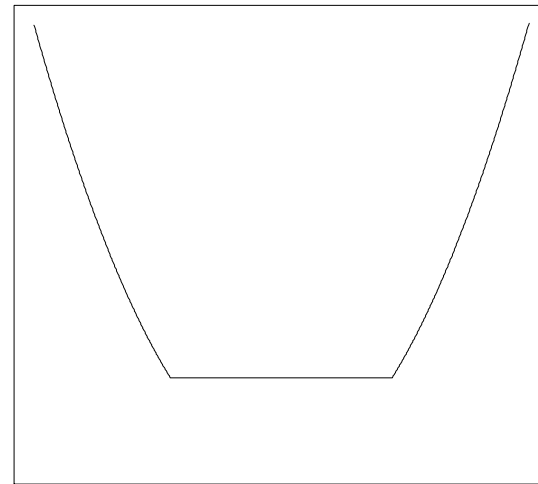
Suppose $\mathbb{X} \in \mathbb{R}^{n \times p}$ and $p > n$

Then $\text{rank}(\mathbb{X}) = n$ and the equation $\mathbb{X}\hat{\beta} = Y$:

- can be solved *exactly* (that is; the training error is 0)
- has an infinite number of solutions



$n > p$



$n < p$

HIGH DIMENSIONAL REGIME: EXAMPLE

