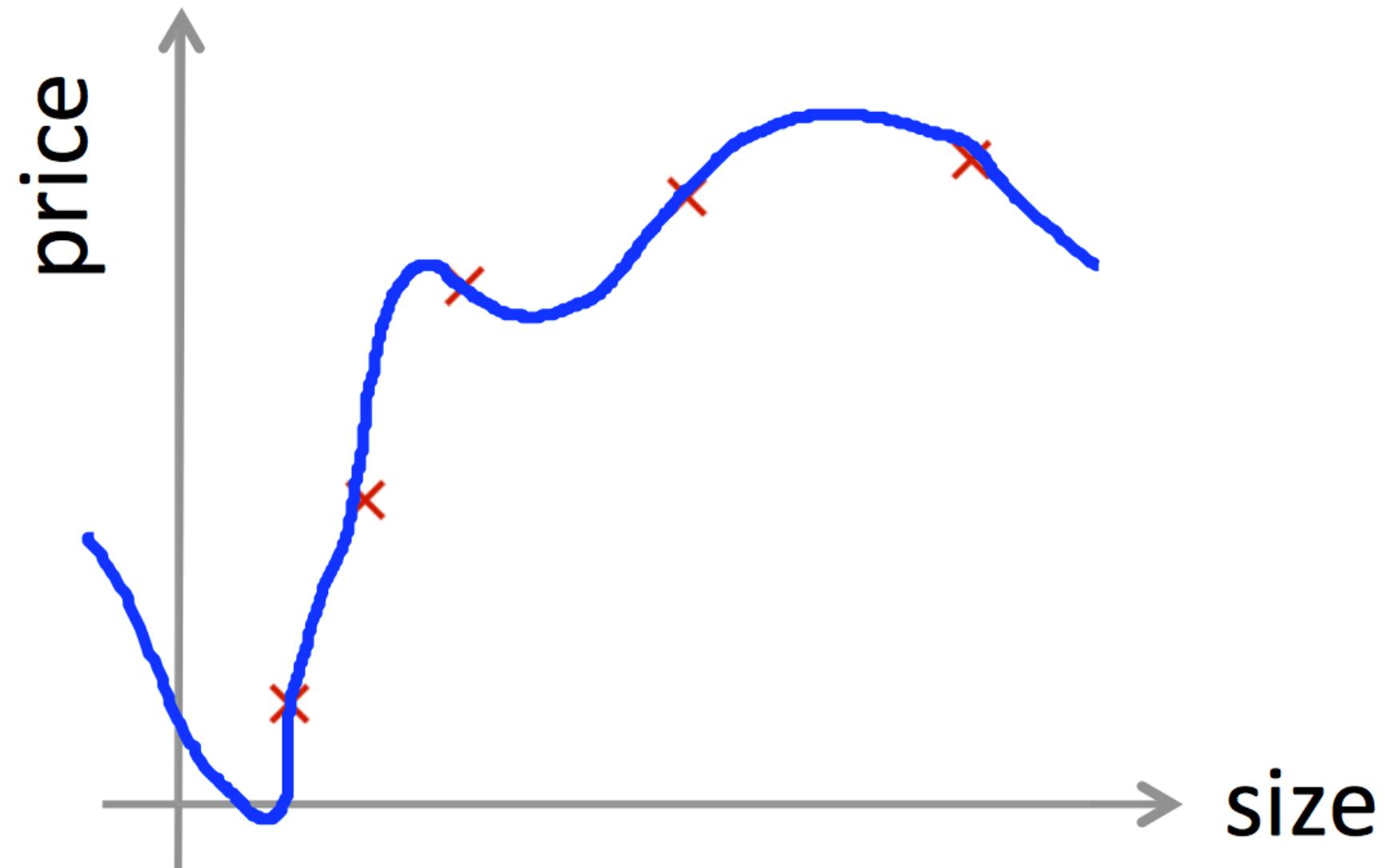


Evaluating your hypothesis



$$\Rightarrow h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

Fails to generalize to new examples not in training set.

