

# CS150A Database

Lu Sun

School of Information Science and Technology  
ShanghaiTech University

Dec. 20, 2022

Today:

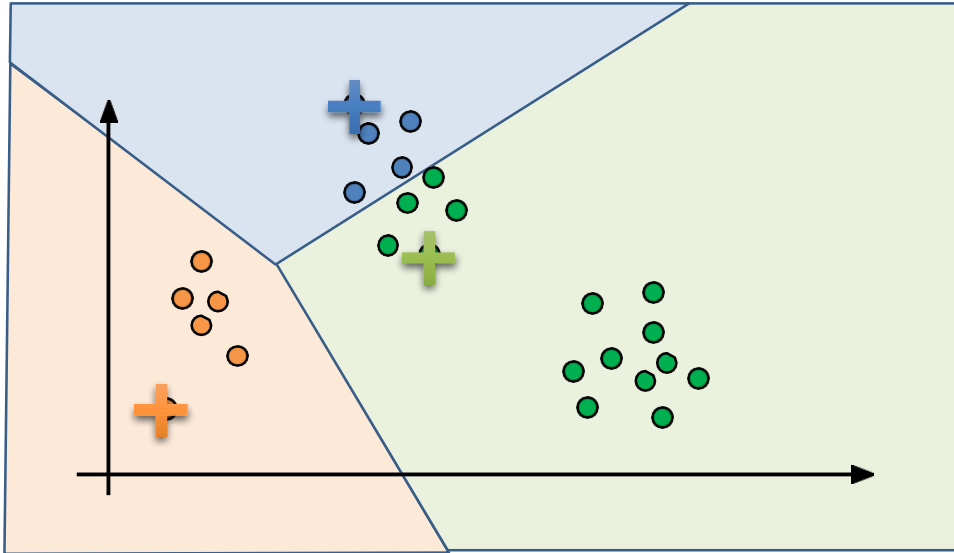
- Analytics and ML in Data Systems:
  - Part 3
  - Linear Regression

Readings:

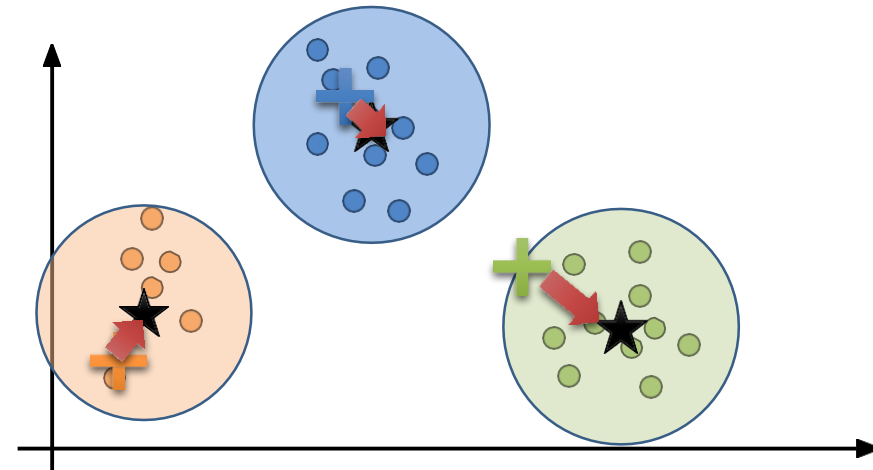
- Database Management Systems (DBMS), Chapters 23&24

# K-Means Clustering

Compute Assignments



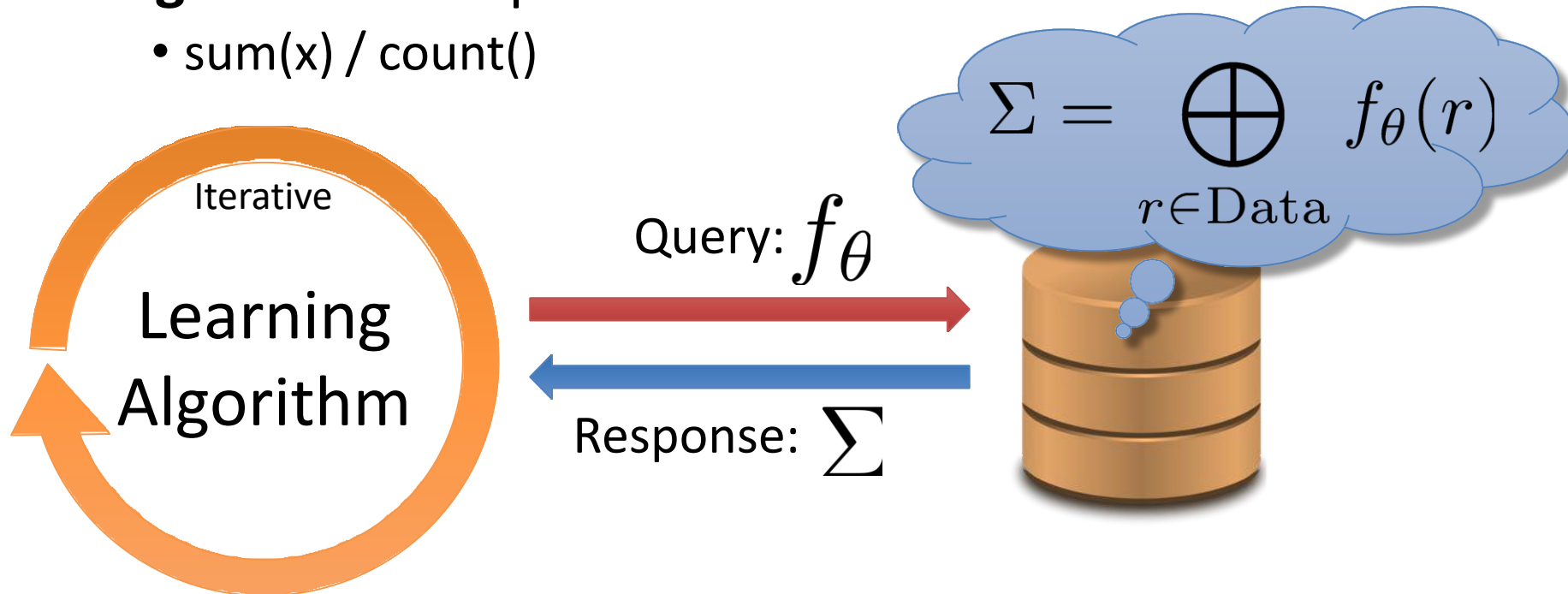
Update Centers



# Statistical Query Pattern

## Common Machine Learning Pattern

- **K-Means Query:** for each old centers compute the sum of the nearest points
- **Response:** sums and counts of points
- **Algorithm:** compute new centers
  - $\text{sum}(x) / \text{count}()$



# Res-A: weighted reservoir sampling

□ **Goal:** Sample  $k$  records from a stream where record  $i$  is included in the sample with probability proportional to  $w_i$

□ **Algorithm:**

- For each record  $i$  draw a uniform random number:

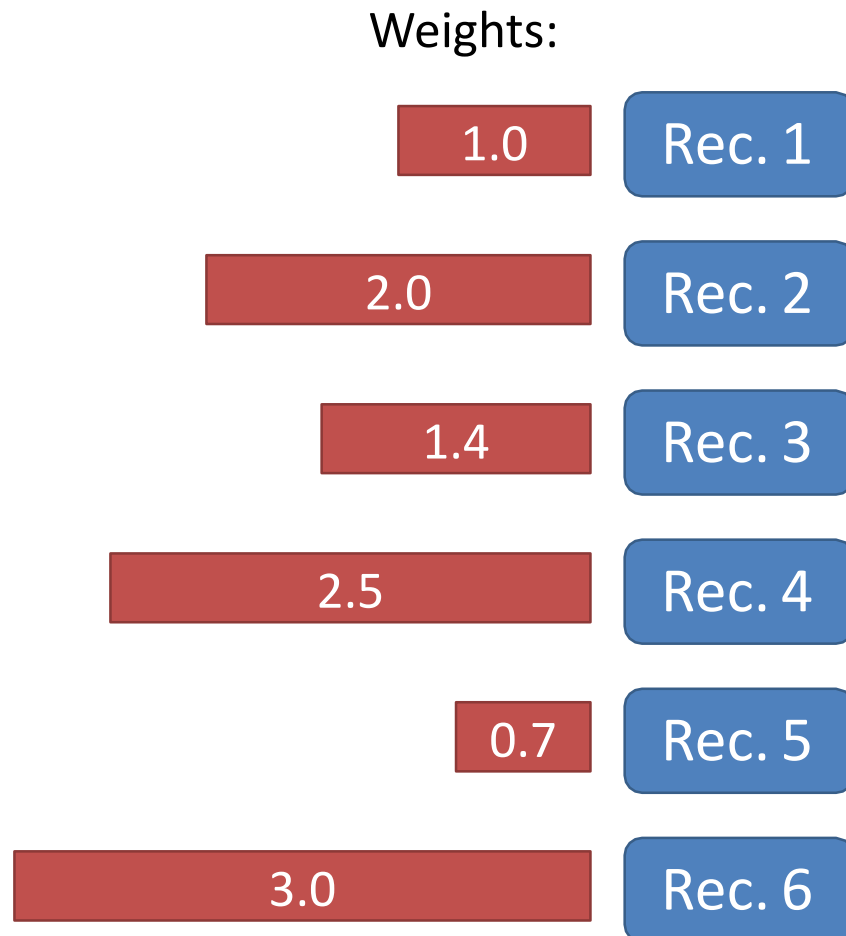
$$u_i \sim \text{Unif}(0, 1)$$

- Select the top- $k$  records ordered by:  $u_i^{1/w_i}$

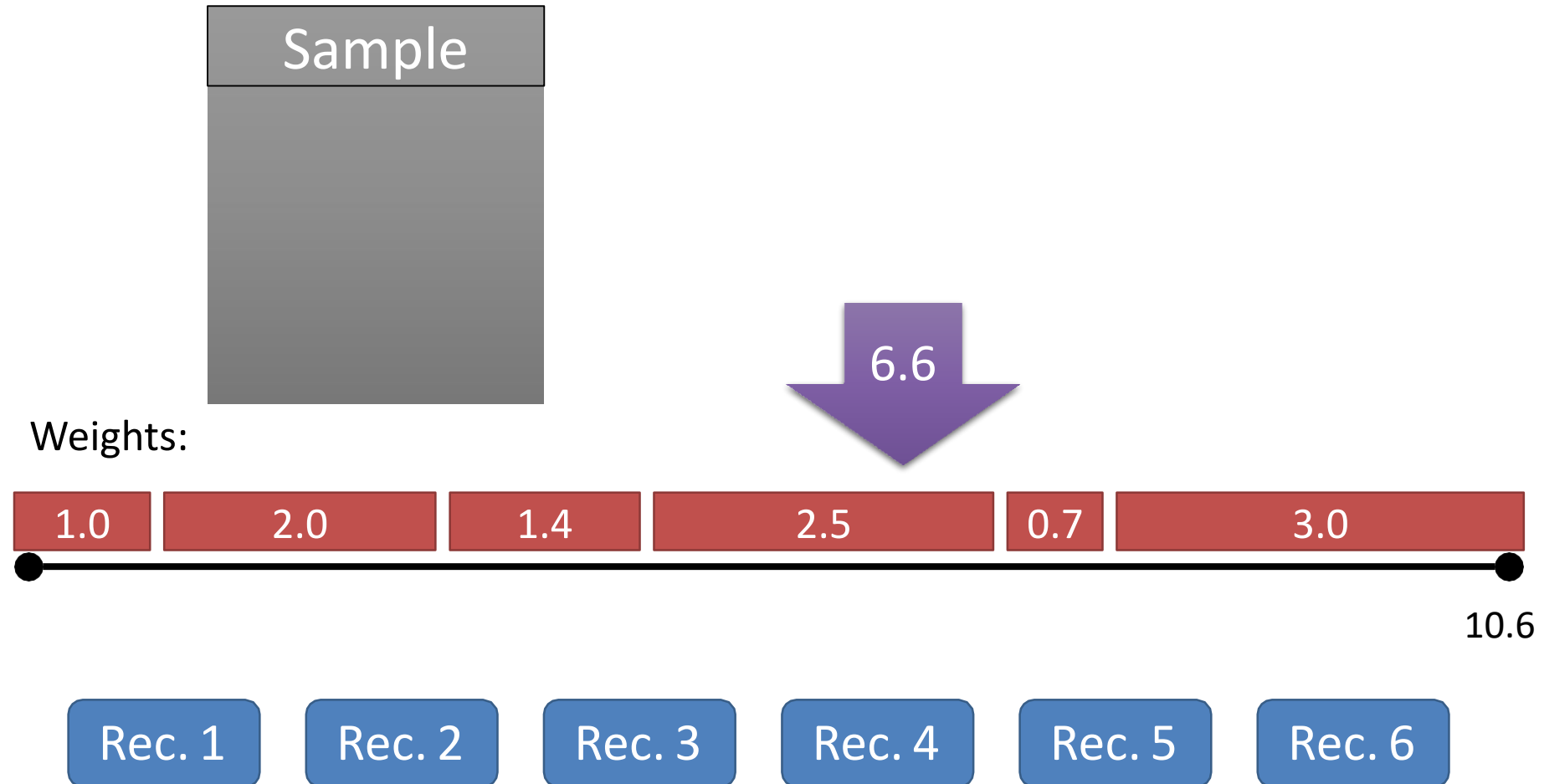
□ **Common ML Pattern?**

- **Query Function:**  $[pow(rand(), 1 / record.w), record]$
- **Agg. Function:** *top-k heap*

