

CS150: Database & Datamining

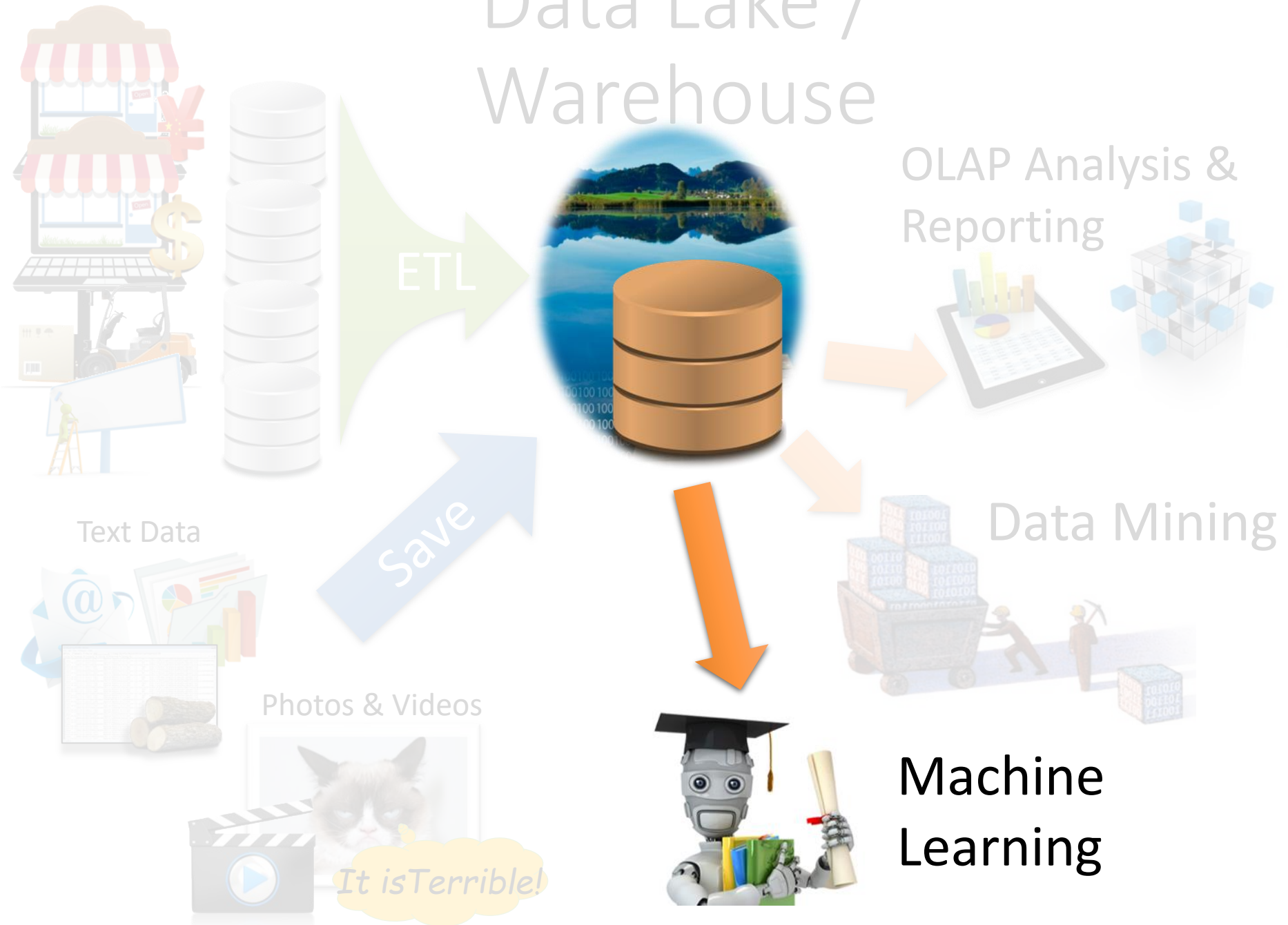
Lecture 21: Analytics & Machine Learning

III

Xuming He
Spring 2019

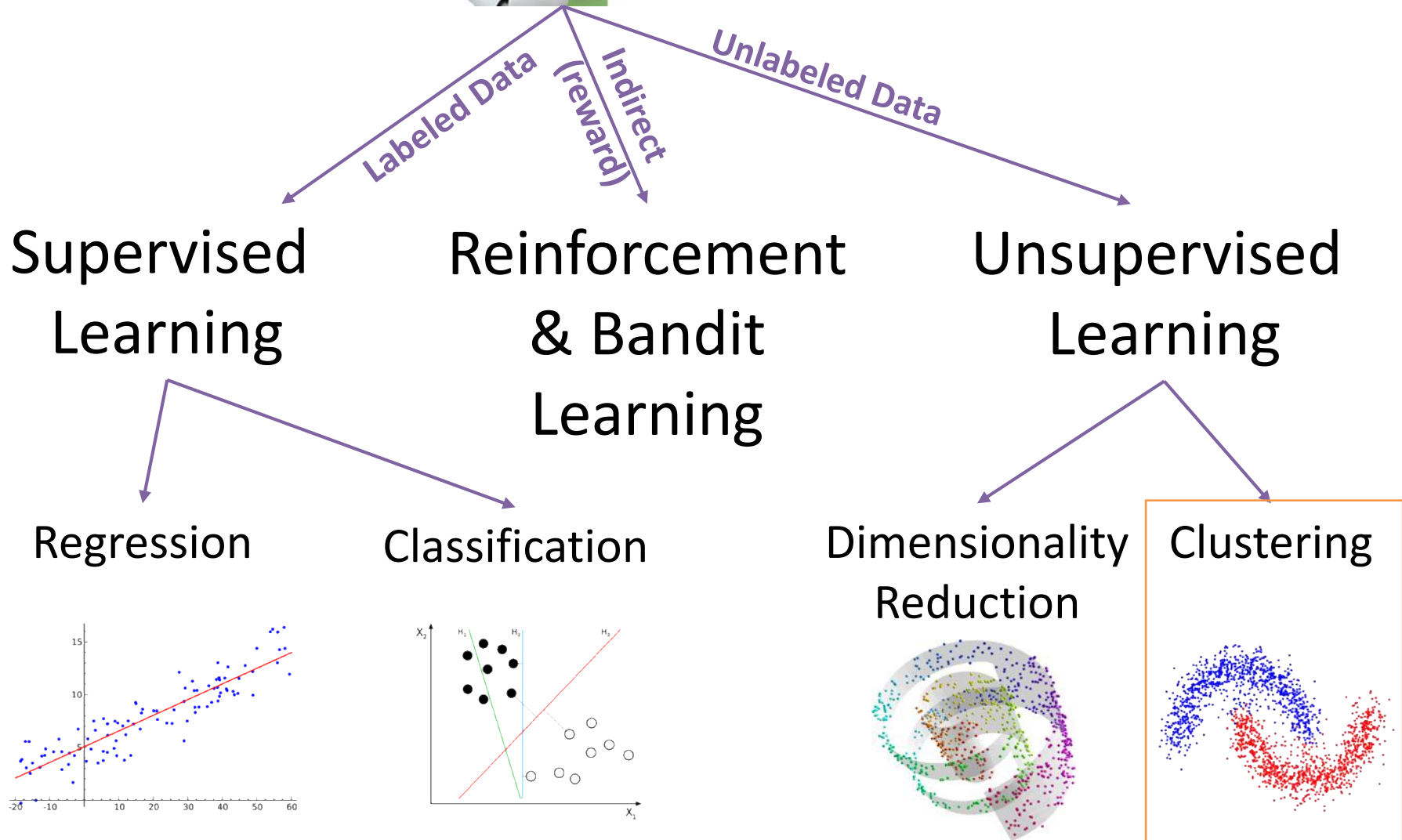
Acknowledgement: Slides are adopted from the Berkeley course CS186 by Joey Gonzalez and Joe Hellerstein, Stanford CS145 by Peter Bailis.

Data Lake / Warehouse



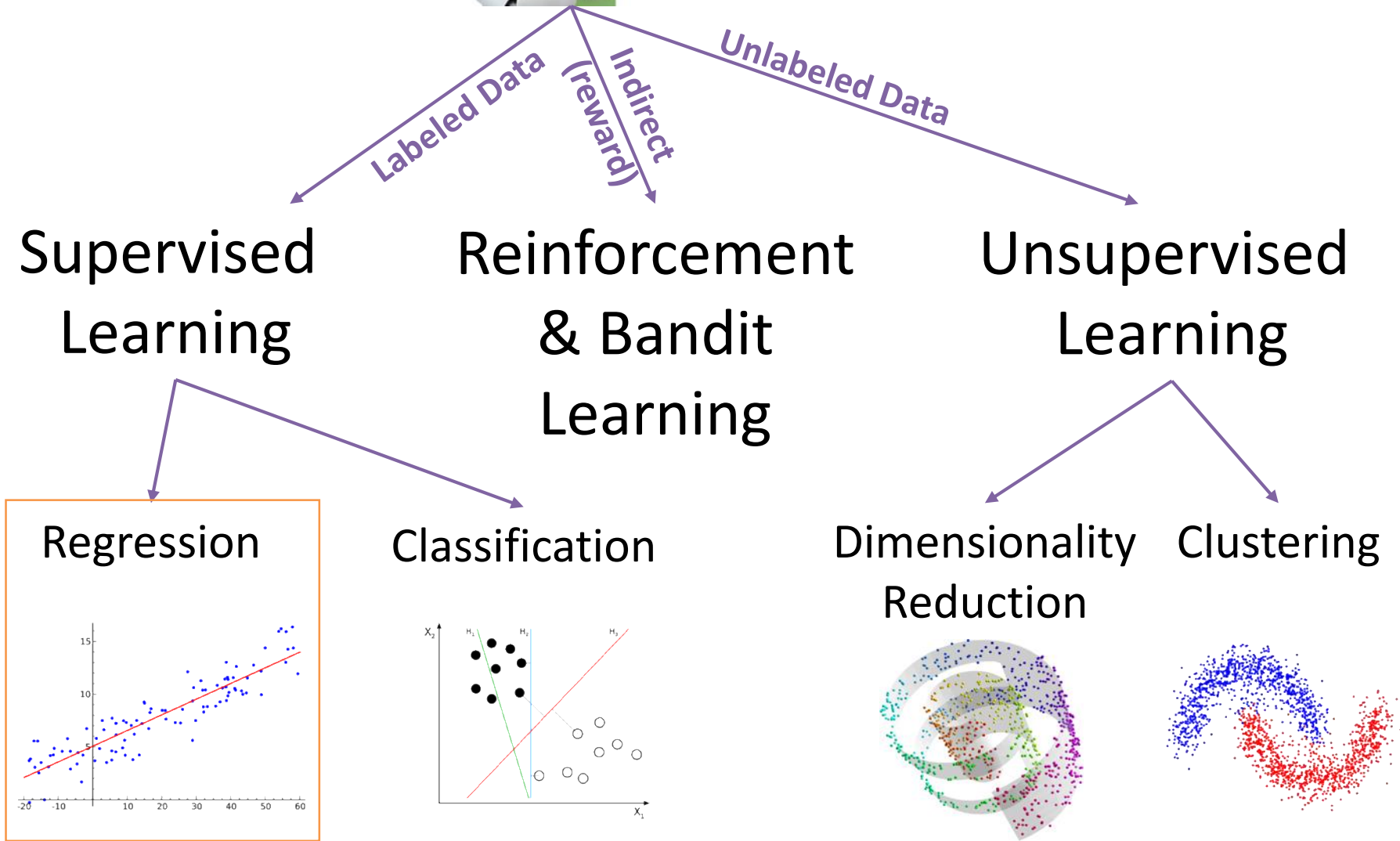


Taxonomy of Machine Learning

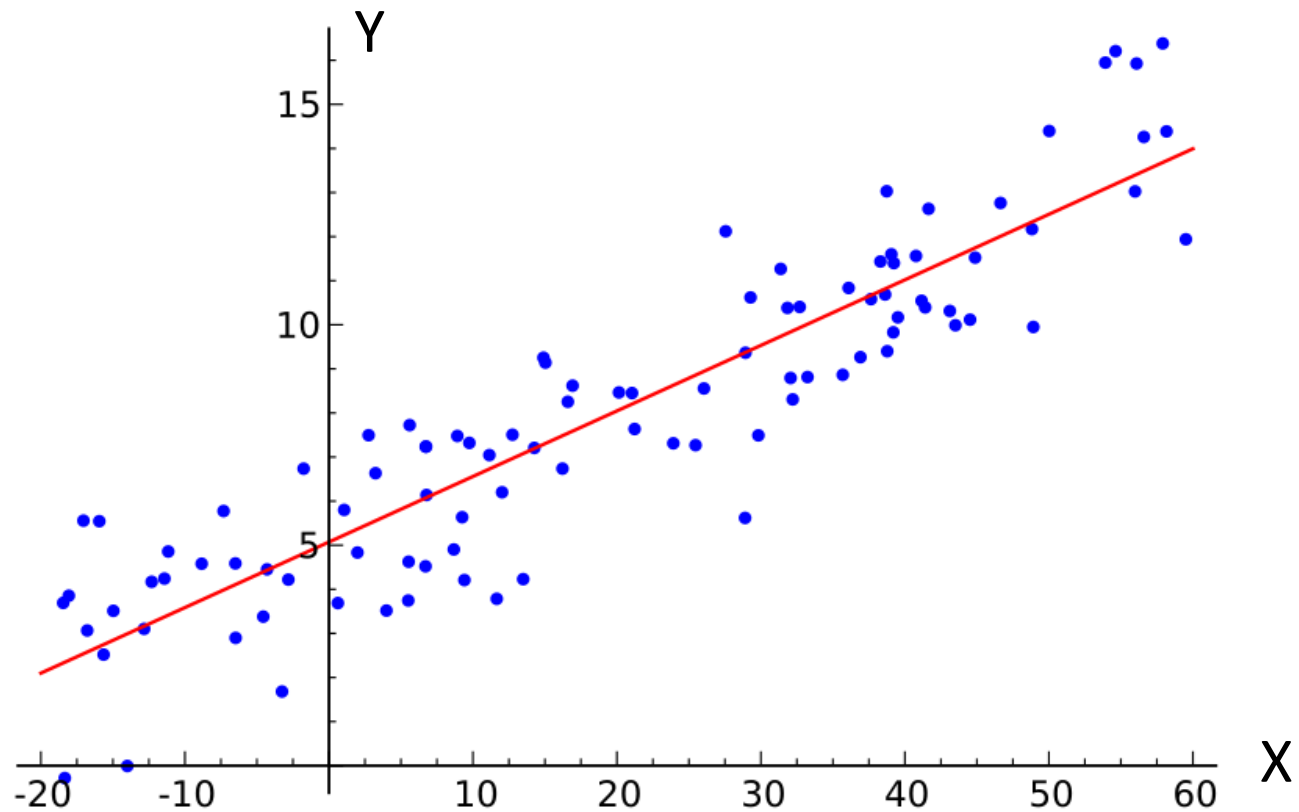




Taxonomy of Machine Learning



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

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

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

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
...

