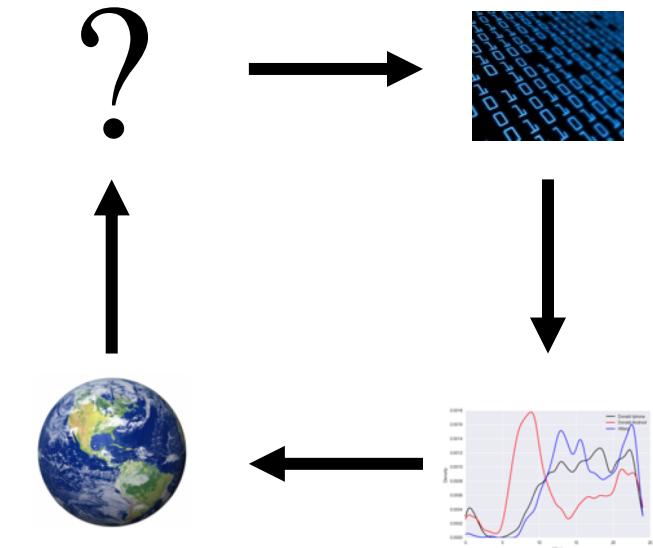


Linear Models & Feature Engineering

Slides by:

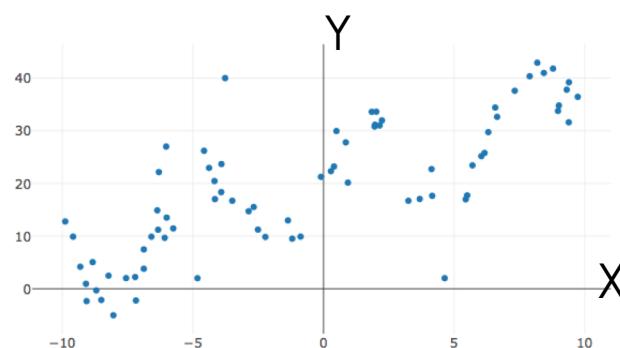
Joseph E. Gonzalez jegonzal@cs.berkeley.edu



Recap

Machine Modeling and Estimation (Learning)

Training Data

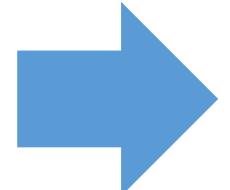


1. Define the model

$$\hat{y} = f_{\theta}(x) = \theta_0 + \theta_1 x$$

2. Choose a loss

$$L(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - f_{\theta}(x_i))^2$$



3. Minimize the loss

$$\hat{\theta} = \arg \min_{\theta} L(\theta)$$

Prediction (Testing)

Sometimes also called inference and scoring

1. Receive a **new** query point

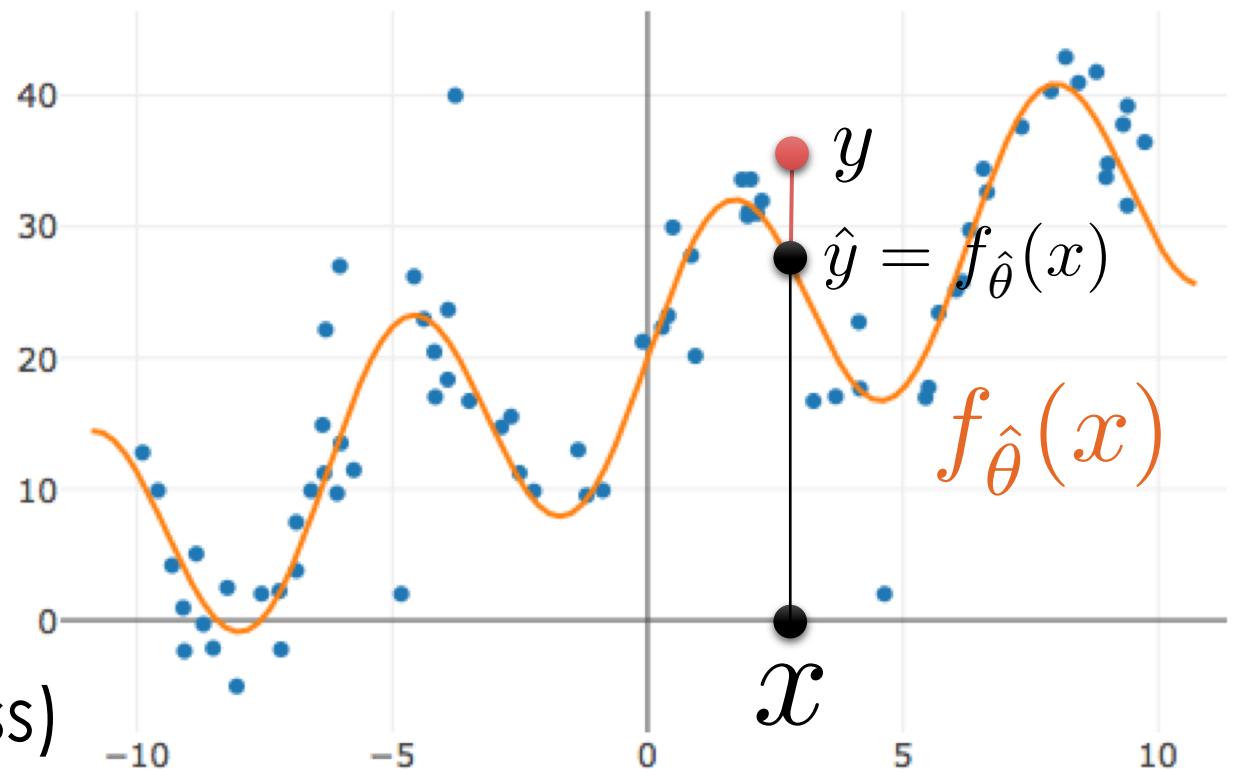
x

2. Make prediction using learned model

$$\hat{y} = f_{\hat{\theta}}(x)$$

3. Test Error (using squared loss)

$$(y - f_{\hat{\theta}}(x))^2 = (y - \hat{y})^2$$



Training Objective

$$\arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n (y_i - f_{\theta}(x_i))^2$$

- Minimize error on training data
 - sample of data from the world
 - estimate of the expected error
- We can compute this directly

Idealized Objective

$$\arg \min_{\theta} \mathbf{E} \left[(y - f_{\theta}(x))^2 \right]$$

- Minimize our expected prediction error over all possible test points
- **Ideal Goal**
 - Can't be computed ... ☹
- But we can analyze it!

Analysis of Squared Error

Quantities in **red** are random variables

Training on a **random sample** of data from the population.

$$(\mathbf{X}_i, Y_i) \sim \mathbf{P}(x, y) \quad \Rightarrow \quad \hat{\theta} = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n (Y_i - f_{\theta}(\mathbf{X}_i))^2$$

Testing at a given query point x and computing **expected squared error**

$$\mathbf{E} \left[(Y - f_{\hat{\theta}}(x))^2 \right]$$

Expectation is taken over all possible Y observations.

Expectation is taken over all possible training datasets

In the last lecture we showed that

$$\mathbf{E} \left[(Y - f_{\hat{\theta}}(x))^2 \right] =$$

Obs. Var. + **(Bias)²** + **Mod. Var.**

Other terminology:

“Noise” + **(Bias)²** + **Variance**

$$\mathbf{E} \left[(\textcolor{red}{Y} - f_{\hat{\theta}}(x))^2 \right] =$$

Assuming 0 mean observation
noise and true function $h(x)$

$$\textcolor{red}{Y} = h(x) + \epsilon$$

$$\mathbf{E} \left[(\textcolor{red}{Y} - h(x))^2 \right] +$$

Obs. Variance
“Noise”

$$(h(x) - \mathbf{E} [f_{\hat{\theta}}(x)])^2 +$$

(Bias)²

$$\mathbf{E} \left[(\mathbf{E} [f_{\hat{\theta}}(x)] - f_{\hat{\theta}}(x))^2 \right]$$

Model Variance

Alternative proof

Courtesy of Allen Shen

Assuming 0 mean observation
noise and true function $h(x)$

$$Y = h(x) + \epsilon$$

$$\mathbf{E} \left[(Y - f_{\hat{\theta}}(x))^2 \right] = \mathbf{E} \left[Y^2 - 2f_{\hat{\theta}}(x)Y + f_{\hat{\theta}}^2(x) \right]$$

Linearity of Expectation

$$= \mathbf{E} [Y^2] - \mathbf{E} [2f_{\hat{\theta}}(x)Y] + \mathbf{E} [f_{\hat{\theta}}^2(x)]$$

Definition of Y

$$= \mathbf{E} [(h(x) - \epsilon)^2] - \mathbf{E} [2f_{\hat{\theta}}(x)Y] + \mathbf{E} [f_{\hat{\theta}}^2(x)]$$

$$\mathbf{E} [(h(x) - \epsilon)^2] = h^2(x) - 2h(x)\mathbf{E} [\epsilon] + \mathbf{E} [\epsilon^2]$$

\Downarrow
 O

\Downarrow
 ∞

Defn' of ϵ

$$\mathbf{E} \left[(Y - f_{\hat{\theta}}(x))^2 \right] = \mathbf{E} \left[Y^2 - 2f_{\hat{\theta}}(x)Y + f_{\hat{\theta}}^2(x) \right]$$

Linearity of Expectation $= \mathbf{E} [Y^2] - \mathbf{E} [2f_{\hat{\theta}}(x)Y] + \mathbf{E} [f_{\hat{\theta}}^2(x)]$

Definition of Y $= \mathbf{E} [(h(x) - \epsilon)^2] - \mathbf{E} [2f_{\hat{\theta}}(x)Y] + \mathbf{E} [f_{\hat{\theta}}^2(x)]$

$$\mathbf{E} [(h(x) - \epsilon)^2] = h^2(x) - 2h(x)\mathbf{E} [\epsilon] + \mathbf{E} [\epsilon^2]$$

\Downarrow

\Downarrow σ^2 Defn' of ϵ

$$= h(x)^2 + \sigma^2 - \mathbf{E} [2f_{\hat{\theta}}(x)Y] + \mathbf{E} [f_{\hat{\theta}}^2(x)]$$

$$\begin{aligned}
\mathbf{E} \left[(Y - f_{\hat{\theta}}(x))^2 \right] &= \mathbf{E} \left[Y^2 - 2f_{\hat{\theta}}(x)Y + f_{\hat{\theta}}^2(x) \right] \\
&= h(x)^2 + \sigma^2 - \mathbf{E} \left[2f_{\hat{\theta}}(x)Y \right] + \mathbf{E} \left[f_{\hat{\theta}}^2(x) \right] \\
&\stackrel{\text{Y is independent of } \theta \text{ (only depends on noise)}}{=} h(x)^2 + \sigma^2 - 2\mathbf{E} \left[f_{\hat{\theta}}(x) \right] \mathbf{E} [Y] + \mathbf{E} \left[f_{\hat{\theta}}^2(x) \right] \\
&= h(x)^2 + \sigma^2 - 2\mathbf{E} \left[f_{\hat{\theta}}(x) \right] \mathbf{E} [h(x) + \epsilon] + \mathbf{E} \left[f_{\hat{\theta}}^2(x) \right] \\
&\quad \text{Definition of Y} \\
&= h(x)^2 + \sigma^2 - 2\mathbf{E} \left[f_{\hat{\theta}}(x) \right] h(x) + \mathbf{E} \left[f_{\hat{\theta}}^2(x) \right]
\end{aligned}$$

Linearity of expectation

Assuming 0 mean observation noise and true function $h(x)$

$$Y = h(x) + \epsilon$$

$$\mathbf{E} \left[(Y - f_{\hat{\theta}}(x))^2 \right] = \mathbf{E} \left[Y^2 - 2f_{\hat{\theta}}(x)Y + f_{\hat{\theta}}^2(x) \right]$$

$$= h(x)^2 + \sigma^2 - 2\mathbf{E} [f_{\hat{\theta}}(x)] h(x) + \mathbf{E} [f_{\hat{\theta}}^2(x)]$$

Definition of Variance

$$\mathbf{Var} [f_{\hat{\theta}}] = \mathbf{E} [f_{\hat{\theta}}^2(x)] - \mathbf{E} [f_{\hat{\theta}}(x)]^2$$

$$= h(x)^2 + \sigma^2 - 2\mathbf{E} [f_{\hat{\theta}}(x)] h(x) + \mathbf{E} [f_{\hat{\theta}}(x)]^2 + \mathbf{Var} [f_{\hat{\theta}}(x)]$$

Rearranging terms

$$= \sigma^2 + h(x)^2 - 2\mathbf{E} [f_{\hat{\theta}}(x)] h(x) + \mathbf{E} [f_{\hat{\theta}}(x)]^2 + \mathbf{Var} [f_{\hat{\theta}}(x)]$$

$$= \sigma^2 + (h(x) - \mathbf{E} [f_{\hat{\theta}}(x)])^2 + \mathbf{Var} [f_{\hat{\theta}}(x)]$$

Bonus study material!