- x_1 = size of house
- x_2 = no. of bedrooms
- x_3 = no. of floors
- x_4 = age of house
- x_5 = average income in neighborhood
- x_6 = kitchen size
- ⋮
- x_{100}

Evaluating a Learning Algorithm

Evaluating your hypothesis

Dataset:

Size	Price
2104	400
1600	330
2400	369
1416	232
3000	540
1985	300
1534	315
1427	199
1380	212
1494	243

70%

Training set

30%

Test set



$(x^{(1)}, y^{(1)})$
$(x^{(2)}, y^{(2)})$
\vdots
\vdots
$(x^{(m)}, y^{(m)})$



$(x_{test}^{(1)}, y_{test}^{(1)})$
$(x_{test}^{(2)}, y_{test}^{(2)})$
\vdots
$(x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$

m_{test} = no.
of test example
 $(x_{test}^{(i)}, y_{test}^{(i)})$

Evaluating a Learning Algorithm

Training/testing procedure for linear regression

- Learn parameter $\underline{\theta}$ from training data (minimizing training error $J(\theta)$)

\uparrow 70%

- Compute test set error:

$$J_{test}(\theta) = \frac{1}{2m_{test}} \sum_{i=1}^{m_{test}} \left(h_{\theta}(x_{test}^{(i)}) - y_{test}^{(i)} \right)^2$$

Training/testing procedure for logistic regression

- Learn parameter θ from training data
- Compute test set error:

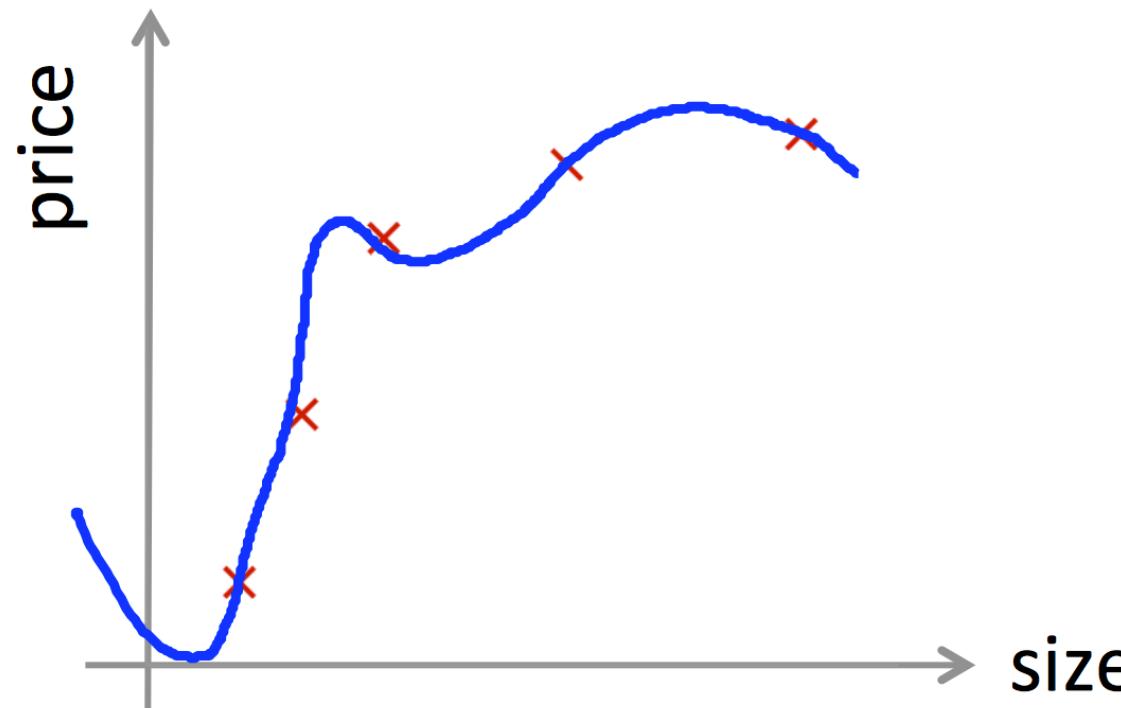
$$J_{test}(\theta) = -\frac{1}{m_{test}} \sum_{i=1}^{m_{test}} y_{test}^{(i)} \log h_{\theta}(x_{test}^{(i)}) + (1 - y_{test}^{(i)}) \log (1 - h_{\theta}(x_{test}^{(i)}))$$

- Misclassification error (0/1 misclassification error):