Vector of
Parameters

Vector of
Features

Real Valued
Observations

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

$y \in \mathbb{R}$ (Real Value)
 $\theta \in \mathbb{R}^p$ (Vector of Parameters)
 $x \in \mathbb{R}^p$ (Vector of Features)
 $\epsilon \in \mathbb{R}$ (Real Value Noise)

Noise Model:

$$\epsilon \sim \mathcal{N}(0, \sigma^2)$$

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

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

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?

“Learning about machine learning is fun.”



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

Vector

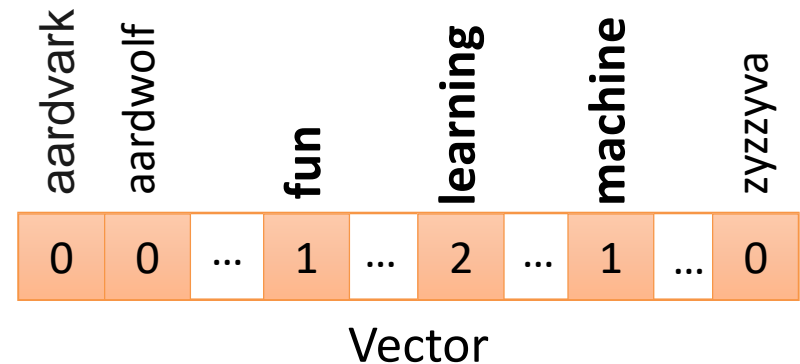
➤ Bag-of-words model

- 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”)

Working with Text Data

➤ Features: Transforming the data

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

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



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
...

$$f_{\theta}(x) := \overbrace{\theta^T}^{\text{Vector of Parameters}} \overbrace{x}^{\text{Vector of Features}}$$

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

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

Finding the Best Parameters

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

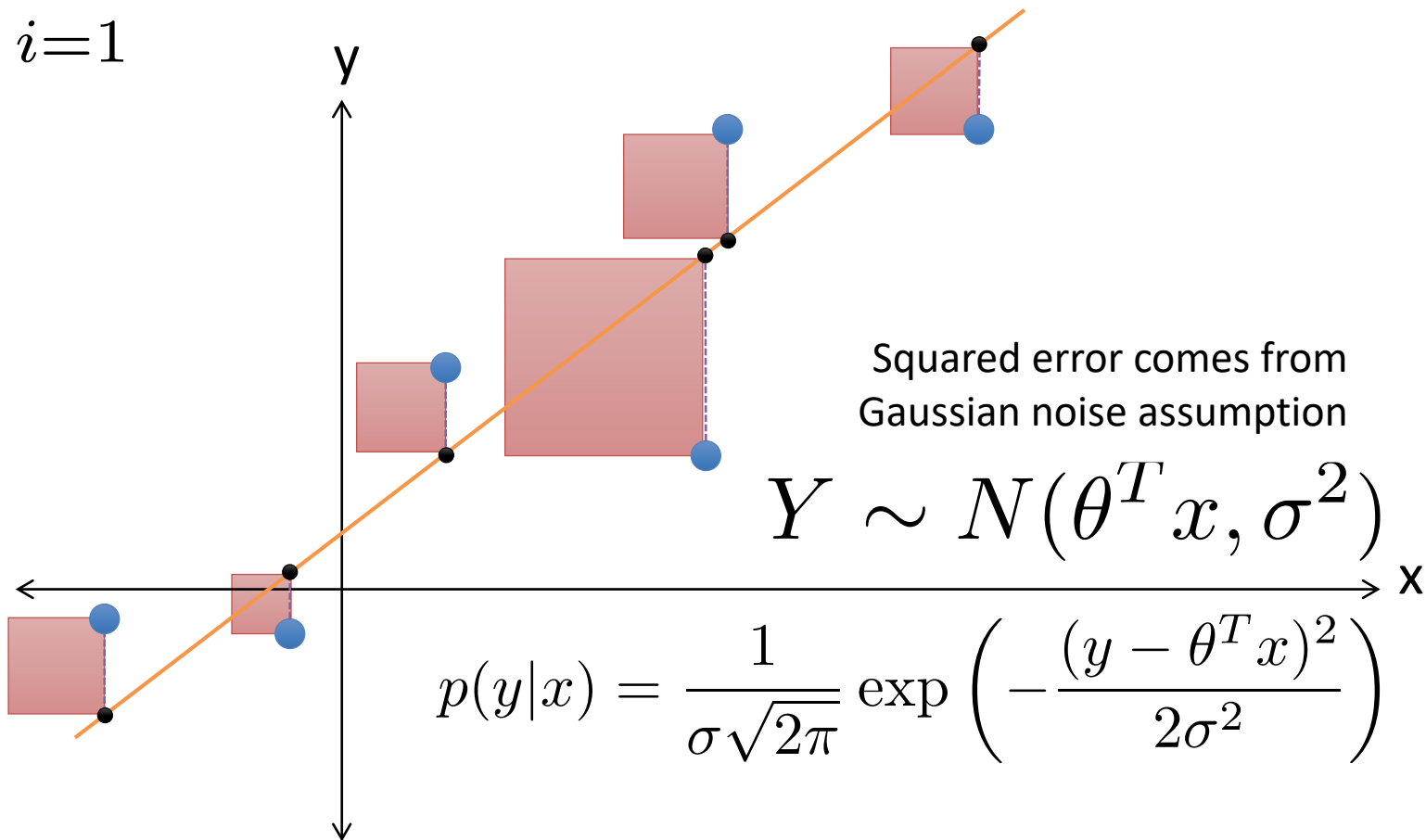
Step 1: define a **Loss Function**: *Average Prediction Error*

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

➤ Difference between **true** (y_i) and **predicted** $f_{\theta}(x_i)$ values

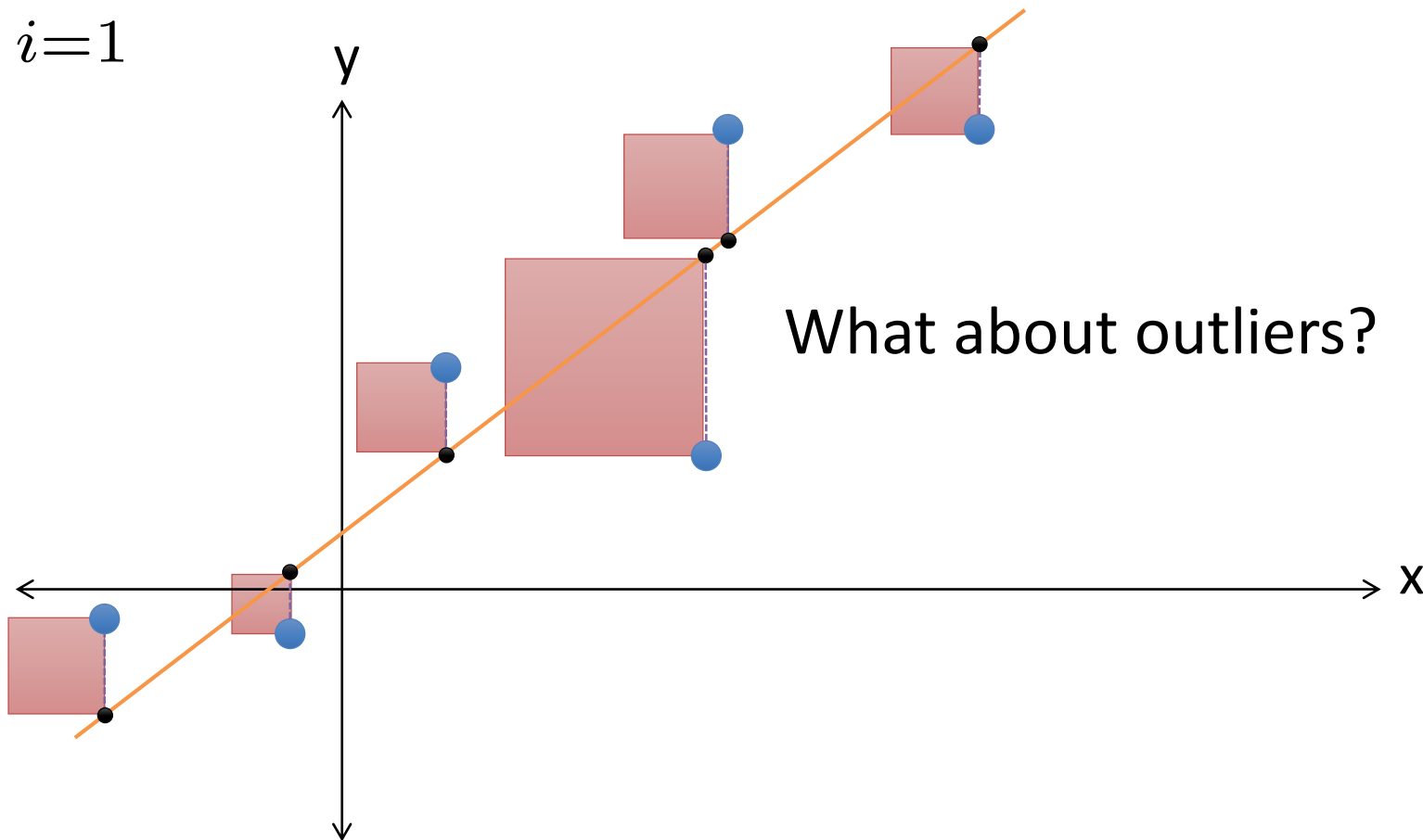
The meaning of Squared Loss (Error)

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



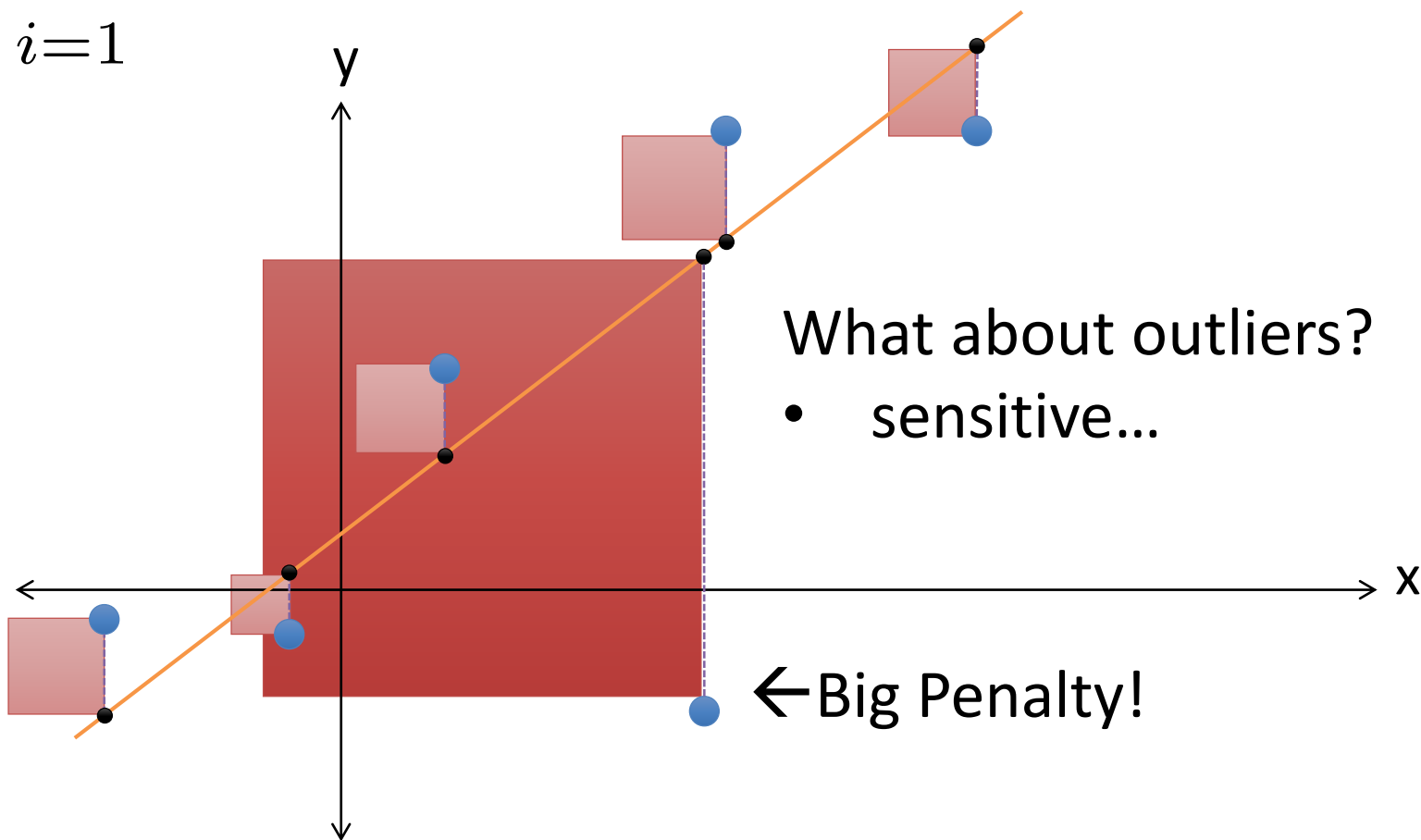
The meaning of Squared Loss (Error)