Summary

$$(\textcolor{red}{X}_i, \textcolor{red}{Y}_i) \sim \mathbf{P}(x, y) \rightarrow \hat{\theta} = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n (\textcolor{red}{Y}_i - f_{\theta}(\textcolor{red}{X}_i))^2$$

Expectation is taken over all possible Y observations.

$$\mathbf{E} \left[(\textcolor{red}{Y} - f_{\hat{\theta}}(x))^2 \right] = \sigma^2 + (h(x) - \mathbf{E} [f_{\hat{\theta}}(x)])^2 + \mathbf{Var} [f_{\hat{\theta}}(x)]$$

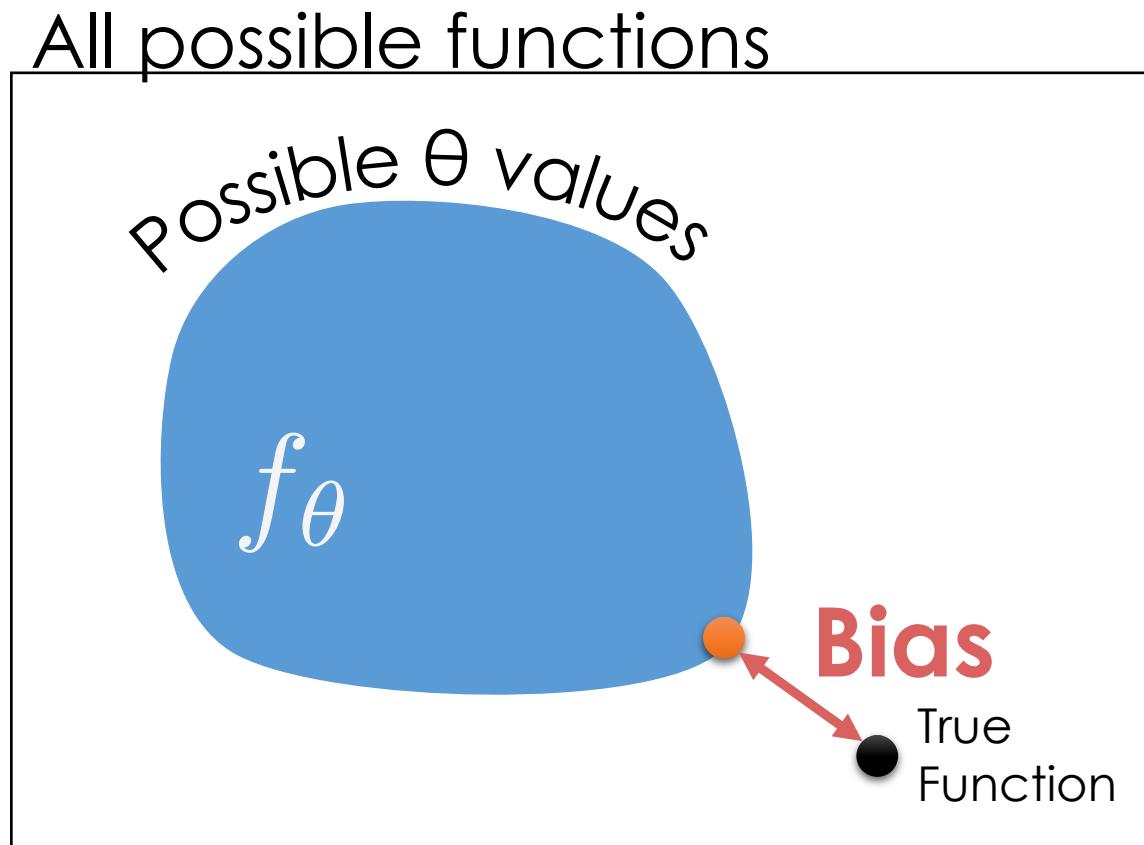
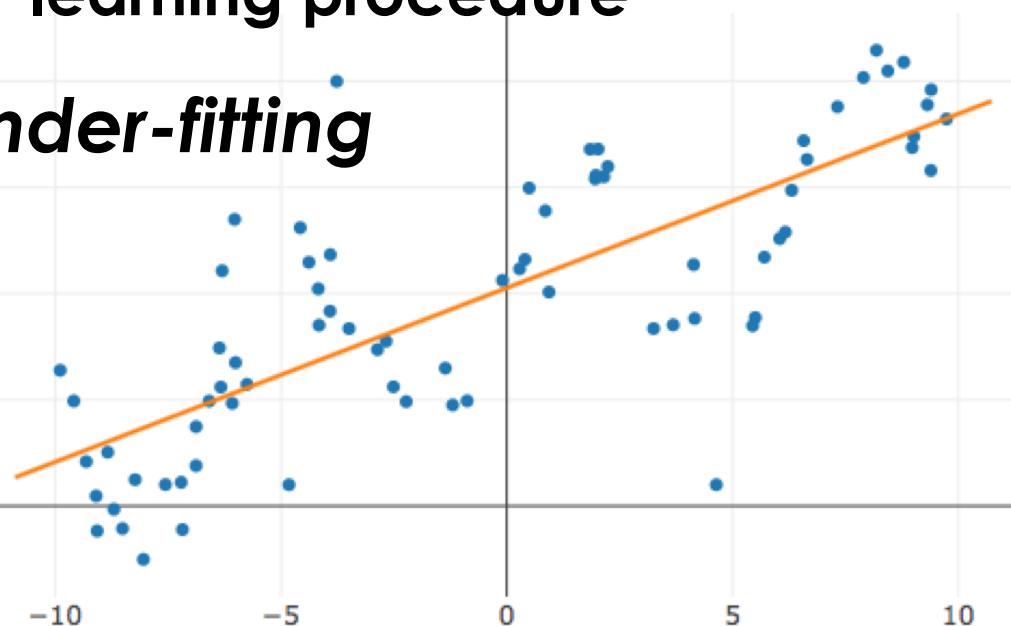
Obs. Var. + (Bias)² + Mod. Var.

Expectation is taken over all possible training datasets

$$\text{Bias} = h(x) - \mathbf{E} [f_{\hat{\theta}}(x)]$$

The expected deviation between the predicted value and the true value

- Depends on both the:
 - choice of f
 - learning procedure
- **Under-fitting**

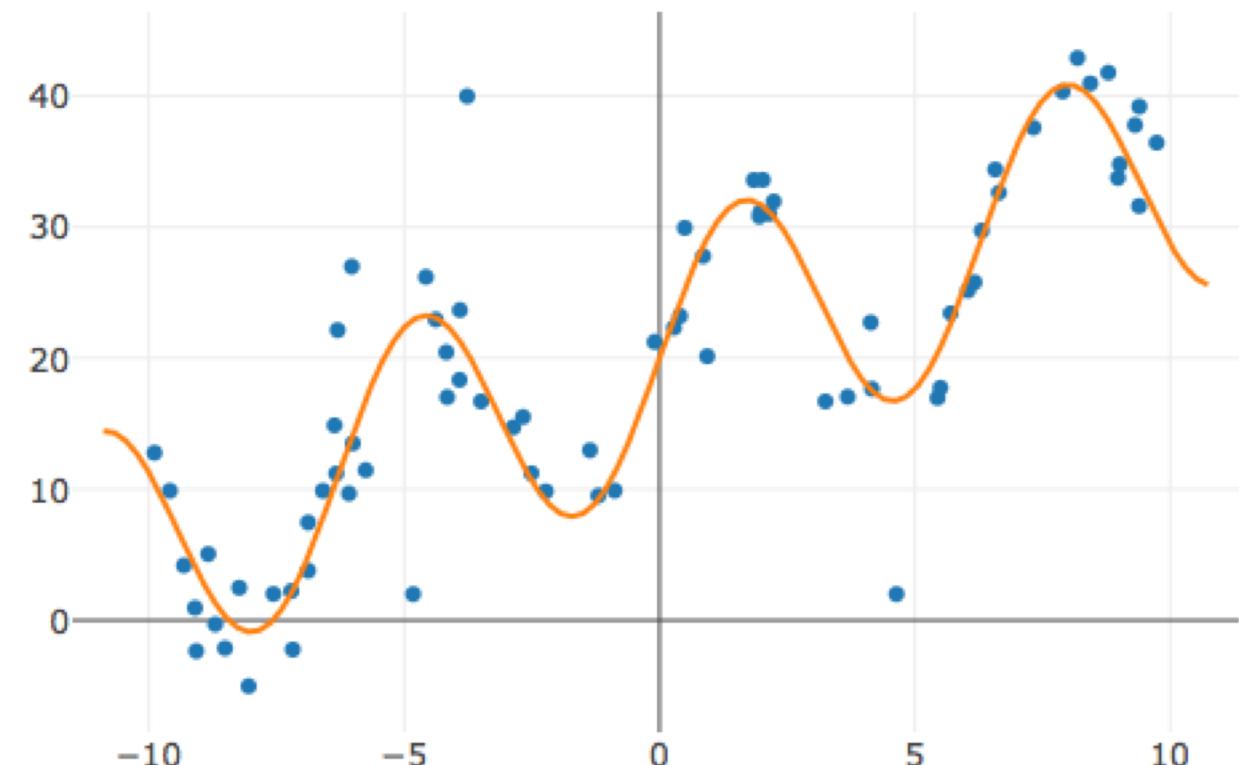


$$\text{Observation Variance} = \mathbb{E}[(Y - h(x))^2] = \sigma^2$$

the variability of the random noise in the process we are trying to model

- measurement variability
- stochasticity
- missing information

**Beyond our control
(usually)**

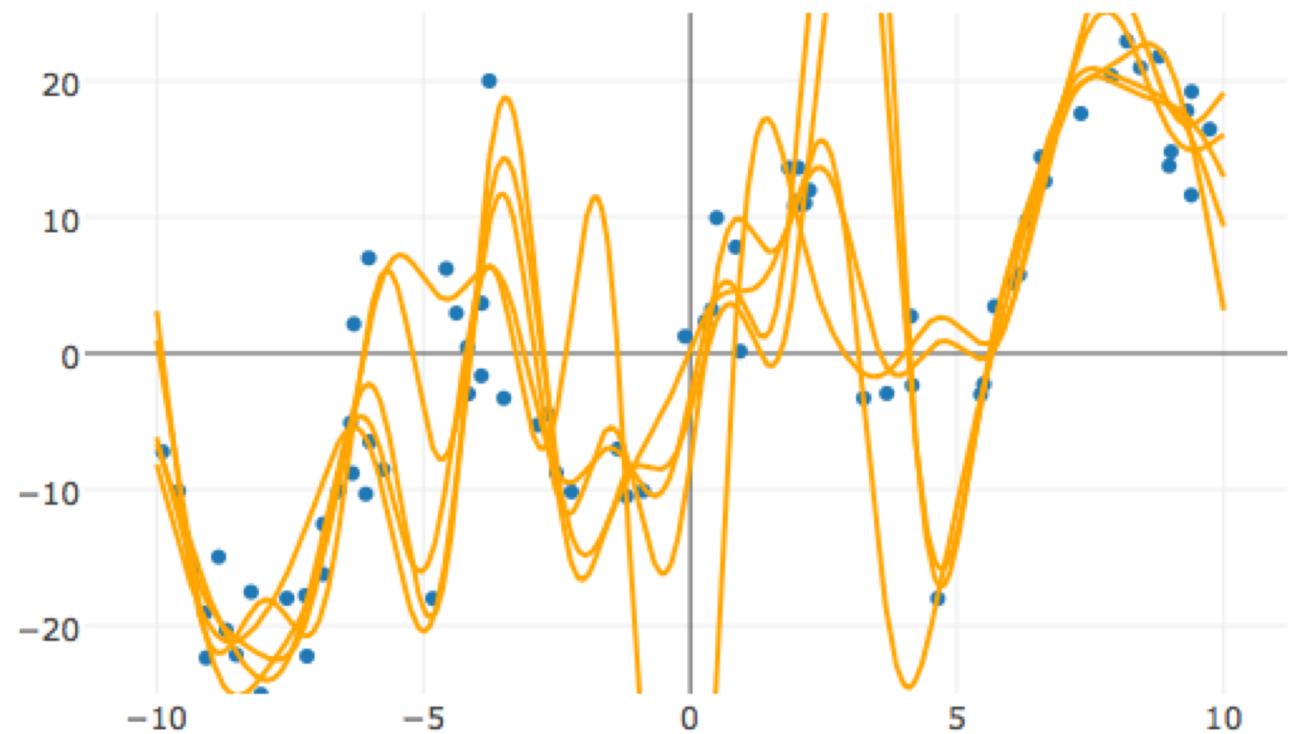


Estimated Model Variance =

$$\text{Var} [f_{\hat{\theta}}(x)] = \mathbf{E} [(f_{\hat{\theta}}(x) - \mathbf{E} [f_{\hat{\theta}}(x)])]$$

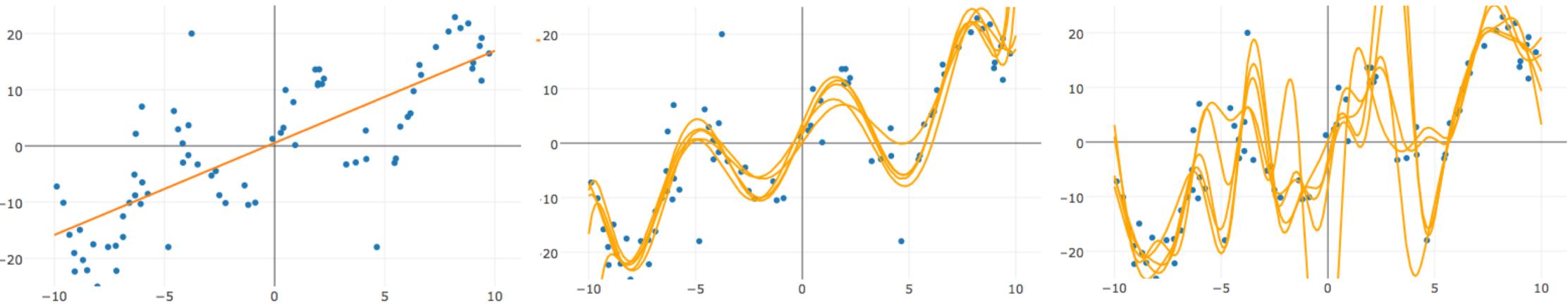
variability in the predicted value across different training datasets