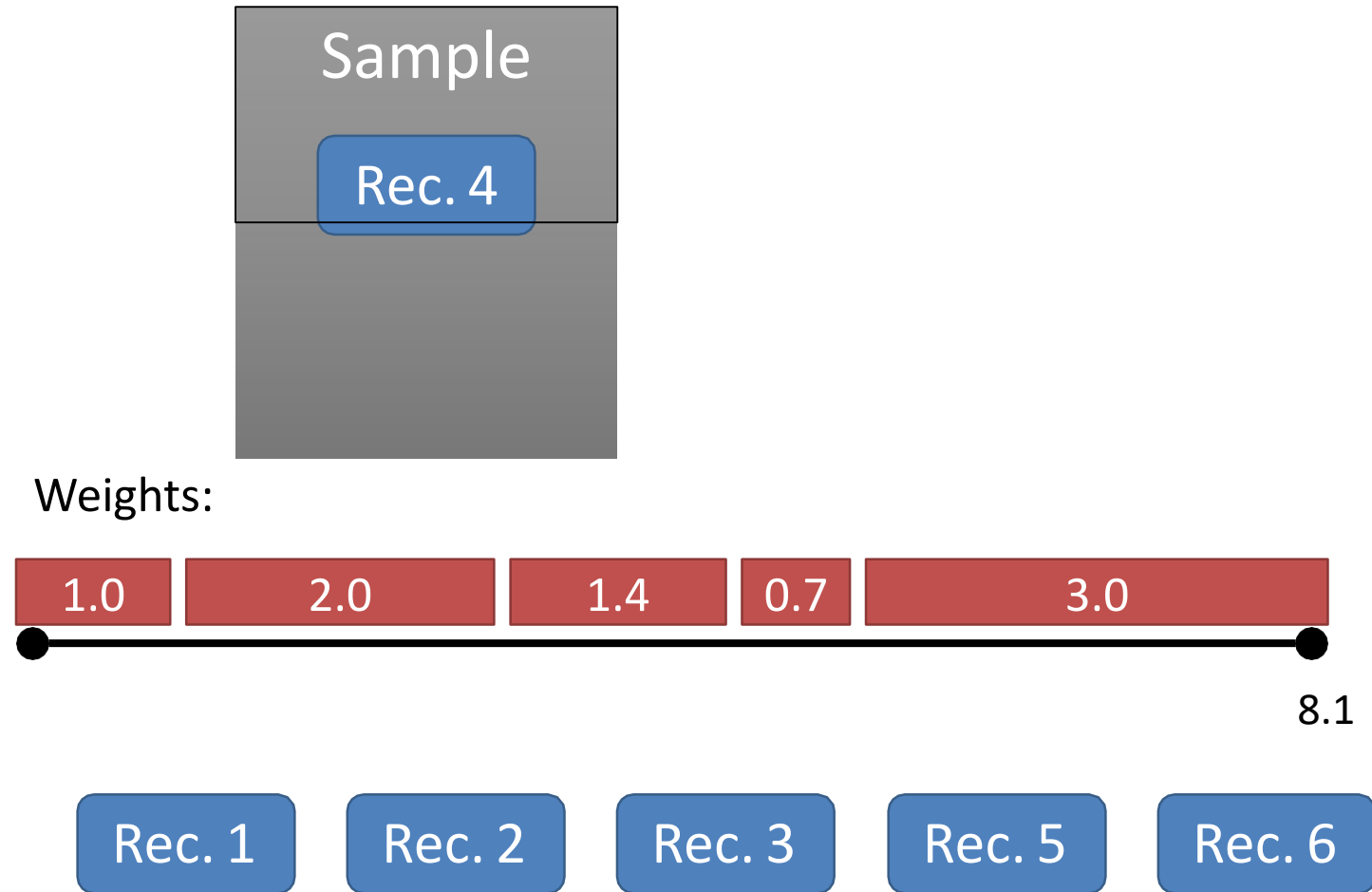
# Illustrating Basic Weighted Sampling



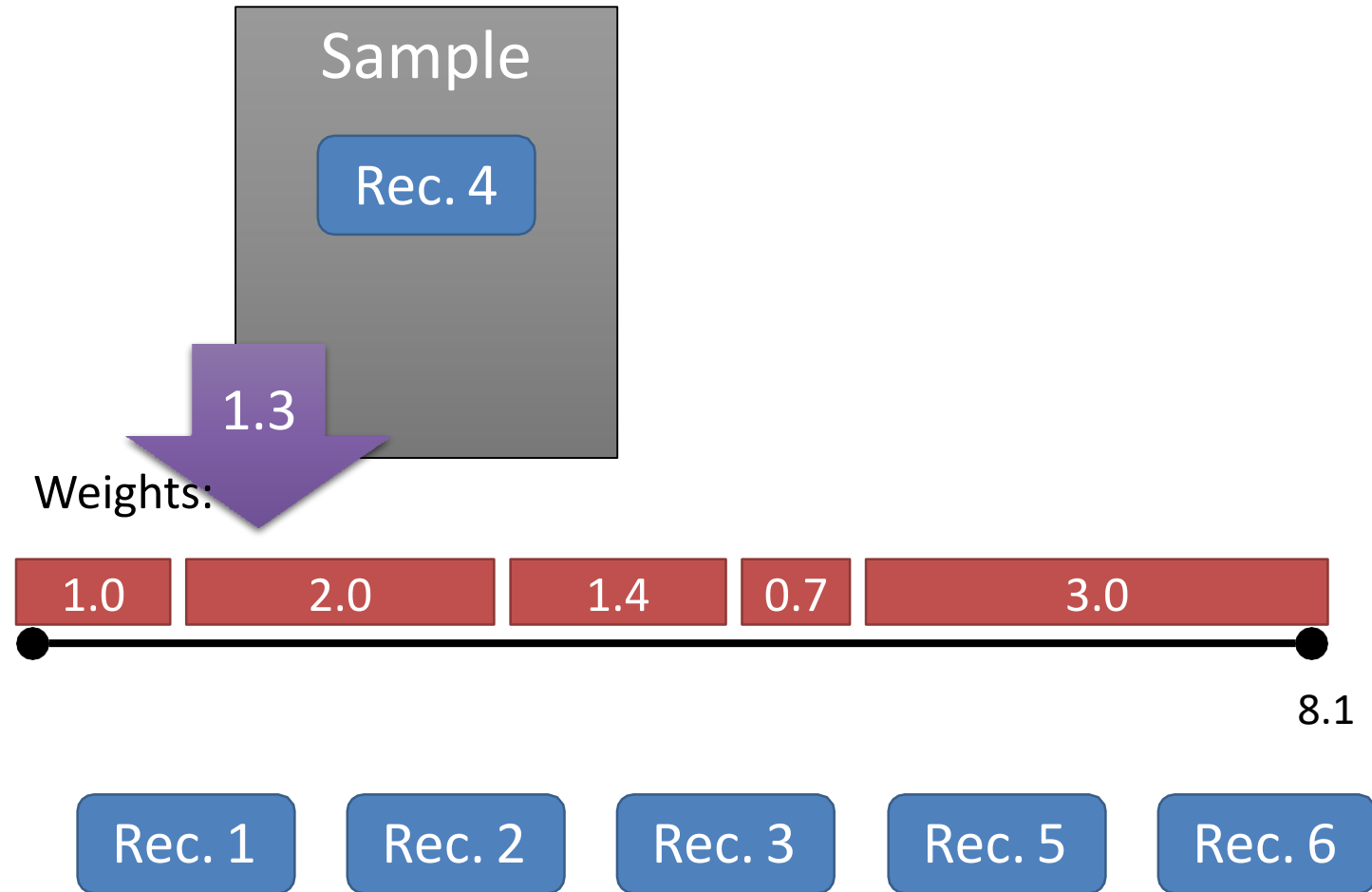
# Illustrating Basic Weighted Sampling



# Illustrating Basic Weighted Sampling

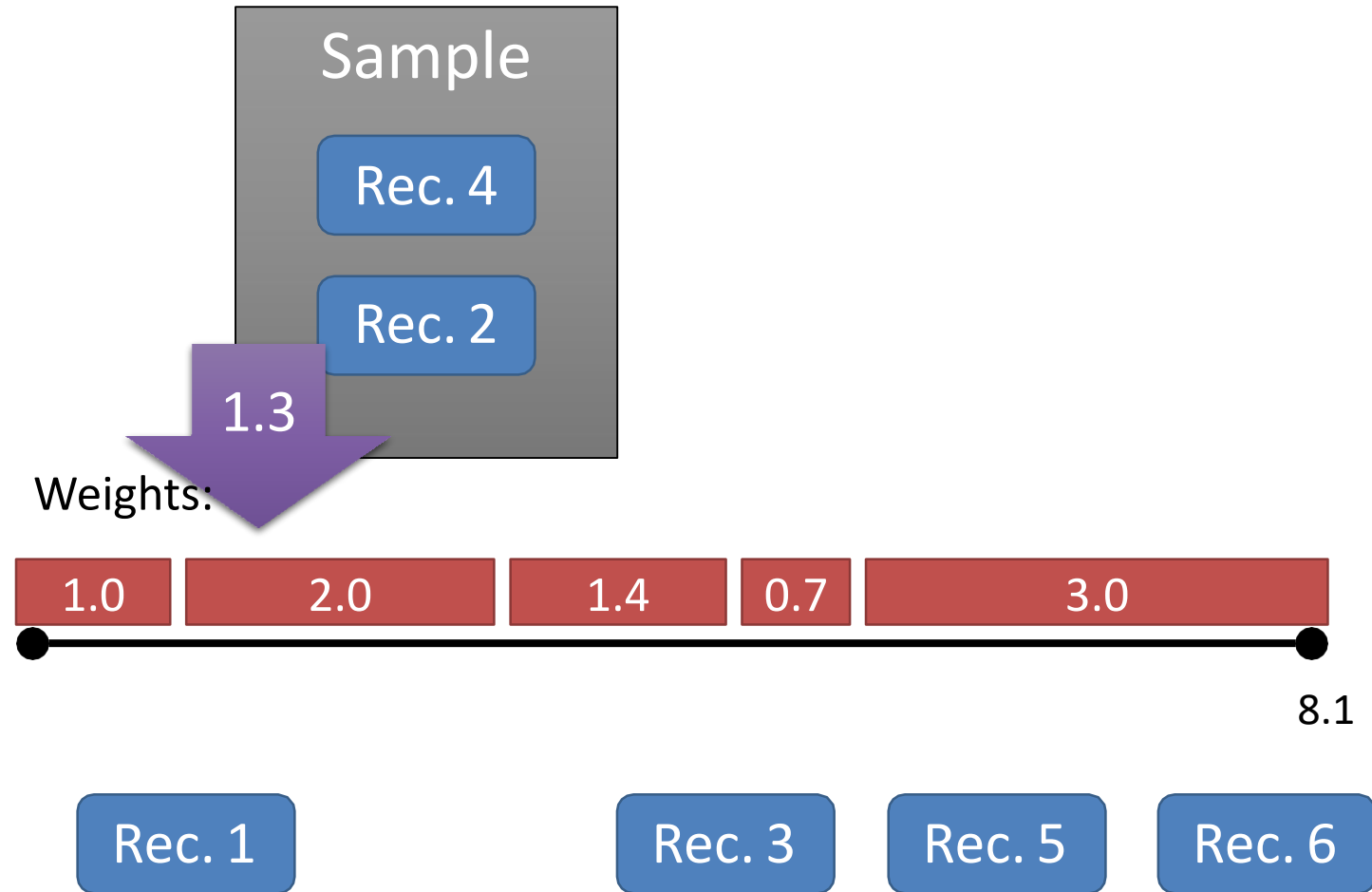


# Illustrating Basic Weighted Sampling

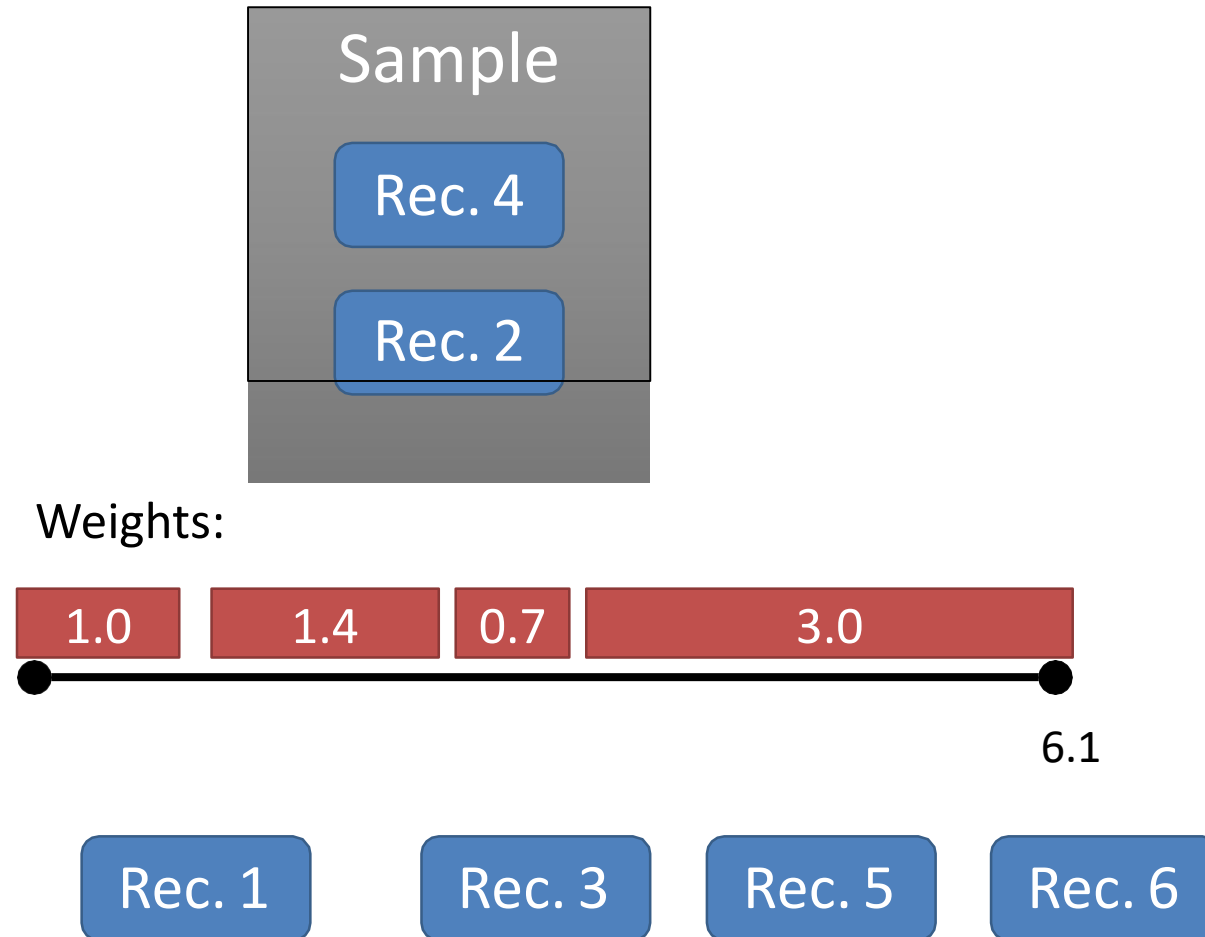




# Illustrating Basic Weighted Sampling



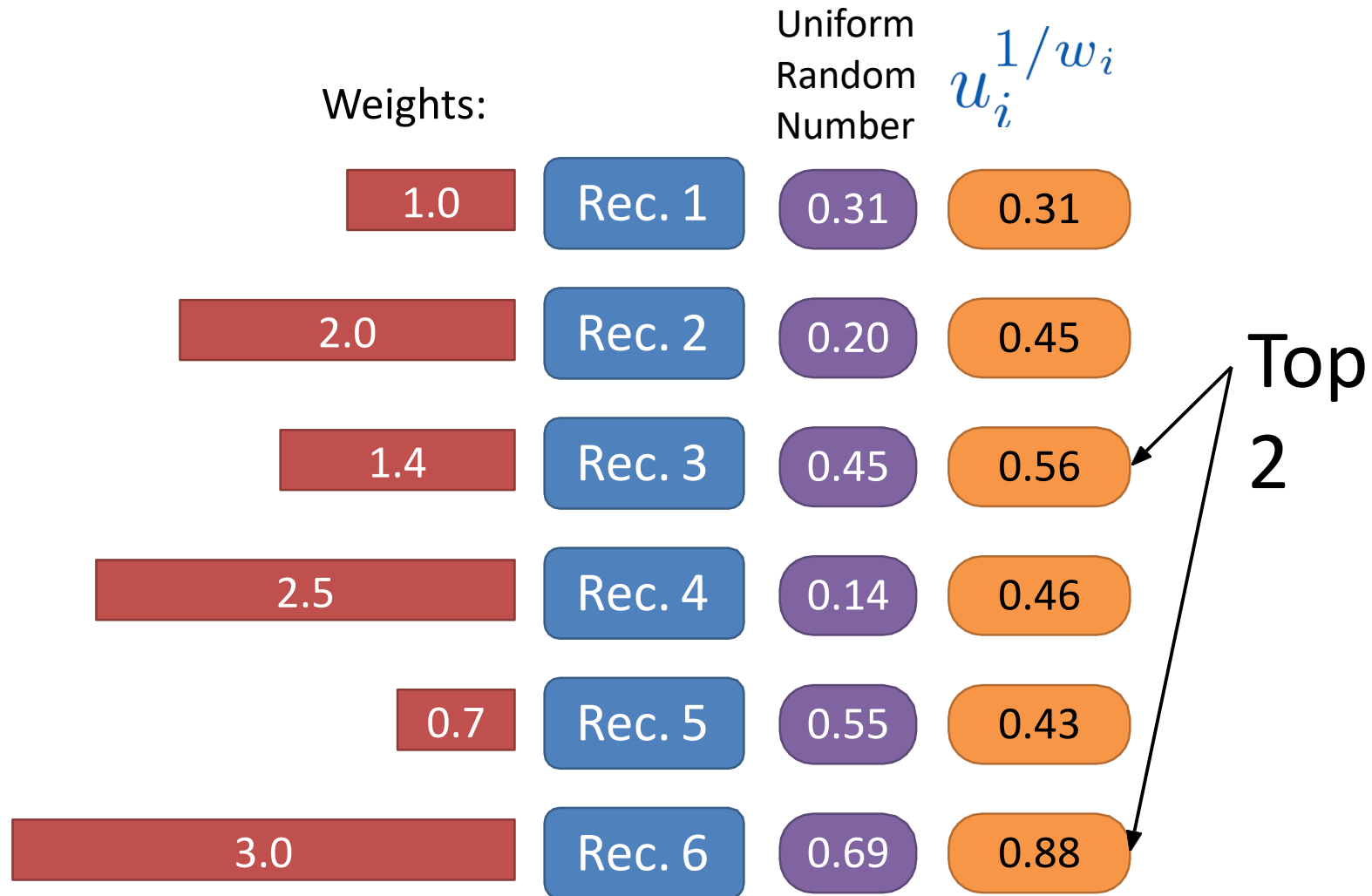
# Illustrating Basic Weighted Sampling



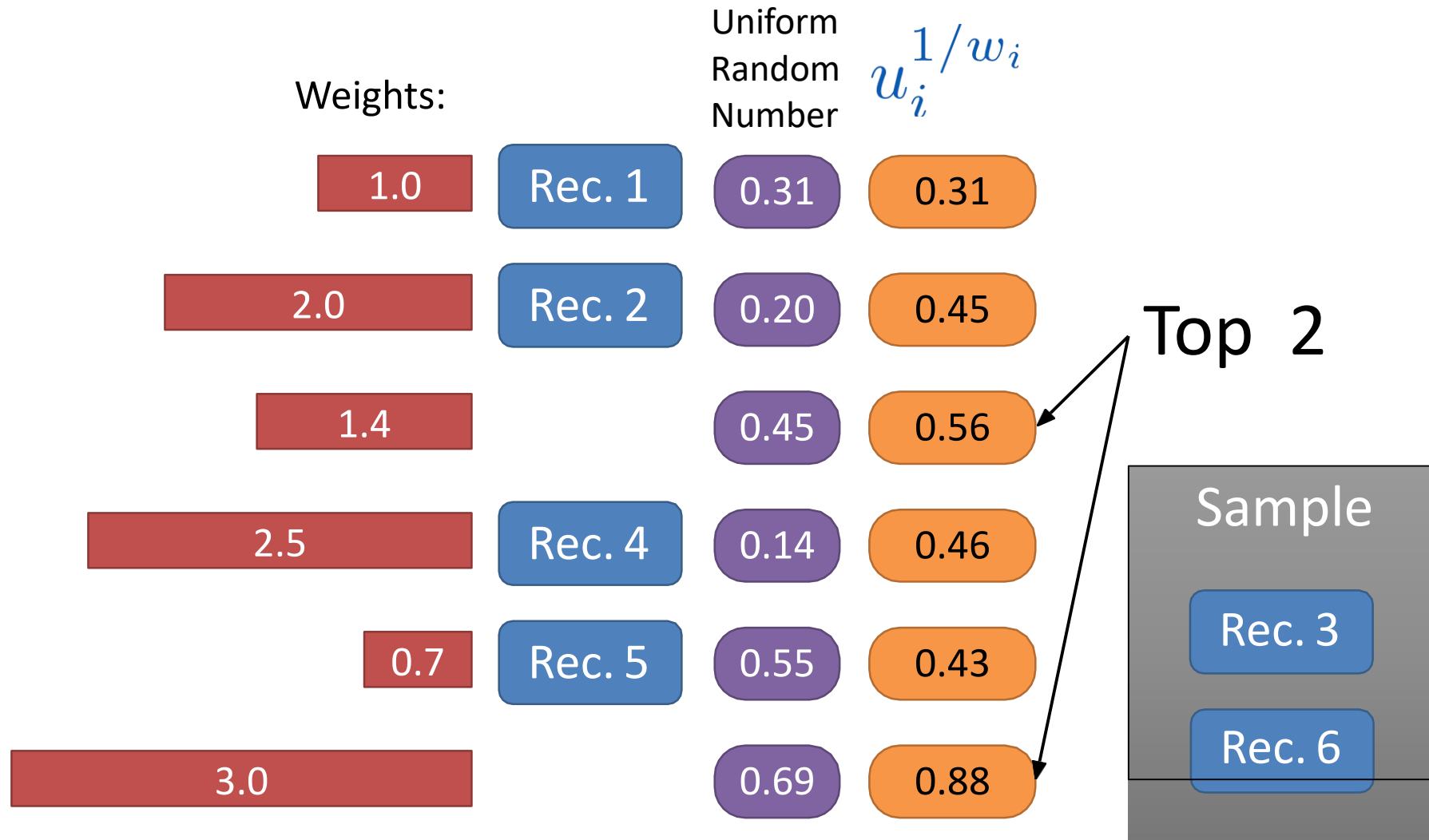
# Illustrating Res-A Algorithm

Weights:		Uniform Random Number	$u_i^{1/w_i}$
1.0	Rec. 1	0.31	0.31
2.0	Rec. 2	0.20	0.45
1.4	Rec. 3	0.45	0.56
2.5	Rec. 4	0.14	0.46
0.7	Rec. 5	0.55	0.43
3.0	Rec. 6	0.69	0.88