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



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**:

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

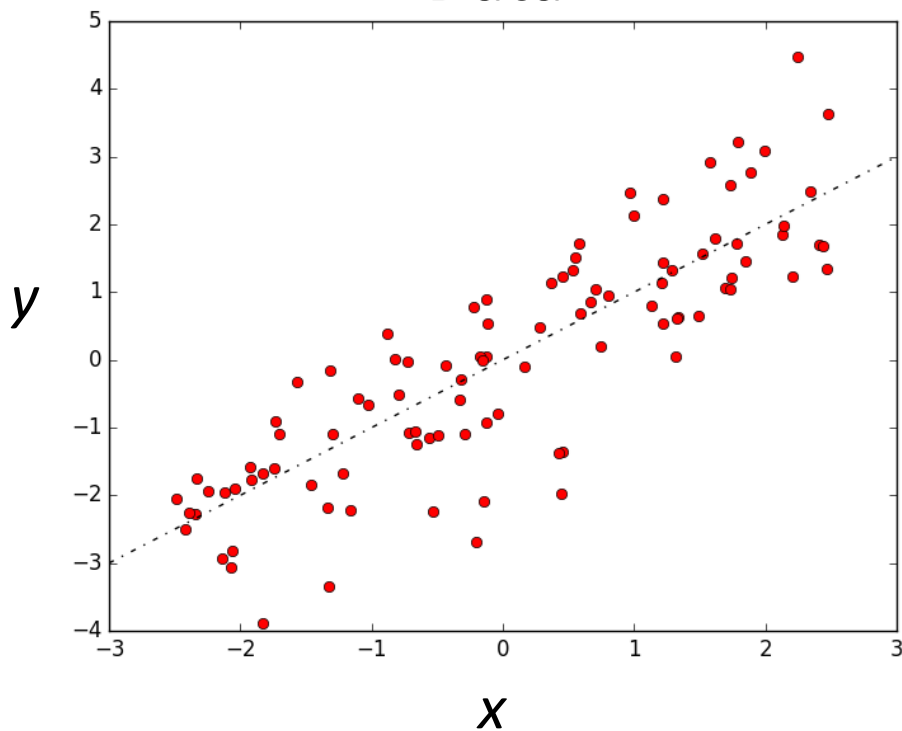
Step 2: Search for best model parameters θ

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

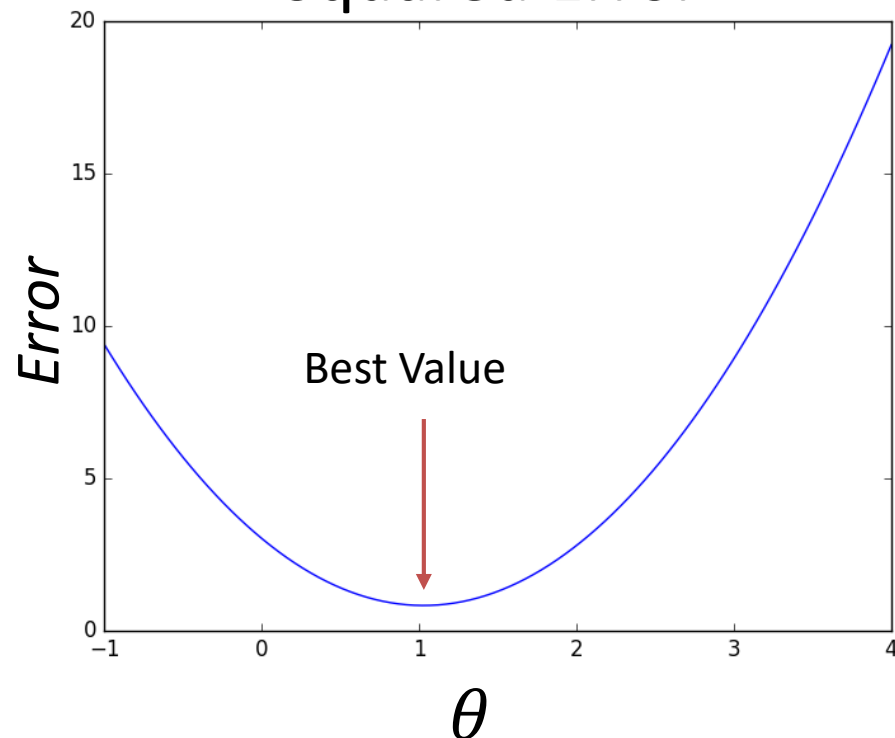
Minimizing the Squared Error

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

Data



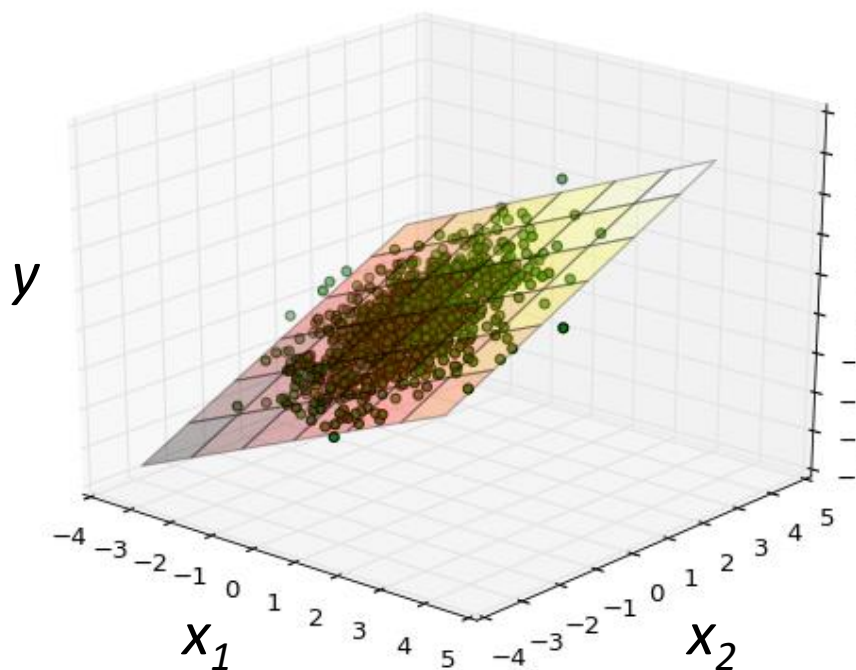
Squared Error



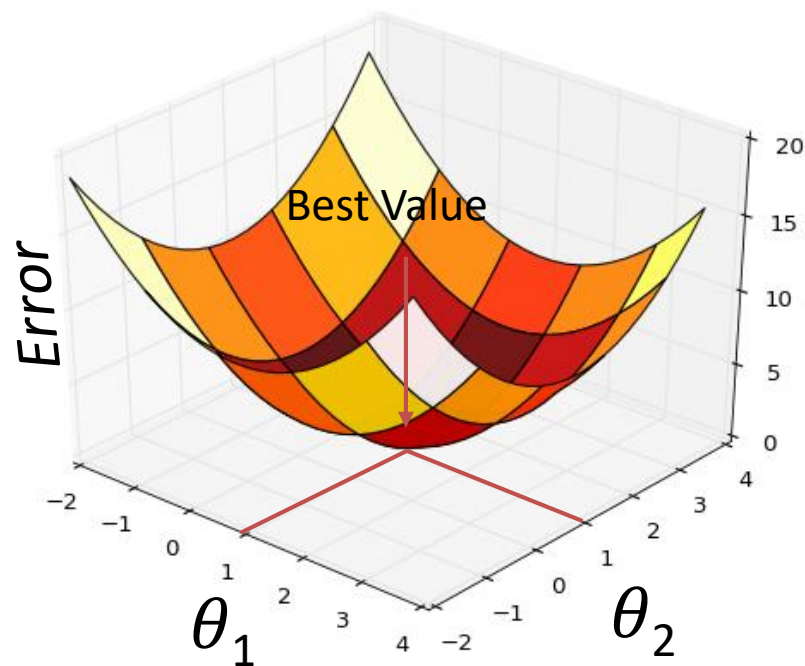
Minimizing the Squared Error

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

Data

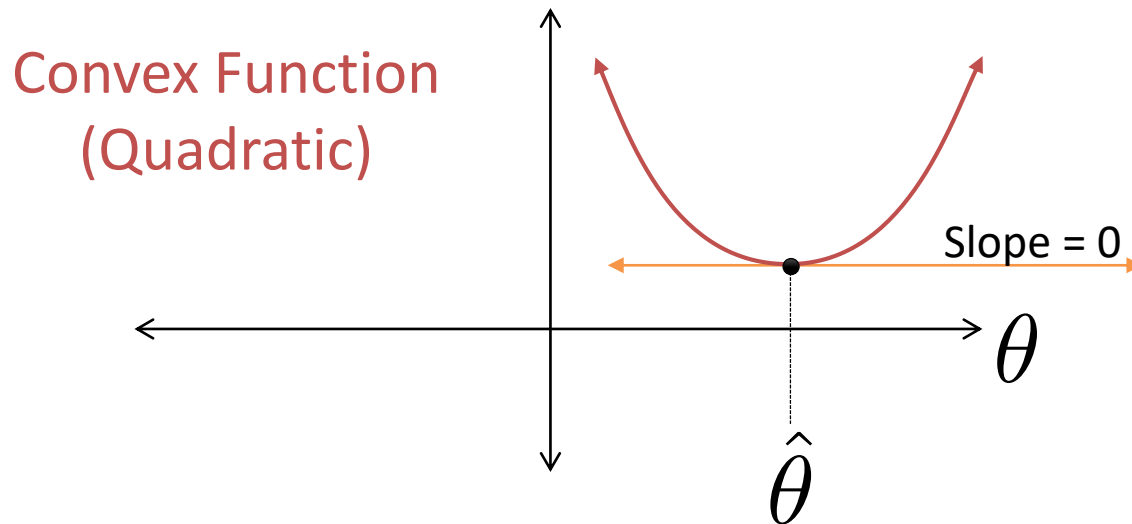


Squared Error



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 θ (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:

<p>Covariate (Design) Matrix</p> $X = \begin{matrix} \underbrace{\hspace{1cm}}_n \left[\begin{array}{c} \text{--- } x_1 \text{ ---} \\ \text{--- } x_2 \text{ ---} \\ \vdots \\ \text{--- } x_n \text{ ---} \end{array} \right] \in \mathbb{R}^{np} \\ \underbrace{\hspace{1cm}}_p \end{matrix}$	<p>Response Vector</p> $Y = \begin{matrix} \underbrace{\hspace{1cm}}_n \left[\begin{array}{c} y_1 \\ y_2 \\ \vdots \\ y_n \end{array} \right] \in \mathbb{R}^n \\ \underbrace{\hspace{1cm}}_1 \end{matrix}$
---	--

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 θ :

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

➤ Normal Equation:

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

➤ Solved using any standard linear algebra library

Examples

➤ Goal: Predict movie rating automatically

复仇者联盟3：无限战争 Avengers: Infinity War (2018)



导演: 安东尼·罗素 / 乔·罗素

编剧: 杰克·科比 / 克里斯托弗·马库斯 / 斯蒂芬·麦克菲利 / 吉姆·斯特林

主演: 小罗伯特·唐尼 / 克里斯·海姆斯沃斯 / 克里斯·埃文斯 / 马克·鲁弗洛 / 乔什·布洛林 / 更多...