- Sensitivity to variation in the training data
- Poor generalization
- **Overfitting**



The Bias-Variance Tradeoff

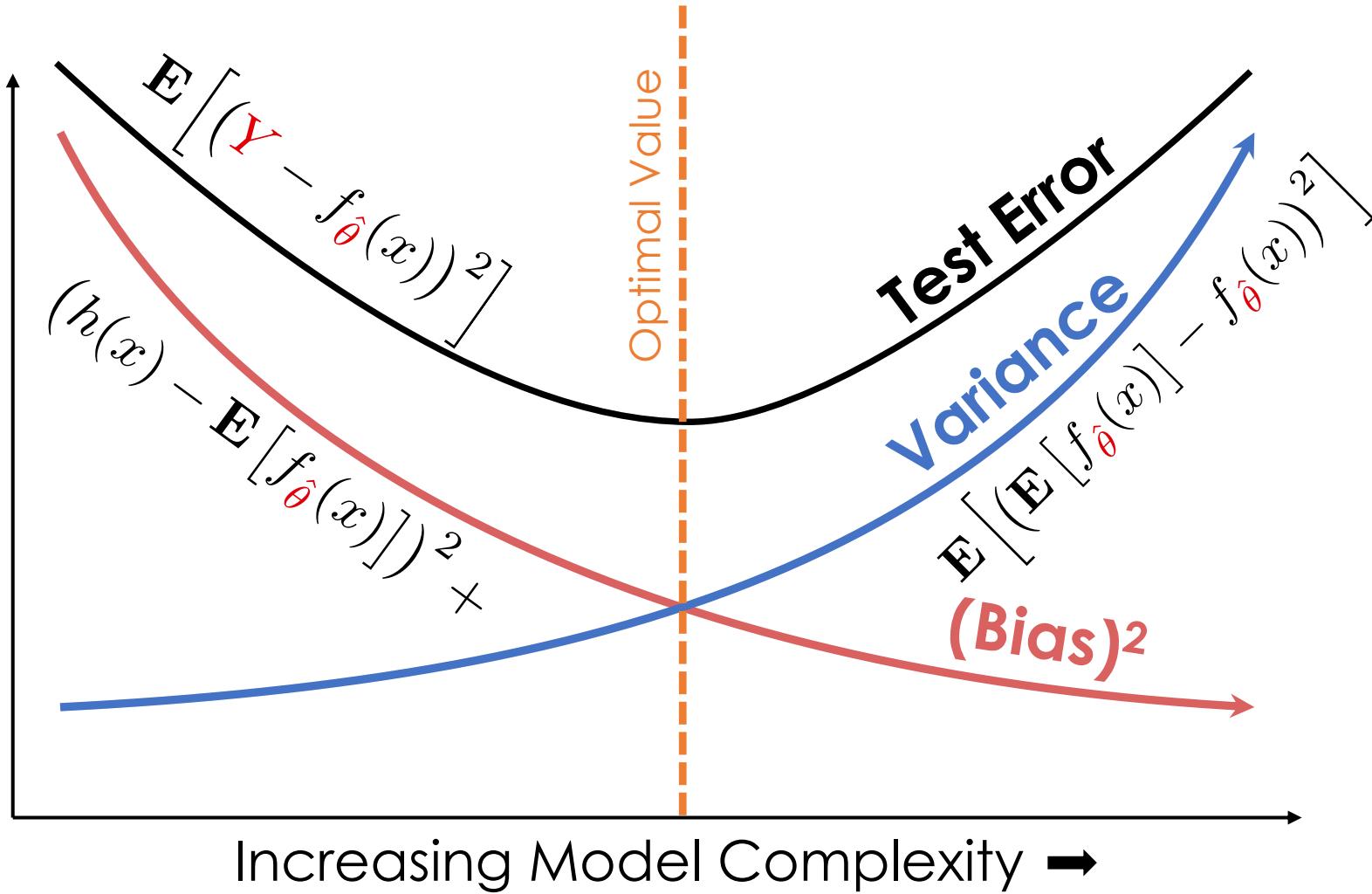
Estimated Model Variance



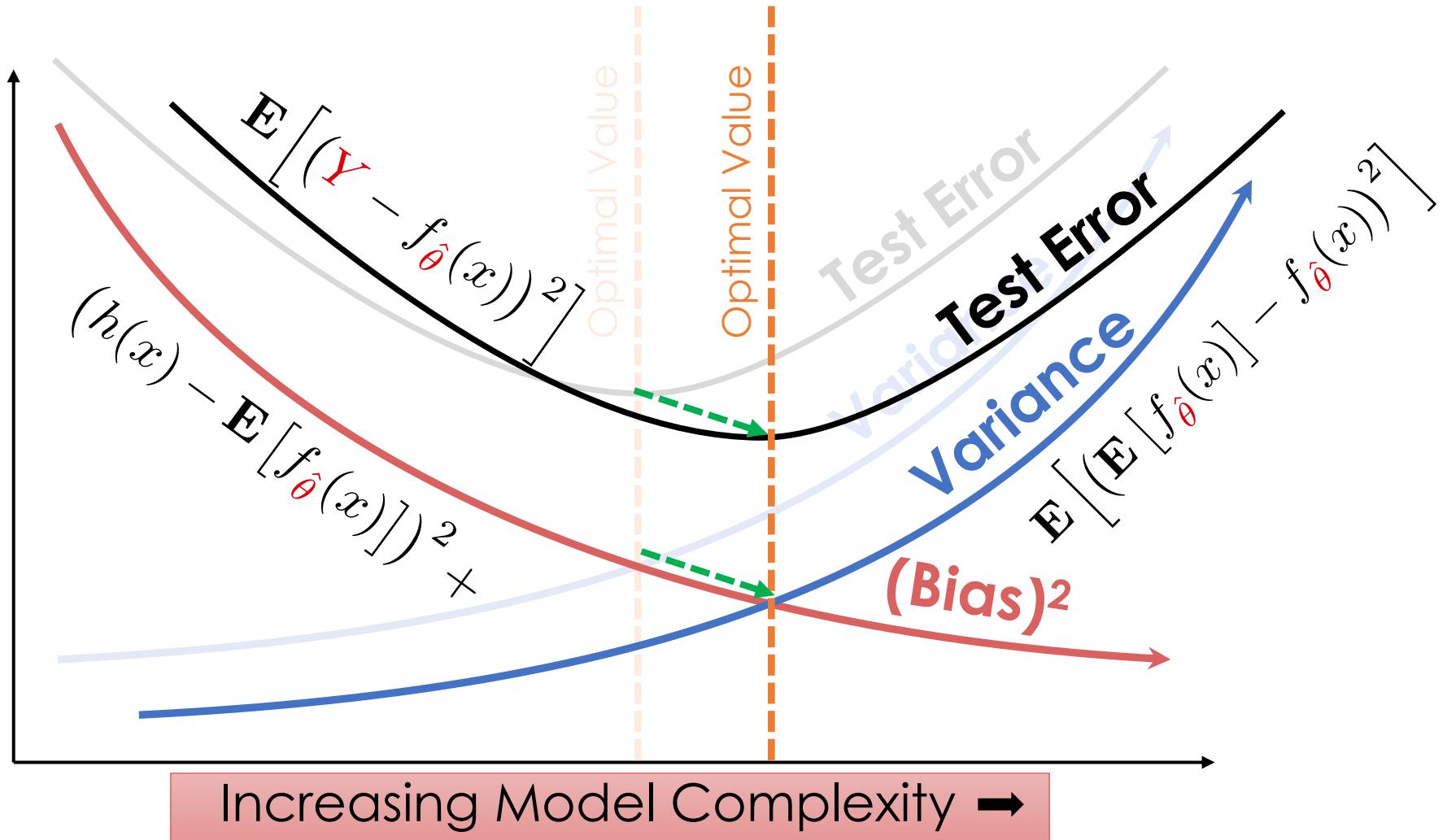
We want to **decrease both bias and variance** but often decreasing one results in an increase in the other.

Bias

Bias Variance Plot



More Data supports More Complexity



Model Complexity

- Roughly: capacity of the model to fit the data
- Many different measures and factors
 - Covered in machine learning class
- Dominant factors in **linear models**
 - Number and types of features
 - Regularization



Start with this

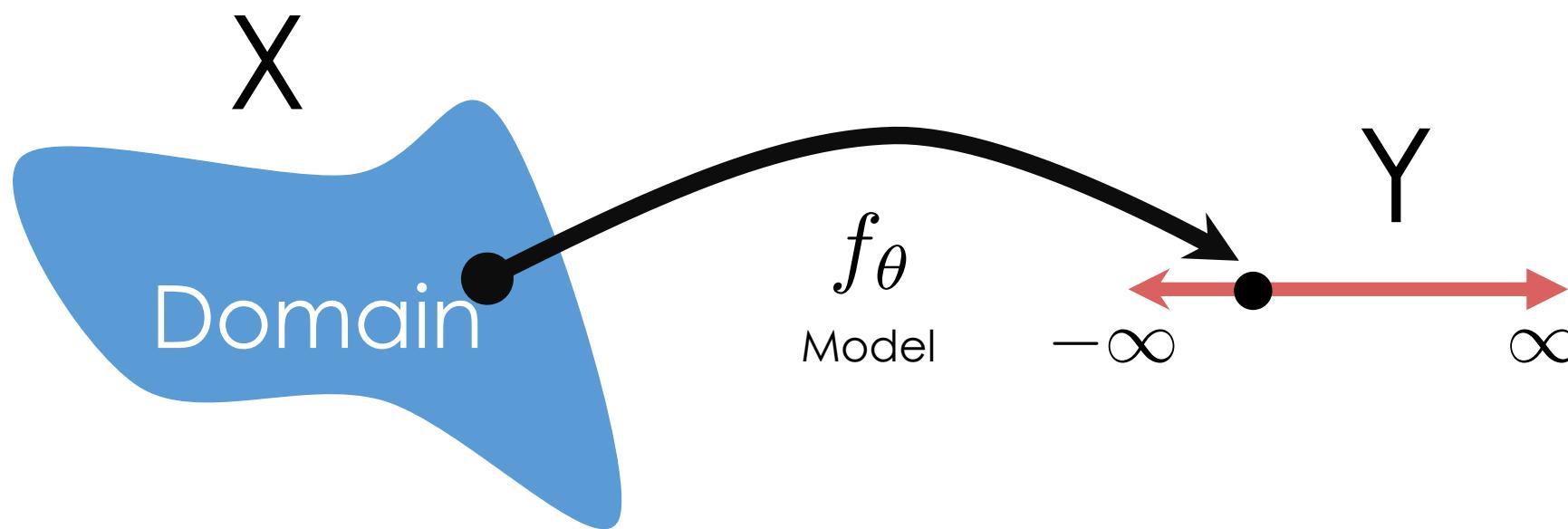
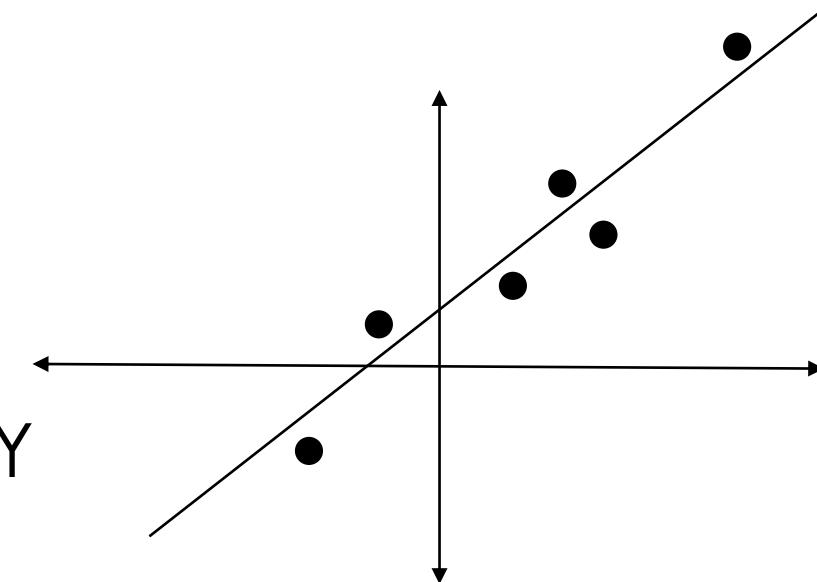


Return to this

Regression and Linear Models

Regression

- Estimating relationship between X and Y
 - Y is a quantitative value
 - We will soon see X can be almost anything ...



Least Squares Linear Regression

One of the most widely used tools in machine learning and data science

Model

$$\hat{y} = f_{\theta}(x) = \sum_{j=1}^d \theta_j \phi_j(x)$$

Linear in the Parameters

Feature Functions

Squared Loss

Loss Minimization

$$\hat{\theta} = \arg \min \frac{1}{n} \sum_{i=1}^n \left(y_i - \sum_{j=1}^d \theta_j \phi_j(x_i) \right)^2$$

We will return to
solving this soon!

Linear Models and Feature Functions

$$\hat{y} = f_{\theta}(x) = \sum_{j=1}^d \theta_j \phi_j(x)$$

Linear in the Parameters

Feature Functions

Designing the feature functions is a big part of machine learning and data science.