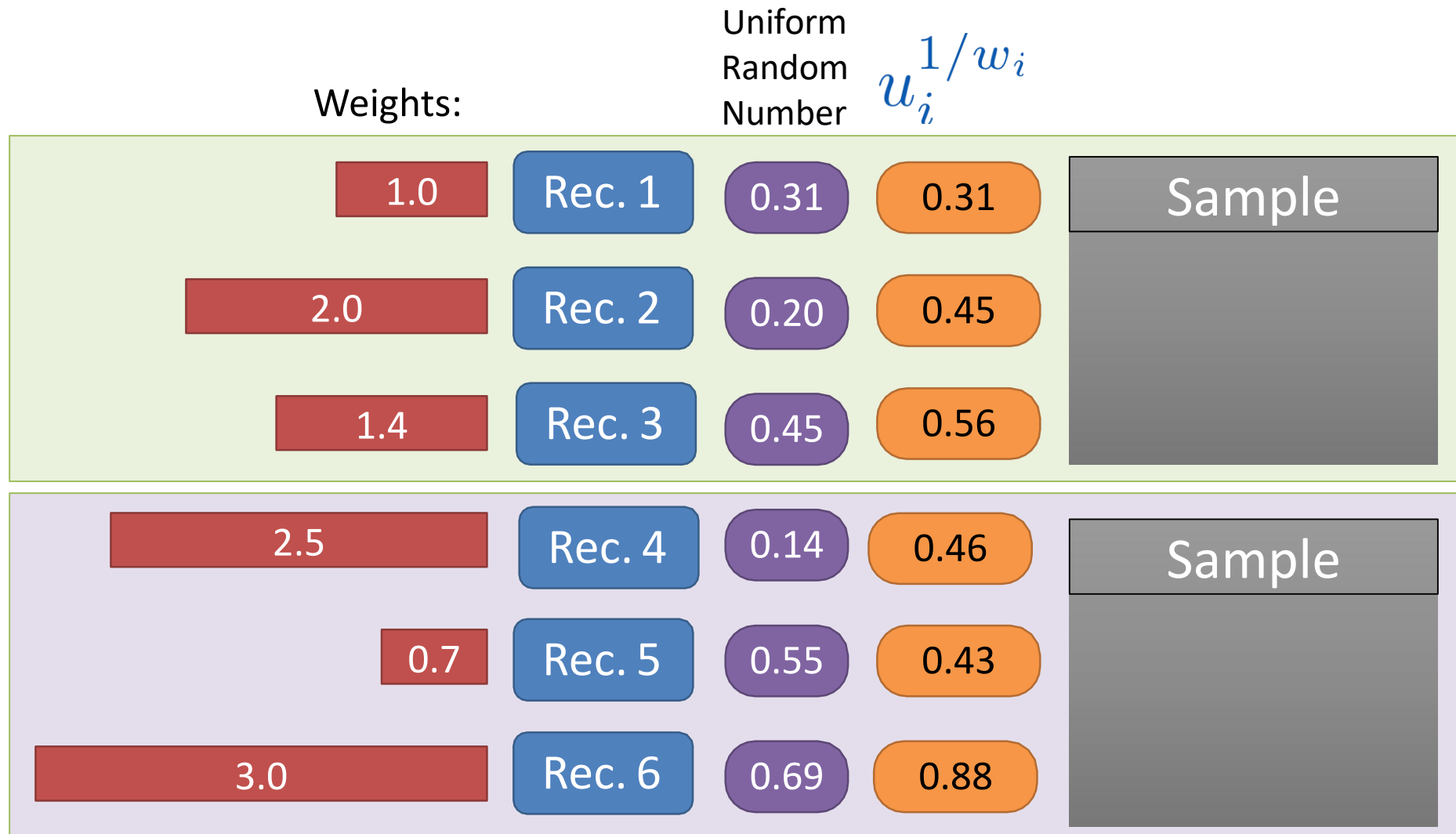
Top 2



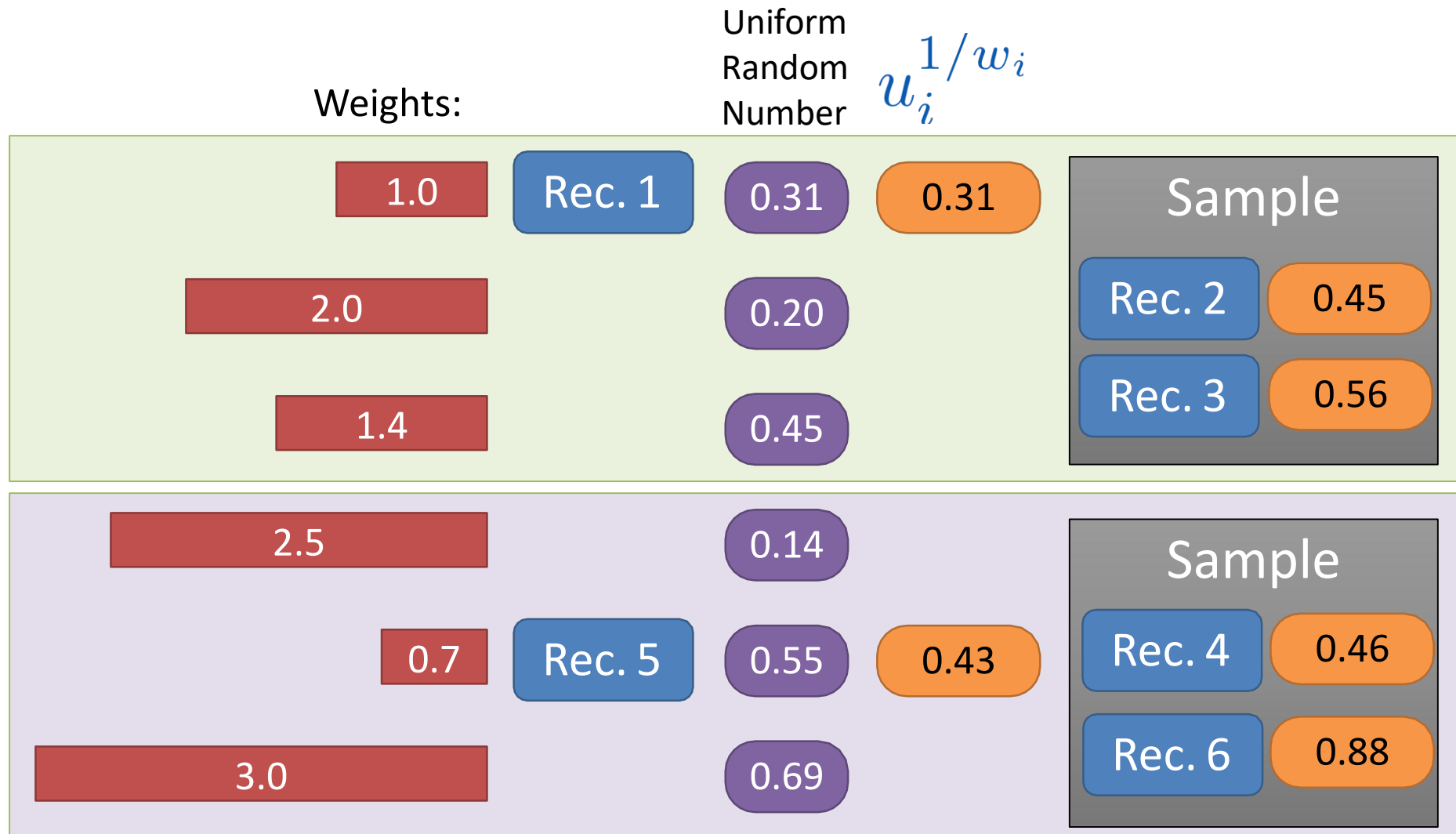
# Illustrating Res-A Algorithm



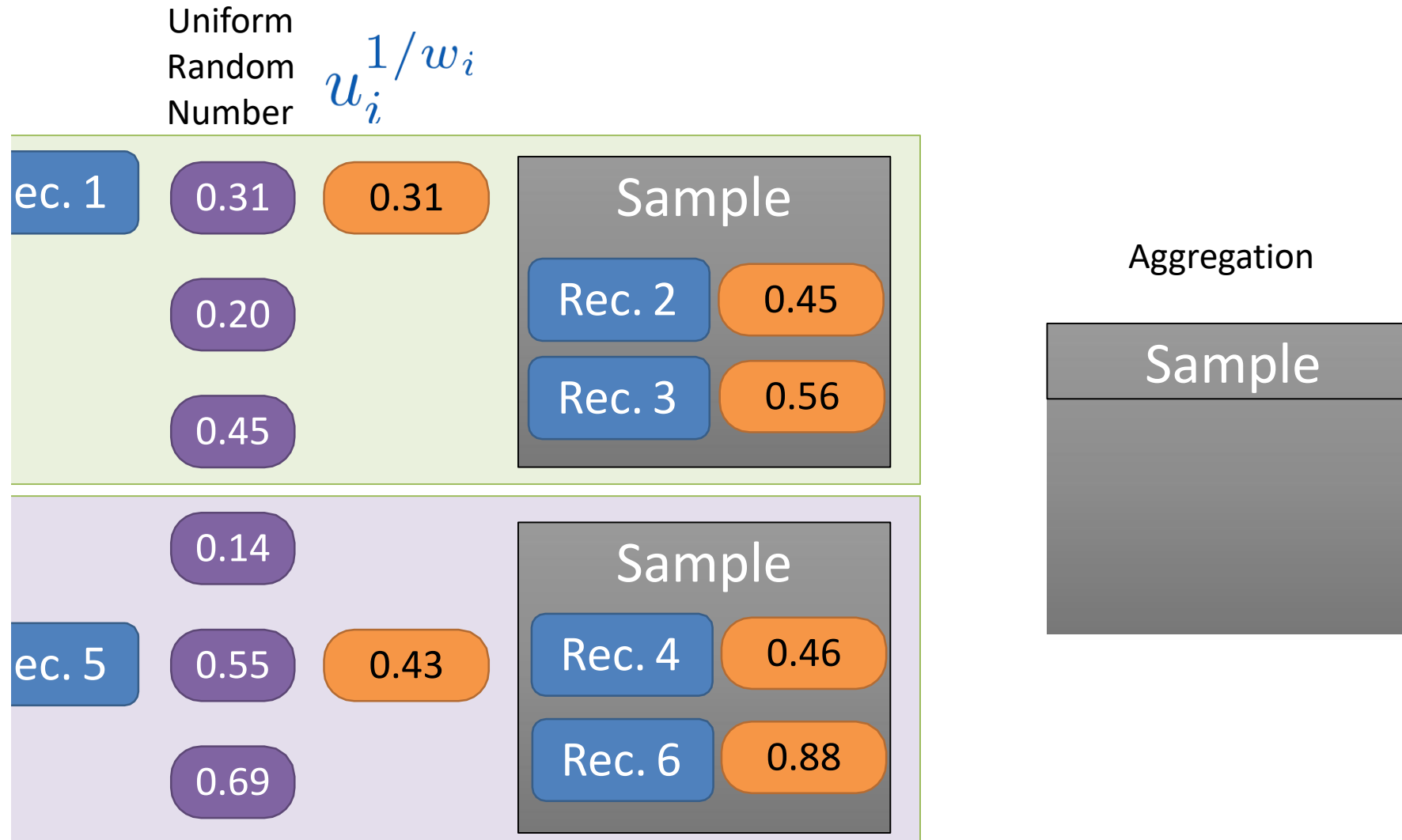
# Parallelizing Res-A Algorithm



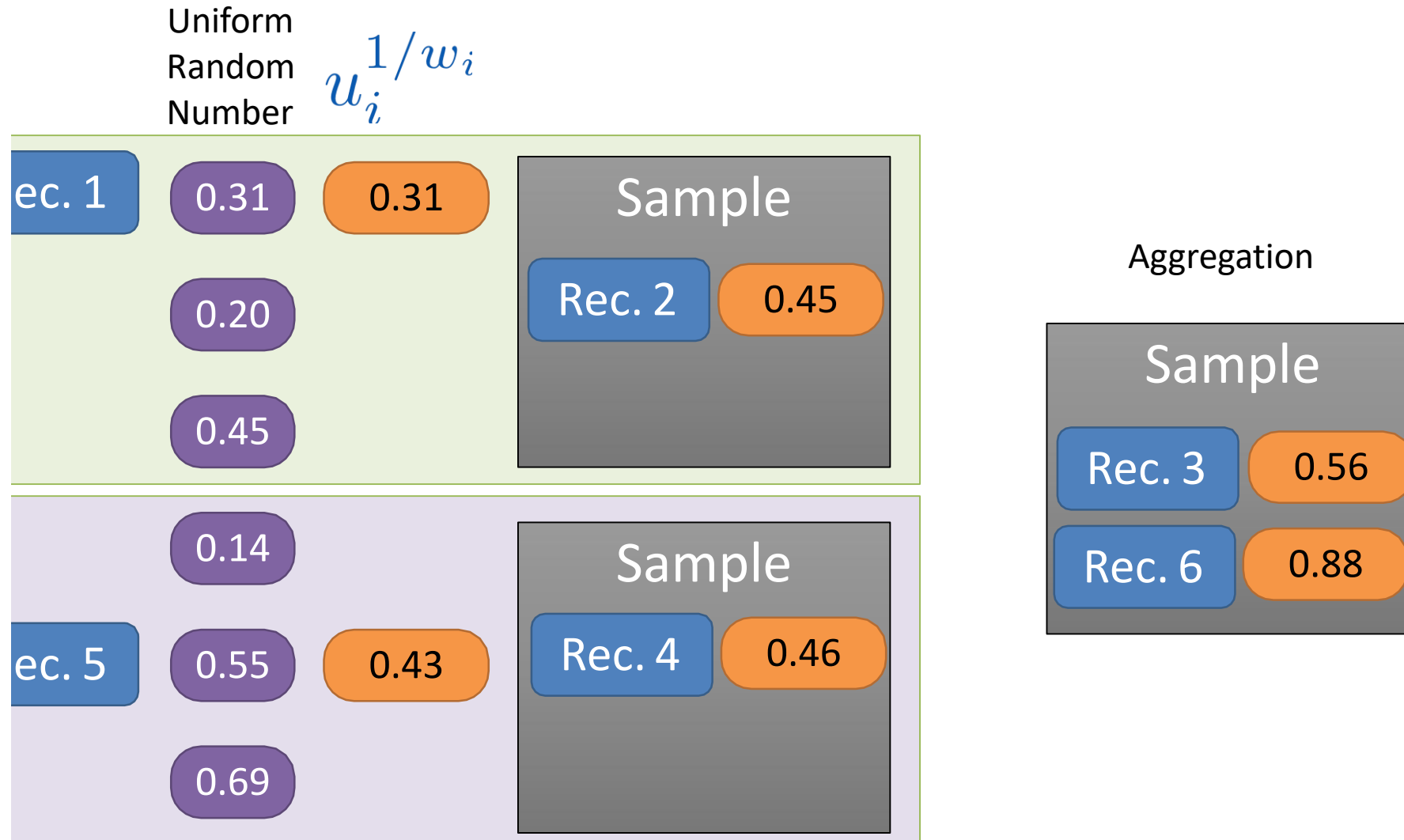
# Parallelizing Res-A Algorithm



# Parallelizing Res-A Algorithm



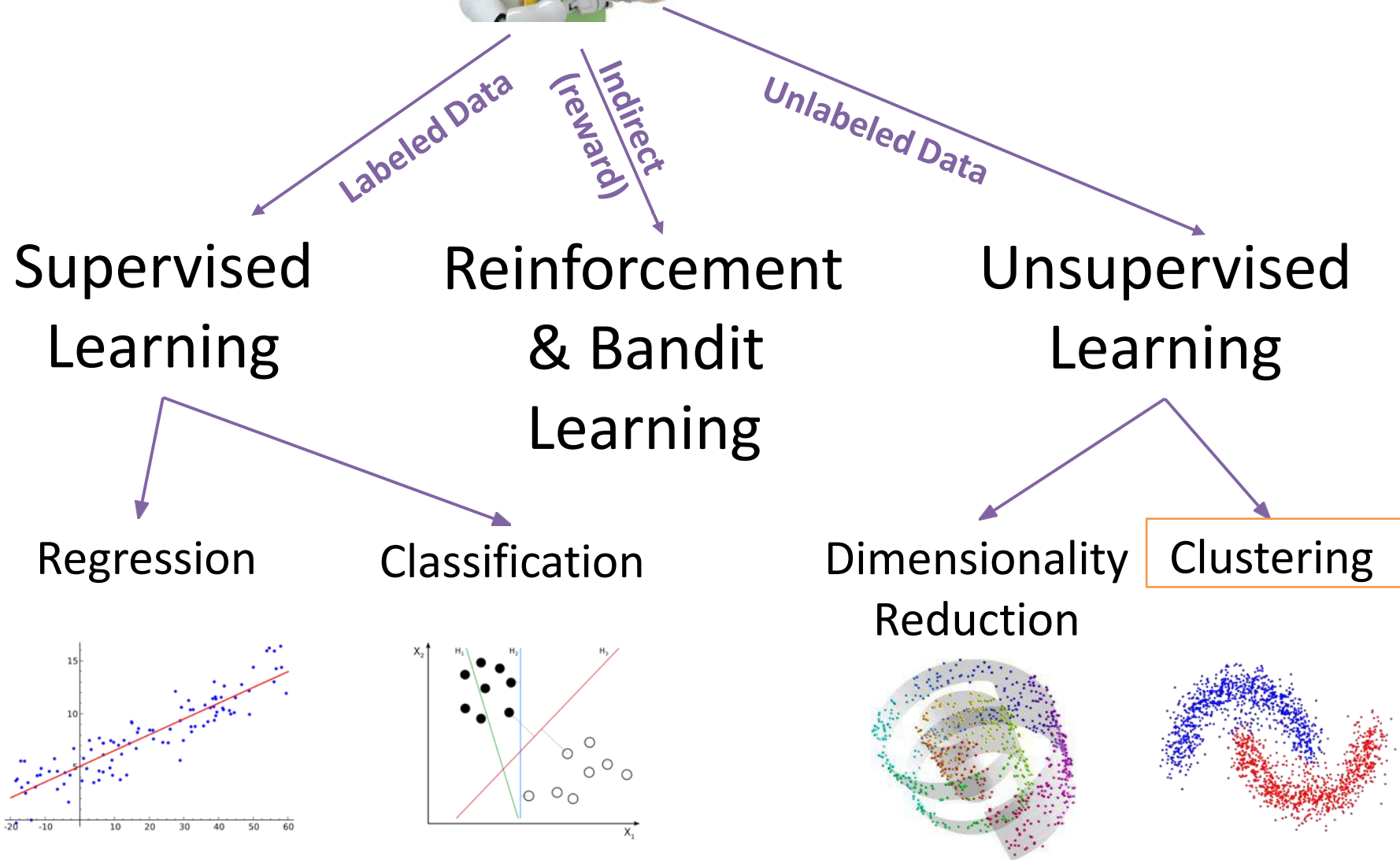
# Parallelizing Res-A Algorithm





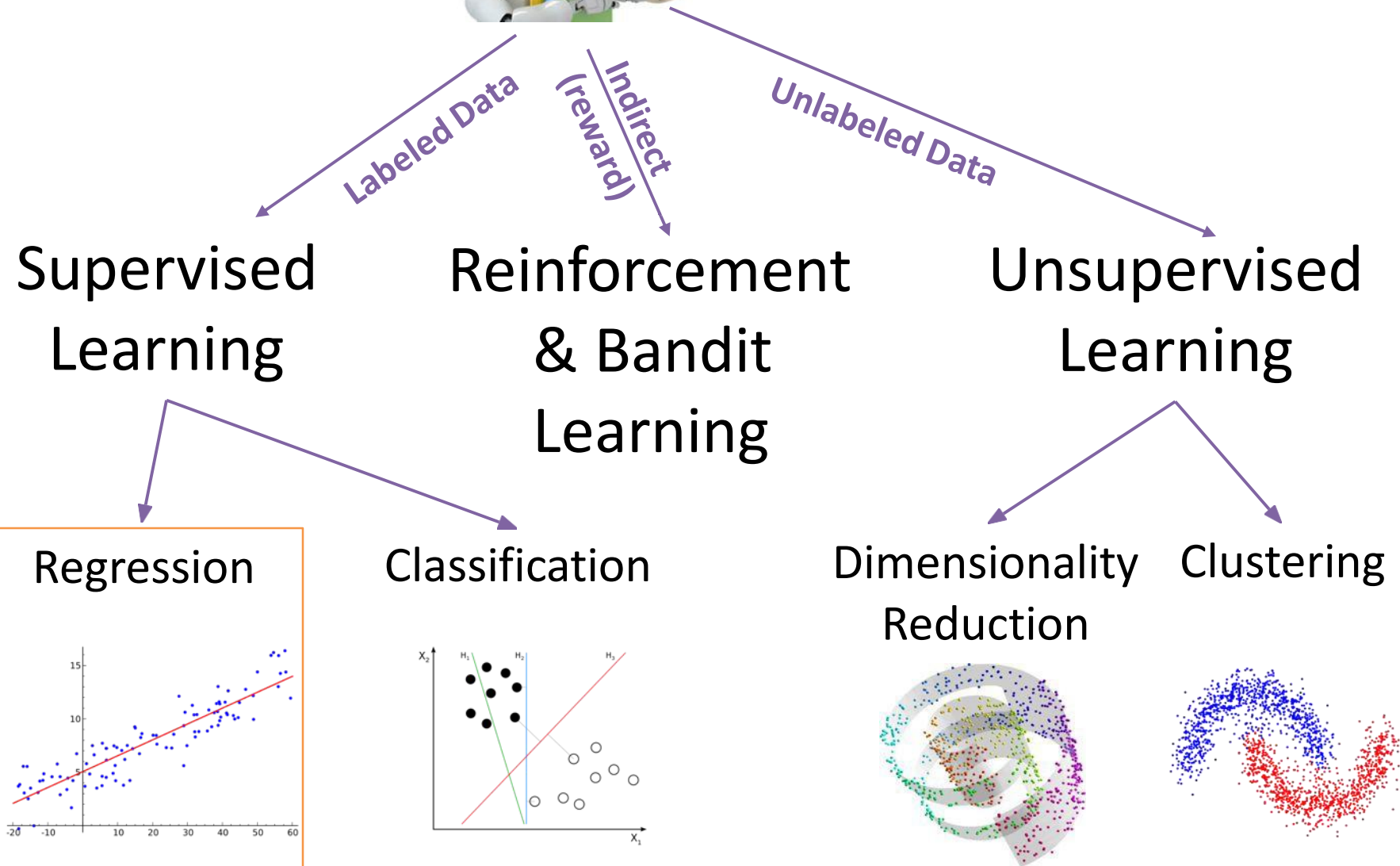


# Taxonomy of Machine Learning

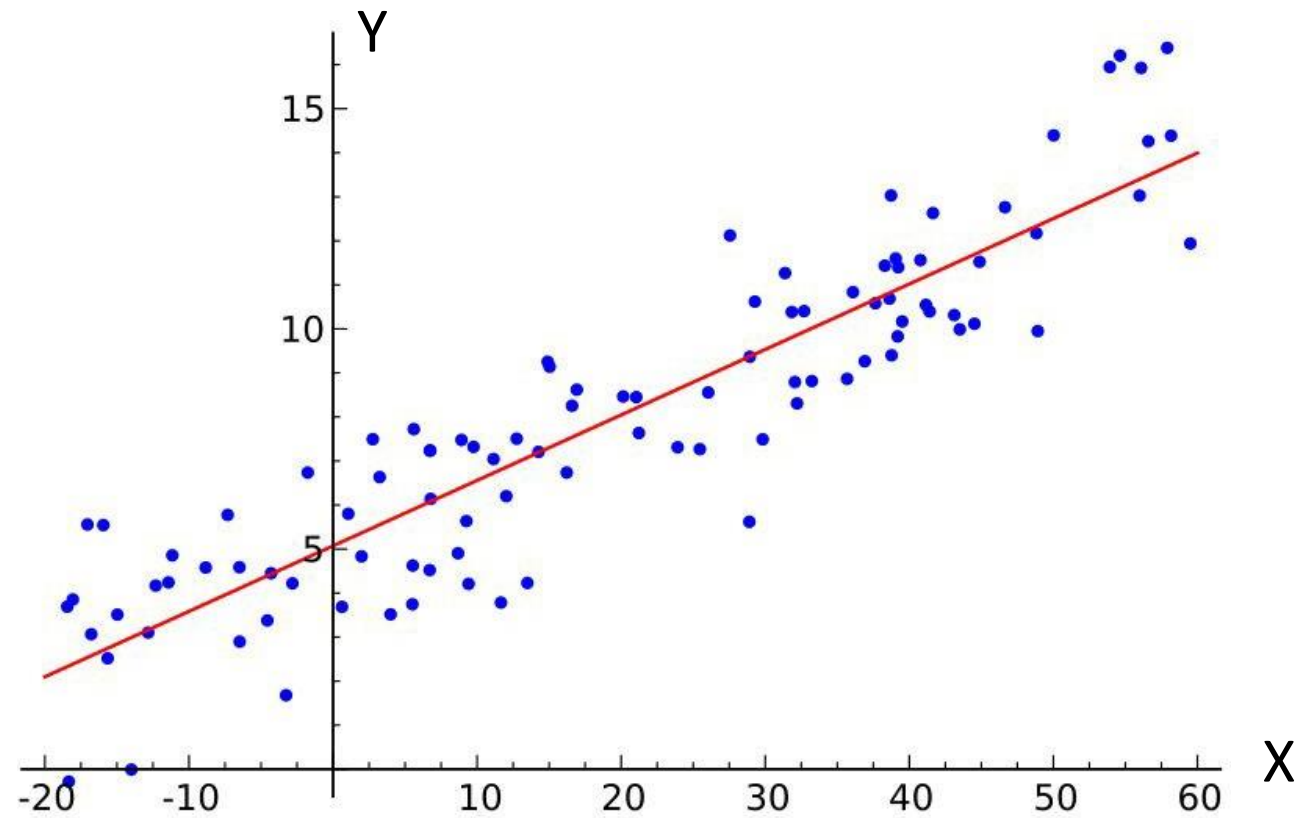




# Taxonomy of Machine Learning



# Simple Linear Regression



# Linear Regression is Powerful

- One of the most widely used techniques
- Fundamental to many larger models
  - Logistic Regression
  - Collaborative filtering
- Easy to interpret
  - e.g., the weights tell us something about the features
    - Positive or negative relationships ...
- Efficient to solve
  - Fast numerical methods
  - Closed form solutions

# The Linear Model

Data:

$x_1$	$x_2$	$y$
1.1	2.7	3.6
4.2	3.2	7.5
9.8	9.2	17
...	...	...