类型: 动作 / 科幻 / 奇幻 / 冒险

官方网站: marvel.com/avengers

制片国家/地区: 美国

语言: 英语

上映日期: 2018-05-11(中国大陆) / 2018-04-23(加州首映) / 2018-04-27(美国)

片长: 150分钟

又名: 复联3 / 复仇者联盟：无限之战(台) / 复仇者联盟3：无尽之战 / Avengers: Infinity War - Part I / The Avengers 3: Part 1

IMDb链接: [tt4154756](https://www.imdb.com/title/tt4154756)

豆瓣评分

8.5

★★★★★
182012人评价

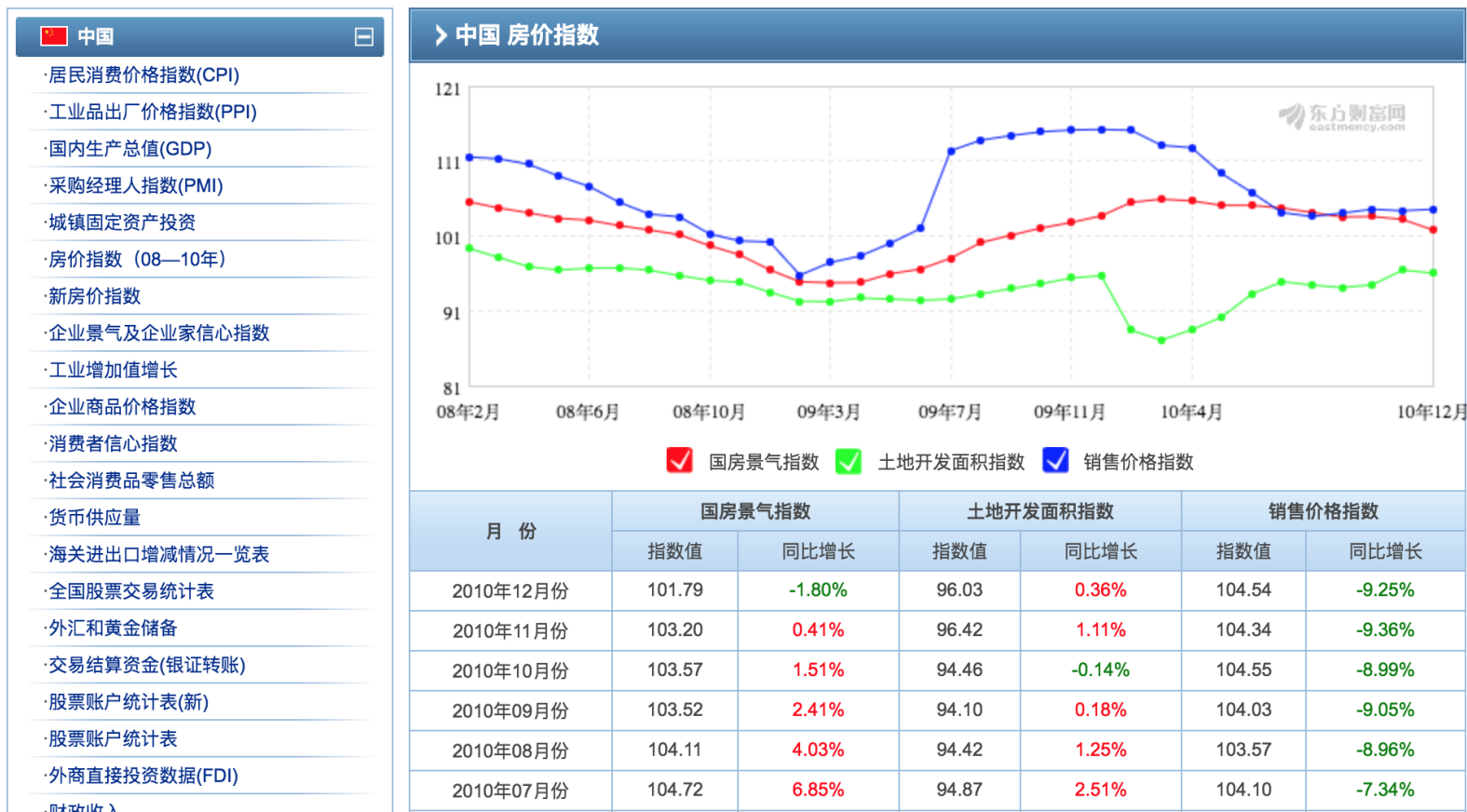
5星	44.8%
4星	38.6%
3星	14.1%
2星	1.8%
1星	0.8%

好于 96% 科幻片

好于 97% 动作片

Examples

➤ Goal: Predict the price of the house



Parallelize the Linear Regression

- Question: How do we estimate the linear model with a large, distributed database?

Recall: Linear Regression

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

➤ Setting equal to zero and solving for θ :

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

➤ Normal Equation:

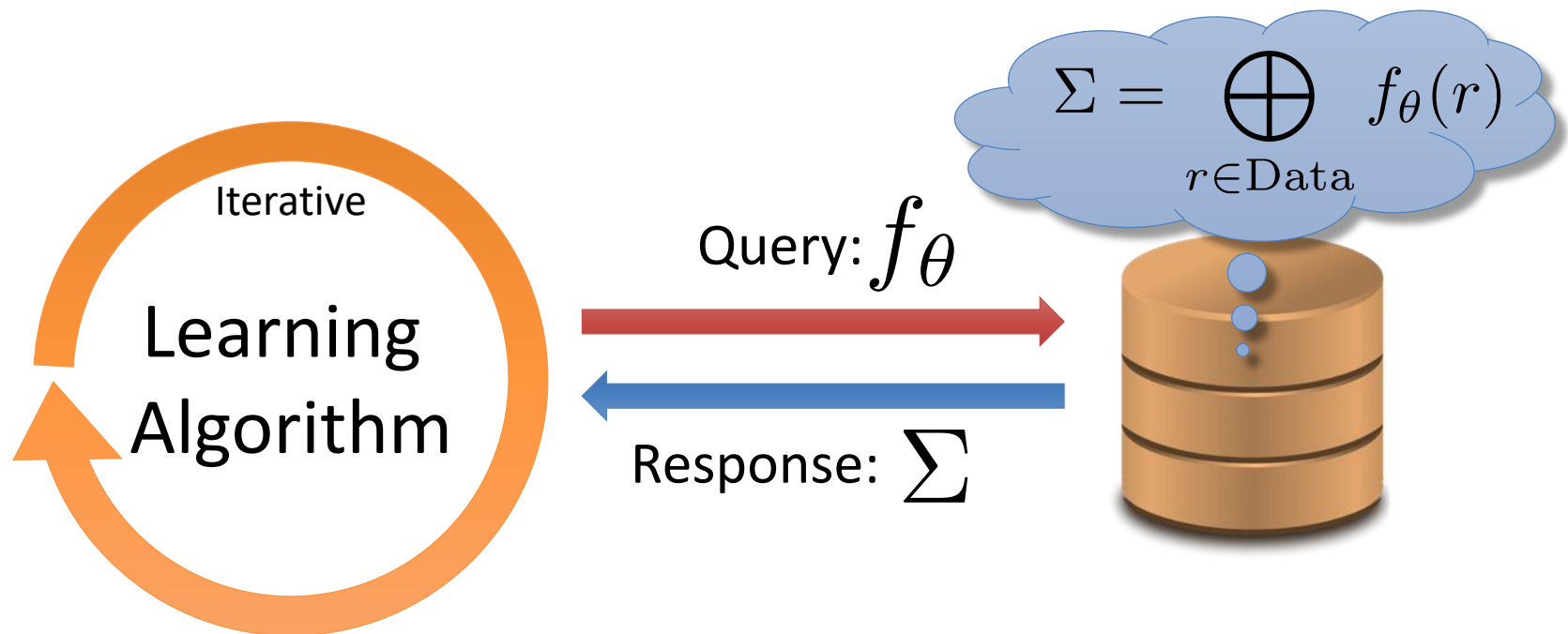
$$\hat{\theta} = (X^T X)^{-1} X^T Y$$

➤ Solved using any standard linear algebra library

Can we compute

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

using the statistical query pattern?



Can we compute

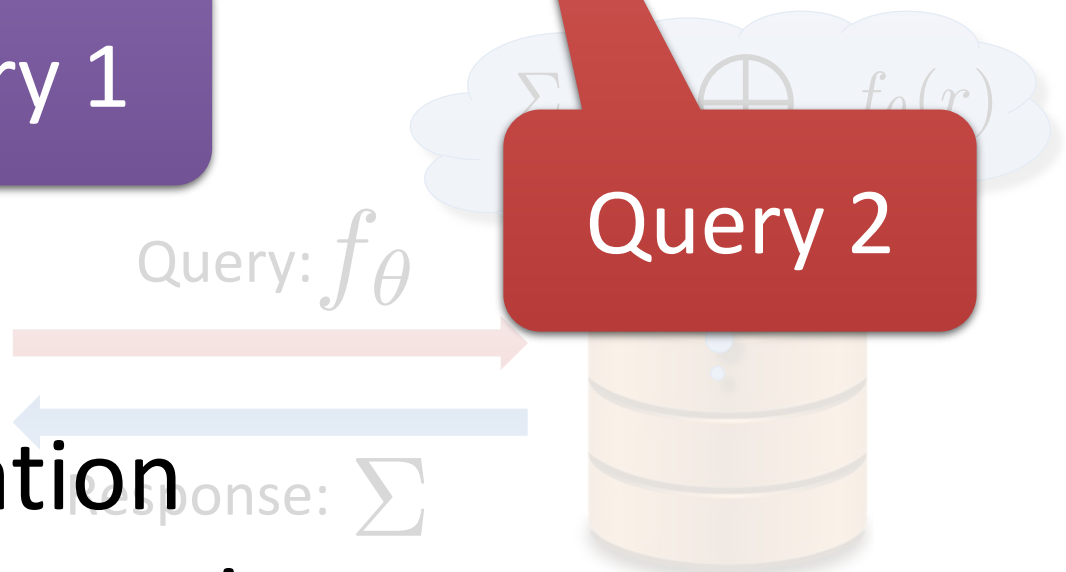
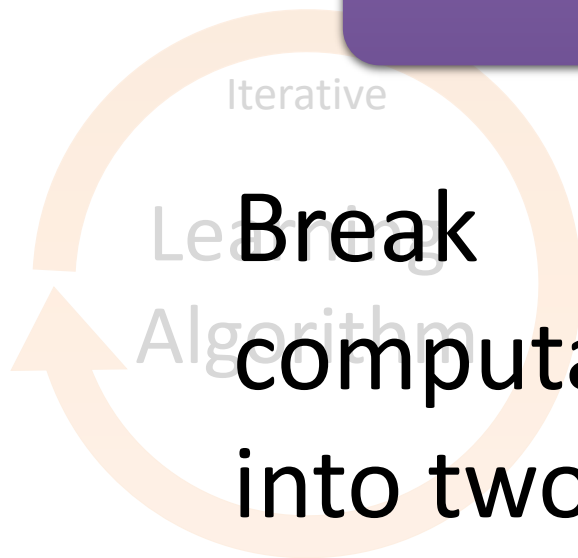
$$\hat{\theta} = (\underline{X^T X})^{-1} \underline{X^T Y}$$

using the statistical query pattern?

Query 1

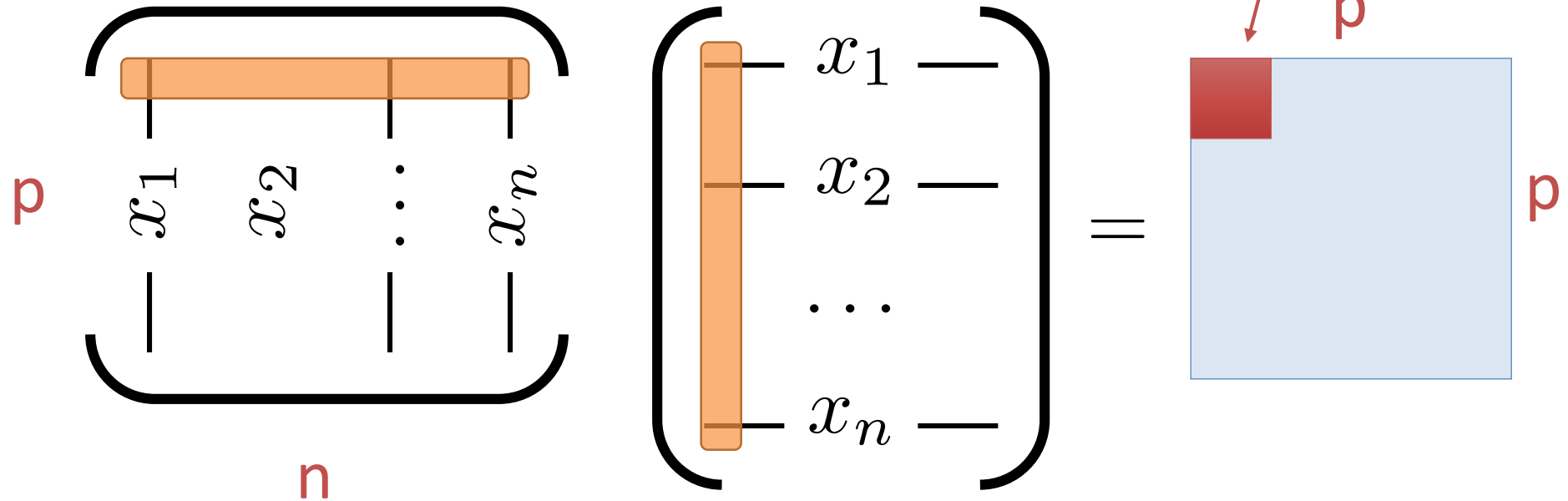
Query 2

Break
computation
into two queries



Computing Query 1

$$X^T X =$$



Computing Query 1

$$X^T X =$$

The diagram illustrates the computation of the matrix element (i,j) of $X^T X$. It shows the dot product of the i -th row of X and the j -th column of X , which is equivalent to the sum of the products of corresponding elements in the i -th row and j -th column.

On the left, a matrix X of size $p \times n$ is shown. The i -th row is highlighted in orange, with elements x_1, x_2, \dots, x_n . The matrix is labeled with p on the left and n on the bottom.

In the middle, a column vector of size n is shown, with elements x_1, x_2, \dots, x_n . The column is highlighted in orange. The matrix is labeled with n on the left.

On the right, a matrix of size $p \times p$ is shown. The element (i,j) is highlighted in red, representing the result of the dot product. The matrix is labeled with p on the left and p on the right.