Feature Functions

- capture domain knowledge
- substantial contribute to expressivity (and complexity)

Linear Models and Feature Functions

$$\hat{y} = f_{\theta}(x) = \sum_{j=1}^d \theta_j \phi_j(x)$$

Linear in the Parameters

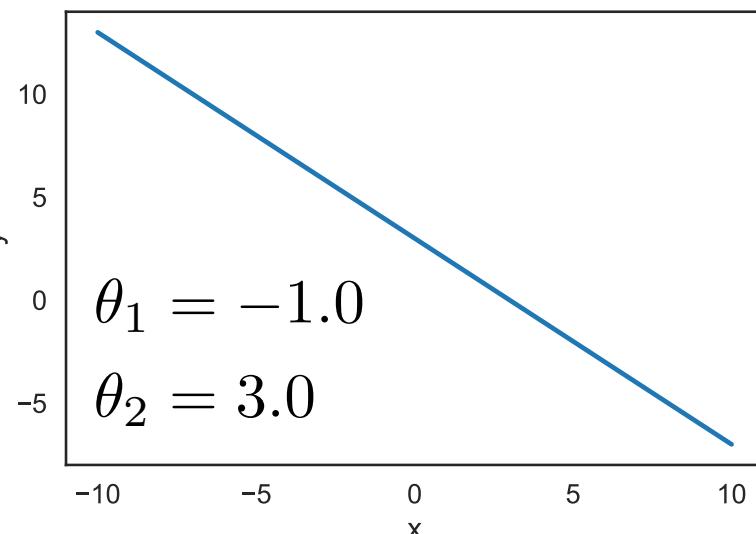
Feature Functions

For Example: Domain: $x \in \mathbb{R}$ Model: $f_{\theta}(x) = \theta_1 x + \theta_2$

Features:

$$\phi_1(x) = x$$

$$\phi_2(x) = 1$$



Adding a “**constant**” feature function $\phi_2(x) = 1$

is a common method to introduce an **offset** (also sometimes called **bias**) term.

Linear Models and Feature Functions

$$\hat{y} = f_{\theta}(x) = \sum_{j=1}^d \theta_j \phi_j(x)$$

Linear in the Parameters

Feature Functions

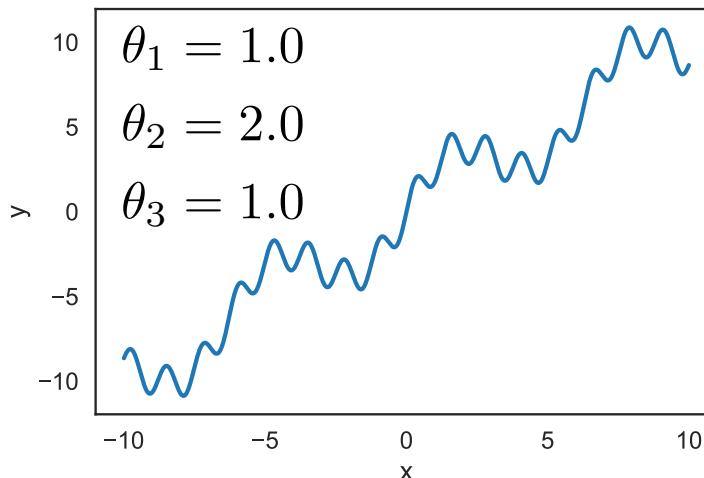
For Example: $x \in \mathbb{R}$ $f_{\theta}(x) = \theta_1 x + \theta_2 \sin(x) + \theta_3 \sin(5x)$

Features:

$$\phi_1(x) = x$$

$$\phi_2(x) = \sin(x)$$

$$\phi_3(x) = \sin(5x)$$



← This is a linear model!

Linear in the parameters

Linear Models and Feature Functions

$$\hat{y} = f_{\theta}(x) = \sum_{j=1}^d \theta_j \phi_j(x)$$

Linear in the Parameters
Feature Functions

For Example: $x \in \mathbb{R}^2$

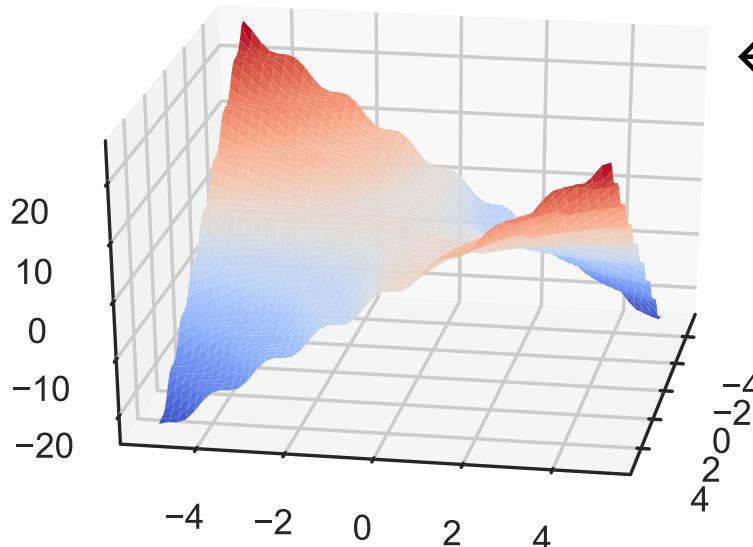
$$f_{\theta}(x) = \theta_1 x_1 x_2 + \theta_2 \cos(x_2 x_1) + \theta_3 \mathbb{I}[x_1 > x_2]$$

Features:

$$\phi_1(x) = x_1 x_2$$

$$\phi_2(x) = \cos(x_2 x_1)$$

$$\phi_3(x) = \mathbb{I}[x_1 > x_2]$$



← This is a linear model!

Linear in the parameters

Linear Models and Feature Functions

$$\hat{y} = f_{\theta}(x) = \sum_{j=1}^d \theta_j \phi_j(x)$$

Linear in the Parameters

Feature Functions

What if x is a record with numbers, text, booleans, etc...

X					Y
uid	age	state	hasBought	review	rating
0	32	NY	True	"Meh."	2.0
42	50	VA	False	Wanted out of box ..."	4.5
57	16	CA	True	It's a tots lit yo ..."	4.1

Answer:

Feature engineering

How do we define ϕ ?

Feature Engineering

Keeping it *Real*

Feature Engineering

- The process of transforming the inputs to a model to improve prediction accuracy.
- A key focus in many applications of data science

- Feature Engineering enables you to:
 - **capture domain knowledge** (e.g., periodicity or relationships between features)
 - **encode non-numeric features** to be used as inputs to models
 - **express non-linear relationships** using linear models

Predict rating from review information

uid	age	state	hasBought	review	rating
0	32	NY	True	"Meh."	2.0
42	50	WA	True	"Worked out of the box ..."	4.5
57	16	CA	NULL	"Hella tots lit yo ..."	4.1

Schema:

```
RatingsData(uid INTEGER, age FLOAT,  
           state STRING, hasBought BOOLEAN,  
           review STRING, rating FLOAT)
```

As a Linear Model?

```
RatingsData(uid INTEGER, age FLOAT,  
state STRING, hasBought BOOLEAN,  
review STRING, rating FLOAT)
```

X =

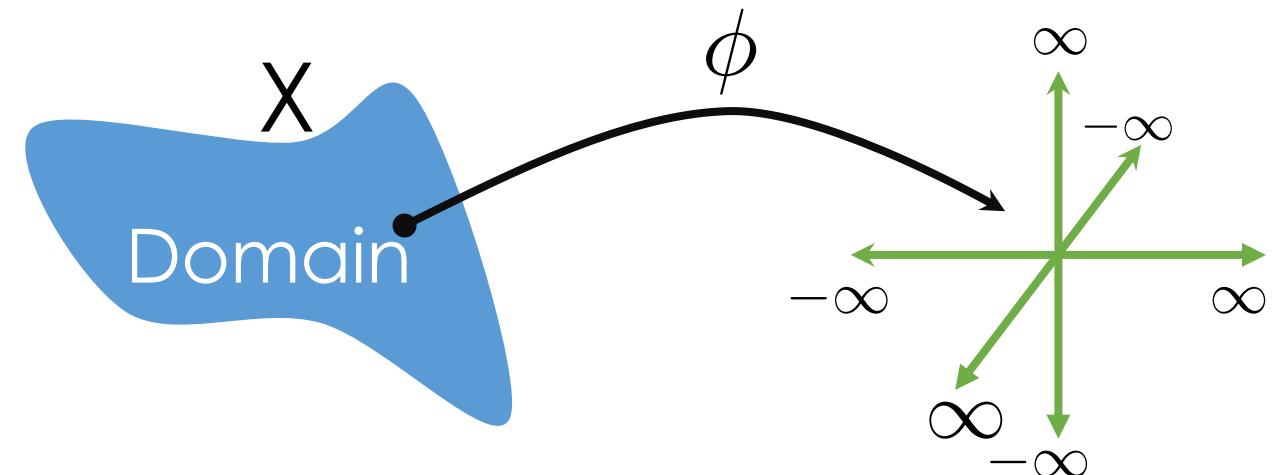
uid	age	state	hasBought	review
0	32	NY	True	"Meh."
42	50	WA	True	"Worked out of the box ..."
57	16	CA	NULL	"Hella tots lit yo ..."

Y =

rating
2.0
4.5
4.1

Can I use X and Y directly
in a linear model

- No! Why?
- Text, Categorical data,
Missing values...



Basic Transformations

- Uninformative features: (e.g., UID)
 - Is this informative (probably not?)
 - **Transformation:** remove uninformative features (why?)
 - Could increase model variance ...
- Quantitative Features (e.g., Age)
 - **Transformation:** May apply non-linear transformations (e.g., log)
 - **Transformation:** Normalize/standardize (more on this later ...)
 - Example: $(x - \text{mean})/\text{stdev}$
- Categorical Features (e.g., State)
 - How do we convert State into meaningful numbers?
 - Alabama = 1 , ..., Utah = 50 ?
 - Implies order/magnitude means something ... we don't want that ...
 - **Transformation:** One-hot-Encode

One Hot Encoding (dummy encoding)

- Transform categorical feature into many binary features:

state	AK	...	CA	...	NY	...	WA	...	WY
NY	0	...	0	...	1	...	0	...	0
WA	0	...	0	...	0	...	1	...	0
CA	0	...	1	...	0	...	0	...	0

$$\phi_1(x) = \mathbb{I}[x \text{ is 'AK'}]$$

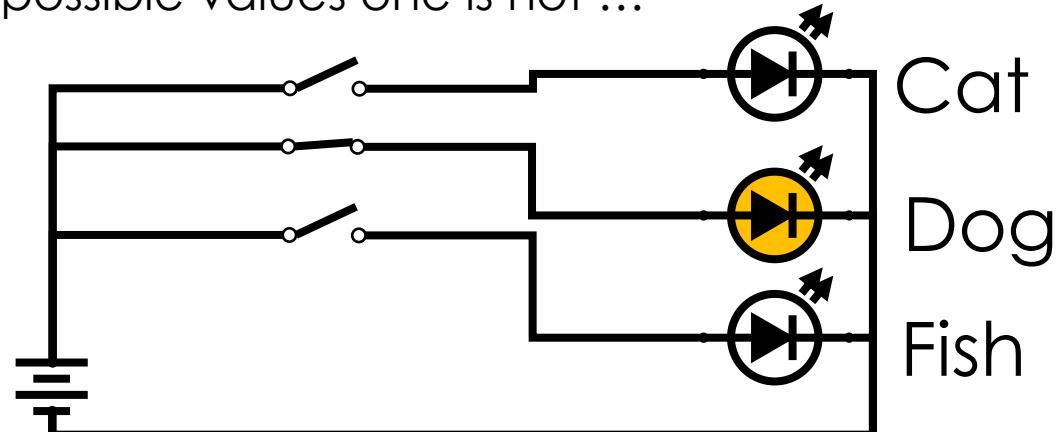
Corresponding
feature
functions

$$\phi_2(x) = \mathbb{I}[x \text{ is 'AL'}]$$

...

$$\phi_{50}(x) = \mathbb{I}[x \text{ is 'WY'}]$$

Origin of the term: multiple “wires” for possible values one is hot ...



See notebook
for example
code.