Evaluating a Learning Algorithm

Model selection and train/validation/test sets

Overfitting example



Model selection

$$d=1 \quad 1. \quad h_{\theta}(x) = \theta_0 + \theta_1 x$$

$$d=2 \quad 2. \quad h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$

$$d=3 \quad 3. \quad h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_3 x^3$$

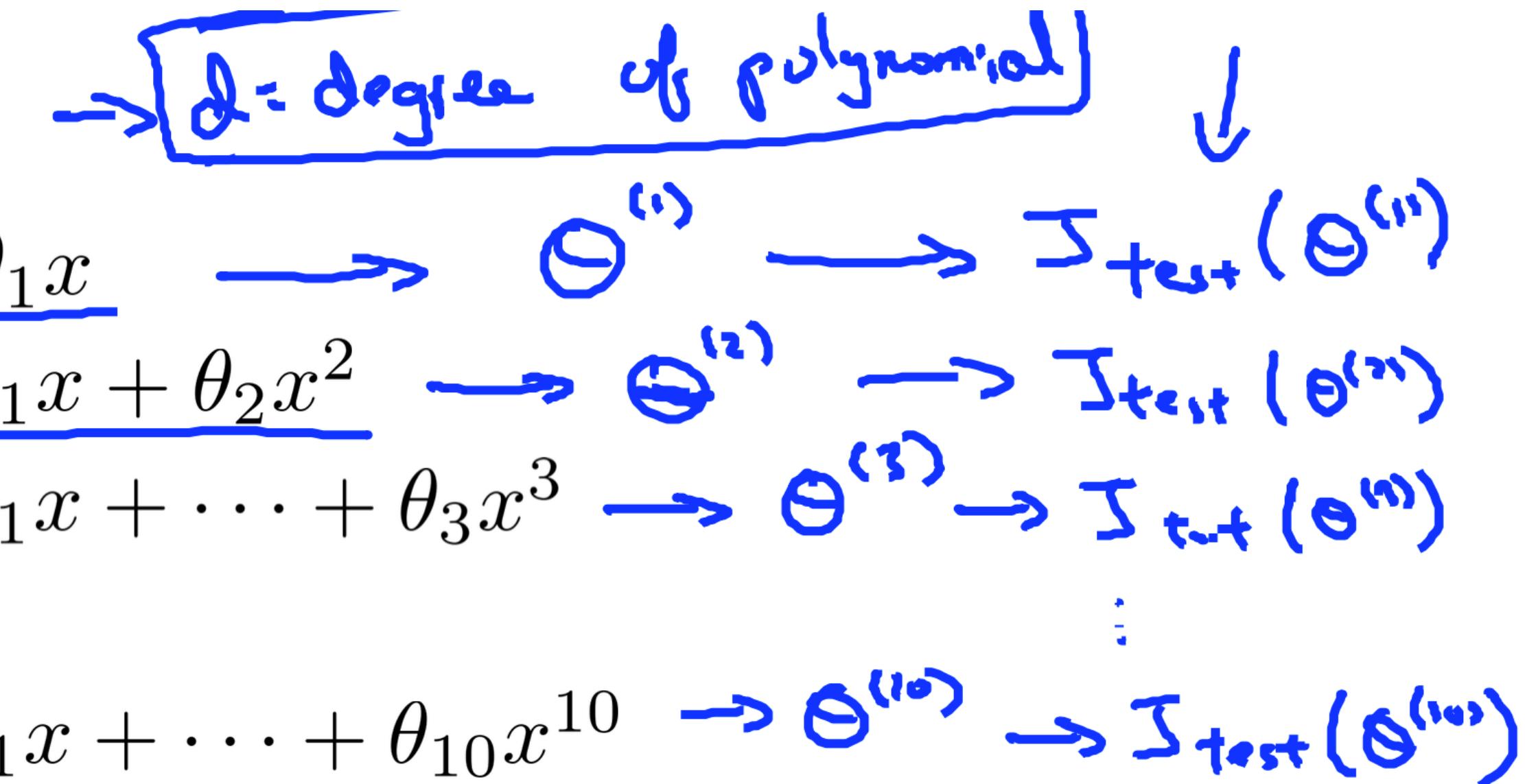
$$\vdots$$

$$d=10 \quad 10. \quad h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x^{10}$$

Choose $\theta_0 + \dots + \theta_5 x^5$

How well does the model generalize? Report test set error $J_{test}(\theta^{(5)})$.