Real Valued Observations

Vector of Parameters

Vector of Features

$$y = \theta^T x + \epsilon$$

Real Value Noise

Linear Combination of Covariates

$$\sum_{i=1}^p \theta_i x_i$$

$$\theta, x \in \mathbb{R}^p$$

Data:

$x_1$	$x_2$	$y$
1.1	2.7	3.6
4.2	3.2	7.5
9.8	9.2	17
...	...	...

Vector of  
Features

$$\mathcal{X} \in \mathbb{R}^p$$

“Real” data doesn’t typically consist of entirely **real** valued features.

# Real Data and Vector Spaces

- What about data with more complex schemas?

$x_1$	$x_2$	Date	prod_id	comment	y
1.1	2.7	8/21/16	7	"the best glider ..."	3.6
4.2	3.2	8/14/16	3	"vacation for two ..."	7.5
9.8	9.2	9/20/16	4	"A special gift for ..."	17
...	...	...	...	...	...

- The math wants the features to be vectors ...
- How do we encode dates, categorical fields, and text?

# Feature Engineering

- A key part of most machine learning applications
- Common tasks:
  - Transforming raw features into **vector representations**
  - Encoding **prior knowledge** (e.g., translating currencies)
  - Transformations that **increase the expressivity** of the model ... (more on this soon)
- Critical to model performance:
  - **engineers compete** to get the best features
- A few standard techniques (that we will cover):
  - one-hot encoding
  - bag-of-words



# Encoding Categorical Data

- How do we represent fields like “Product Category”
- **Proposal 1:** *Enumerate categories*
  - Sports = 1, Furniture = 2, Clothing = 3, Shoes = 4, ...
  - Store field number as a feature
  - **Implications:**
    - **similarity:** sports is closer to Furniture than shoes
    - **magnitude:** larger values ?
  - Not typically used (unless there are two categories ...)

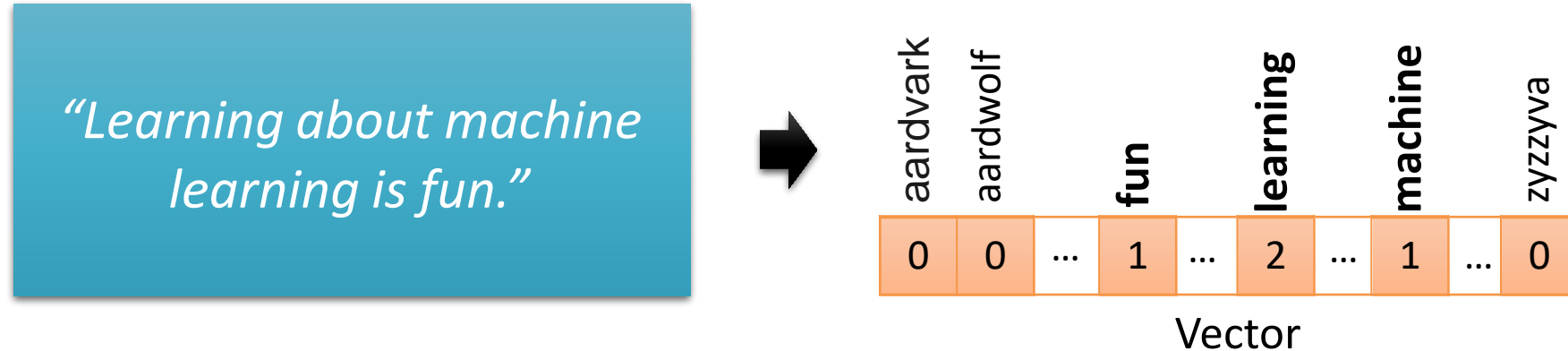
# One-hot encoding

- How do we represent fields like “Product Category”
- **Proposal 1:** *Enumerate categories*
  - Not typically used (unless there are two categories ...)
- **Proposal 2:** *Encode as binary vectors:*
  - Very commonly used and built-in to many packages
  - Enumerate all possible product categories ( $m$ )
  - Add  $m$  additional features to the record:
  - Put a one in the feature corresponding to the product category and a zero everywhere else.

0	0	0	0	1	0	0
Sports	Furniture	Clothing	Shoes	Electronics	Music	Books

# Working with Text Data

- How do we convert text to vectors?



- Bag-of-words model
  - Transform emails into  $d$ -dimensional vectors
    - $d$  is the number of unique words in the language (big!)
  - Each entry is number of occurrences of that word
  - **Sparse**: Most words don't occur in most emails
  - **Remove Stop-Words**: common words that provide little information (e.g., "is", "about")

If all you had was this  
vector could you tell  
what the passage is  
about?

aardvark	aardwolf		fun		learning		machine		zyzzyva
0	0	...	1	...	2	...	1	...	0

Vector



# The Linear Model

Vector of  
Parameters

Vector of  
Features

$$f_{\theta}(x) := \theta^T x$$

Data:

$x_1$	$x_2$	$y$
1.1	2.7	3.6
4.2	3.2	7.5
9.8	9.2	17
...	...	...

- Encode data is real valued vectors
- **Next:** find the optimal value for  $\theta$ 
  - How?

$$\theta, x \in \mathbb{R}^p$$

# Finding the Best Parameters

Model:  $f_{\theta}(x) := \theta^T x$

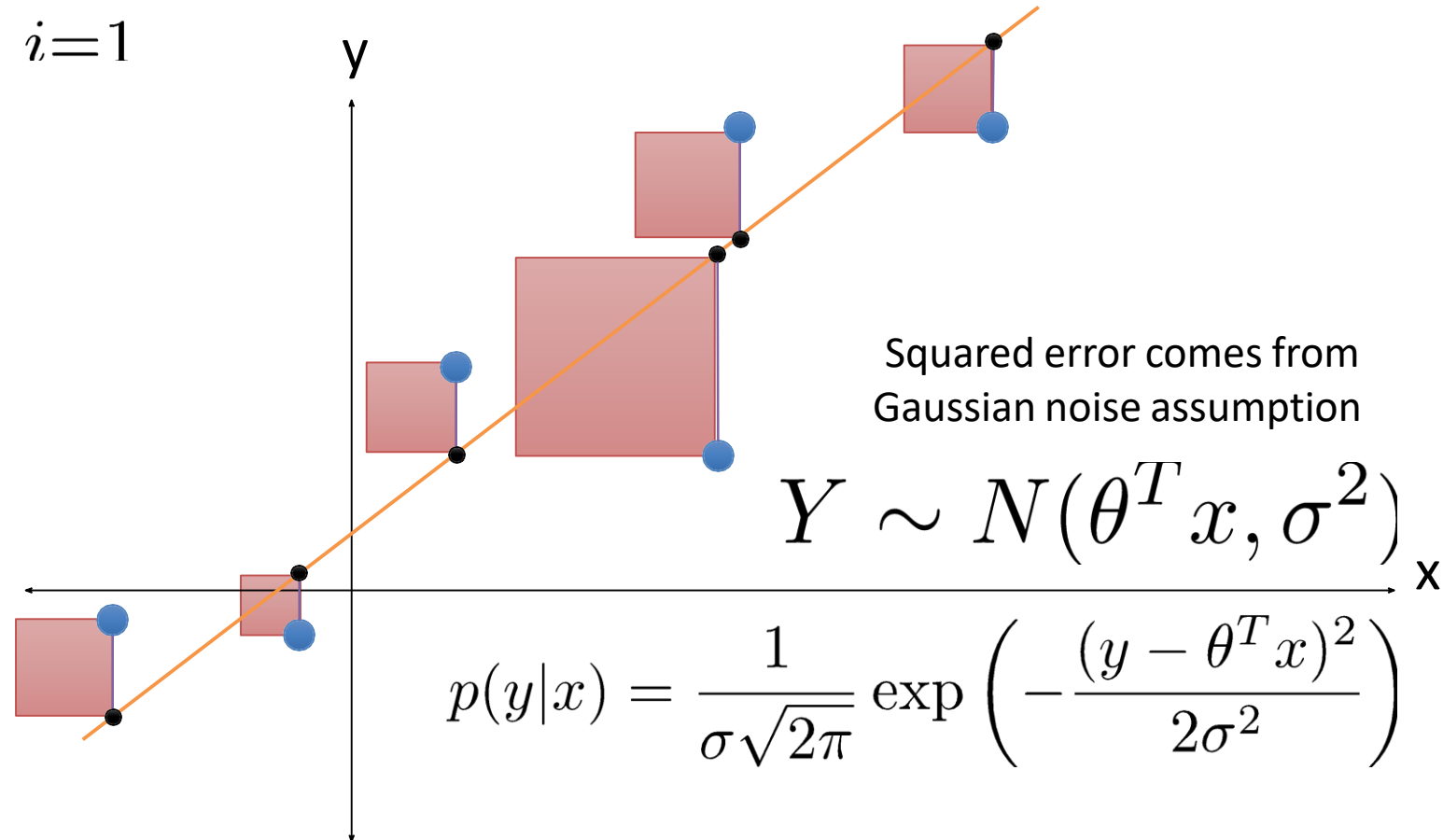
**Step 1:** define a **Loss Function:** *Average Prediction Error*

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

□ Difference between **true**  $(\mathbf{y}_i)$  and **predicted**  $\mathbf{f}_{\theta}(\mathbf{x}_i)$  values

# The meaning of Squared Loss (Error)

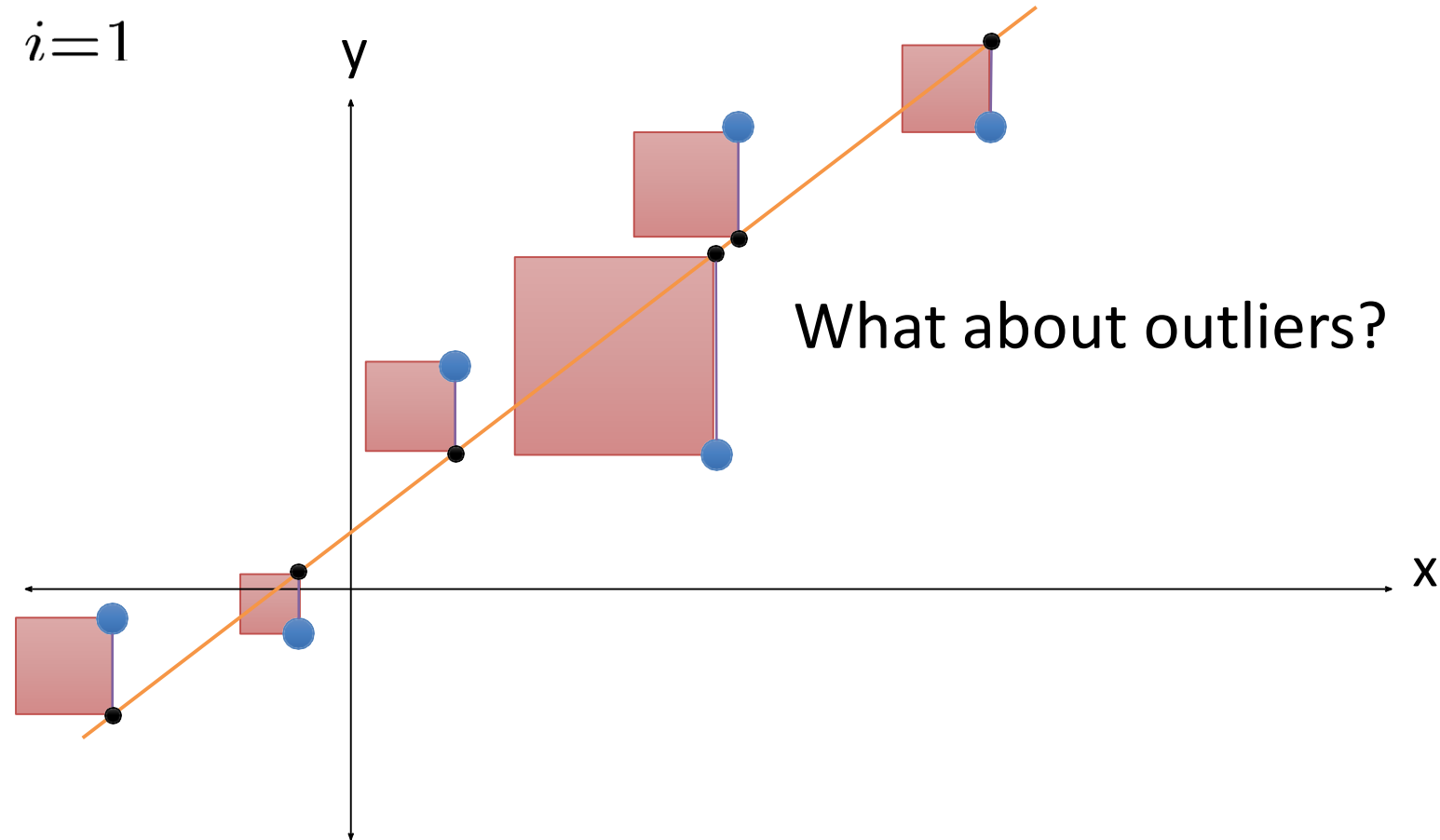
$$\frac{1}{n} \sum_{i=1}^n (y_i - \theta^T x_i)^2$$





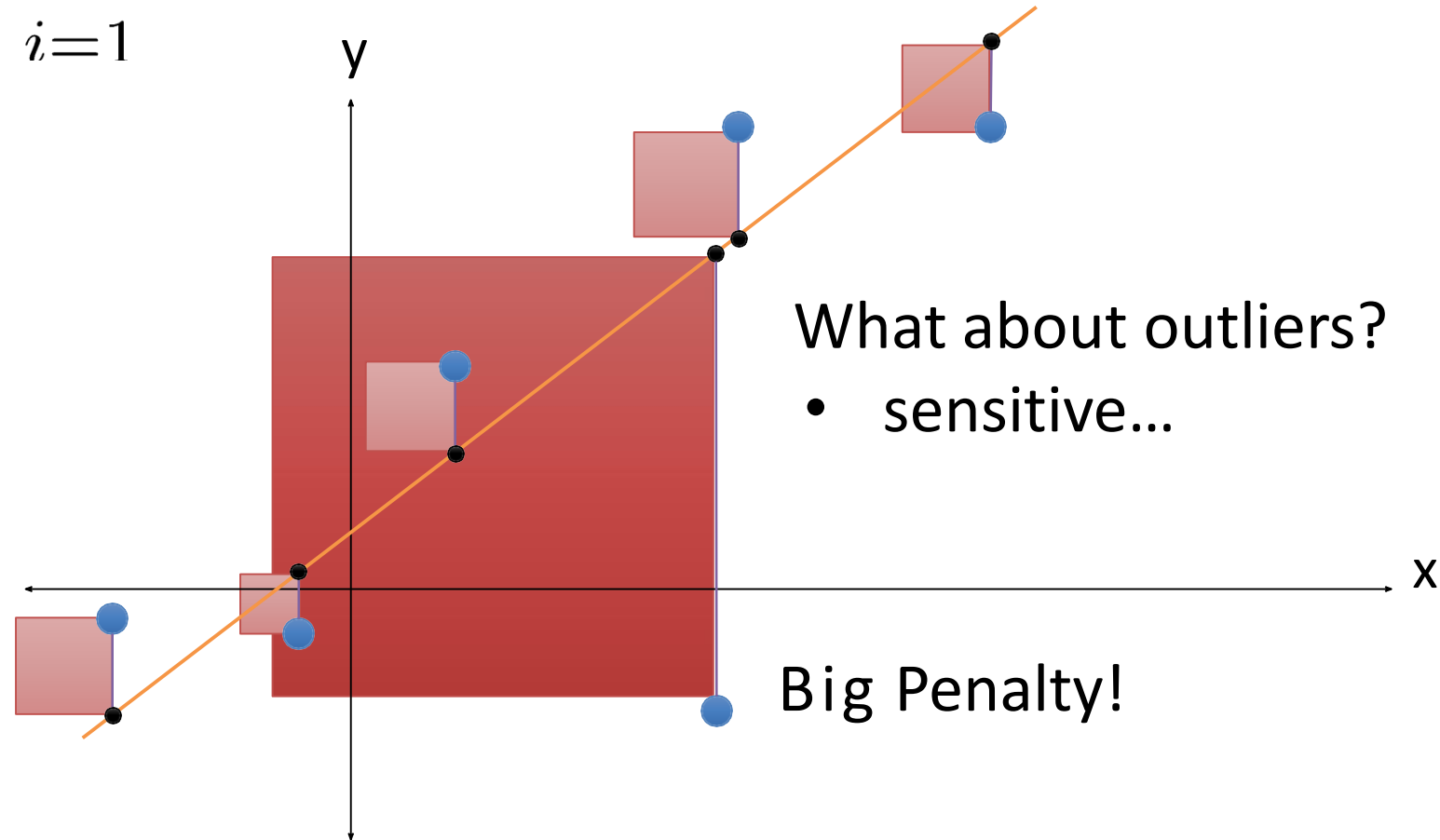
# The meaning of Squared Loss (Error)

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



# The meaning of Squared Loss (Error)

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



# Finding the Best Parameters

Model:  $f_{\theta}(x) := \theta^T x$

**Step 1:** define a **Loss Function**:

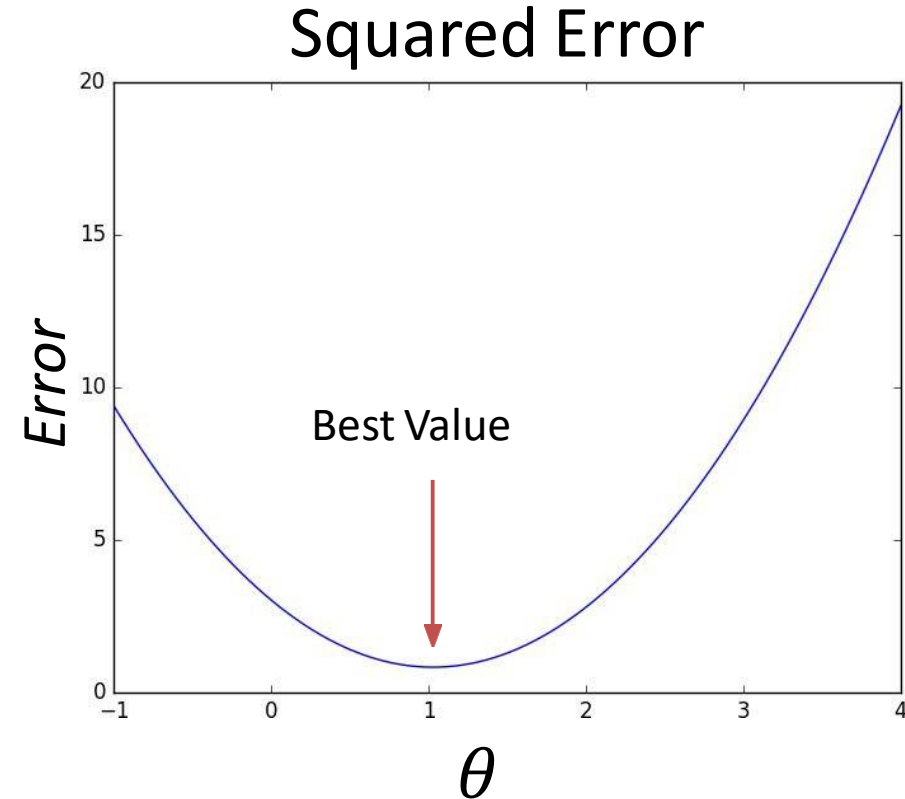
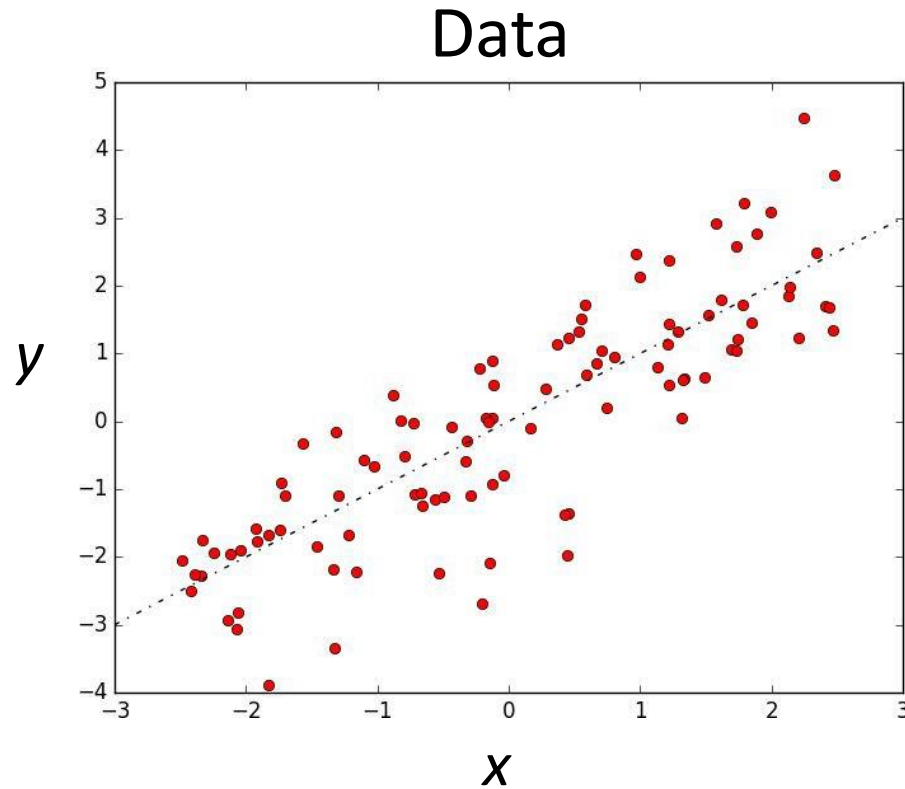
$$\frac{1}{n} \sum_{i=1}^n (y_i - f_{\theta}(x_i))^2$$