Encoding Missing Values

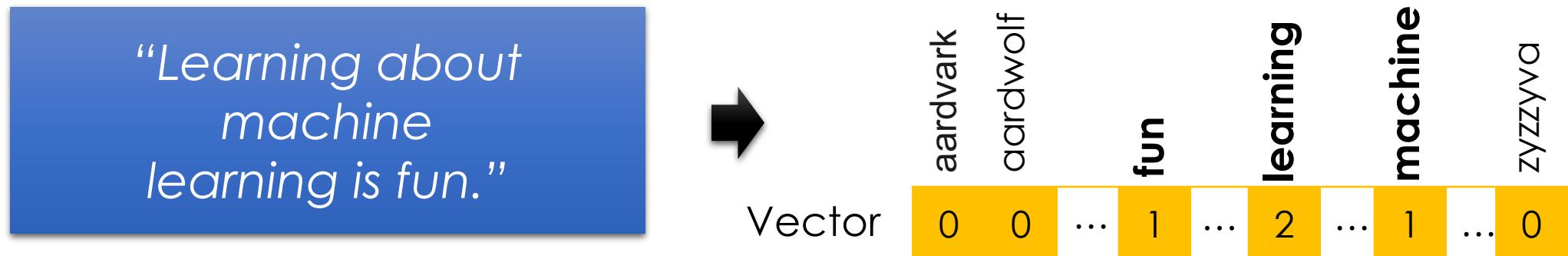
- Missing values in **Quantitative Data**
 - Try to impute (estimate) missing values... (tricky)
 - Substitute the sample mean
 - Try more sophisticated algorithms to predict the missing value ...
 - Add a binary field called “missing_col_name”. (why?)
 - Sometimes missing data is signal!
- Missing values in **Categorical Data**
 - Add an addition category called “missing_col_name”
 - Some Boolean values can be converted into
 - True => +1, False => -1, Missing => 0

Encoding categorical data

- Categorical Data → One-hot encoding:

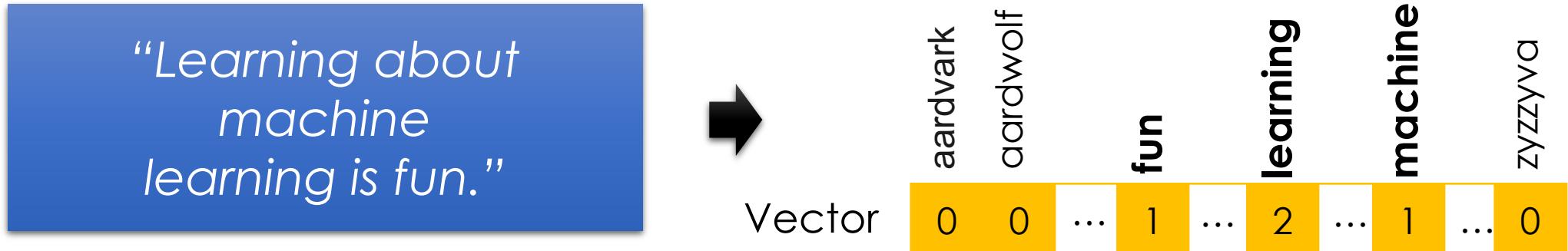
state	AL	...	CA	...	NY	...	WA	...	WY
NY	0	...	0	...	1	...	0	...	0
WA	0	...	0	...	0	...	1	...	0
CA	0	...	1	...	0	...	0	...	0

- Text Data
 - Bag-of-words & N-gram models



Bag-of-words Encoding

- Generalization of one-hot-encoding for a string of text:



- Encode text as a long vector of word counts (Issues?)
 - Long = millions of columns → typically high dimensional and very sparse
 - Word order information is lost... (is this an issue?)
 - New unseen words at prediction (test) time → drop them ...
- A **bag** is another term for a **multiset**: an unordered collection which may contain multiple instances of each element.
- **Stop words**: words that do not contain significant information
 - Examples: the, in, at, or, on, a, an, and ...
 - Typically removed

I made this art piece
in graduate school

Do you see the stop word?

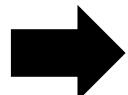
There used to be a
dustbin and broom
... but the janitors
got confused ...



N-Gram Encoding

- Sometimes word order matters:

*The book was not well
written but I did enjoy it.*



*The book was well written
but I did not enjoy it.*

- How do we capture word order in a “vector” model?
 - N-Gram: “Bag-of- sequences-of-words”

Removed
stop words

The book was well written but I did not enjoy it.

book well

well written

written not

not enjoy

2-Gram Encoding

aardvark airlines

apple pen

book well

not enjoy

well written

written not

zyzzyva sf

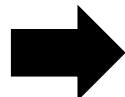
Vector



N-Gram Encoding

- Sometimes word order matters:

*The book was not well
written but I did enjoy it.*



*The book was well written
but I did not enjoy it.*

- How do we capture word order in a “vector” model?
 - N-Gram: “Bag-of- sequences-of-words”
- Issues:
 - Can be very sparse (many combinations occur only once)
 - Many combinations will only occur at prediction time → drop ..
 - Often use hashing approximation:
 - Increment counter at **hash(“not enjoy”)** collisions are okay

Feature Transformations to Capture Domain Knowledge

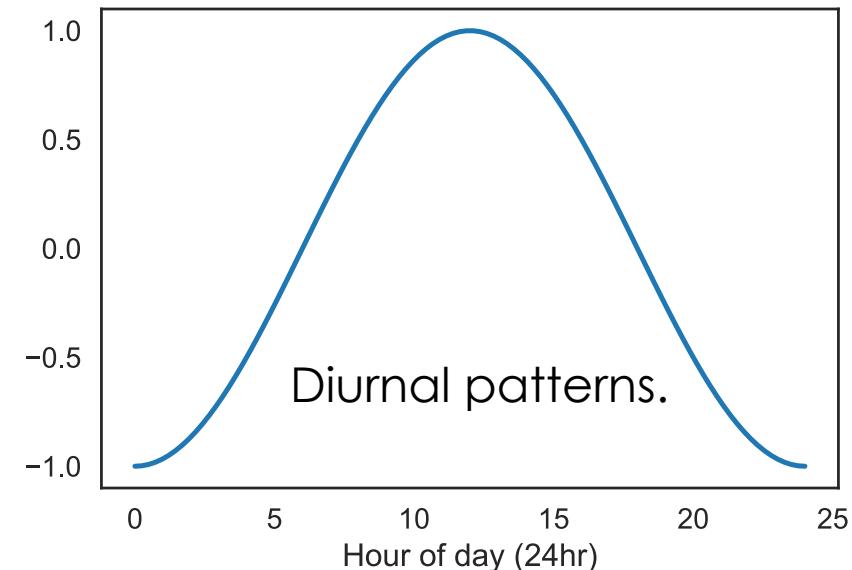
- Feature functions capture domain knowledge by introducing **additional information** from other sources **and/or combining features**

Could do a database lookup

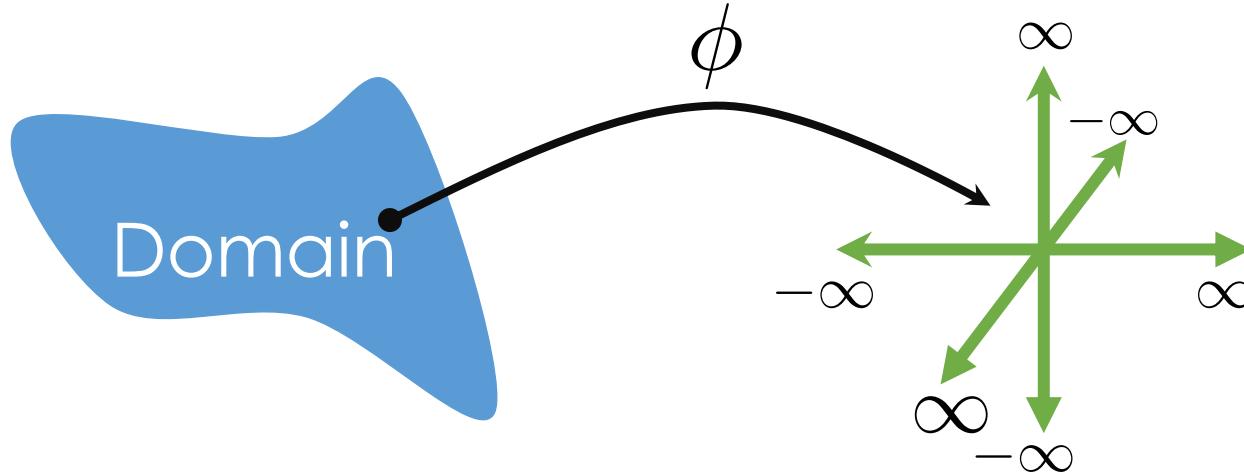
$$\phi_i(x) = \text{isWinter}(x_{\text{date}}, x_{\text{location}})$$

- Encoding non-linear patterns

$$\phi_i(x) = \cos\left(\frac{x_{\text{hour}}}{12}\pi + \pi\right)$$



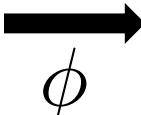
The Feature Matrix Φ



X DataFrame

uid	age	state	hasBought	review
0	32	NY	True	"Meh."
42	50	WA	True	"Worked out of the box ..."
57	16	CA	NULL	"Hella tots lit..."

$$\Phi \in \mathbb{R}^{n \times d}$$



AK	...	NY	...	WY	age	hasBought	hasBought missing
0	...	1	...	0	32	1	0
0	...	0	...	0	50	1	0
0	...	0	...	0	16	0	1

Entirely **Quantitative** Values

The Feature Matrix Φ

AK	...	NY	...	WY	age	hasBought	hasBought missing
0	...	1	...	0	32	1	0
0	...	0	...	0	50	1	0
0	...	0	...	0	16	0	1

Entirely **Quantitative** Values

$$\Phi \in \mathbb{R}^{n \times d} = \phi(X) = \begin{matrix} n \\ \left[\begin{array}{c} \phi(x^{(1)}) \\ \phi(x^{(2)}) \\ \dots \\ \phi(x^{(n)}) \end{array} \right] \\ d \end{matrix}$$

DataFrame

Rows of the Φ matrix correspond to records.

Columns of the Φ matrix correspond to features.

Making Predictions

$$\Phi \in \mathbb{R}^{n \times d} = \phi(X) = \begin{array}{c} \text{DataFrame} \\ \left[\begin{array}{c} \phi(x^{(1)}) \\ \phi(x^{(2)}) \\ \dots \\ \phi(x^{(n)}) \end{array} \right] \end{array}$$

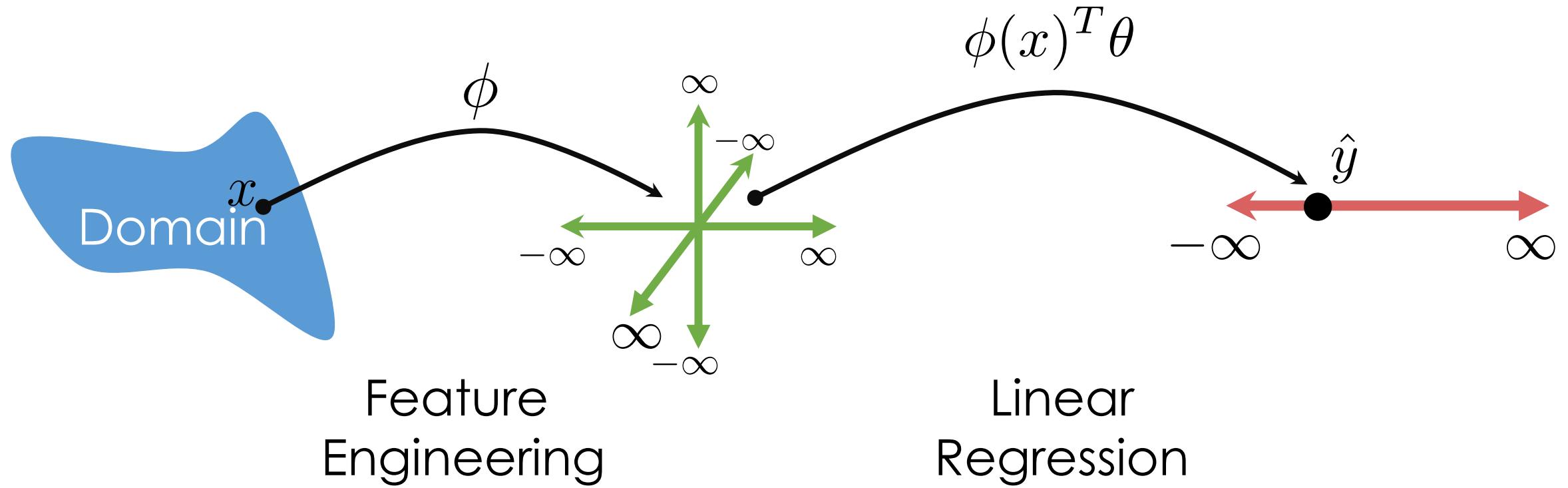
Rows of the Φ matrix correspond to records.

Columns of the Φ matrix correspond to features.

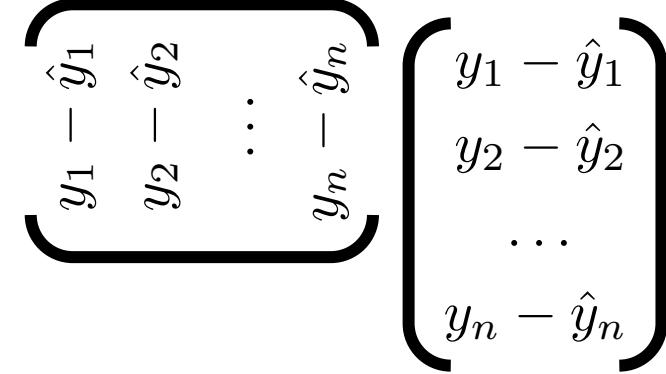
Prediction

$$\hat{Y} = f_{\hat{\theta}}(X) = \Phi \hat{\theta} = \left[\begin{array}{c} \phi(x^{(1)}) \\ \phi(x^{(2)}) \\ \dots \\ \phi(x^{(n)}) \end{array} \right] \left| \begin{array}{c} \hat{\theta} \end{array} \right| = \left[\begin{array}{c} \hat{y}^{(1)} \\ \hat{y}^{(2)} \\ \dots \\ \hat{y}^{(n)} \end{array} \right]$$

Summary of Notation



Optimizing the Loss (Bonus Material)

$$\begin{aligned} L(\theta) &= \frac{1}{n} \sum_{i=1}^n \left(y_i - \sum_{j=1}^d \theta_j \phi_j(x_i) \right)^2 = (Y - \hat{Y})^T (Y - \hat{Y}) \\ &= \frac{1}{n} (Y - \Phi\theta)^T (Y - \Phi\theta) \\ &= \frac{1}{n} (Y^T Y - 2Y^T \Phi\theta + \theta^T \Phi^T \Phi\theta) \end{aligned}$$


Taking the Gradient of the loss

Optimizing the Loss (Bonus Material)

Deriving the Normal Equation

$$L(\theta) = \frac{1}{n} (Y^T Y - 2Y^T \Phi\theta + \theta^T \Phi^T \Phi\theta)$$

Taking the Gradient of the loss

$$\nabla_{\theta} L(\theta) = -\frac{2}{n} \Phi^T Y + \frac{2}{n} \Phi^T \Phi\theta$$

Rule 1

Rule 2

Setting the gradient equal to 0 and solving for θ :

$$0 = -\frac{2}{n} \Phi^T Y + \frac{2}{n} \Phi^T \Phi\theta \longrightarrow$$

$$\hat{\theta} = (\Phi^T \Phi)^{-1} \Phi^T Y$$

“Normal Equation”

Useful Matrix Derivative Rules:

$$(1) \nabla_{\theta} (A\theta) = A^T$$

$$(2) \nabla_{\theta} (\theta^T A\theta) = A\theta + A^T\theta$$

The Normal Equation $\hat{\theta} = (\Phi^T \Phi)^{-1} \Phi^T Y$

$$\hat{\theta} = \left(\begin{array}{c|c} \Phi^T & \Phi \end{array} \right)^{-1} \left(\begin{array}{c|c|c} d & \Phi^T & Y \\ \hline n & 1 \end{array} \right)$$

Note: For inverse to exist Φ needs to be full column rank.