Problem: $J_{test}(\theta^{(5)})$ is likely to be an optimistic estimate of generalization error. I.e. our extra parameter (d = degree of polynomial) is fit to test set.



$\Theta_0, \Theta_1, \dots$

Evaluating a Learning Algorithm

Model selection and train/validation/test sets

Evaluating your hypothesis

Dataset:

Size	Price
2104	400
1600	330
2400	369
1416	232
3000	540
1985	300
1534	315
1427	199
1380	212
1494	243

Annotations on the table:

- Top row: "60%" with a bracket under the first six rows.
- Middle row: "20%" with a bracket under the next two rows.
- Bottom row: "20%" with a bracket under the last two rows.
- Arrows point from the "Training set" bracket to the top section of the table, from the "Cross validation (cv)" bracket to the middle section, and from the "test set" bracket to the bottom section.
- Handwritten labels:
 - "Training set" next to the first bracket.
 - "Cross validation (cv)" next to the second bracket.
 - "test set" next to the third bracket.
 - $m_{cv} = \text{no. of cv examples}$ next to the middle section.
 - $(x_{cv}^{(i)}, y_{cv}^{(i)})$ next to the first example in the middle section.
 - m_{test} next to the bottom section.

Train/validation/test error

Training error:

$$\Rightarrow J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

Cross Validation error:

$$\Rightarrow J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_\theta(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$$

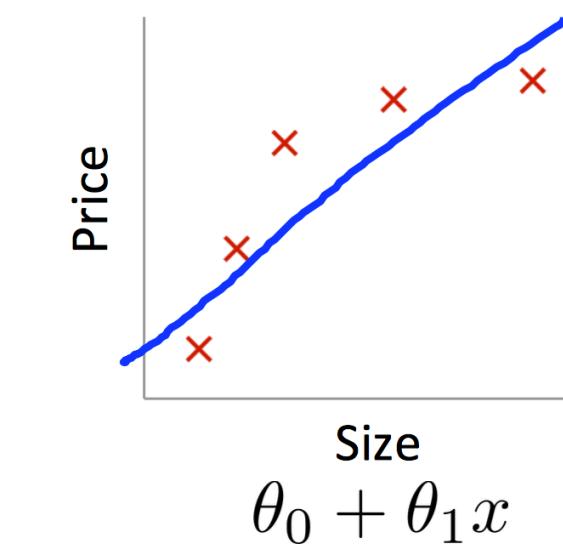
Test error:

$$\Rightarrow J_{test}(\theta) = \frac{1}{2m_{test}} \sum_{i=1}^{m_{test}} (h_\theta(x_{test}^{(i)}) - y_{test}^{(i)})^2$$