**Step 2:** Search for best model parameters  $\theta$

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

# Minimizing the Squared Error

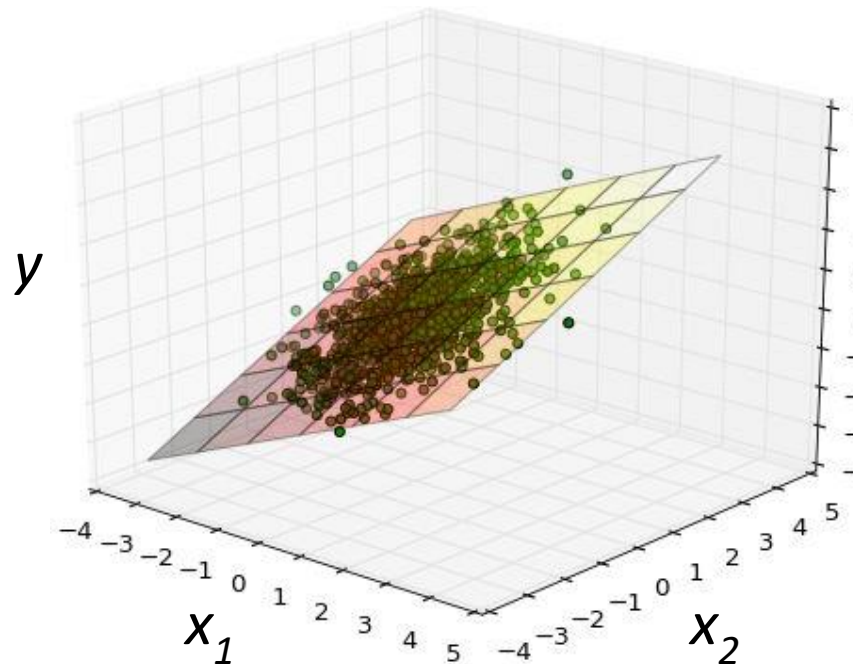
$$\hat{\theta} = \arg \min_{\theta \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n (y_i - \theta^T x_i)^2$$



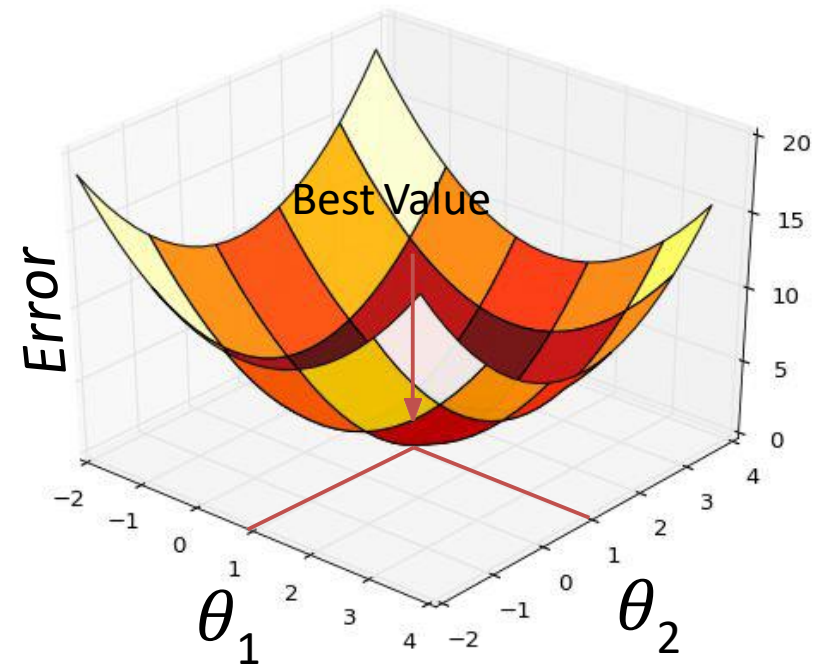
# Minimizing the Squared Error

$$\hat{\theta} = \arg \min_{\theta \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n (y_i - \theta^T x_i)^2$$

Data

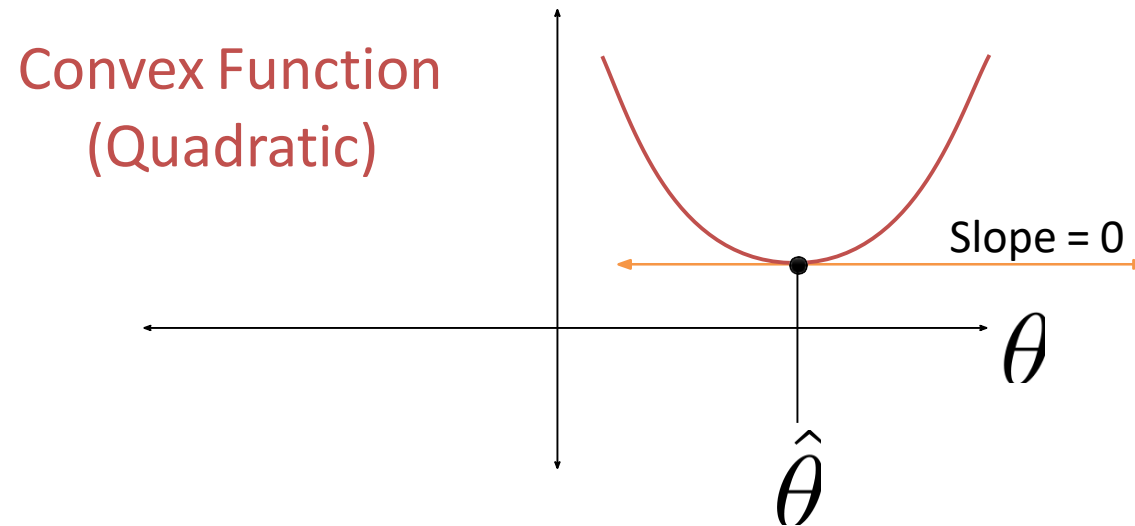


Squared Error



# Minimizing the Squared Error

$$\hat{\theta} = \arg \min_{\theta \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n (y_i - \theta^T x_i)^2$$



- Take the gradient and set it equal to zero

# Minimizing the Squared Error

$$\hat{\theta} = \arg \min_{\theta \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n (y_i - \theta^T x_i)^2$$

□ Taking the gradient

$$\begin{aligned} \nabla_{\theta} \frac{1}{n} \sum_{i=1}^n (y_i - \theta^T x_i)^2 &= -2 \frac{1}{n} \sum_{i=1}^n (y_i - \theta^T x_i) x_i && \text{Chain Rule} \\ &= -2 \frac{1}{n} \sum_{i=1}^n y_i x_i + 2 \frac{1}{n} \sum_{i=1}^n (\theta^T x_i) x_i \end{aligned}$$

□ Setting equal to zero and solving for  $\theta$  (sys. Linear eq.)

$$\sum_{i=1}^n (\theta^T x_i) x_{ij} = \sum_{i=1}^n y_i x_{ij} \quad \forall j \in \{1, \dots, d\}$$

Easier in matrix form ...

# Writing the data in Matrix form

□ Represent data  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$  as:

$$\begin{array}{cc} \text{Covariate (Design)} & \text{Response} \\ \text{Matrix} & \text{Vector} \\ X = \begin{array}{c} \overbrace{\left( \begin{array}{c} x_1 \\ \text{---} \\ x_2 \\ \text{---} \\ \vdots \\ \text{---} \\ x_n \end{array} \right)}^n & \in \mathbb{R}^{np} \end{array} & Y = \begin{array}{c} \overbrace{\left( \begin{array}{c} y_1 \\ y_2 \\ \vdots \\ y_n \end{array} \right)}^n & \in \mathbb{R}^n \end{array} \end{array}$$



# Minimizing the Squared Error

$$\hat{\theta} = \arg \min_{\theta \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n (y_i - \theta^T x_i)^2$$

□ Setting equal to zero and solving for  $\theta$ :

$$\sum_{i=1}^n (\theta^T x_i) x_i = \sum_{i=1}^n y_i x_i \Rightarrow X^T X \theta = X^T y$$

□ Normal Equation:

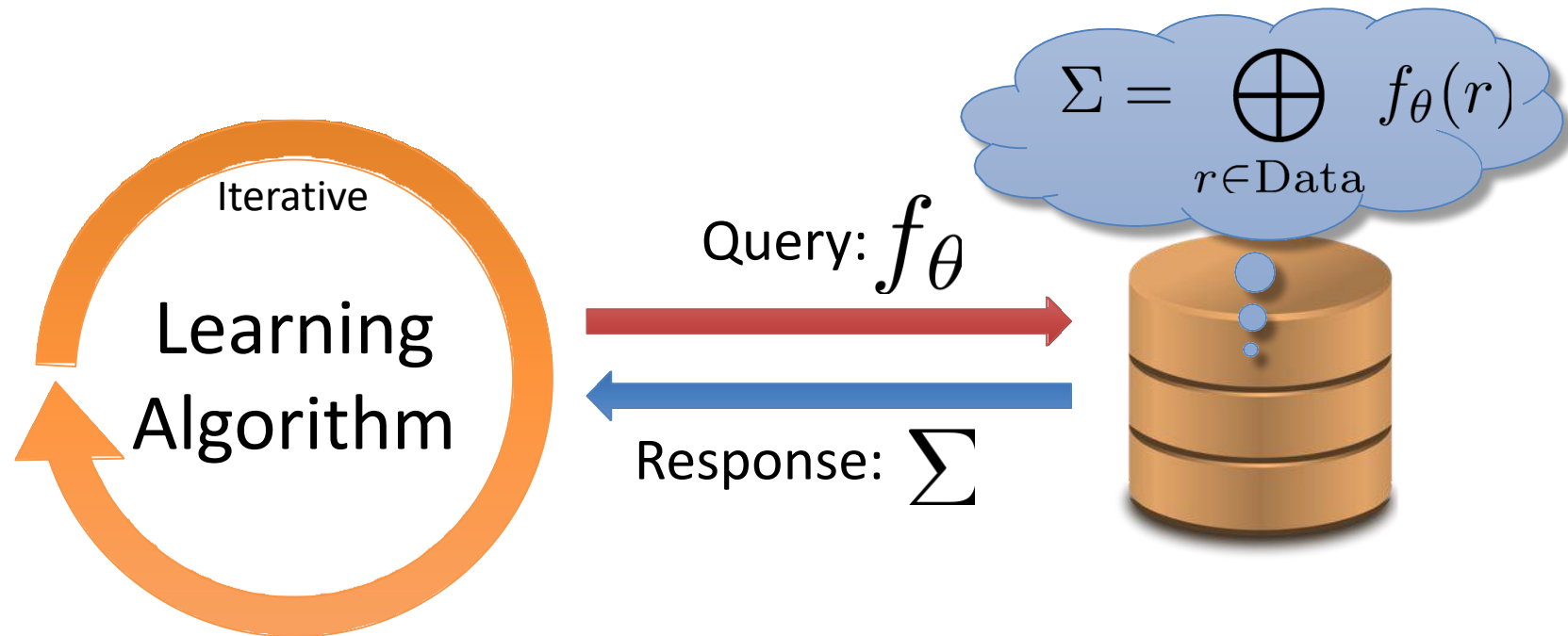
$$\hat{\theta} = (X^T X)^{-1} X^T Y$$

□ Solved using any standard linear algebra library

# Can we compute

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

## using the statistical query pattern?



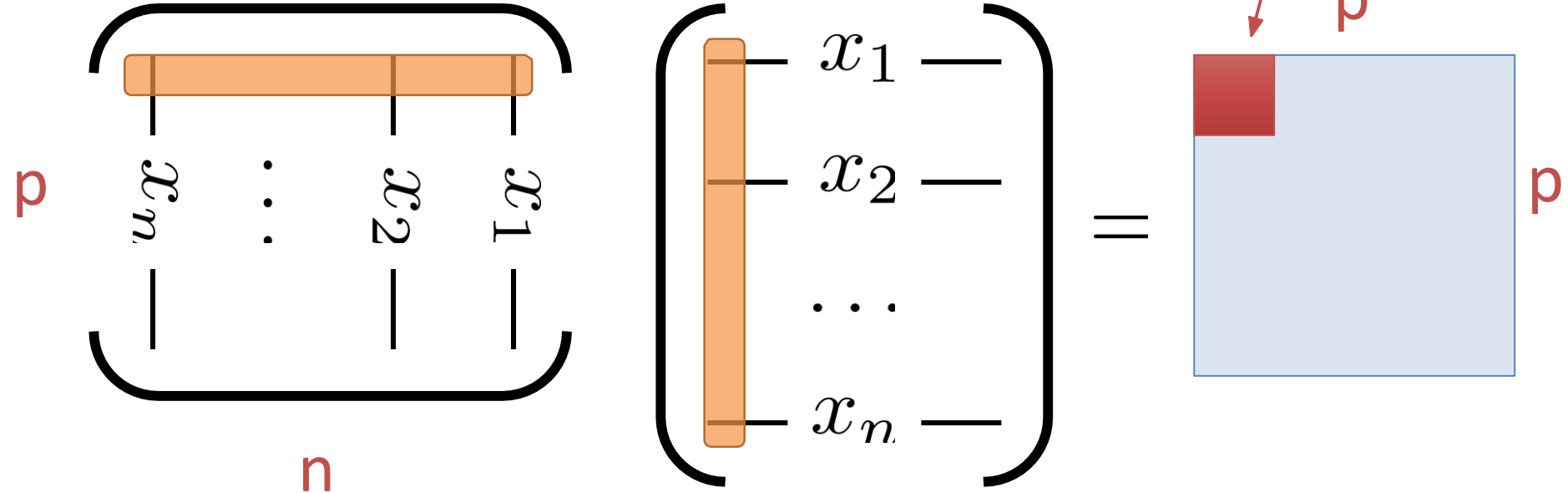
$$\hat{\theta} = (\underline{X^T X})^{-1} \underline{X^T Y}$$

Query 1

Query 2

Break  
computation  
into two queries

# Computing Query 1

$$X^T X =$$


The diagram illustrates the computation of the first column of the matrix product  $X^T X$ . The matrix  $X$  is shown with dimensions  $p$  by  $n$ . The first column of  $X$  is highlighted in orange. The product  $X^T X$  is shown as a sum of outer products of the columns of  $X$ . The first term in the sum is shown as a  $p$  by  $p$  matrix with a red square in the top-left corner, representing the outer product of the first column of  $X$  with itself.

# Computing Query 1

$$X^T X =$$

The diagram illustrates the computation of the matrix product  $X^T X$ . The first matrix is a  $p \times n$  matrix with columns  $x_1, \dots, x_n$ . The second matrix is an  $n \times n$  matrix with rows  $x_1, x_2, \dots, x_n$ . The result is a  $p \times p$  matrix with element  $(i,j)$  highlighted in red.

The summation formula is shown as:

$$\sum_{k=1}^n X_{k,i} X_{k,j}$$

The result matrix is shown as a  $p \times p$  matrix with element  $(i,j)$  highlighted in red.

# Computing Query 1

□ Compute the row-wise some:

$$X^T X = \sum_{i=1}^n x_i x_i^T$$

$x_1$	$x_2$	$y$
1.1	2.7	3.6
4.2	3.2	7.5
9.8	9.2	17
...	...	...

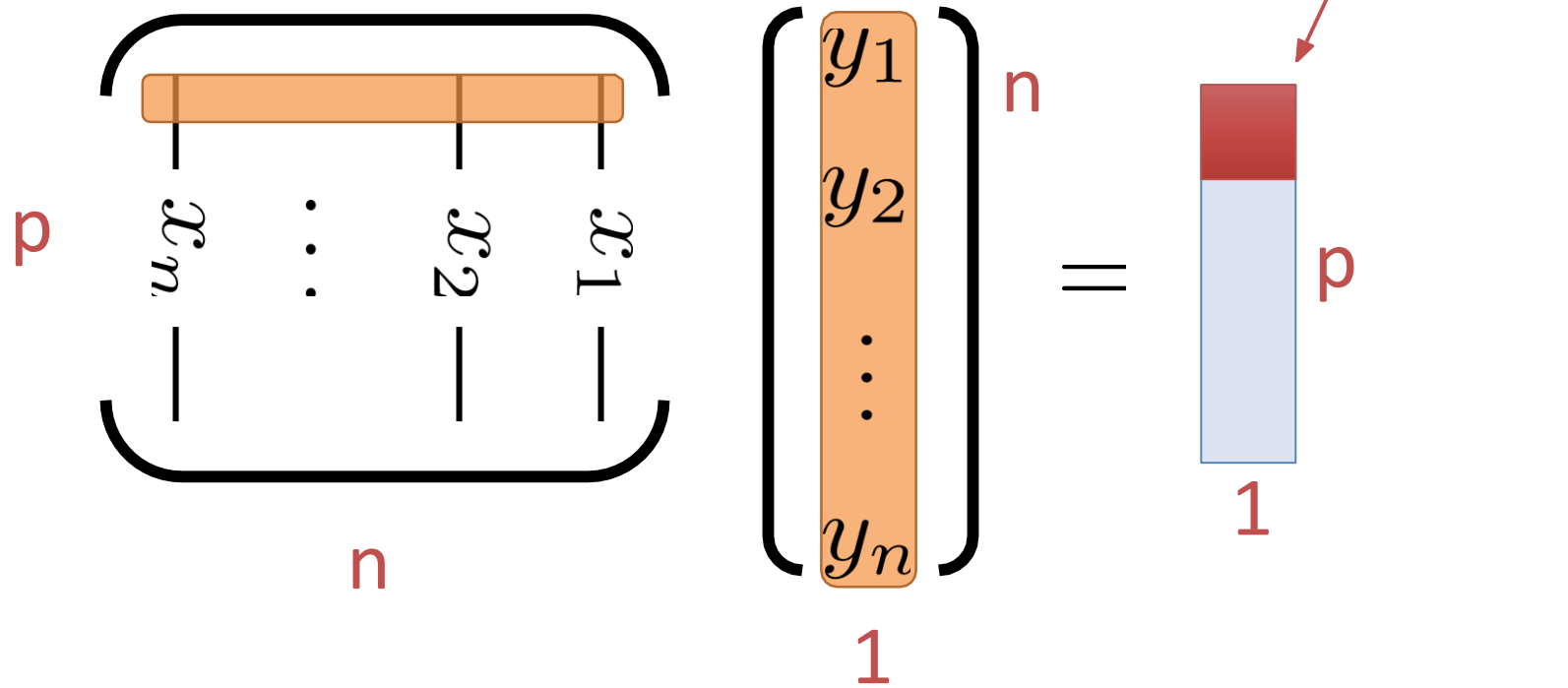
- **MapFunction(x)**: computes  $p$  by  $p$  outer product:  $x x^T$
- **ReduceFunction**: matrix sum:

□ Pure SQL Expression:

```
SELECT
    sum(x1*x1) AS c11, sum(x1*x2) AS c12,
    sum(x2*x1) AS c21, sum(x2*x2) AS c22
FROM data
```

# Computing Query 2

$$X^T y =$$



# Computing Query 2

$$X^T y =$$

The diagram illustrates the computation of the dot product  $X^T y$ . It shows a matrix  $X$  of size  $p \times n$  (with  $p$  rows and  $n$  columns) and a vector  $y$  of size  $n \times 1$ . The matrix  $X$  is represented by a large bracket with  $n$  columns and  $p$  rows. A specific row  $i$  is highlighted in orange, with elements  $x_{i1}, x_{i2}, \dots, x_{in}$ . The vector  $y$  is represented by a large bracket with  $n$  rows and  $1$  column, with elements  $y_1, y_2, \dots, y_n$ . The dot product of row  $i$  and column  $y$  is shown as a result of the matrix multiplication, represented by a vertical bar of size  $1 \times p$  with a red segment of size  $p$ . The sum  $\sum_{k=1}^n X_{k,i} y_k$  is indicated by a red arrow pointing to the red segment.