→ cannot have co-linear features

This can be addressed by adding regularization ...

**In practice we will use regression software
(e.g., scikit-learn) to estimate θ**

Geometric Derivation (Bonus Material)

- Examine the column spaces:

Columns space of Φ

$$\Phi = \begin{bmatrix} & & \\ \vdots & \vdots & \vdots \\ \Phi^{(1)}, \Phi^{(2)}, \dots, \Phi^{(d)} \\ \vdots & \vdots & \vdots \end{bmatrix} \in \mathbb{R}^{n \times d}$$
$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \in \mathbb{R}^n$$

The diagram illustrates the geometric interpretation of the linear model $Y = \Phi \beta$. The matrix Φ is shown as a vertical stack of d columns, each labeled $\Phi^{(1)}, \Phi^{(2)}, \dots, \Phi^{(d)}$. The height of this stack is indicated by a red bracket labeled n at the top and d at the bottom. To the right, the vector Y is shown as a vertical stack of n entries, labeled y_1, y_2, \dots, y_n , with a red bracket at the bottom labeled 1 .

- Linear model $\rightarrow Y$ is a linear combination of columns Φ

Columns space of Φ

$$\Phi = \begin{bmatrix} | & | & | \\ \Phi^{(1)}, \Phi^{(2)}, \dots, \Phi^{(d)} \\ | & | & | \end{bmatrix} \in \mathbb{R}^{n \times d}$$

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \in \mathbb{R}^n$$

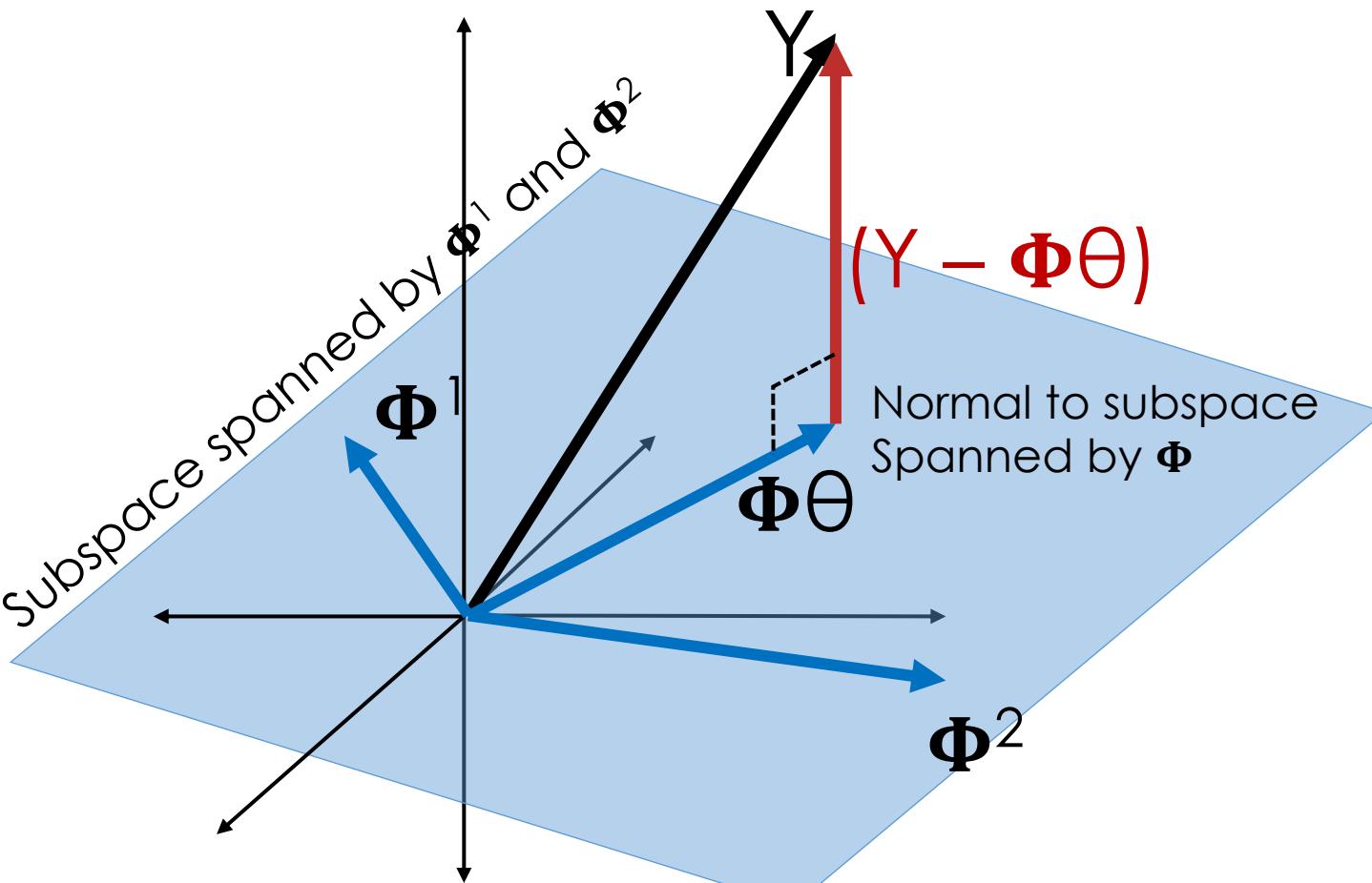
n d 1

➤ Linear model $\rightarrow Y$ is a linear combination of columns Φ

$$Y \approx \hat{Y} = \Phi \hat{\theta} \quad \Rightarrow \quad \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{bmatrix} \approx \begin{bmatrix} | & | & | \\ \Phi^{(1)}, \Phi^{(2)}, \dots, \Phi^{(d)} \\ | & | & | \end{bmatrix} \quad \mid \quad \hat{\theta}$$

$$Y \approx \hat{Y} = \Phi \hat{\theta} \rightarrow \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{bmatrix} \approx \begin{bmatrix} | & | & | \\ \Phi^{(1)}, \Phi^{(2)}, \dots, \Phi^{(d)} \\ | & | & | \end{bmatrix} \hat{\theta}$$

➤ \hat{Y} is in the subspace spanned by the columns of Φ



Definition of orthogonal

$$0 = \Phi^T (Y - \Phi\theta)$$

$$0 = \Phi^T Y - \Phi^T \Phi \theta$$

$$\Phi^T \Phi \theta = \Phi^T Y \quad \text{Normal Equations}$$

$$\theta = (\Phi^T \Phi)^{-1} \Phi^T Y$$

"Normal Equation"

Least Squares Regression in Practice

- Use optimized software packages
 - Address numerical issues with matrix inversion
- Incorporate some form of regularization
 - Address issues of collinearity
 - Produce more robust models
- We will be using scikit-learn:
 - http://scikit-learn.org/stable/modules/linear_model.html
 - See Homework 6 for details!

Daily Quiz

<http://bit.ly/ds100-sp18-lin>

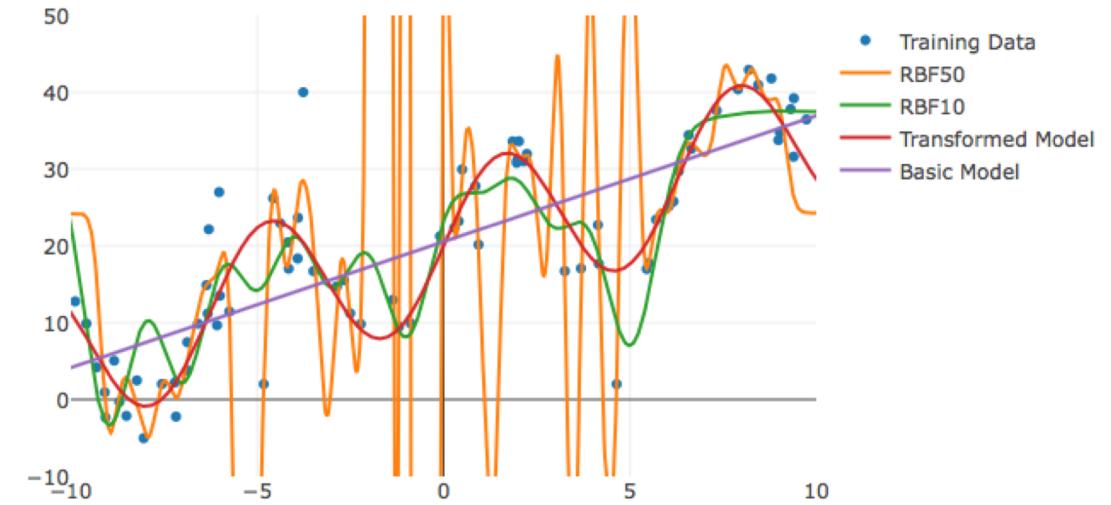
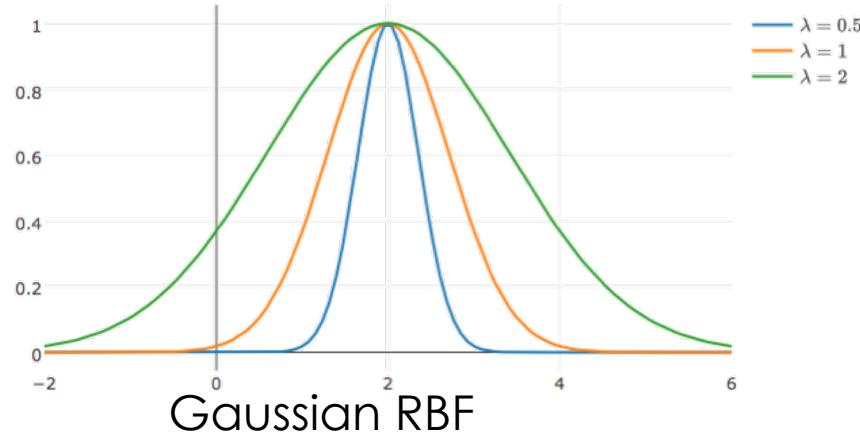
- Which of the following are **not** linear models:

- (A) $f_{\theta}(x) = \theta_1 x^2 + \frac{\theta_2}{\exp(x)}$
- (B) $f_{\theta}(x) = \theta_1 + \theta_2 \sin(\theta_3 x)$
- (C) $f_{\theta}(x) = \theta_1^2 + \theta_2(\theta_3 - x)$
- (D) $f_{\theta}(x) = \theta_1 \exp\left(\frac{(x - \mu)^2}{\sigma^2}\right) + \theta_2$

Notebook Demo

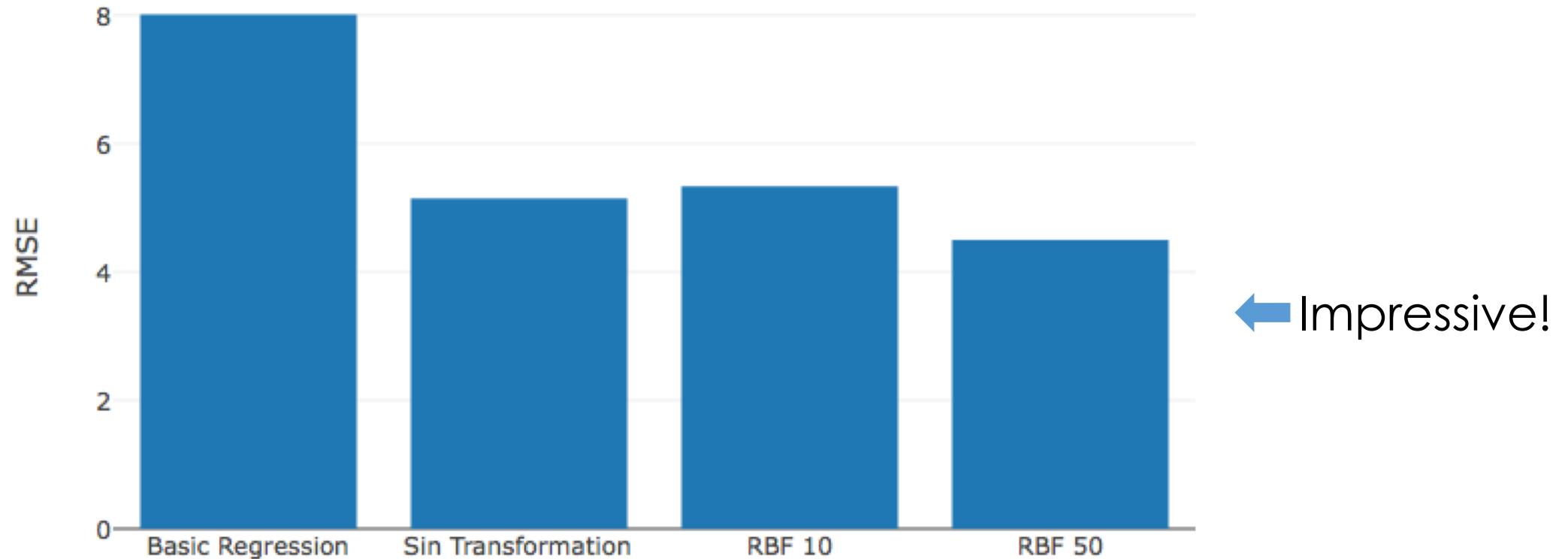
- **Generic Features:** increase model expressivity
- **Gaussian Radial Basis Functions:**

$$\phi_{\lambda_i, \mu_i}(x) = \exp\left(-\frac{\|x - \mu_i\|_2^2}{\lambda_i}\right)$$