The equation is represented as:

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

Computing Query 1

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

➤ Compute the row-wise some:

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

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

➤ Pure SQL Expression:

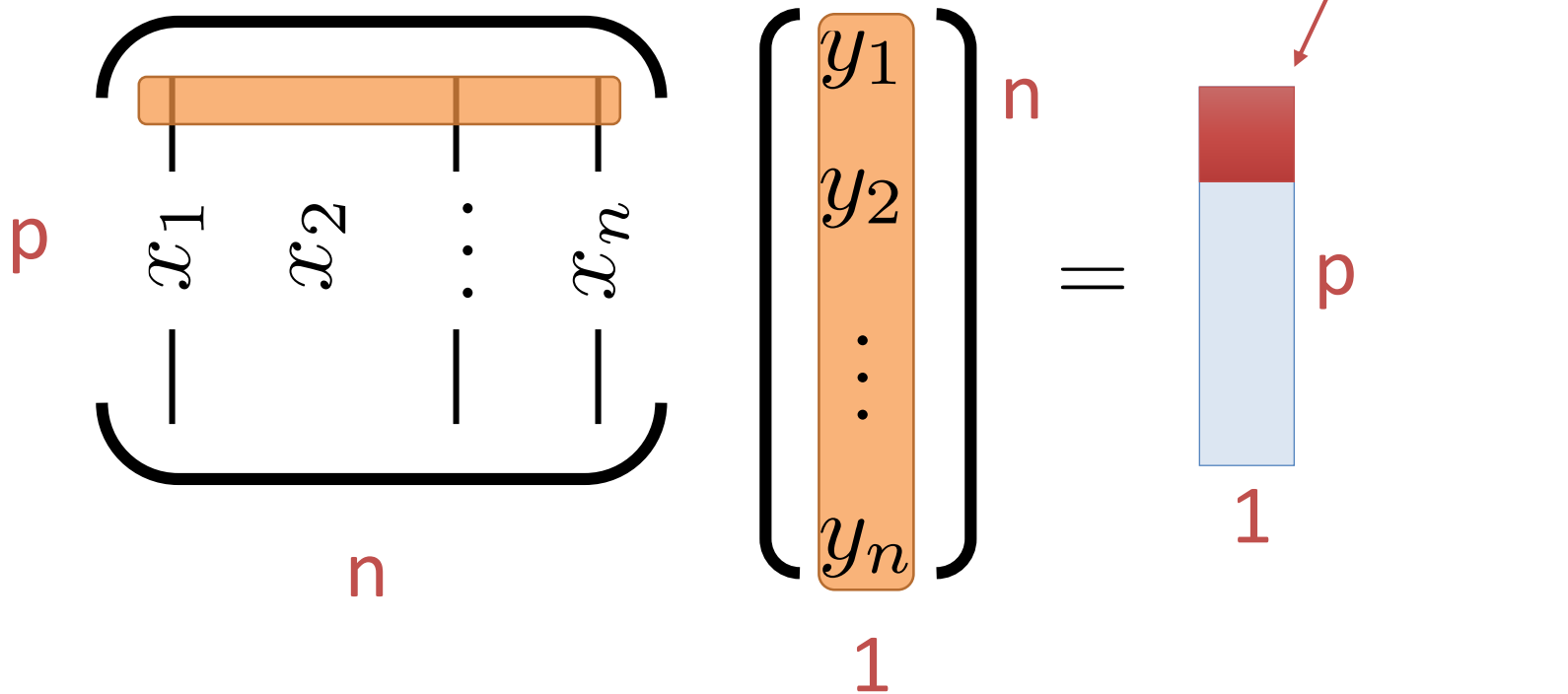
SELECT

sum($x_1 * x_1$) AS c_{11} , sum($x_1 * x_2$) AS c_{12} ,
~~sum($x_2 * x_1$) AS c_{21}~~ , sum($x_2 * x_2$) AS c_{22}

FROM data

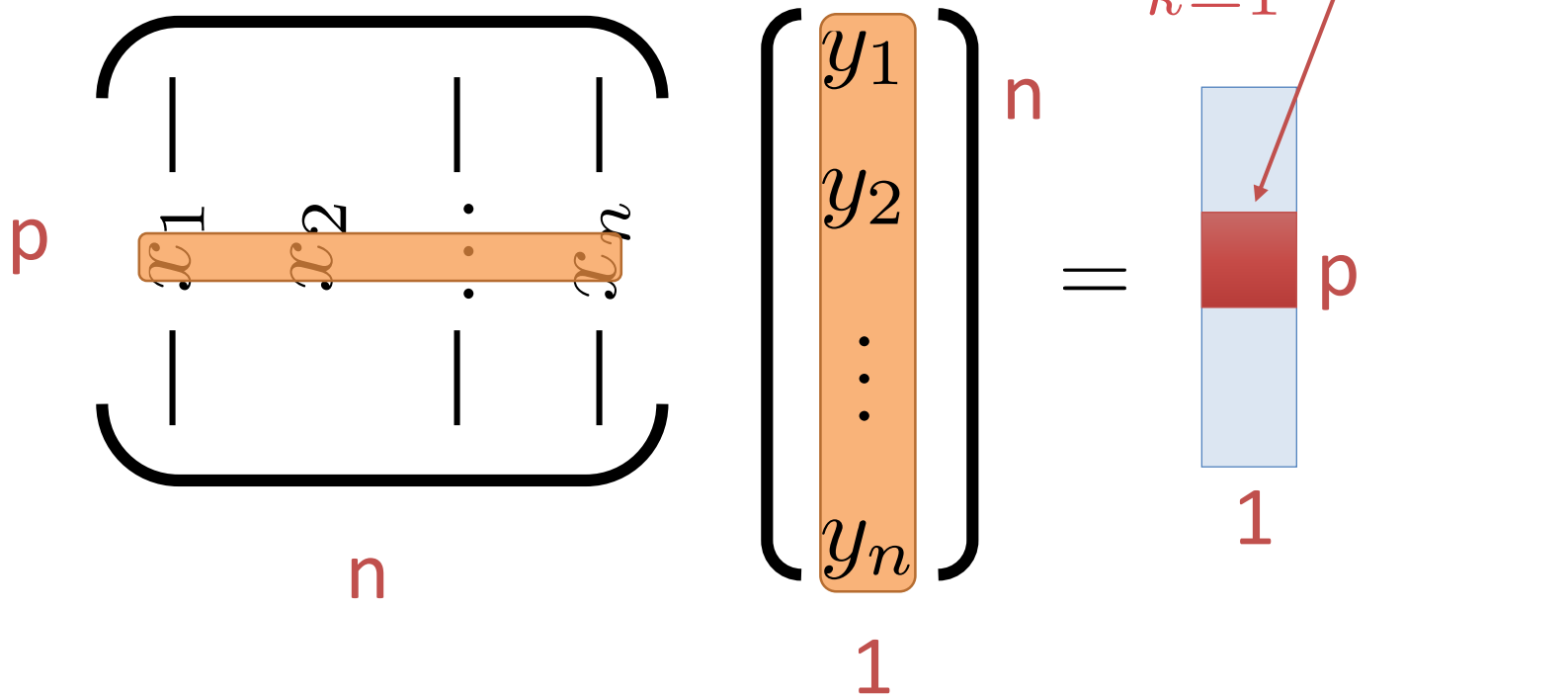
Computing Query 2

$$X^T y =$$



Computing Query 2

$$X^T y =$$



Computing Query 2

➤ Compute the row-wise sum:

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

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: xy
- **ReduceFunction**: vector sum