# Computing Query 2

- Compute the row-wise sum:

$$X^T y = \sum_{i=1}^n x_i y$$

$x_1$	$x_2$	$y$
1.1	2.7	3.6
4.2	3.2	7.5
9.8	9.2	17
...	...	...

- **MapFunction(x)**: computes  $p$  by 1 vector:  $x y$
- **ReduceFunction**: vector sum

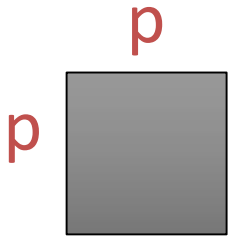
- Pure SQL Expression:


```
SELECT
    sum( $x_1$ * $y$ ) AS  $d_1$ , sum( $x_2$ * $y$ ) AS  $d_2$ 
FROM data
```

# Least Squares Regression using the Statistical Query Pattern

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

□ In database compute sums:

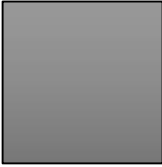
  $C = X^T X = \sum_{i=1}^n x_i x_i^T$   $O(np^2)$

  $d = X^T y = \sum_{i=1}^n x_i y_i$   $O(np)$

□ On client compute:

$$\hat{\theta} = C^{-1} d \quad O(p^3)$$

# What if $p$ is large?

$p$    $C = X^T X = \sum_{i=1}^n x_i x_i^T \quad O(np^2)$

... could be expensive ...

$$\hat{\theta} = C^{-1}d \quad O(p^3)$$

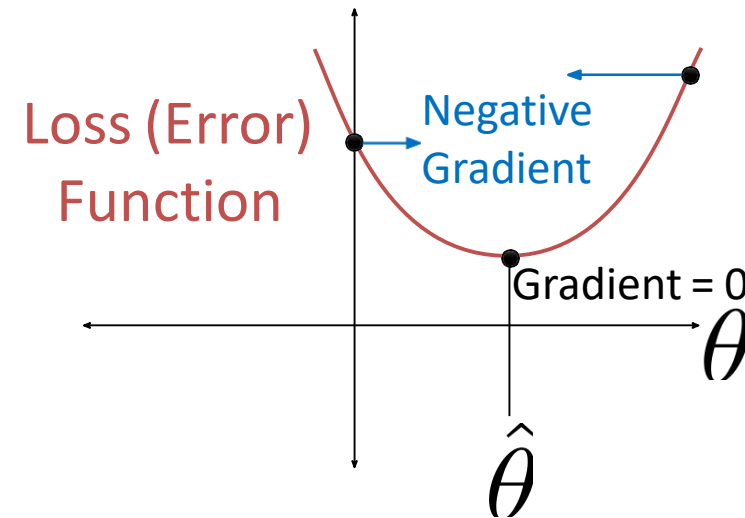
- Rather than directly solving:

$$\hat{\theta} = \arg \min_{\theta \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n L(y_i, \theta^T x_i)$$

- Instead we compute the gradient of the loss:

$$\begin{aligned} G(\theta; X, y) &= \nabla_{\theta} \frac{1}{n} \sum_{i=1}^n L(y_i, \theta^T x_i) = \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} L(y_i, \theta^T x_i) \\ &= \frac{1}{n} \sum_{i=1}^n (y - \theta^T x_i) x_i \end{aligned}$$

- **Big Idea:** Negative gradient points in the direction of steepest descent



# Gradient Descent Algorithm

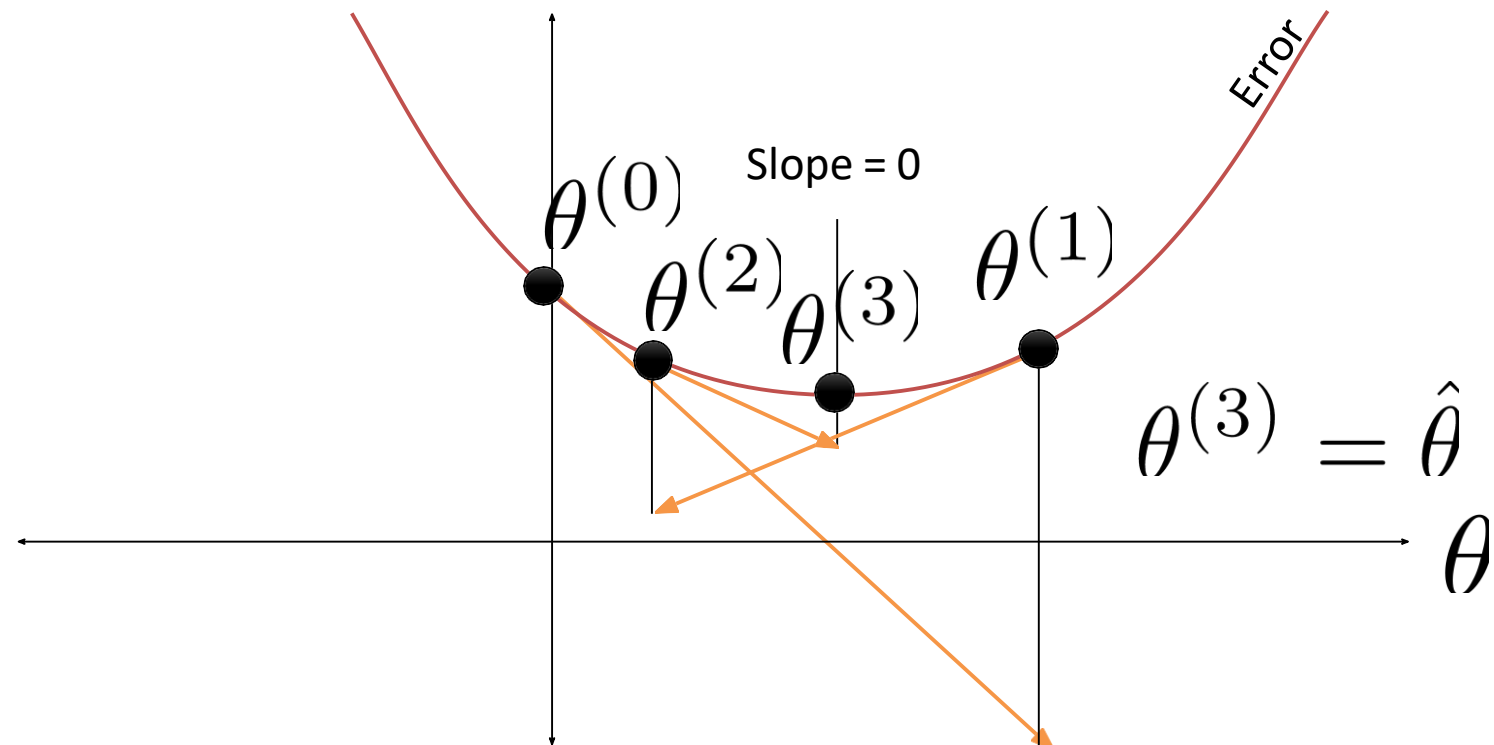
$t \leftarrow 0$

$\theta^{(0)} \leftarrow \text{Vec}(\theta)$

while (not converged):

$\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta * G(\theta; X, Y)$

$t \leftarrow t + 1$



# Gradient Descent Algorithm

$t \leftarrow 0$

$\theta^{(0)} \leftarrow \text{Vec}(\theta)$


while (not converged):

$\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta * G(\theta; X, y)$

$t \leftarrow t + 1$

□ Does this fit the statistical query pattern

- Yes! Only dependence on data is:

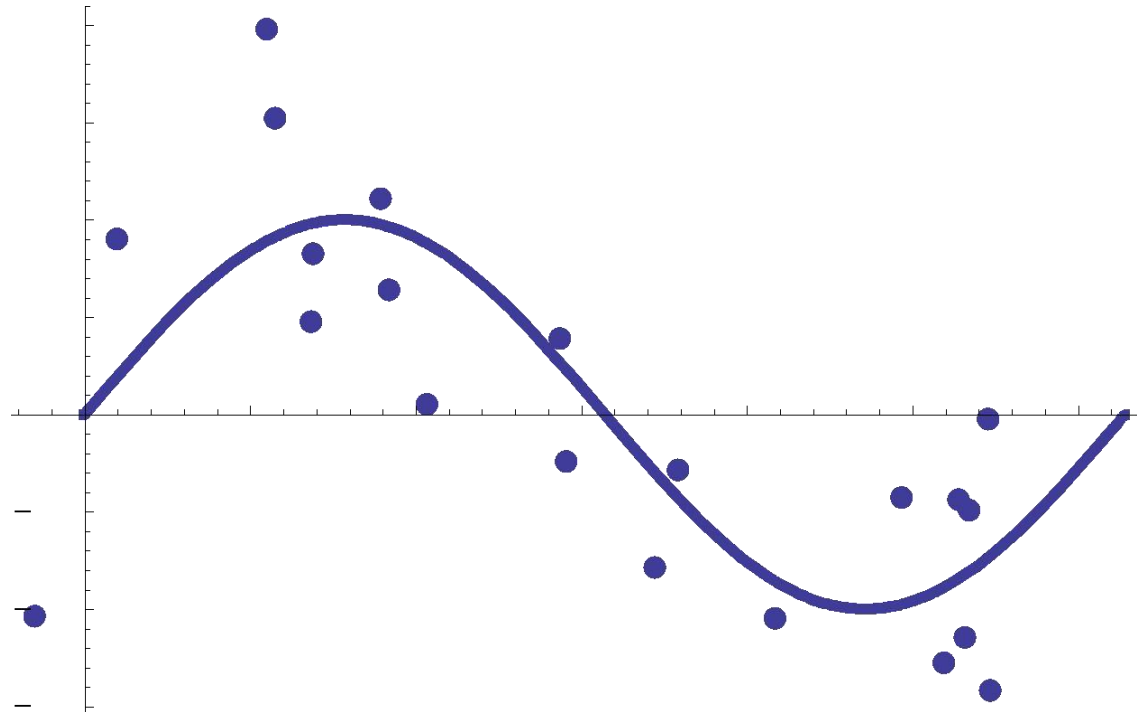
Size?   $G(\theta; X, y) = \frac{1}{n} \sum_{i=1}^n (y - \theta^T x_i) x_i$  Complexity?  $O(np)$

- Can we go even faster?

- **Stochastic Gradient Descent (SGD)**: Approximate the gradient by sampling data (typically several hundred records per query).

# Fitting Non-linear Data

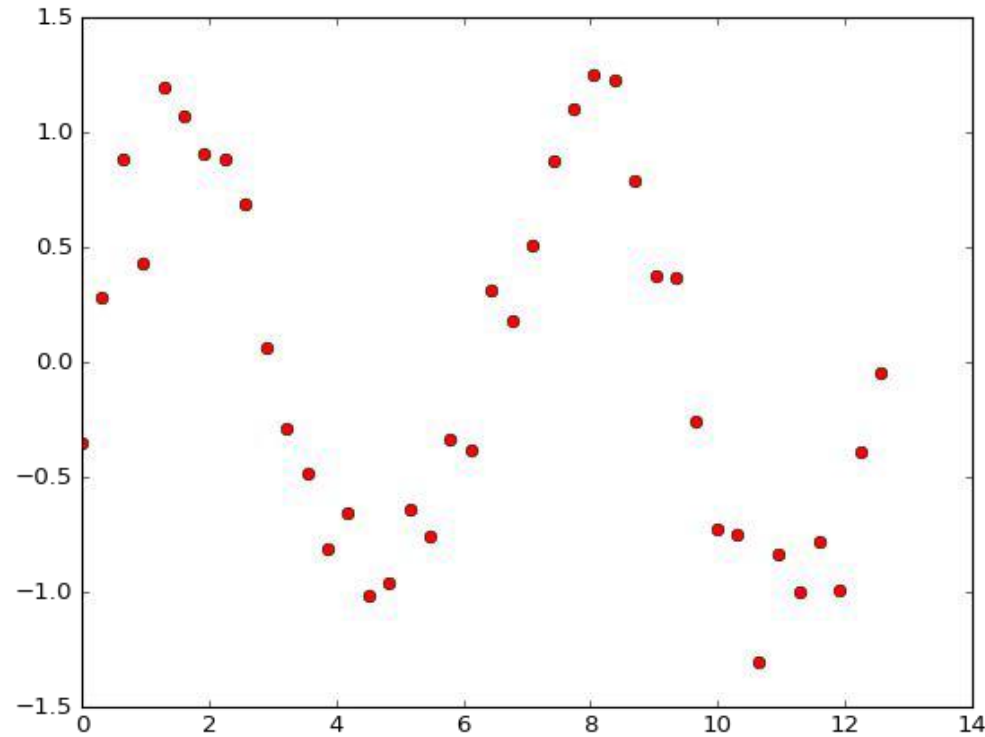
□ What if Y has a non-linear response?



□ Can we still use a linear model?

# Fitting Non-linear Data

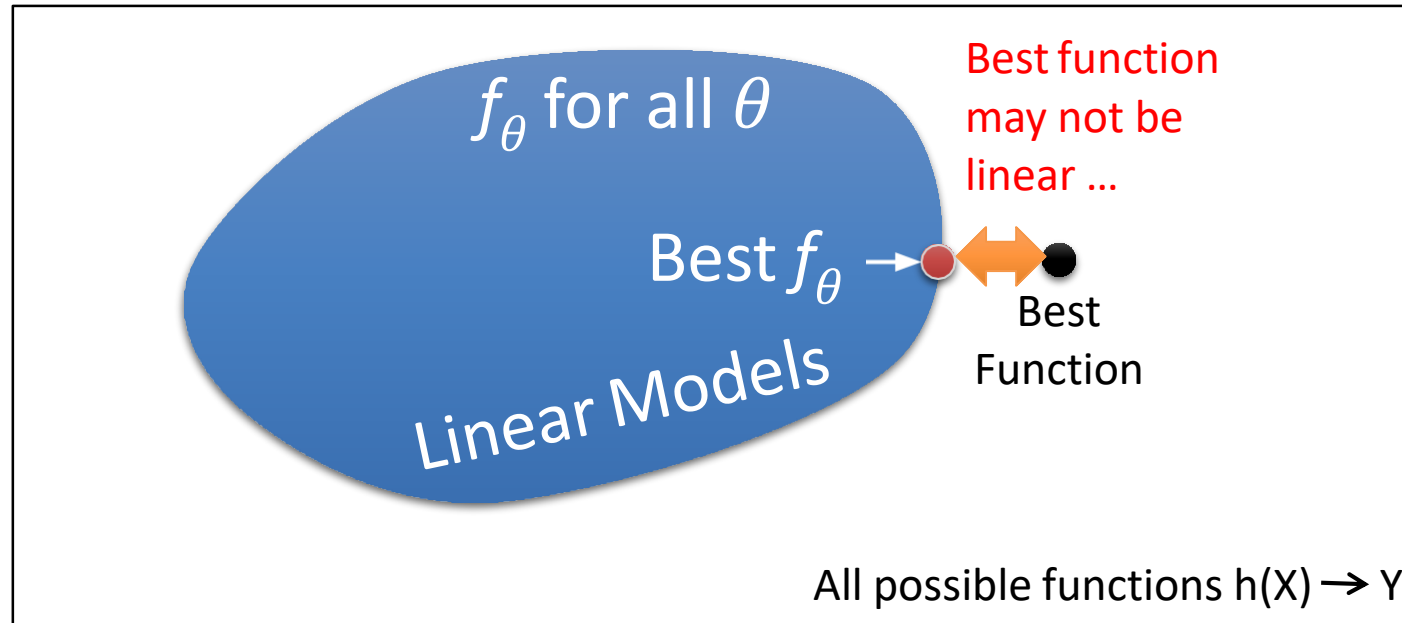
□ What if Y has a non-linear response?



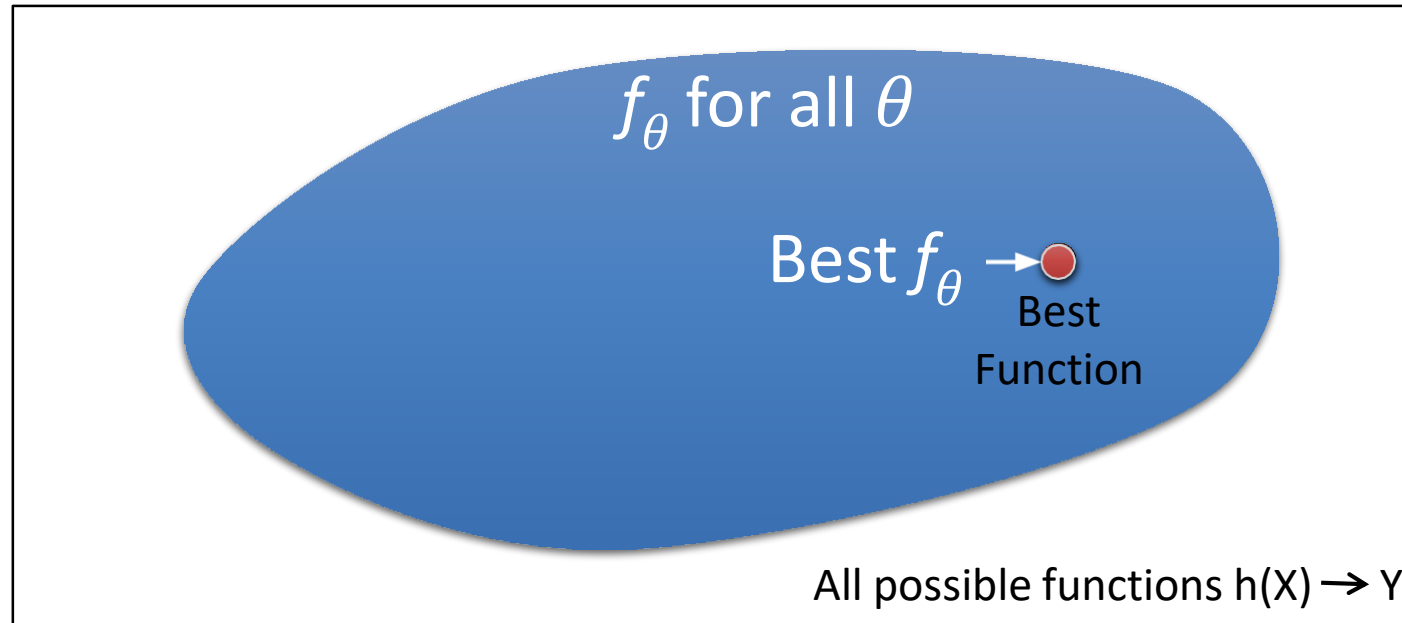
□ Can we still use a linear model?