Dianosing Bias vs. Variance

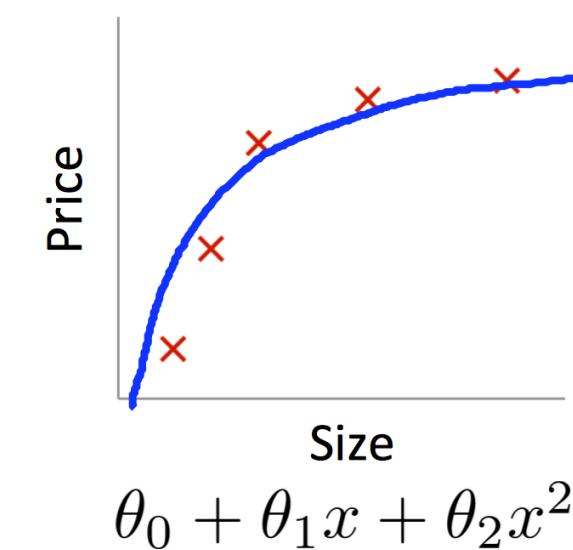
Bias vs. Variance

Bias/variance



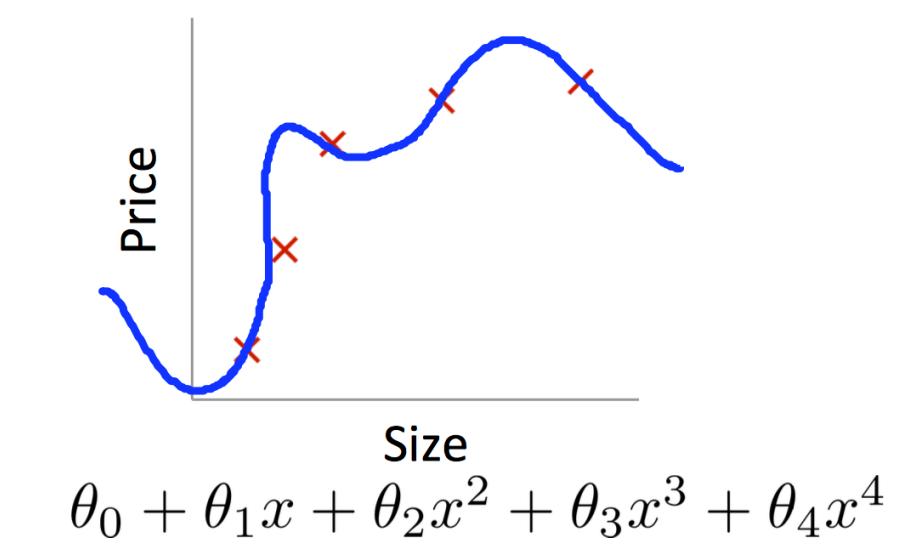
High bias
(underfit)

$d=1$



"Just right"