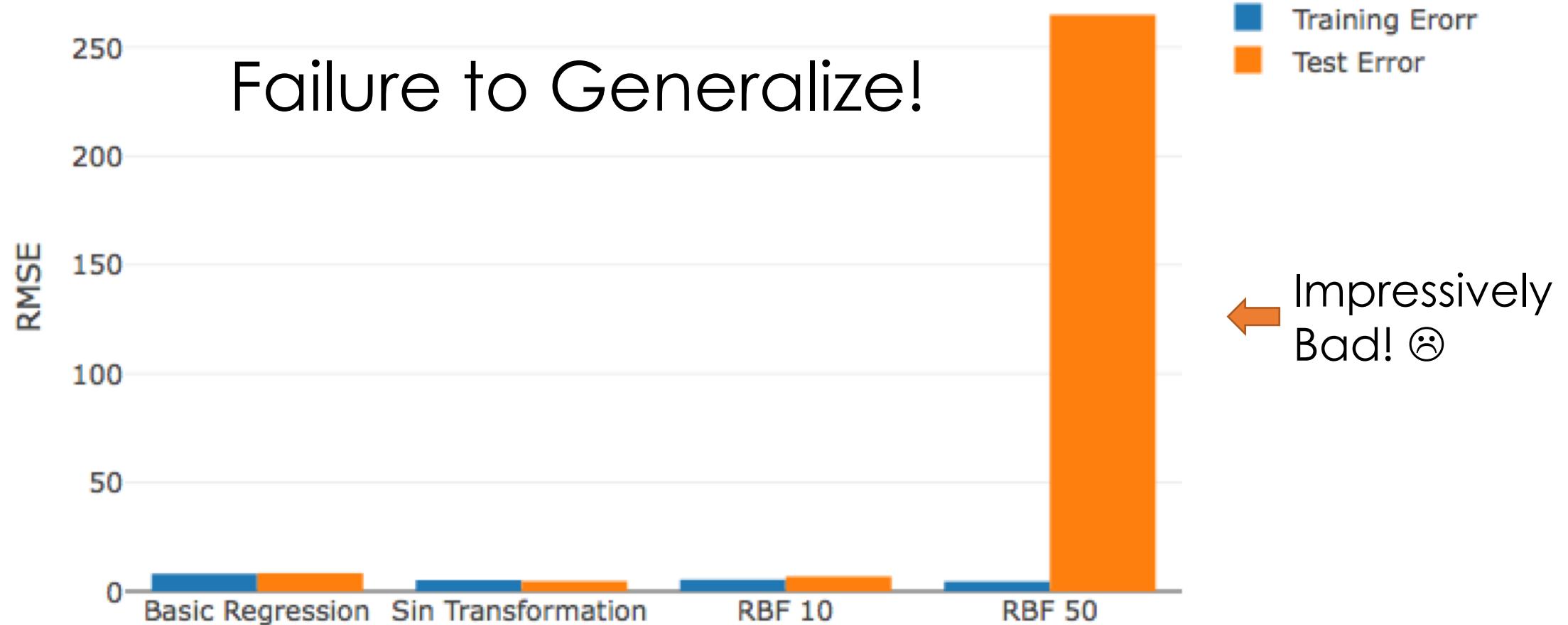


Training Error

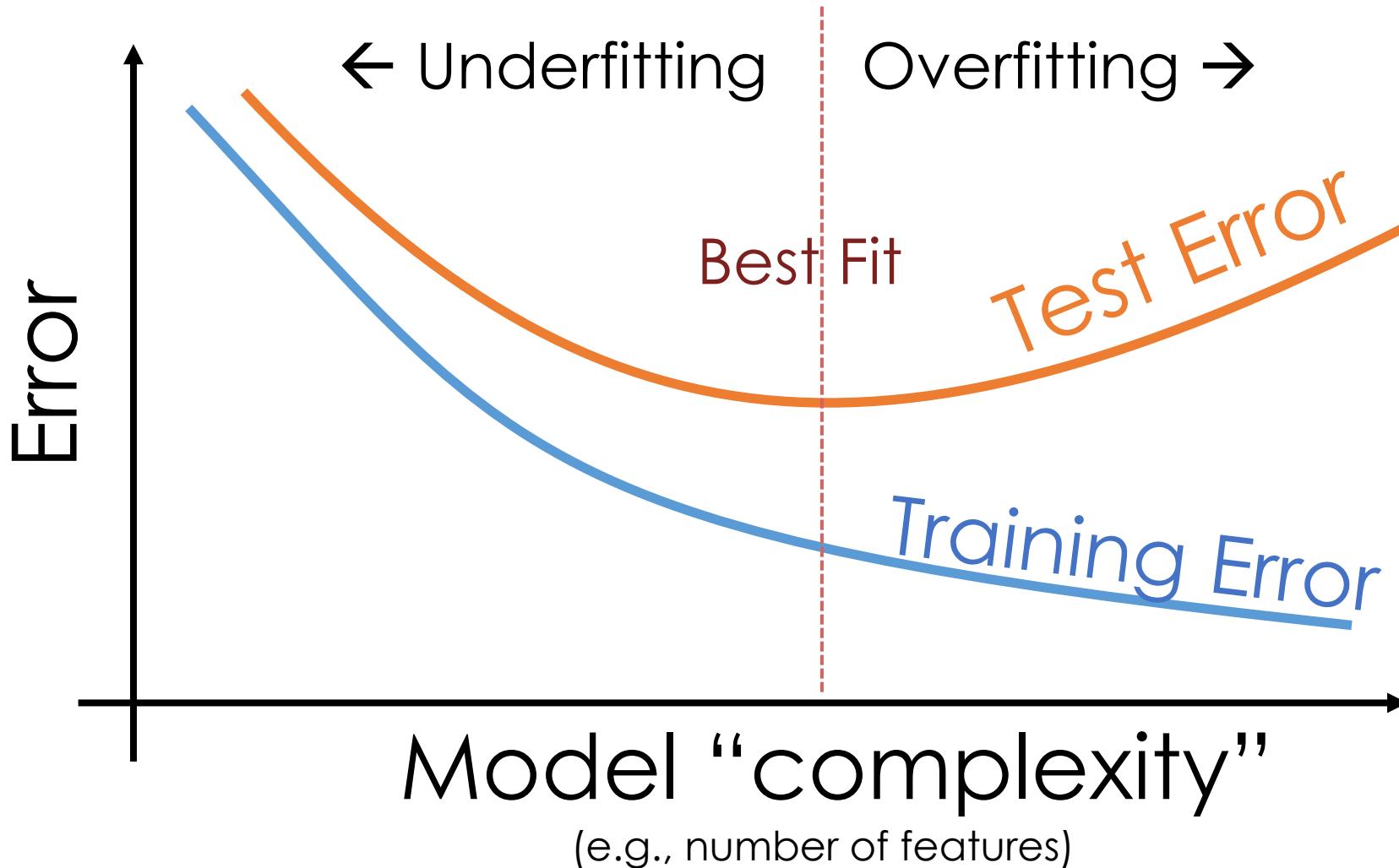
Loss Comparison



Training vs Test Error

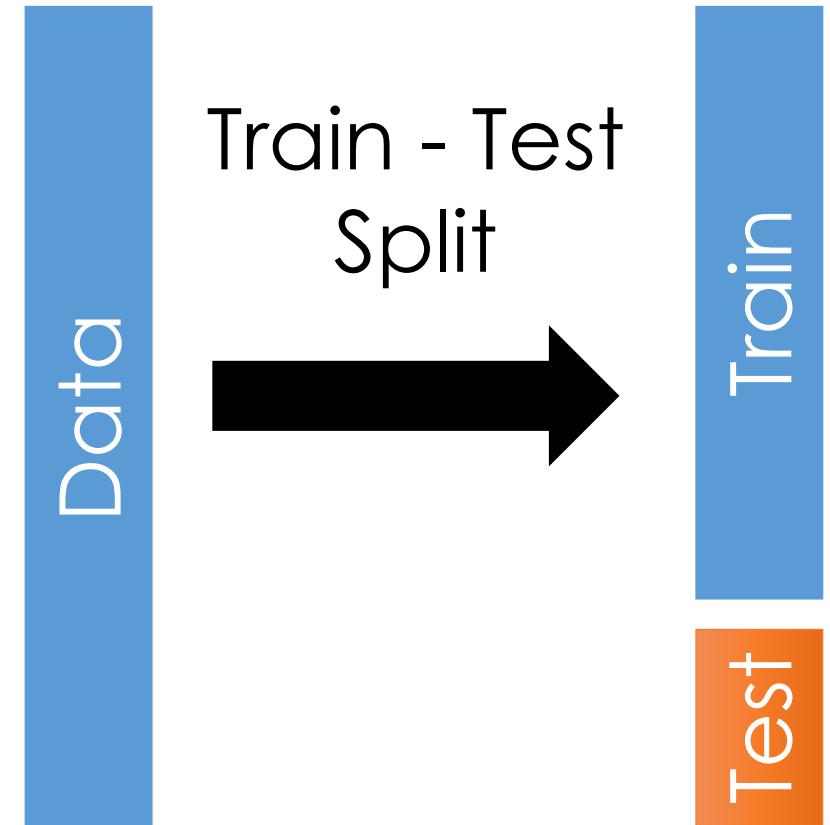


Training vs Test Error



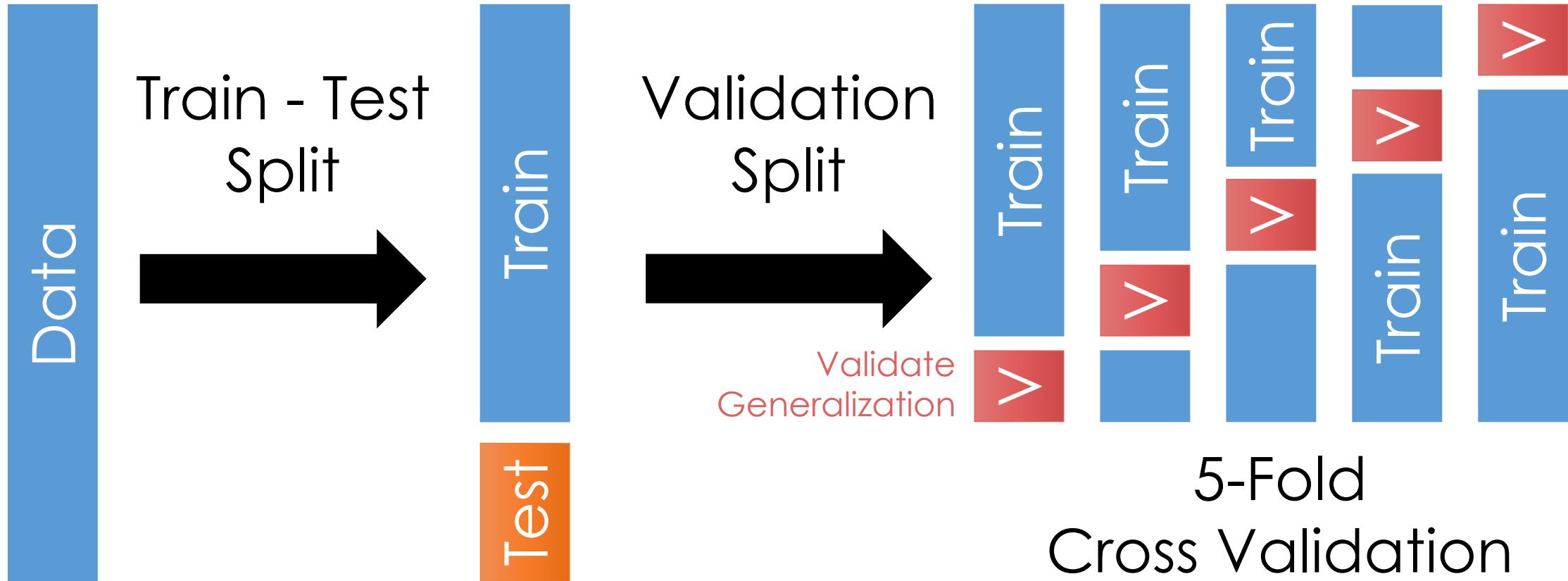
Generalization: The Train-Test Split

- **Training Data:** used to fit model
- **Test Data:** check generalization error
- How to split?
 - Randomly, Temporally, Geo...
 - Depends on application (usually randomly)
- What size? (90%-10%)
 - Larger training set → more complex models
 - Larger test set → better estimate of generalization error
 - Typically between 75%-25% and 90%-10%



You can only use the test dataset once after deciding on the model.

Generalization: Validation Split



Cross validation simulates multiple train test-splits on the training data.