# Finding the Best Parameters



# Finding the Best Parameters



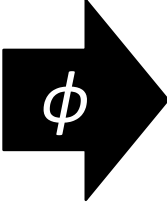
# Feature Engineering

□ By applying non-linear transformation  $\phi$ :

$$\phi : \mathbb{R}^p \rightarrow \mathbb{R}^k$$

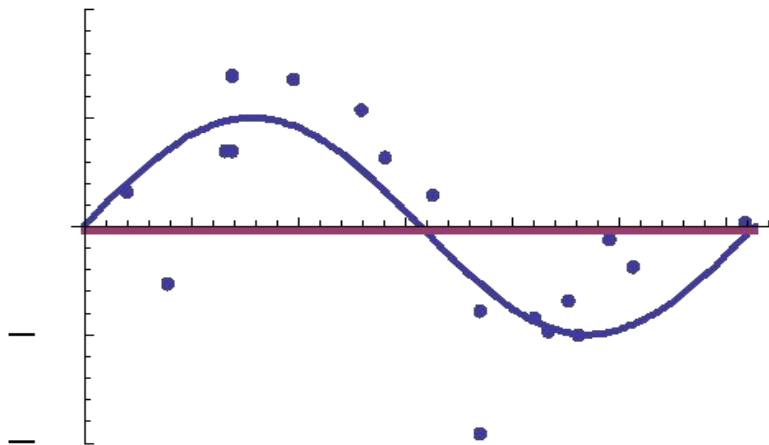
- Example:

$$\phi(x) = \{1, x, x^2, \dots, x^k\}$$

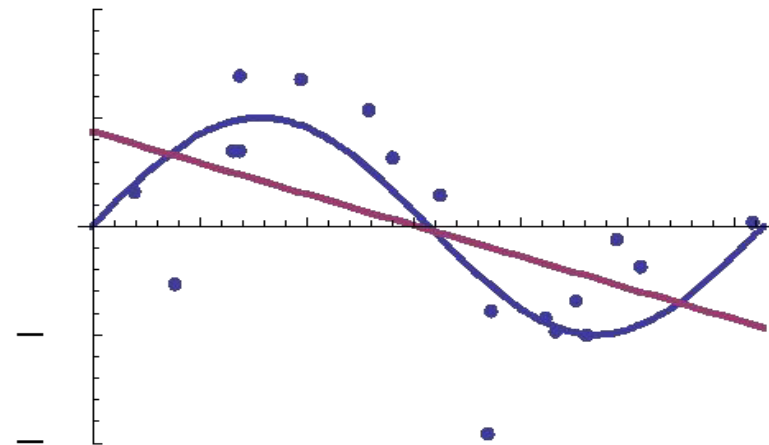
$x_1$	$x_2$	$y$		$x_1$	$x_2$	$x_1 * x_1$	$x_2 * x_2$	$x_1 * x_2$	$y$
1.1	2.7	3.6		1.1	2.7	1.21	7.29	2.97	3.6
4.2	3.2	7.5		4.2	3.2	17.64	10.24	13.44	7.5
9.8	9.2	17		9.8	9.2	90.04	84.64	90.16	17

Under-fitting

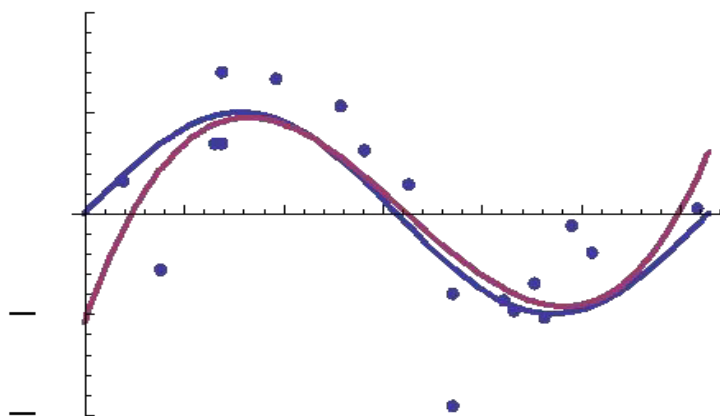
$$\{ 1 \}$$



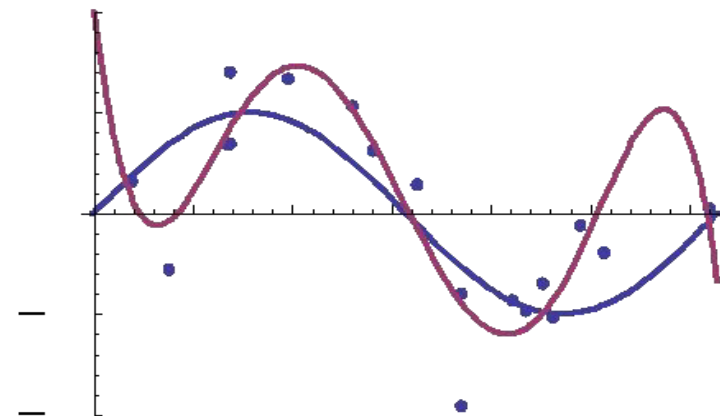
$$\{ 1, x \}$$



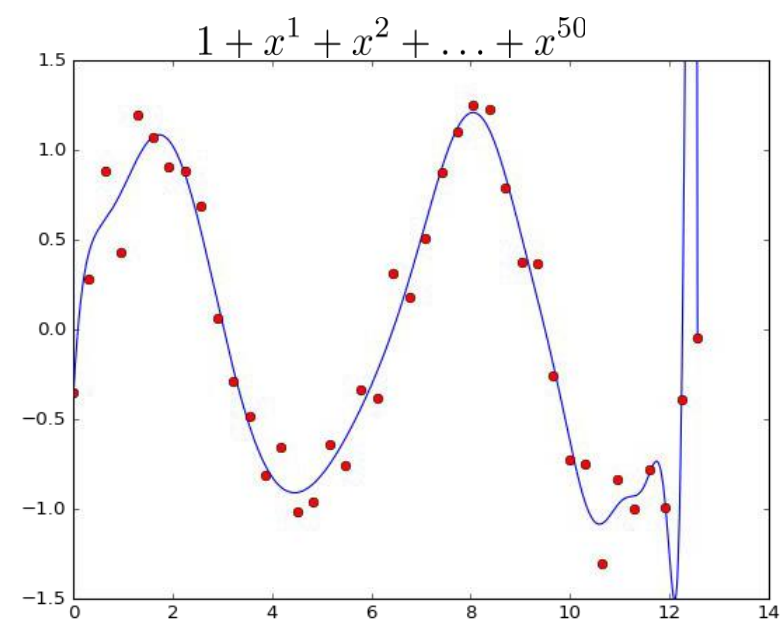
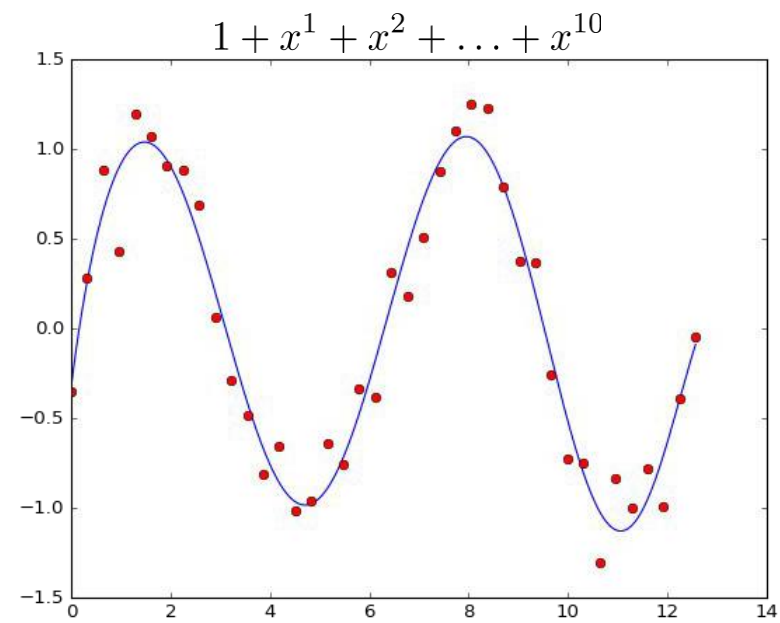
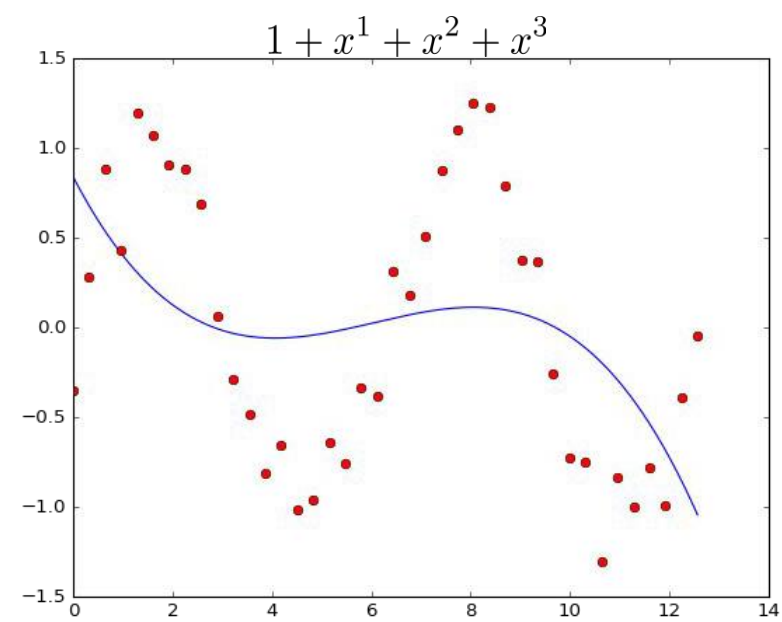
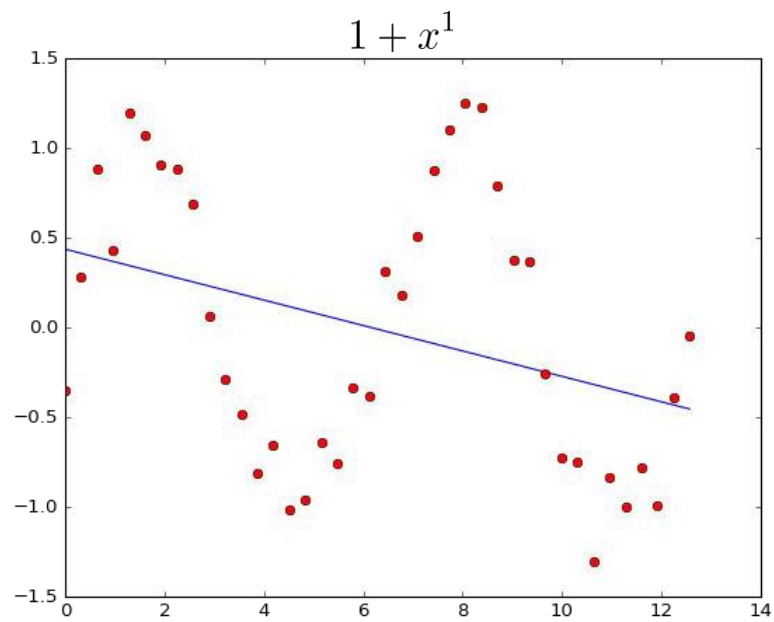
$$\{ 1, x, x^2, x^3 \}$$



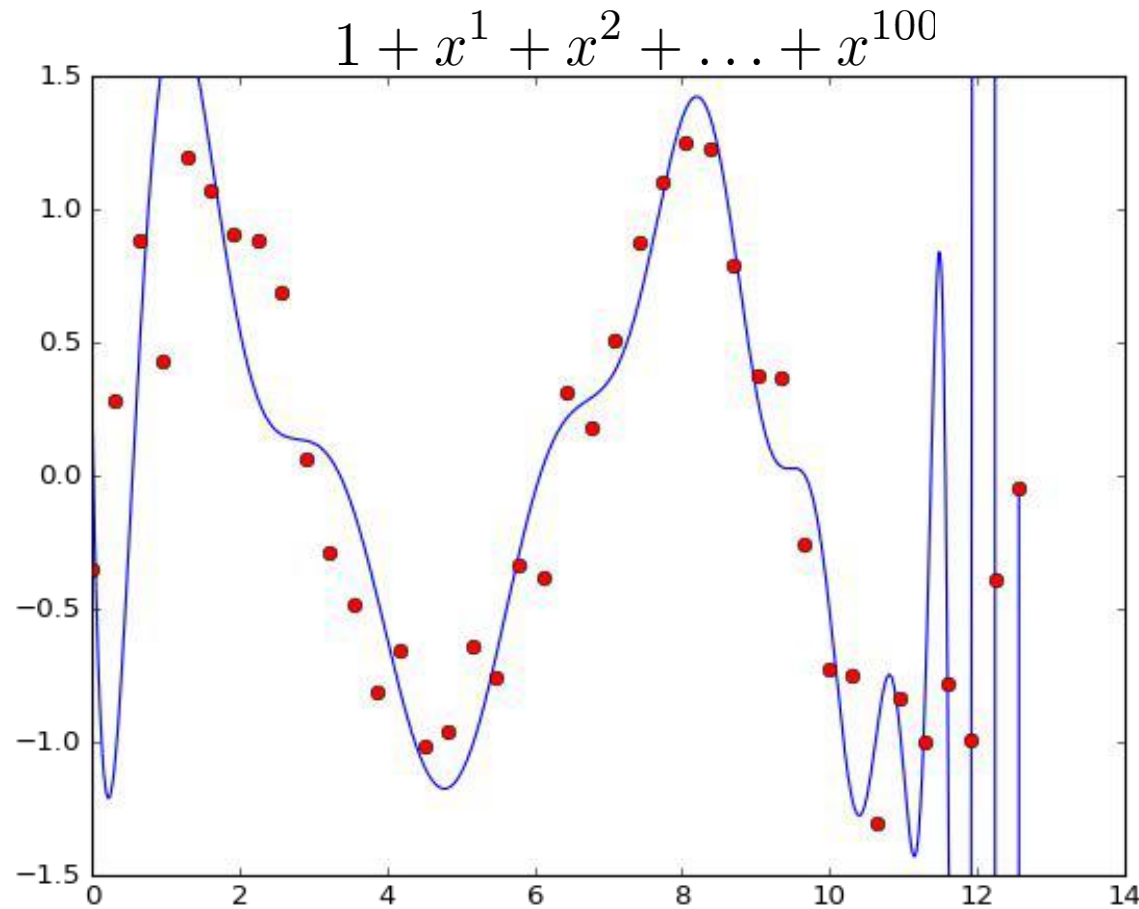
$$\{ 1, x, x^2, x^3, x^4, x^5 \}$$



Over-fitting

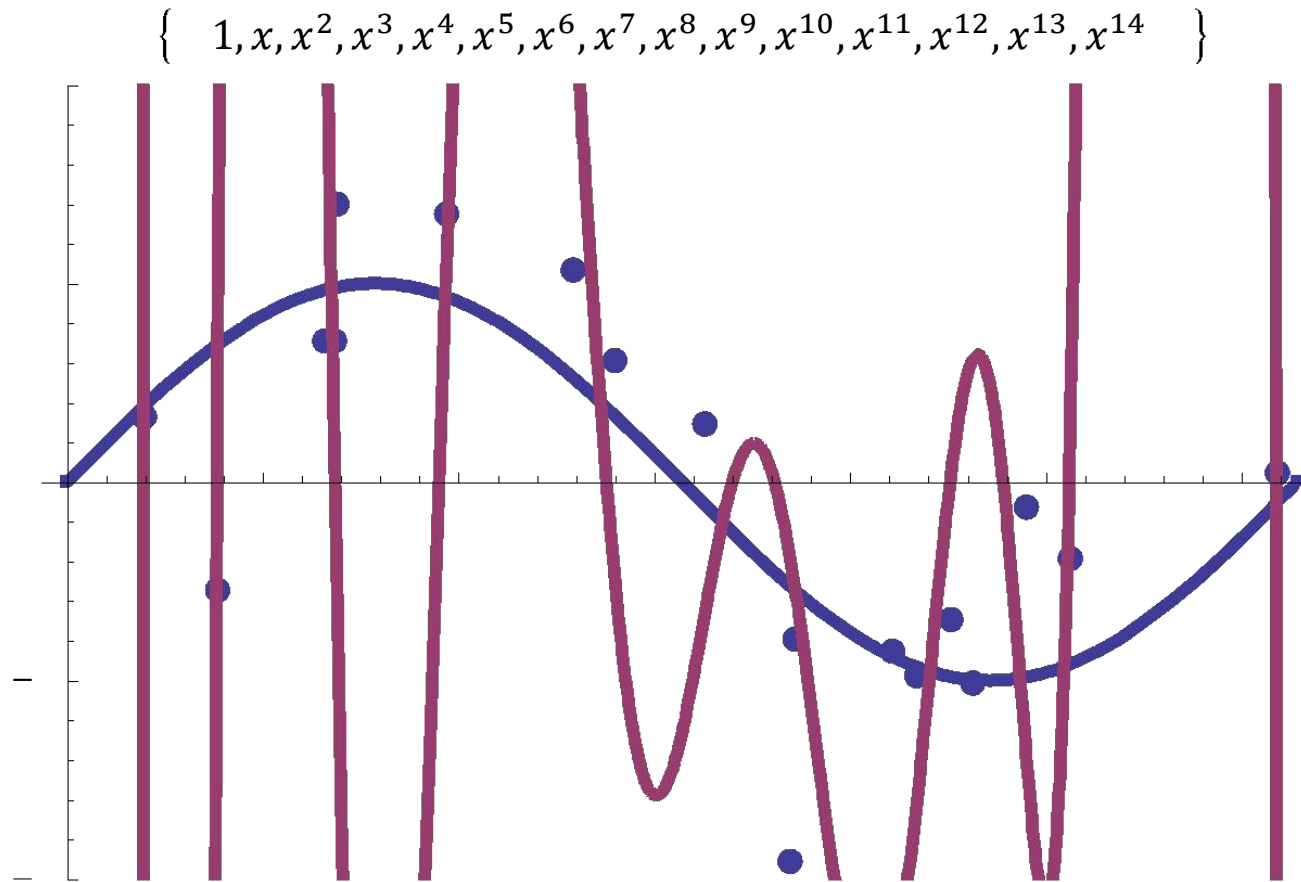


# Over-fitting!



- Errors on training data are small
- But errors on new points are likely to be large

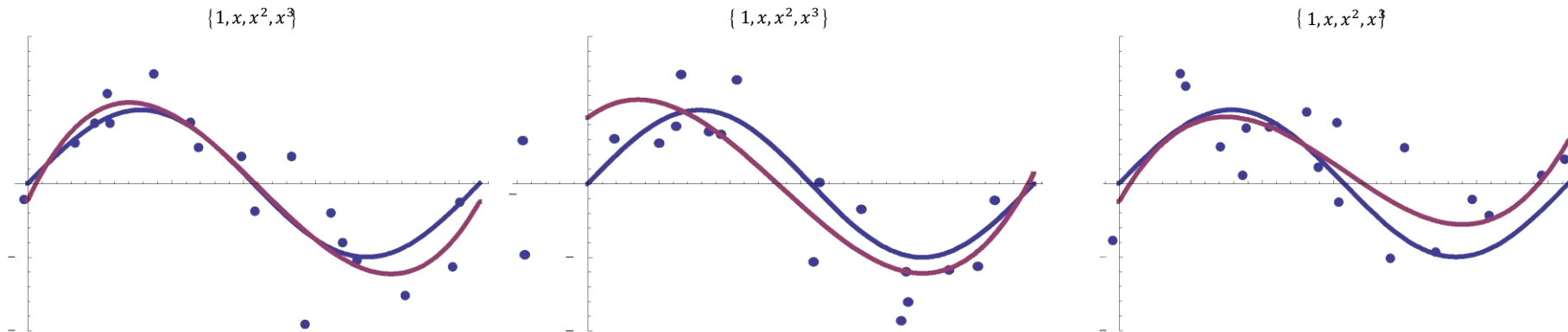
# Really Over-fitting!