$d=2$



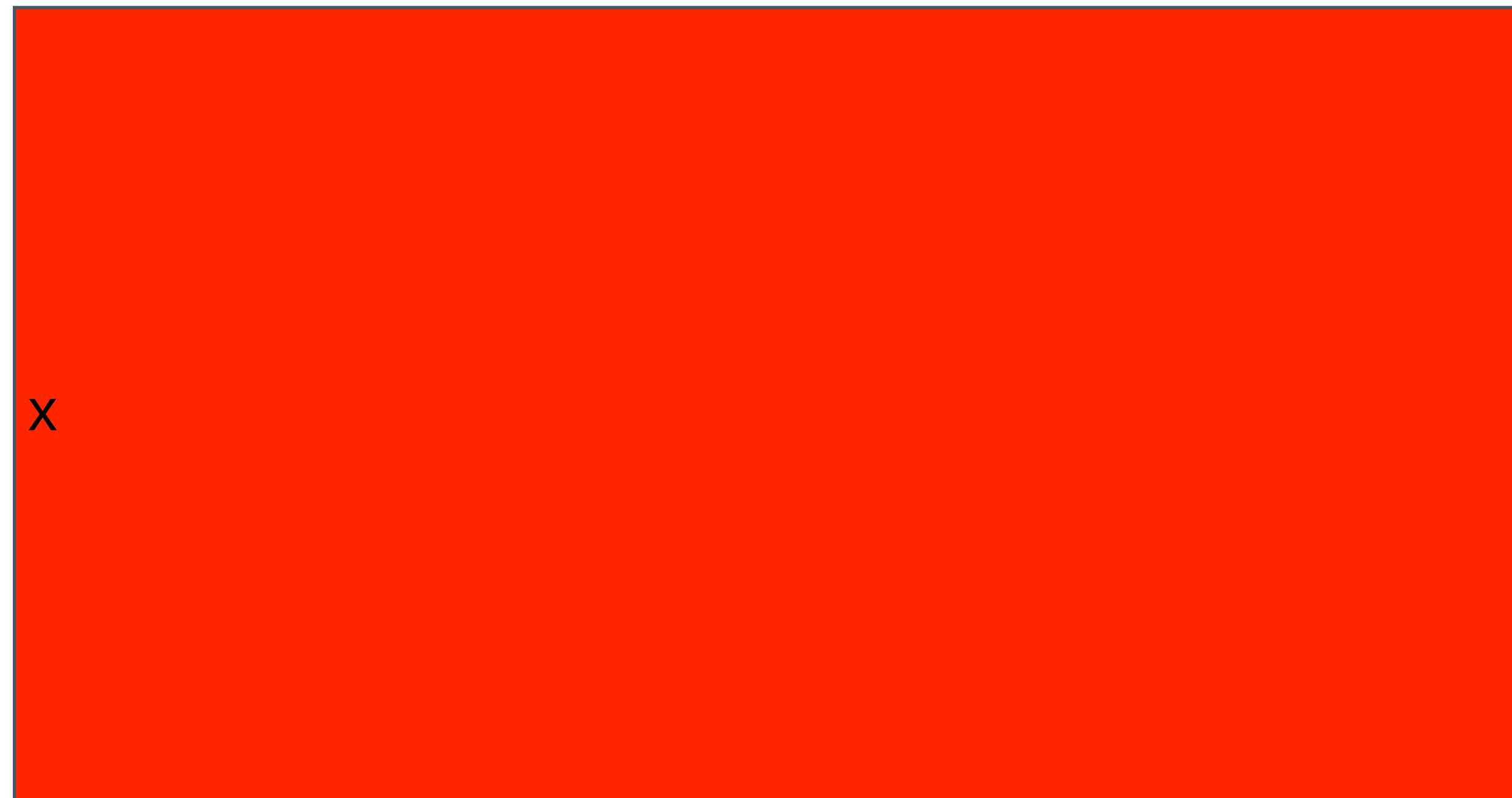
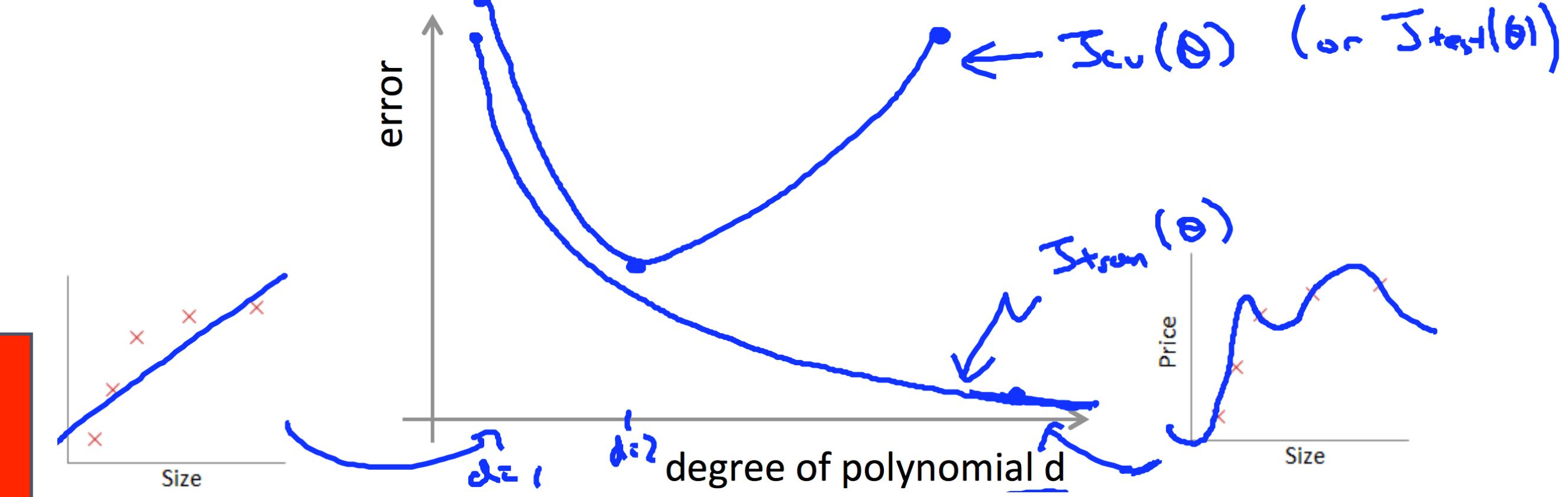
High variance
(overfit)

$d=4$

Bias/variance

Training error: $J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$

Cross validation error: $J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_\theta(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$
(or $J_{test}(\theta)$)