Returning to Regularization

Regularization

Parametrically Controlling the Model Complexity

- Tradeoff:
 - Increase bias
 - Decrease variance



Basic Idea of Regularization

Fit the Data

Penalize
Complex Models

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n \text{Loss}(y_i, f_{\theta}(x_i)) + \lambda \mathbf{R}(\theta)$$

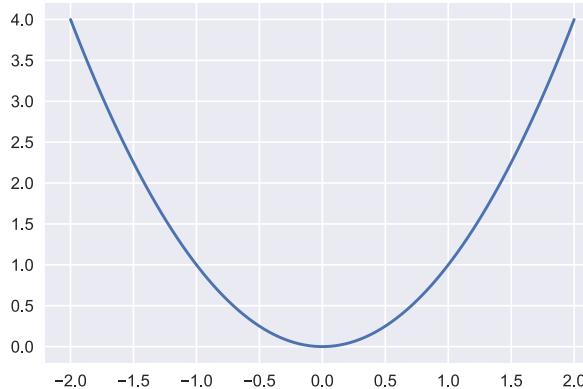
Regularization
Parameter

- How should we define $\mathbf{R}(\theta)$?
- How do we determine λ ?

Common Regularization Functions

Ridge Regression
(L2-Reg)

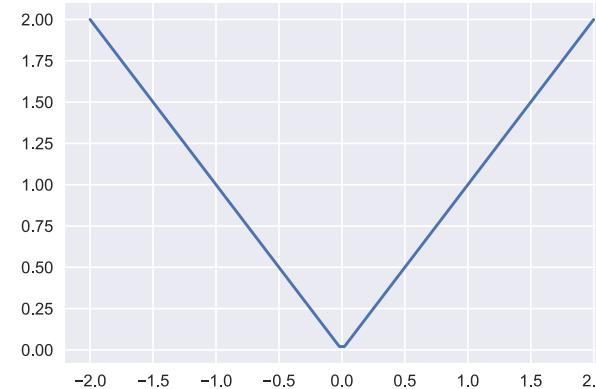
$$R_{\text{Ridge}}(\theta) = \sum_{i=1}^d \theta_i^2$$



- Distributes weight across related features (robust)
- Analytic solution (easy to compute)
- Does not encourage sparsity → small but non-zero weights.

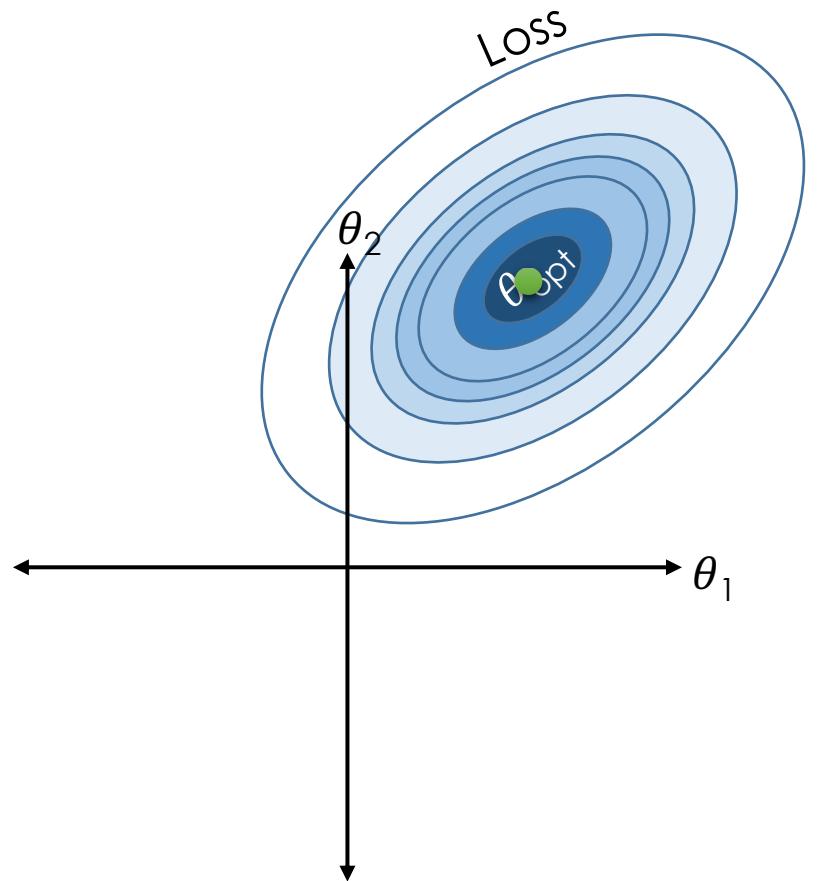
LASSO
(L1-Reg)

$$R_{\text{Lasso}}(\theta) = \sum_{i=1}^d |\theta_i|$$

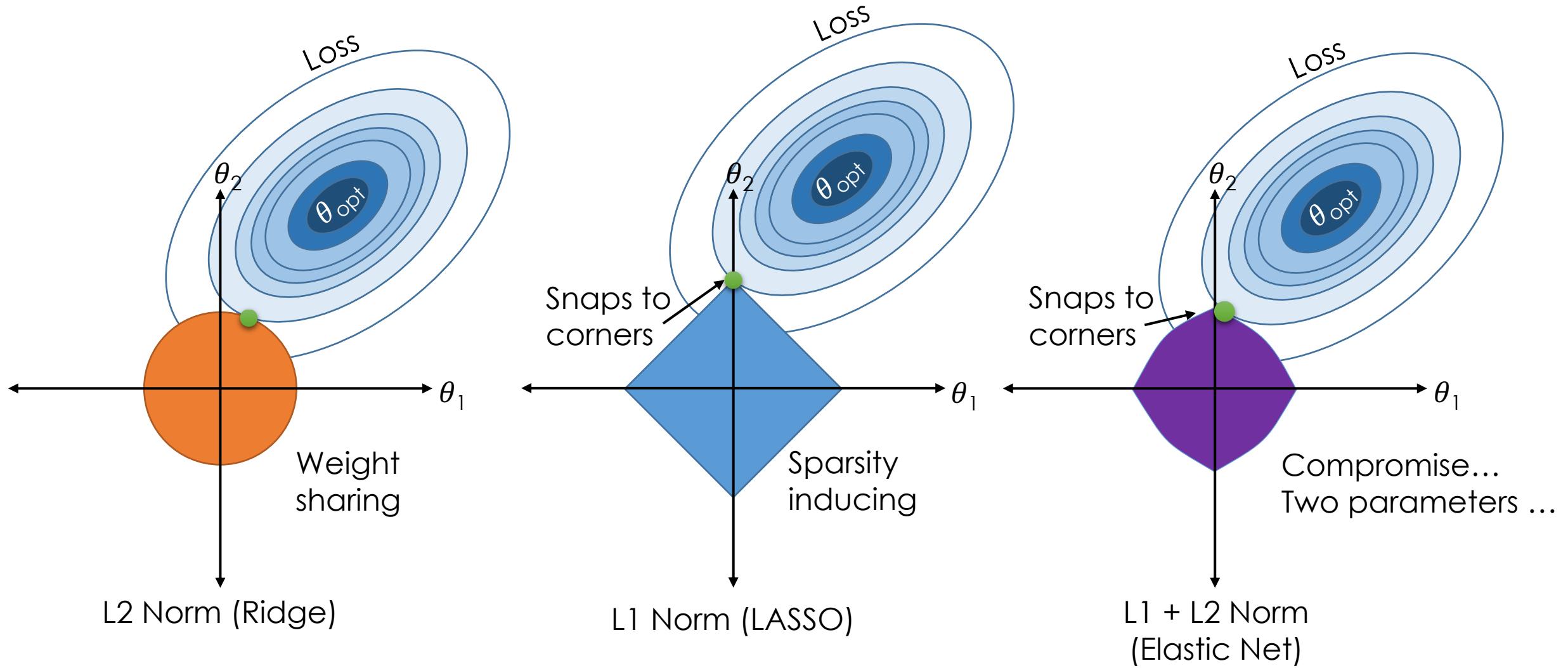


- **Encourages sparsity** by setting weights = 0
 - Used to select informative features
- Does not have an analytic solution → numerical methods

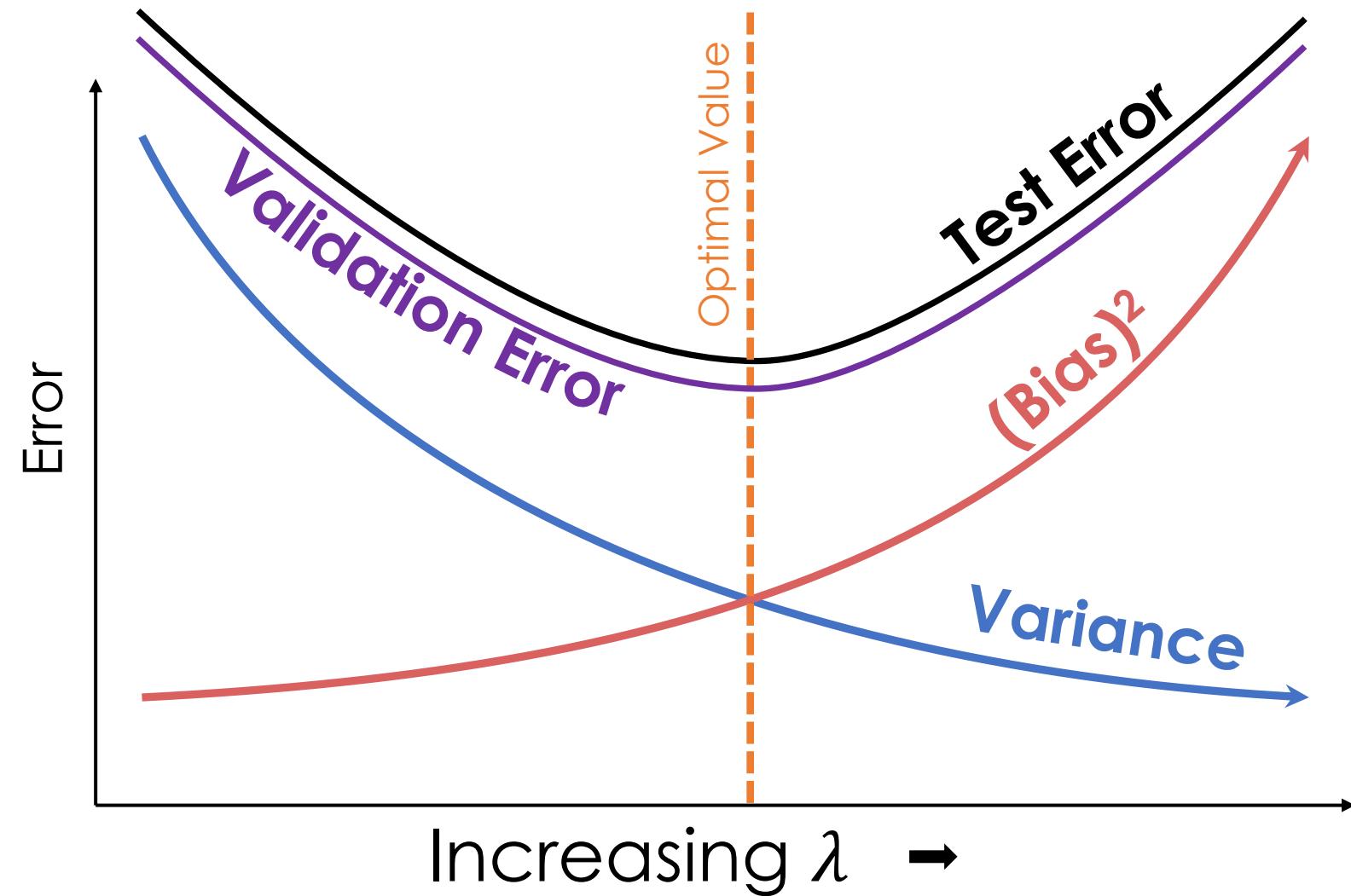
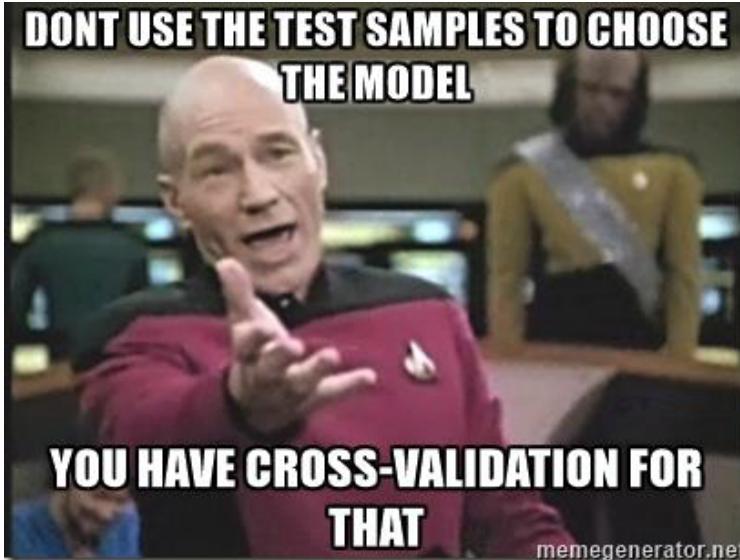
Regularization and Norm Balls



Regularization and Norm Balls



Determining the Optimal λ



- Value of λ determines bias-variance tradeoff
 - Larger values \rightarrow more regularization \rightarrow more bias \rightarrow less variance
- Determined through cross validation

Standardization and the Intercept Term

Height = θ_1 age_in_seconds + θ_2 weight_in_tons

Small

Large

➤ Regularization penalized dimensions equally

➤ **Standardization**

- Ensure that each dimensions has the same scale
- centered around zero

Standardization

For each dimension k :

$$z_k = \frac{x_k - \mu_k}{\sigma_k}$$

➤ **Intercept Terms**

- Typically don't regularize intercept term
- Center y values (e.g., subtract mean)

Regularization and High-Dimensional Data

Regularization is often used with high-dimensional data