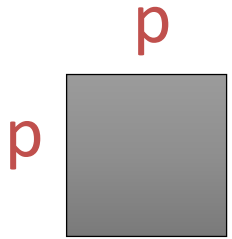
➤ Pure SQL Expression:


```
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:

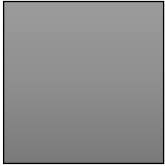
$$\hat{\theta} = C^{-1} d$$


$O(p^3)$

Least Squares Regression using the Statistical Query Pattern

What if p is large?

➤ In database compute sums:

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

1  ... could be expensive ... $O(np)$

➤ On client compute:

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

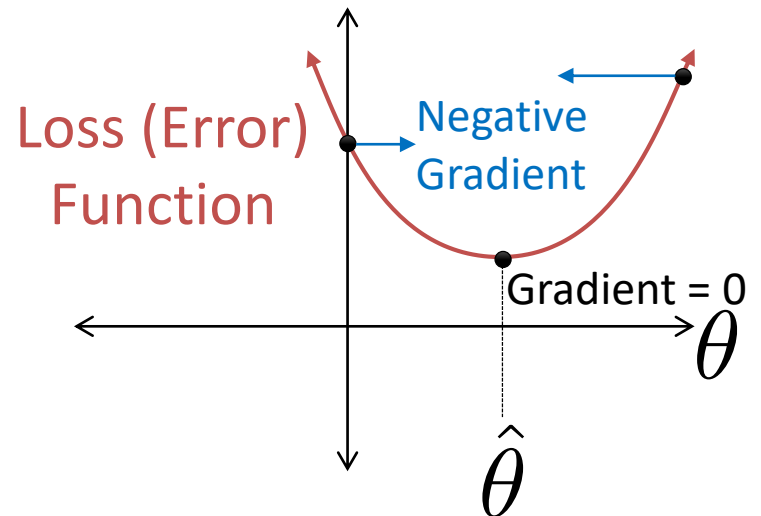
➤ Rather than directly solving:

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

➤ Instead we compute the gradient of the loss:

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

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



Gradient Descent Algorithm

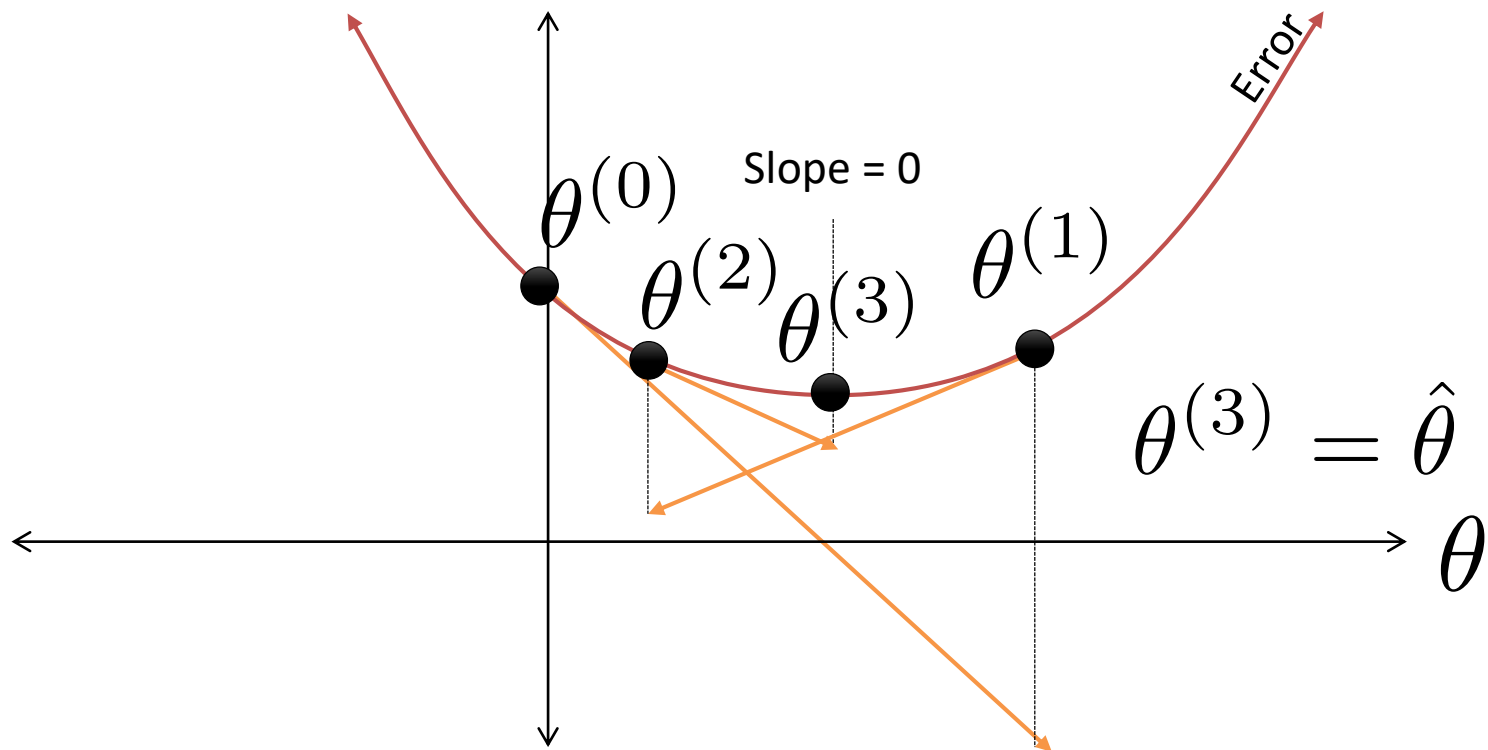
$t \leftarrow 0$

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

while (not converged):

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

$t \leftarrow t + 1$



Gradient Descent Algorithm

$t \leftarrow 0$

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


while (not converged):

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

$t \leftarrow t + 1$

➤ Does this fit the statistical query pattern

- Yes! Only dependence on data is:

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

- Can we go even faster?
 - **Stochastic Gradient Descent (SGD)**: Approximate the gradient by sampling data (typically several hundred records per query).

Stochastic Gradient Descent

- Update the parameters for each training case in turn, according to its own gradients

$$t \leftarrow 0$$

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

while (not converged):

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

$$t \leftarrow t + 1$$