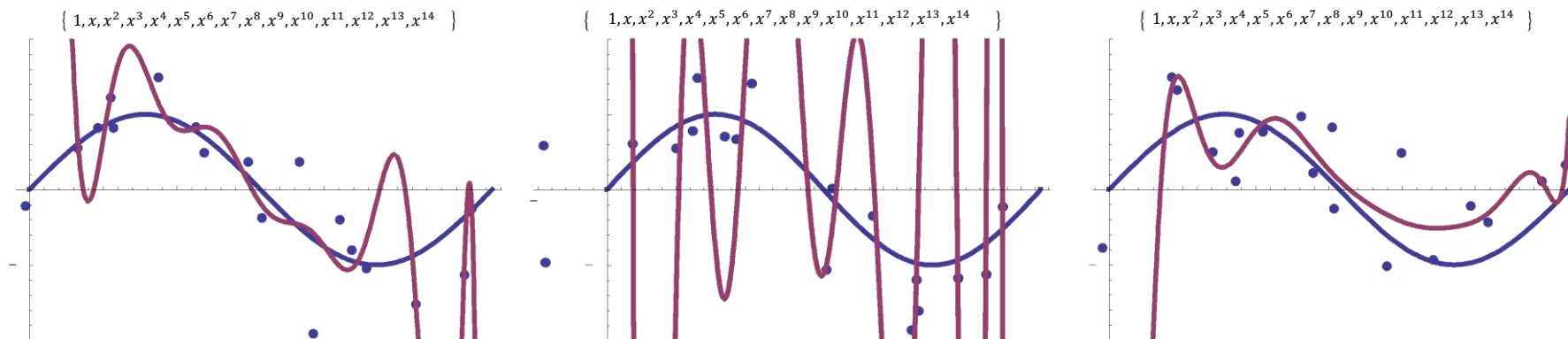
- Errors on training data are small
- But errors on new points are likely to be large

# What if I train on different data?

Simple Model  $\rightarrow$  Low Variability



Complex Model  $\rightarrow$  High Variability





# Bias-Variance Tradeoff

- So far we have minimized the **training error**: the error on the training data.
  - low training error does not guarantee good expected performance (due to over-fitting)
- We would like to reason about the **test error**

## Theorem:

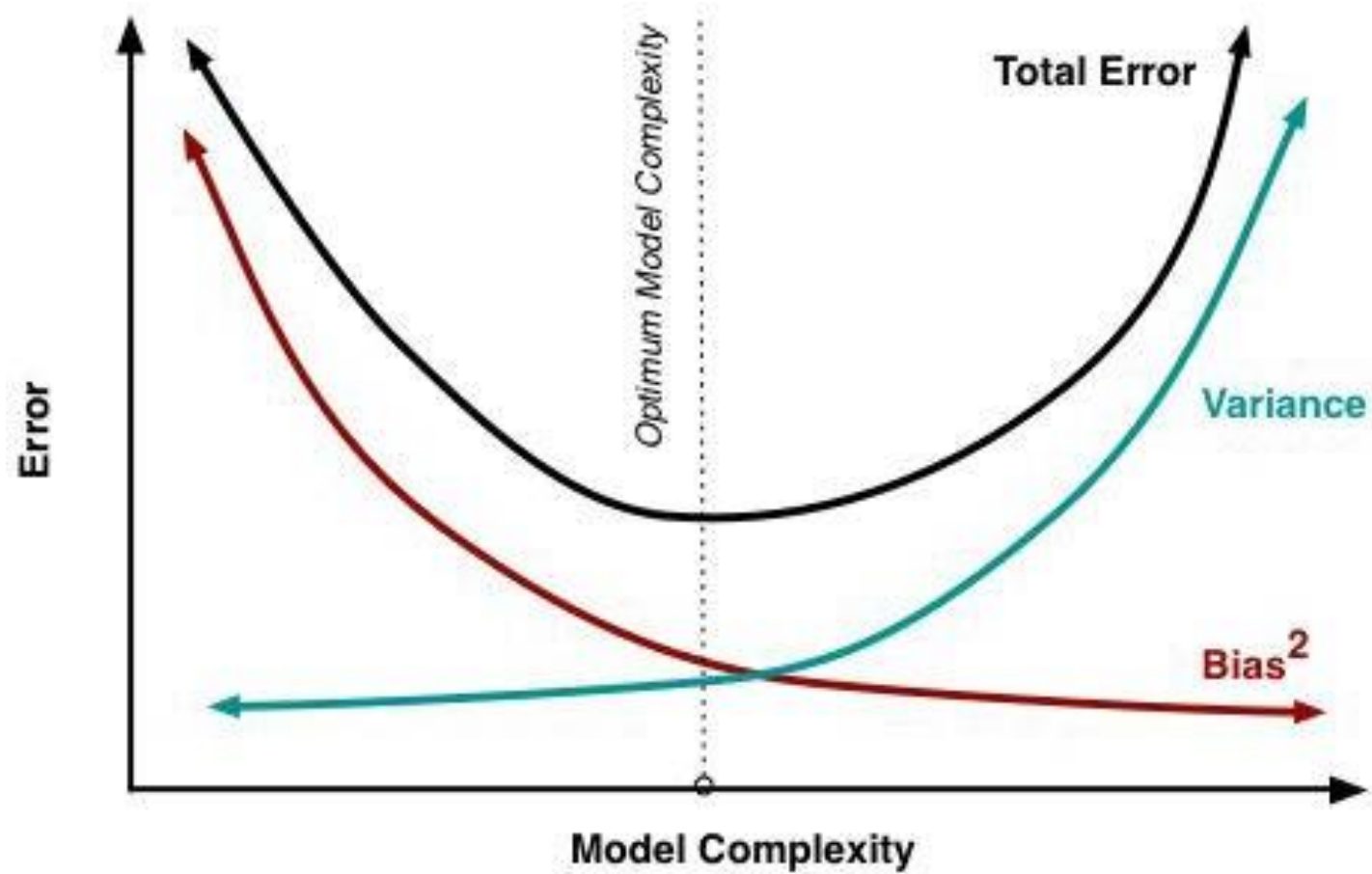
$$\text{Test Error} = \text{Noise} + \text{Bias}^2 + \text{Variance}$$

Noisy data has inherent error (measurement error)

Error due to models inability to fit the data. (**Under Fitting**)

Error due to inability to estimate model parameters. (**Over-fitting**)

# Bias Variance Plot



# Regularization to Reduce Over-fitting

- High dimensional models tend to over-fit
  - How could we “favor” lower dimensional models?

- **Solution Intuition:**

- Too many features over-fitting

$$f(x) = \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_d x_d$$

- Force many of the  $\theta_i \approx 0$  (e.g.,  $i > 2$ ) (“effectively fewer features”)

$$\begin{aligned} f(x) &= \theta_1 x_1 + \theta_2 x_2 + 0x_3 + \dots + 0x_d \\ &= \theta_1 x_1 + \theta_2 x_2 \end{aligned}$$

- Keeping weights close to zero **reduces variance**

# Regularization to Reduce Over-fitting

- We can add a regularization term:

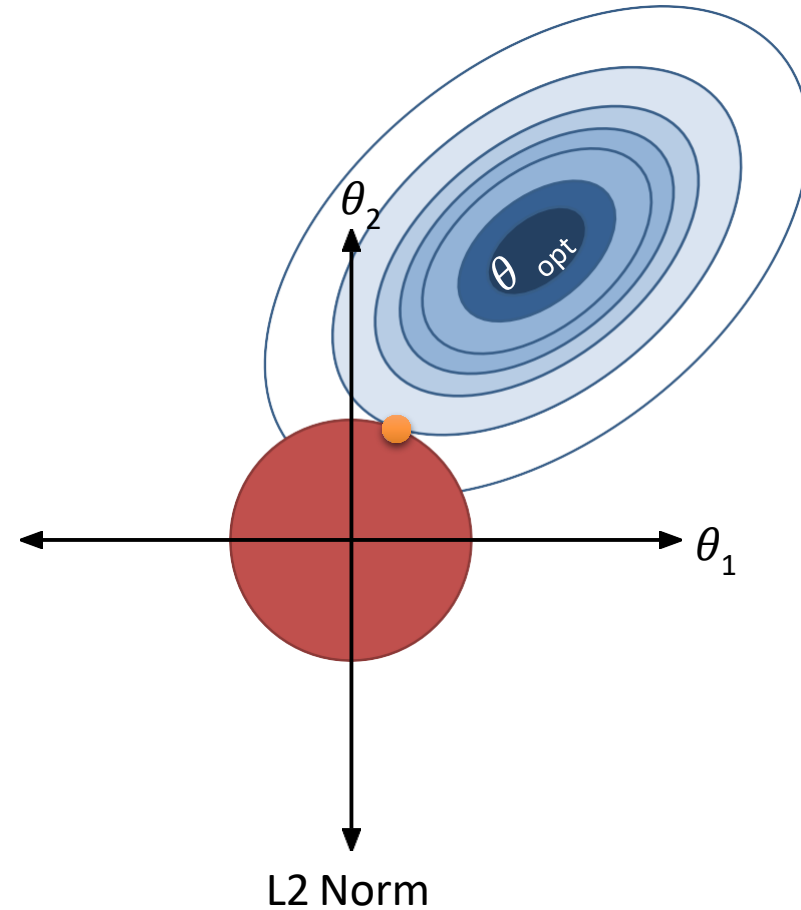
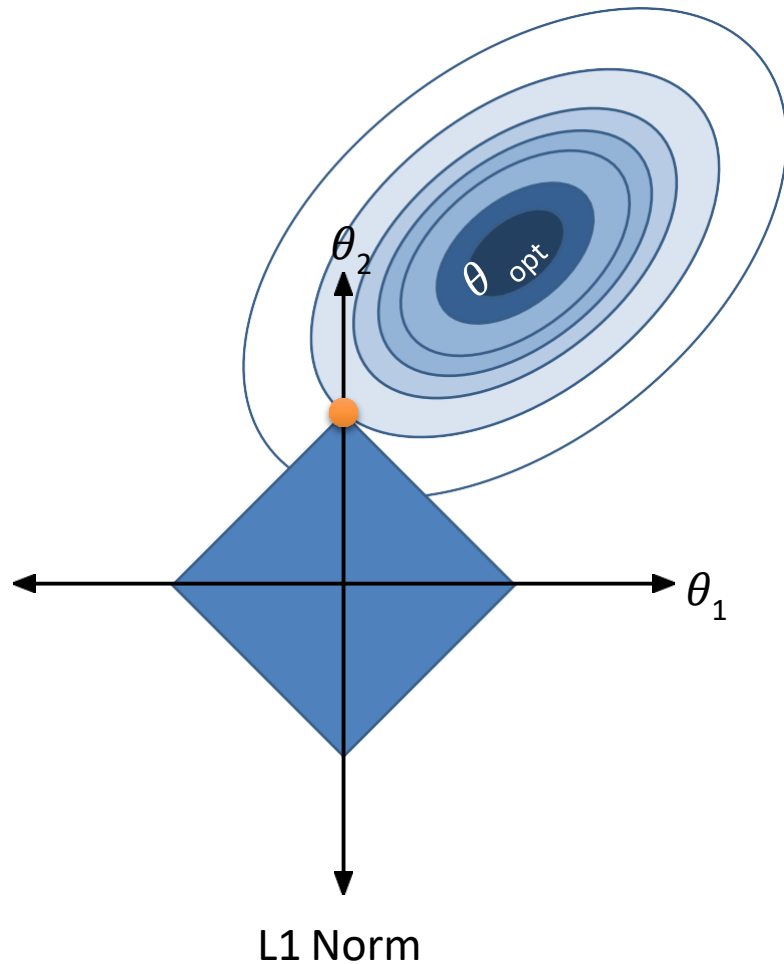
$$\hat{\theta} = \arg \min_{\theta \in \mathbb{R}^p} \quad \frac{1}{n} \sum_{i=1}^n \overbrace{(y_i - \theta^T x_i)^2}^{\text{Squared Loss}} + \underbrace{\lambda R(\theta)}_{\substack{\text{Regularization} \\ \text{Function} \\ \text{Parameter}}}$$

- Common of Regularization Functions:

$$\begin{array}{ll} \text{Ridge (L2-Reg)} & R_{\text{Ridge}}(\theta) = \sum_{i=1}^d \theta_i^2 \\ \text{Regression} & \end{array} \quad \begin{array}{ll} \text{Lasso} & R_{\text{Lasso}}(\theta) = \sum_{i=1}^d |\theta_i| \\ \text{(L1-Reg)} & \end{array}$$

- Encourage small parameter values
- The parameter  $\lambda$  determines amount of reg.
  - Larger  $\rightarrow$  more reg.  $\rightarrow$  lower variance  $\rightarrow$  higher bias

# Regularization and Norm Balls



# Regularization to Reduce Over-fitting

- We can add a regularization term:

$$\hat{\theta} = \arg \min_{\theta \in \mathbb{R}^p} \quad \frac{1}{n} \sum_{i=1}^n \overbrace{(y_i - \theta^T x_i)^2}^{\text{Squared Loss}} + \underbrace{\lambda R(\theta)}_{\substack{\text{Regularization} \\ \text{Function}}} \quad \underbrace{\lambda}_{\text{Regularization Parameter}}$$

- Solving the regularized problem:
  - Closed form solution for Ridge regression (L2):

$$\hat{\theta}_{\text{Ridge}} = (X^T X + \lambda I)^{-1} X^T Y$$

- Iterative methods for Lasso (L1):
  - Stochastic gradient ...

- How do we choose  $\lambda$ ?

# Picking The Regularization Parameter $\lambda$

- **Proposal:** Minimize **training** error

$$\arg \min_{\theta \in \mathbb{R}^p, \lambda \geq 0} \frac{1}{n} \sum_{i=1}^n (y_i - \theta^T x_i)^2 + \lambda R(\theta)$$

- Trivial solution  $\rightarrow \lambda = 0$

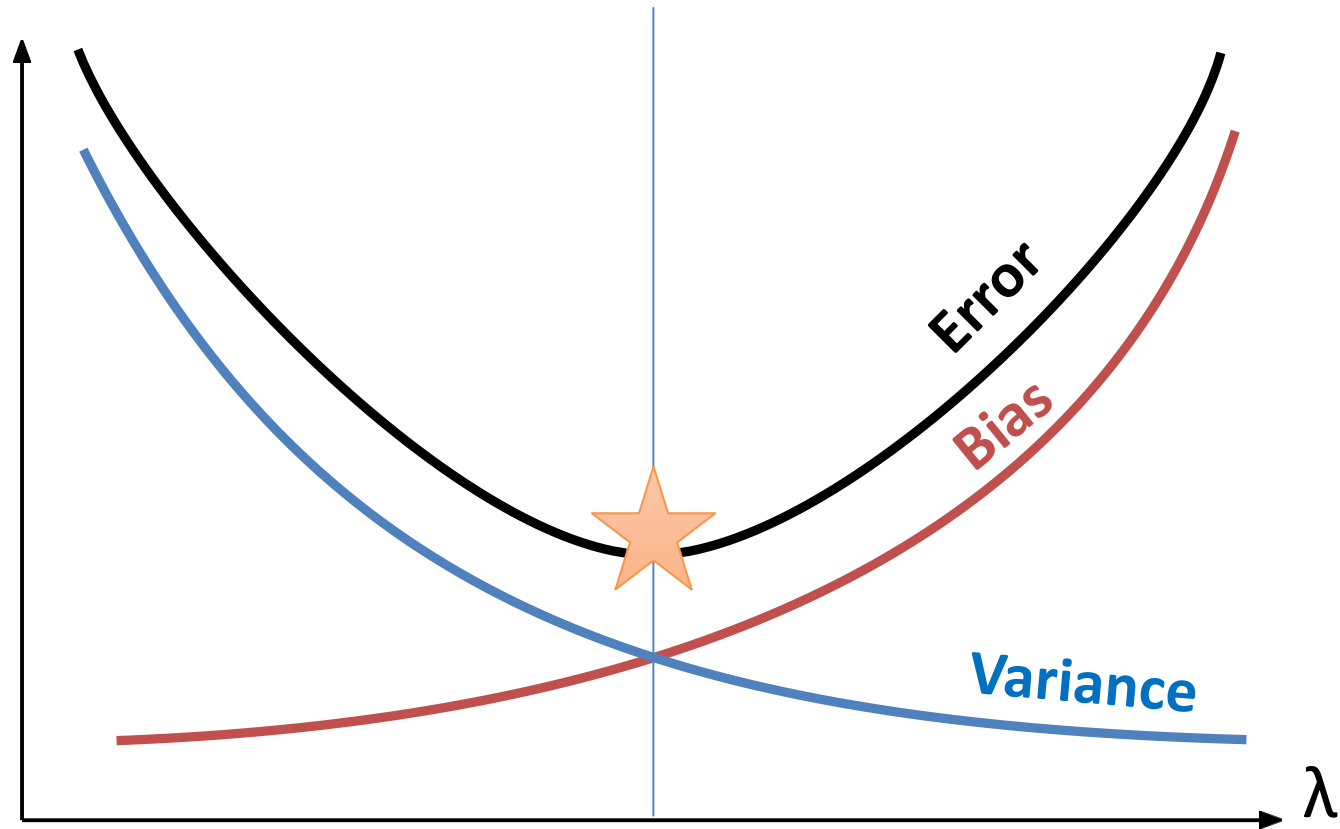
- Intuition we want to minimize **test** error

- **Test error:** error on unseen data

- **2<sup>nd</sup> Proposal:** Split training data into training and evaluation sets

- For a range of  $\lambda$  values compute optimal  $\theta_\lambda$  using only the reduced training set
- Evaluate  $\theta_\lambda$  on the separate evaluation set and select the  $\lambda$  with the lowest error

# Bias Variance Plot

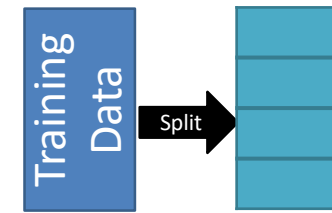


Optimal  $\lambda$

Larger  $\lambda \rightarrow$  Reduced  
Model Complexity



# K-Fold Cross Validation



- Split training data into K-equally sized parts
  - In practice K is relatively small (e.g., 5)
- For each part train on the other k-1 parts and compute the error on that part:



- Compute the average test error over held out parts
- Select reg. param. that minimizes average test error

# Regularization to Reduce Over-fitting

- High dimensional models tend to over-fit
  - How could we “favor” lower dimensional models?

- **Solution Intuition:**

- Too many features over-fitting

$$f(x) = \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_d x_d$$

- Force many of the  $\theta_i \approx 0$  (e.g.,  $i > 2$ ) (“effectively fewer features”)

$$\begin{aligned} f(x) &= \theta_1 x_1 + \theta_2 x_2 + 0x_3 + \dots + 0x_d \\ &= \theta_1 x_1 + \theta_2 x_2 \end{aligned}$$

- Keeping weights close to zero **reduces variance**

# Regularization to Reduce Over-fitting

- We can add a regularization term:

$$\hat{\theta} = \arg \min_{\theta \in \mathbb{R}^p} \quad \underbrace{\frac{1}{n} \sum_{i=1}^n (y_i - \theta^T x_i)^2}_{\text{Squared Loss}} + \underbrace{\lambda R(\theta)}_{\substack{\text{Regularization} \\ \text{Function} \\ \text{Regularization} \\ \text{Parameter}}}$$

- Common of Regularization Functions:

$$\begin{array}{ll} \text{Ridge (L2-Reg)} & R_{\text{Ridge}}(\theta) = \sum_{i=1}^d \theta_i^2 \\ \text{Regression} & \end{array} \quad \begin{array}{ll} \text{Lasso} & R_{\text{Lasso}}(\theta) = \sum_{i=1}^d |\theta_i| \\ \text{(L1-Reg)} & \end{array}$$

- Encourage small parameter values
- The parameter  $\lambda$  determines amount of reg.
  - Larger  $\rightarrow$  more reg.  $\rightarrow$  lower variance  $\rightarrow$  higher bias

# Summary of Regression



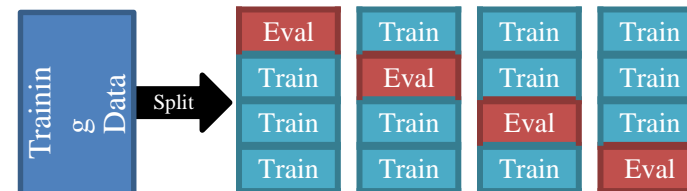
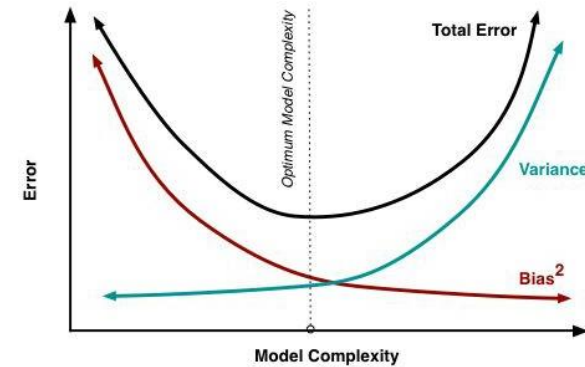
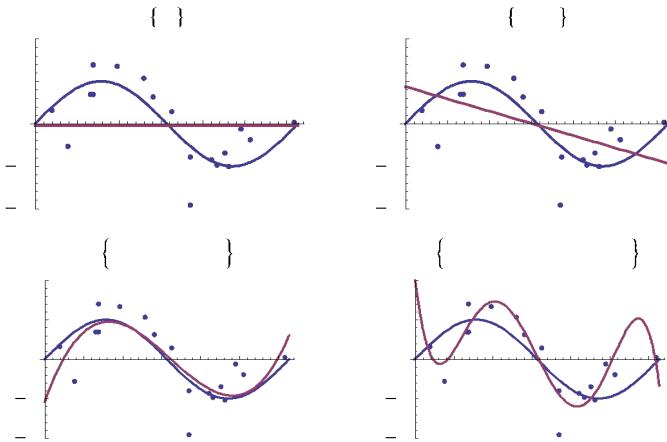
0 0 0 0 1  
 Sports Furniture Clothing Shoes Electronics

$p$

$$C = X^T X = \sum_{i=1}^n x_i x_i^T$$

1

$$d = X^T y = \sum_{i=1}^n x_i y_i$$





# Taxonomy of Machine Learning