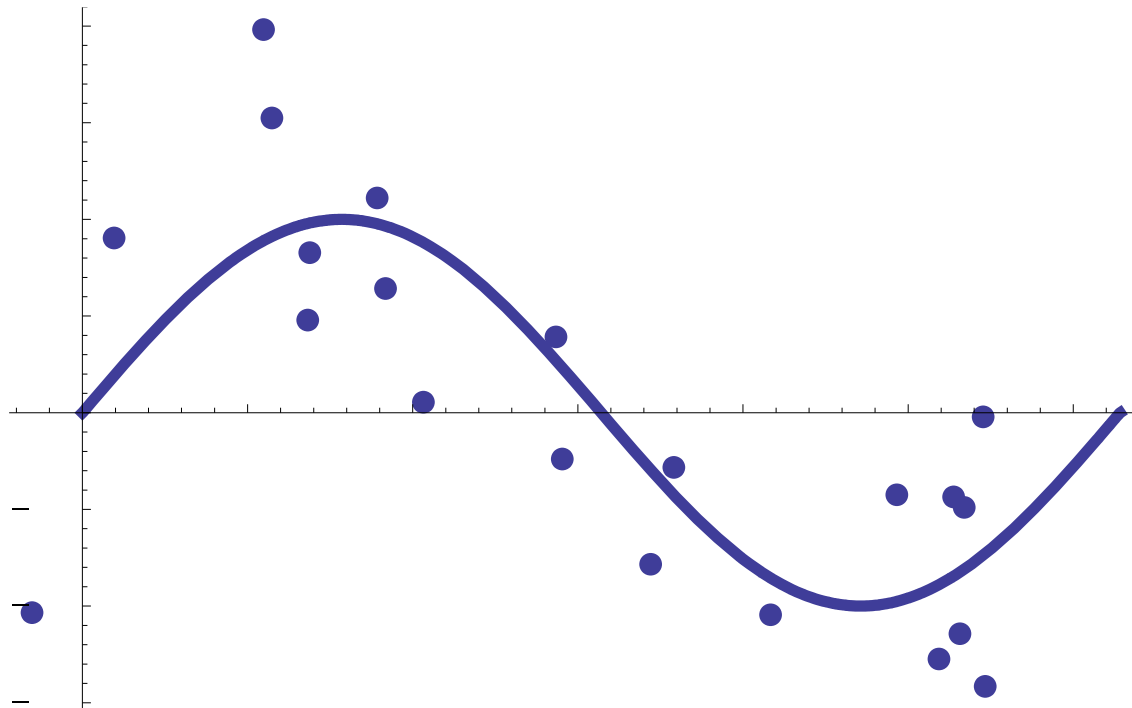
$$G(\theta; X, y) = (\mathbf{y}_i - \theta^T x_i) x_i \quad \text{Complexity?}$$

$$\text{Stepsize}^{(t)} = \frac{1}{t + 1} \quad O(p)$$

Advanced Topics in Linear Regression

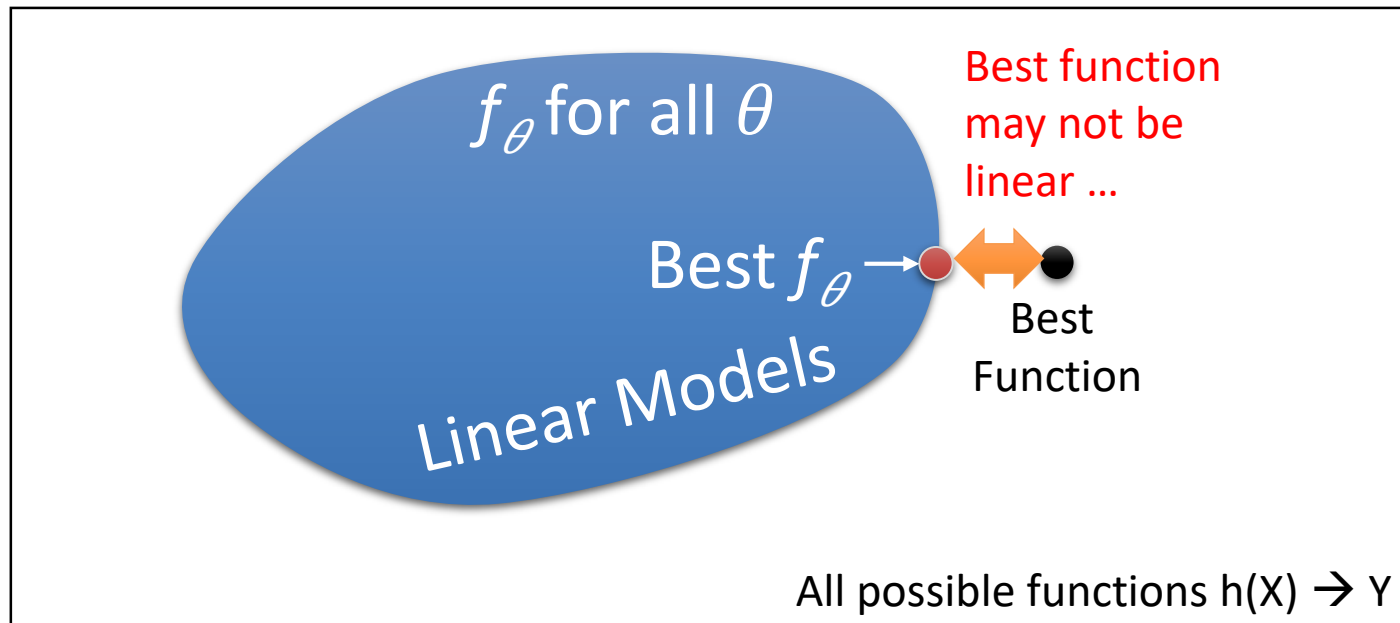
Fitting Non-linear Data

➤ What if Y has a non-linear response?

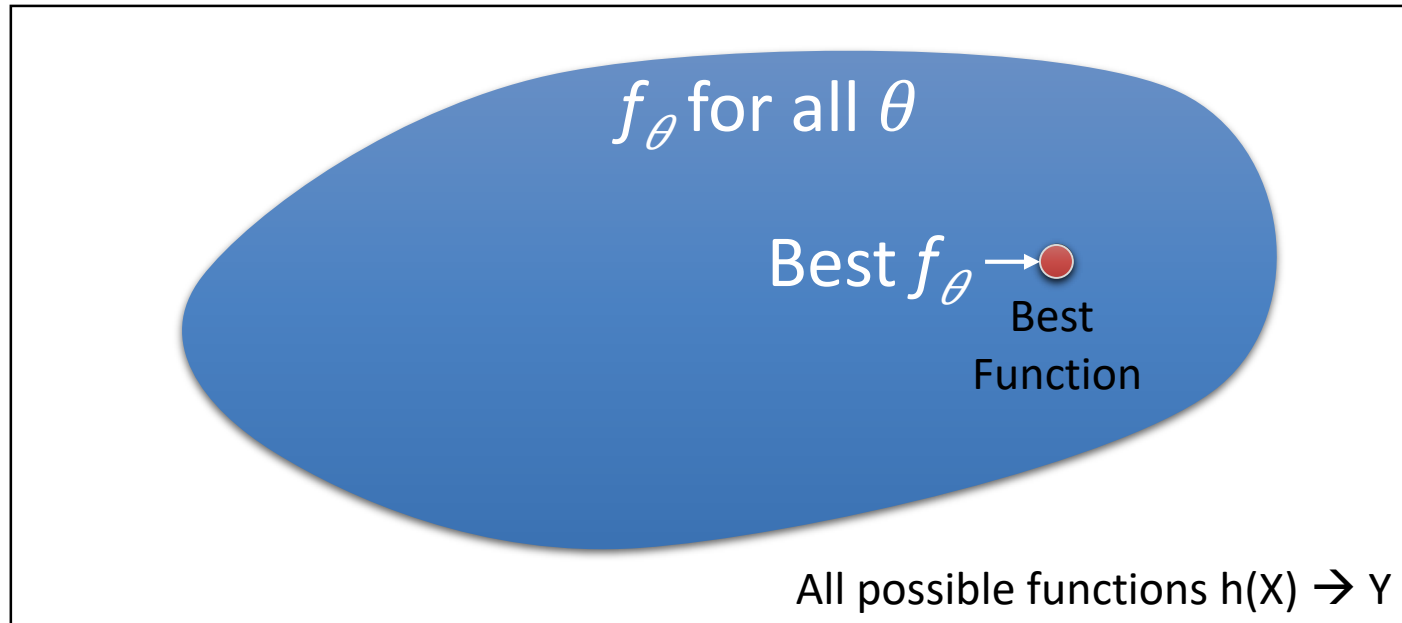


➤ Can we still use a linear model?

Finding the Best Parameters



Finding the Best Parameters



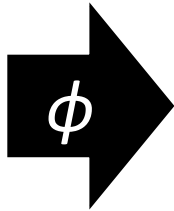
Feature Engineering

➤ By applying non-linear transformation ϕ :

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

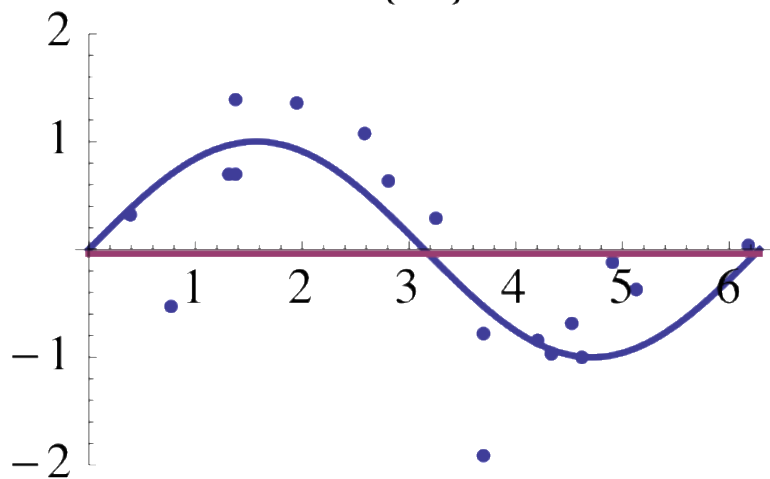
- Example:

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

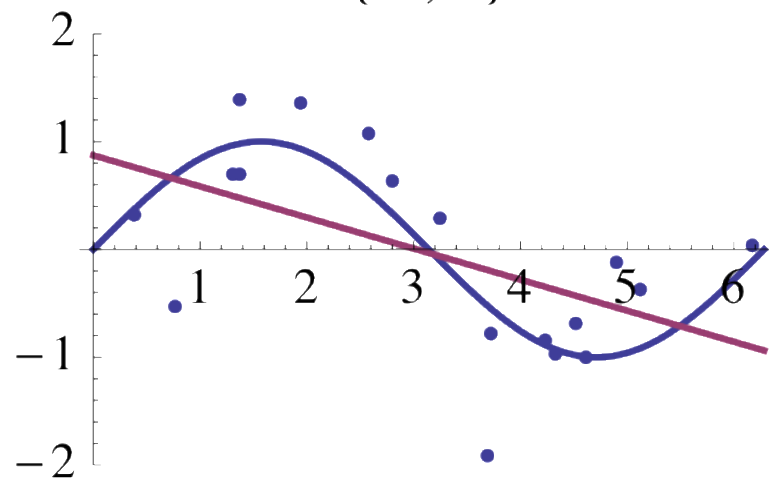
x_1	x_2	y		x_1	x_2	$x_1 * x_1$	$x_2 * x_2$	$x_1 * x_2$	y
1.1	2.7	3.6		1.1	2.7	1.21	7.29	2.97	3.6
4.2	3.2	7.5		4.2	3.2	17.64	10.24	13.44	7.5
9.8	9.2	17		9.8	9.2	90.04	84.64	90.16	17

Under-fitting

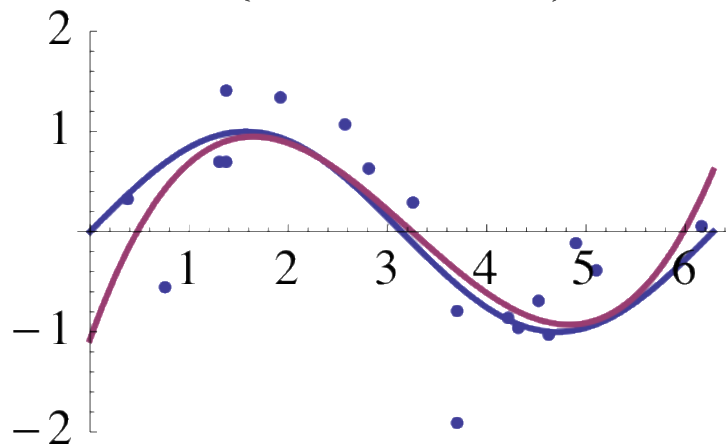
$\{1.\}$



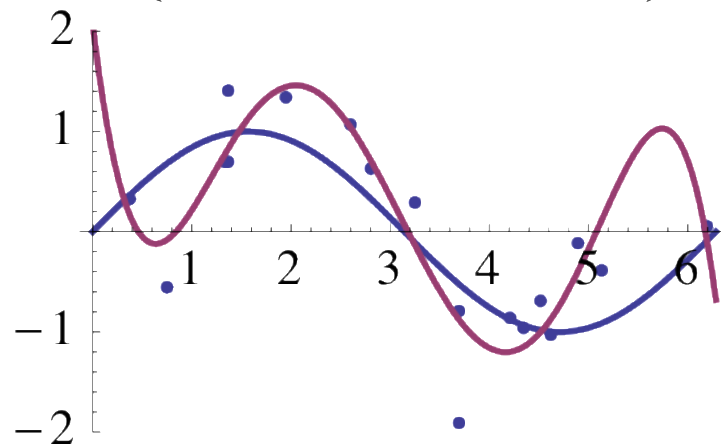
$\{1., x\}$



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

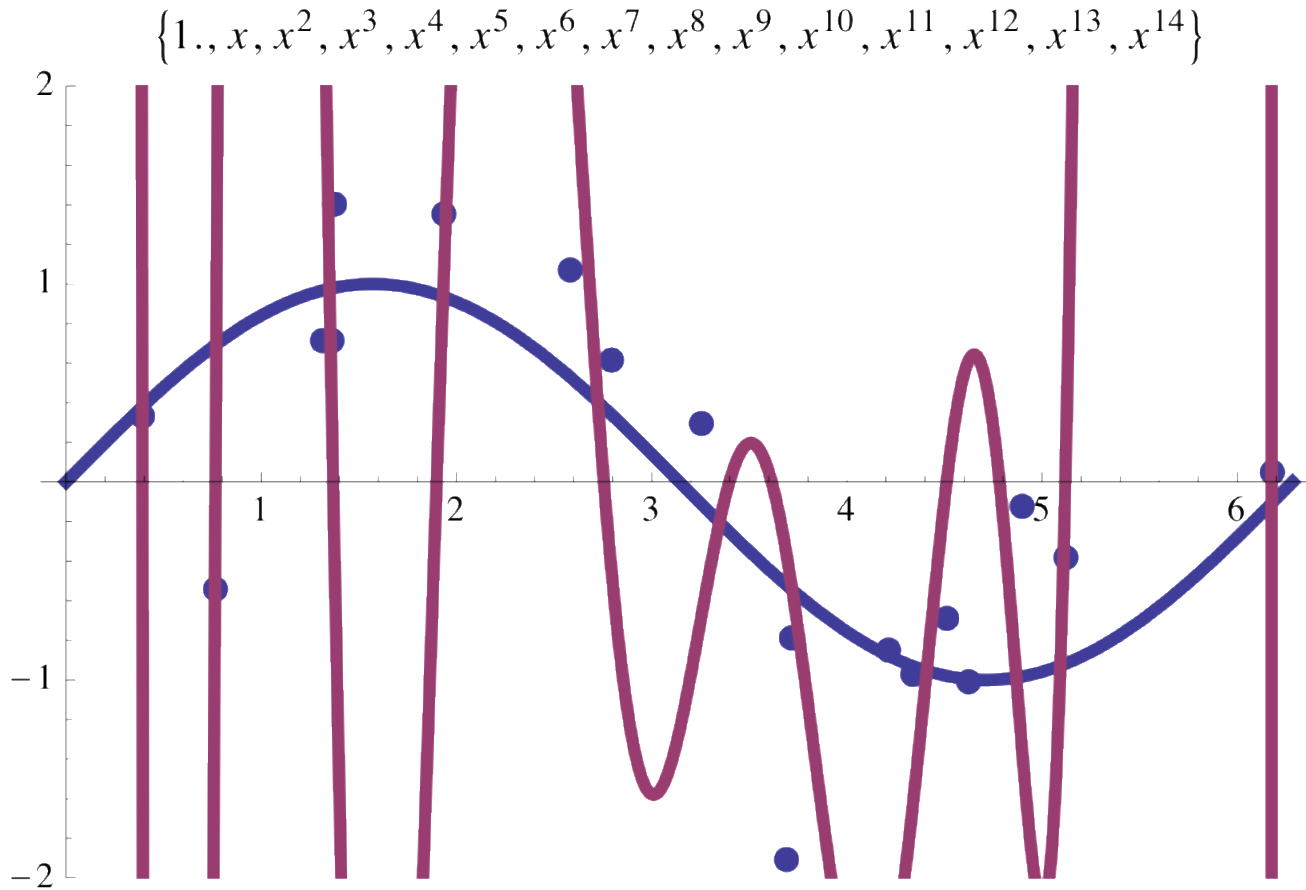


$\{1., x, x^2, x^3, x^4, x^5\}$



Over-fitting

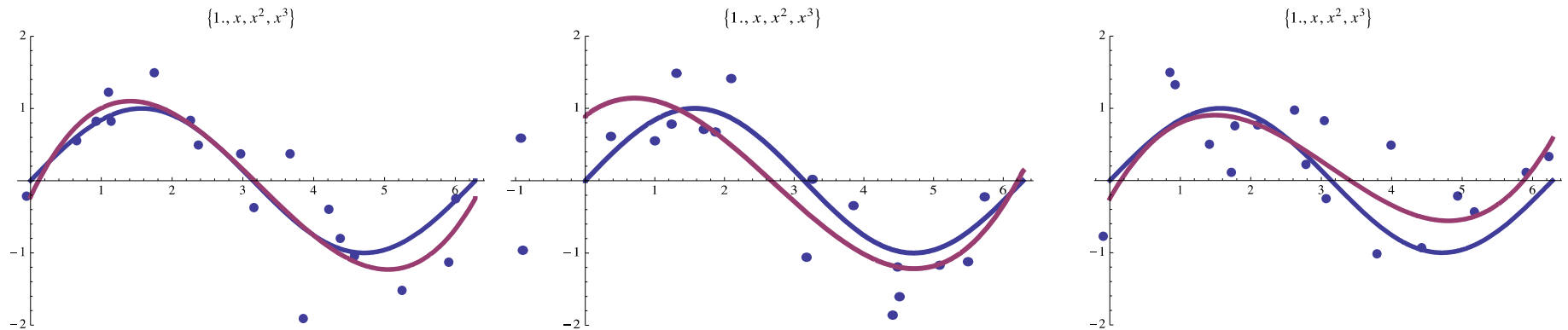
Really Over-fitting!



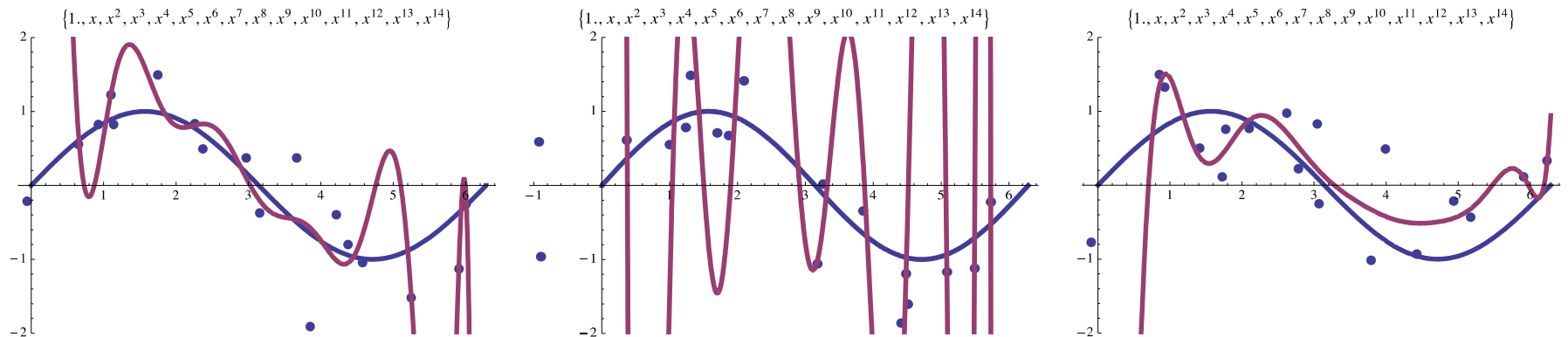
- Errors on training data are small
- But errors on new points are likely to be large

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:

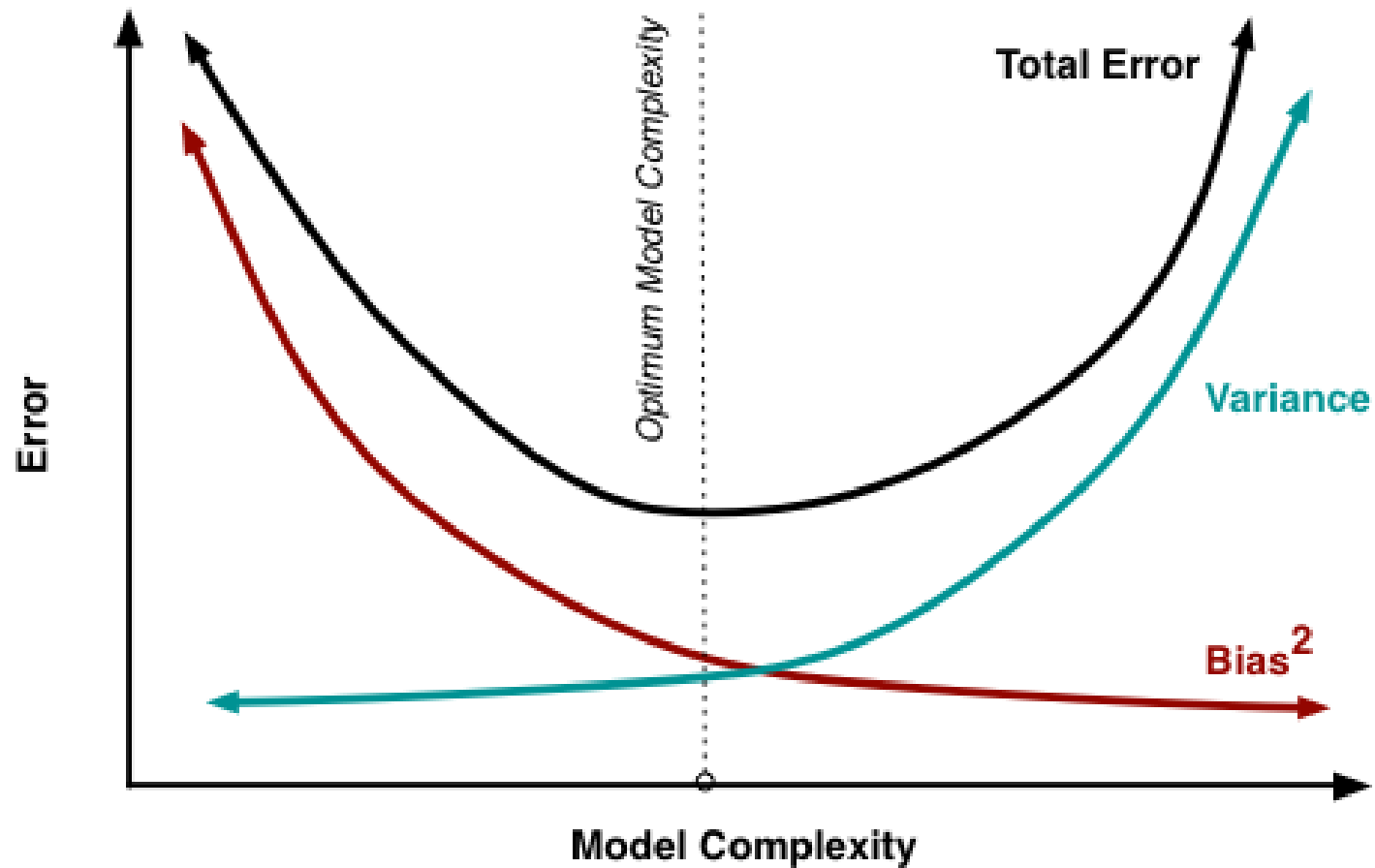
$$\text{Test Error} = \text{Noise} + \text{Bias}^2 + \text{Variance}$$

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:

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

➤ Common of Regularization Functions:

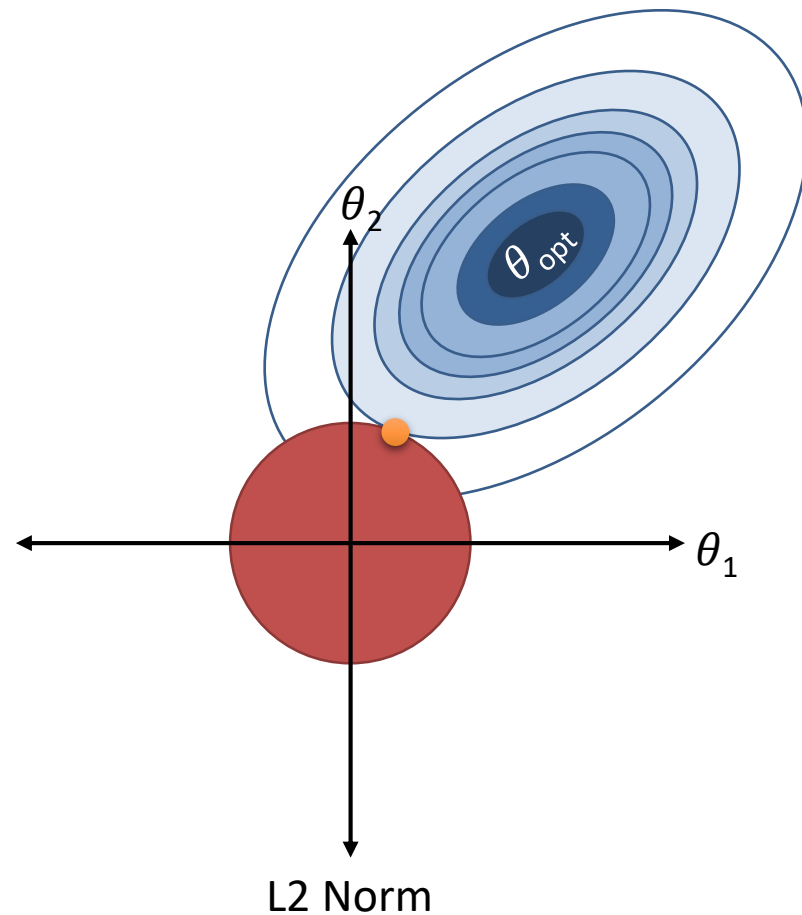
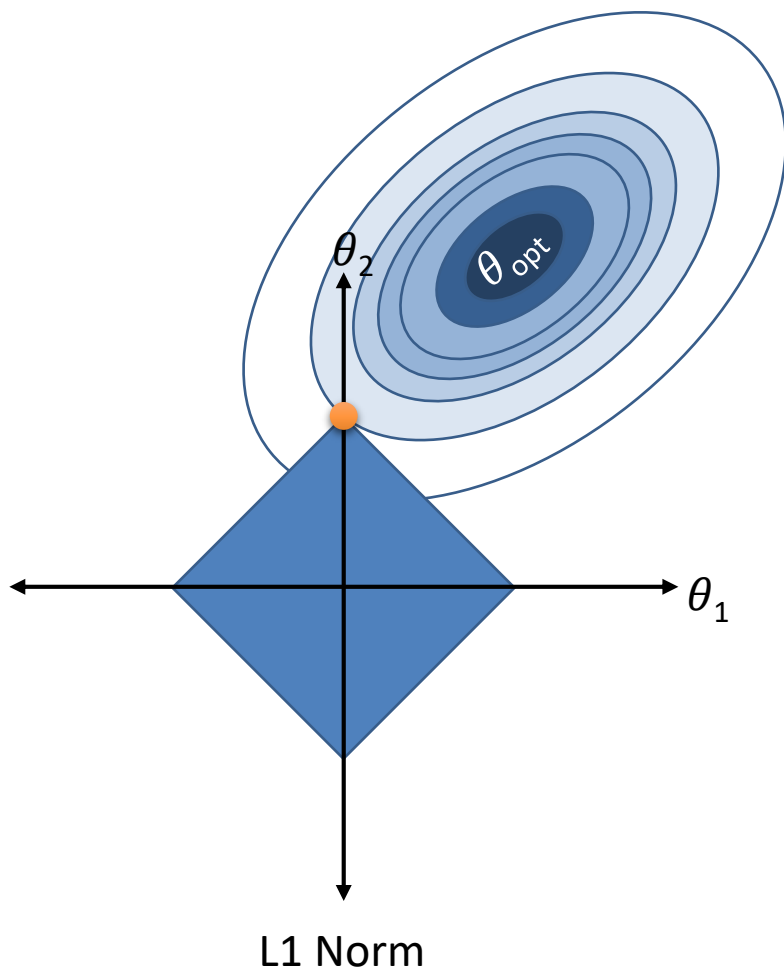
$$\begin{array}{ll} \text{Ridge (L2-Reg)} & R_{\text{Ridge}}(\theta) = \sum_{i=1}^d \theta_i^2 \\ \text{Regression} & \end{array} \quad \begin{array}{ll} \text{Lasso} & R_{\text{Lasso}}(\theta) = \sum_{i=1}^d |\theta_i| \\ \text{(L1-Reg)} & \end{array}$$

- Encourage small parameter values

➤ The parameter λ 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:

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

➤ 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 λ ?

Picking The Regularization Parameter λ

➤ **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 $\rightarrow \lambda = 0$

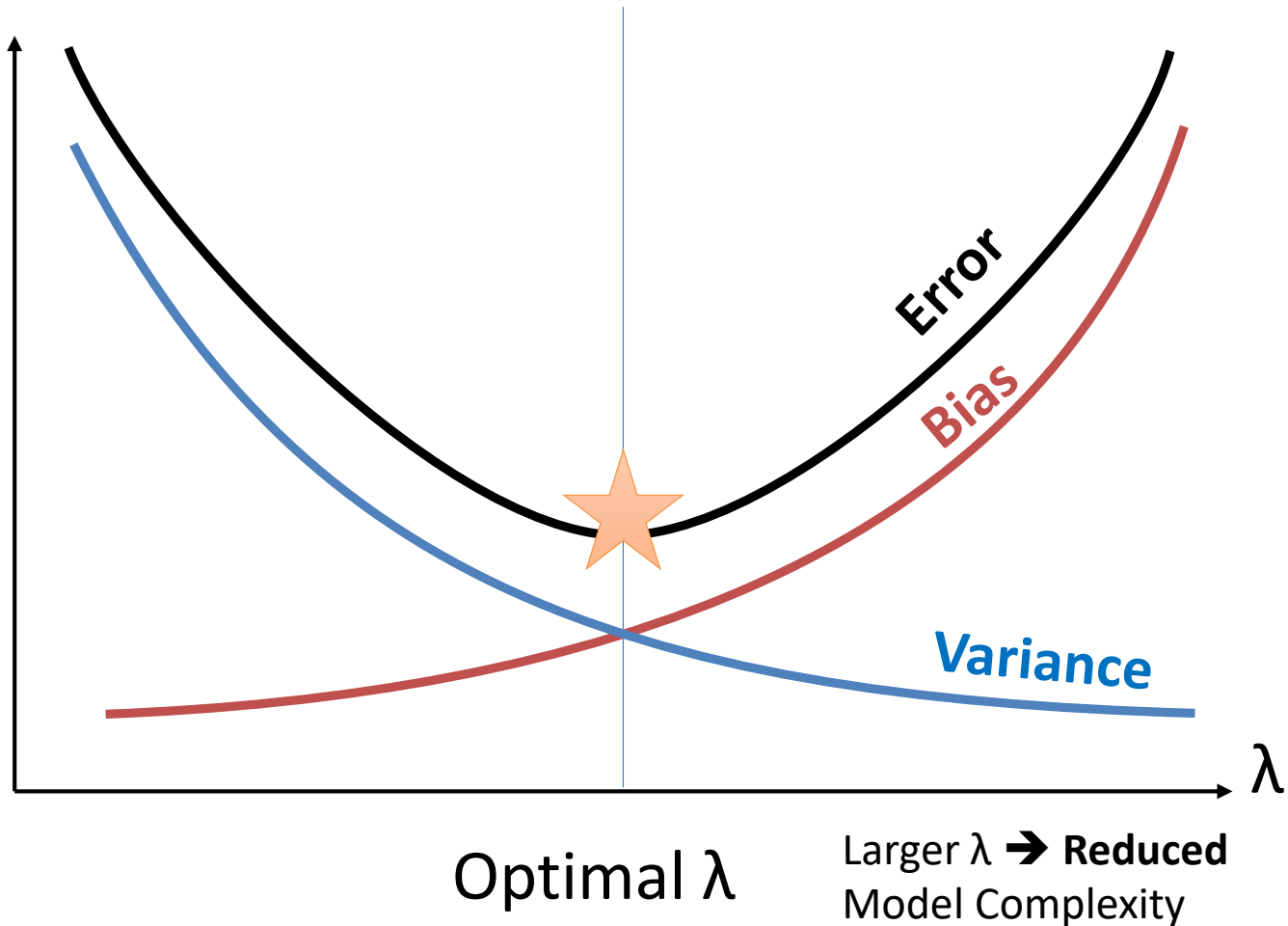
➤ Intuition we want to minimize **test** error

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

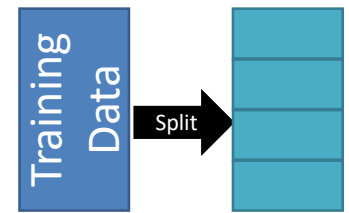
➤ **2nd Proposal:** Split training data into training and evaluation sets

- For a range of λ values compute optimal θ_λ using only the reduced training set
- Evaluate θ_λ on the separate evaluation set and select the λ with the lowest error

Bias Variance Plot



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

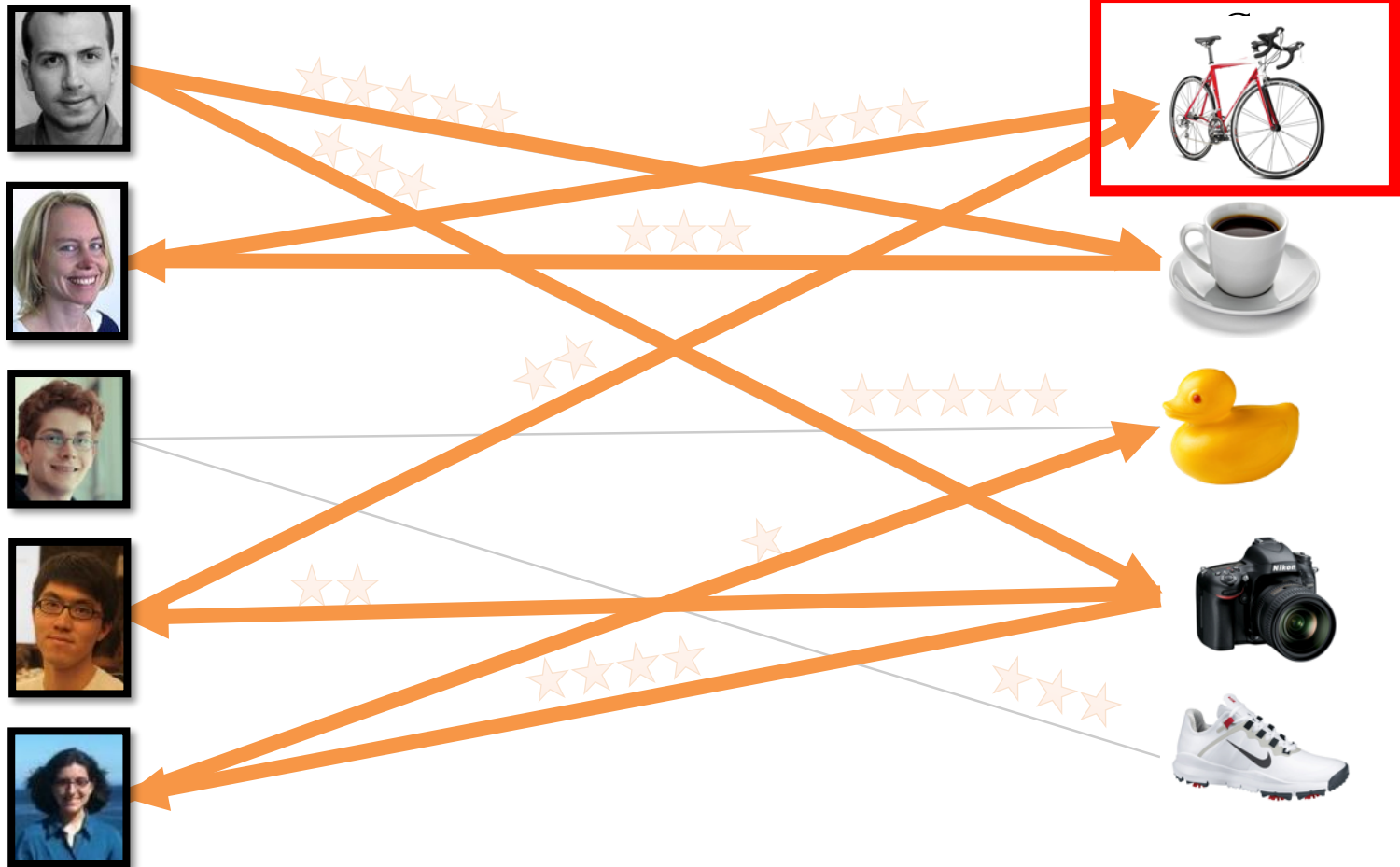
Using Regression for Content Recommendation

Recommending Products

Users

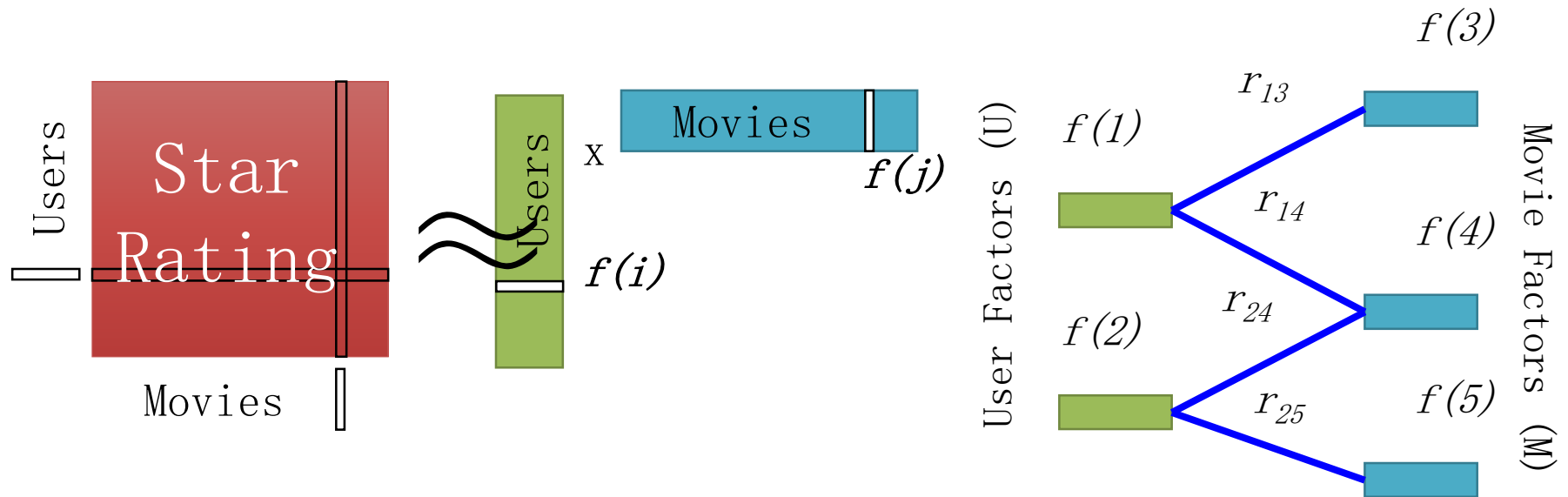
Ratings

Item



Recommending Products

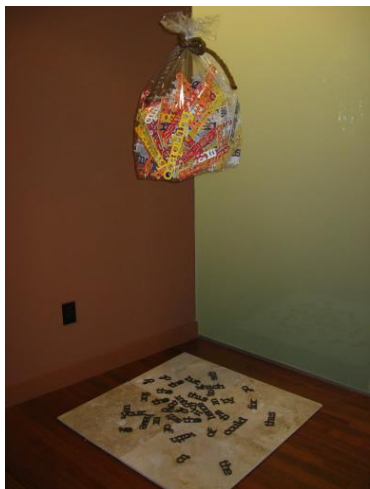
Low-Rank Matrix Factorization:



Iterate:

$$f[i] = \arg \min_{w \in \mathbb{R}^d} \sum_{j \in \text{Nbrs}(i)} (r_{ij} - w^T f[j])^2 + \lambda ||w||_2^2$$

Summary of Regression



0	0	0	0	1
Sports	Furniture	Clothing	Shoes	Electronics

$$\begin{matrix} p \\ p \end{matrix} \begin{matrix} p \\ 1 \end{matrix} \begin{matrix} \square \\ \square \end{matrix} \quad \begin{matrix} C = X^T X = \sum_{i=1}^n X_i X_i^T \\ d = X^T y = \sum_{i=1}^n X_i y_i \end{matrix}$$