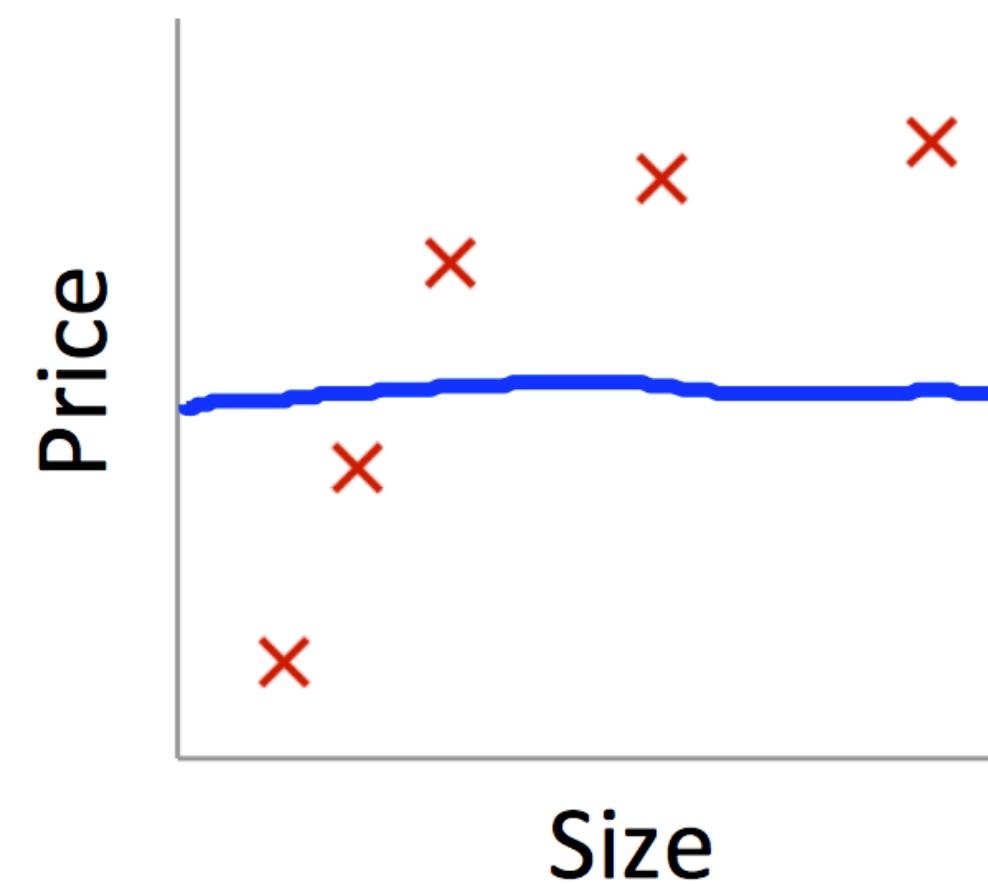


Bias vs. Variance

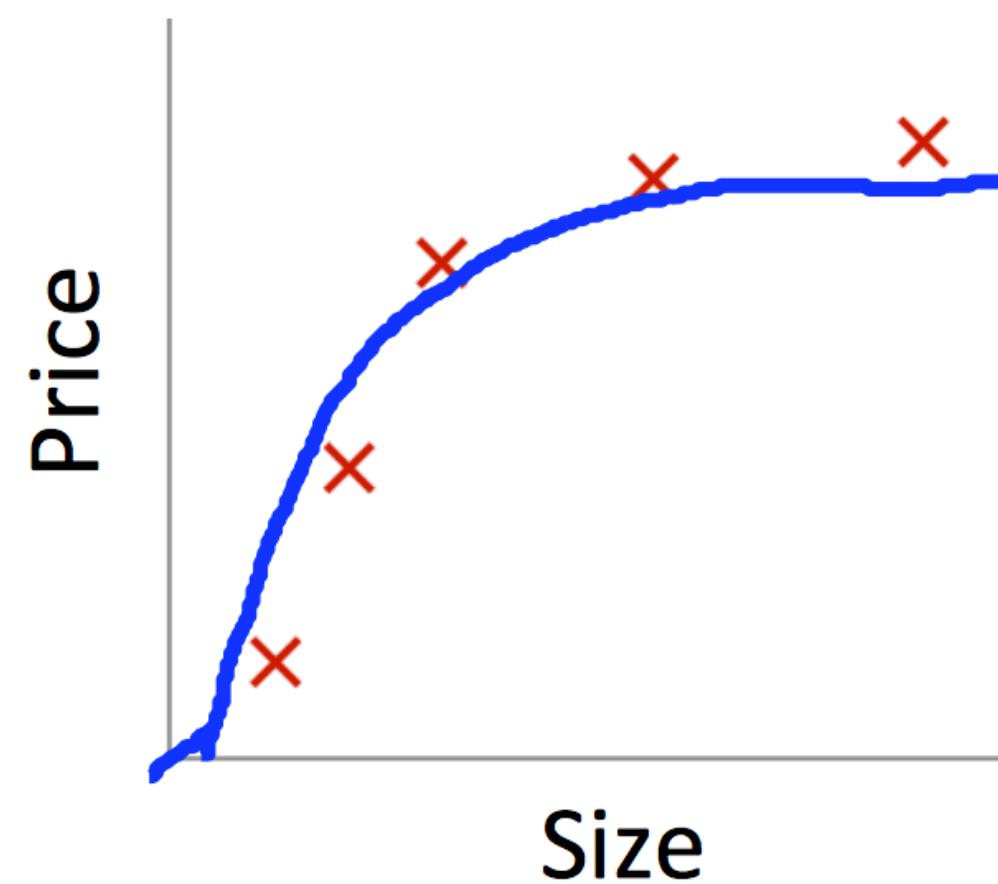
Linear regression with regularization

Model:
$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

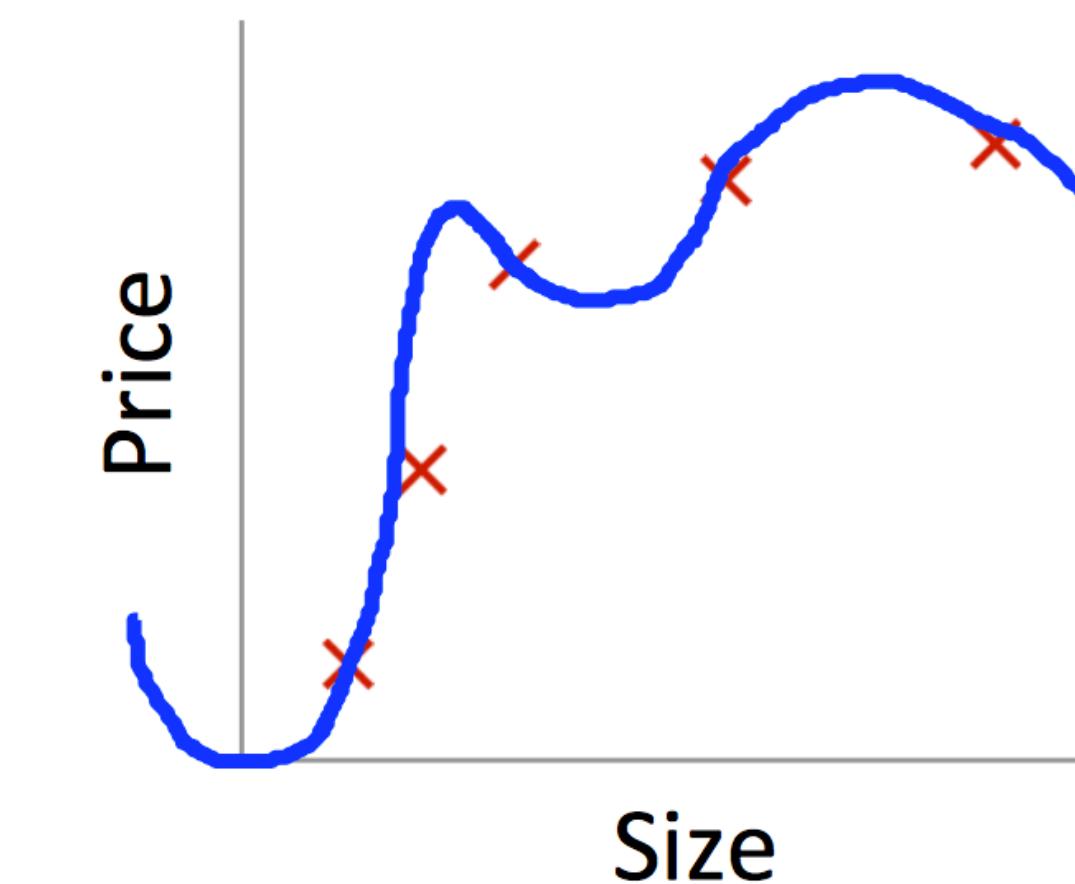
$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^m \theta_j^2$$



Large λ ←
→ High bias (underfit)
 $\lambda = 10000$. $\theta_1 \approx 0, \theta_2 \approx 0, \dots$
 $h_{\theta}(x) \approx \theta_0$



Intermediate λ ←
“Just right”



→ Small λ
 High variance (overfit)
 $\rightarrow \lambda = 0$

Bias vs. Variance

Choosing the regularization parameter λ

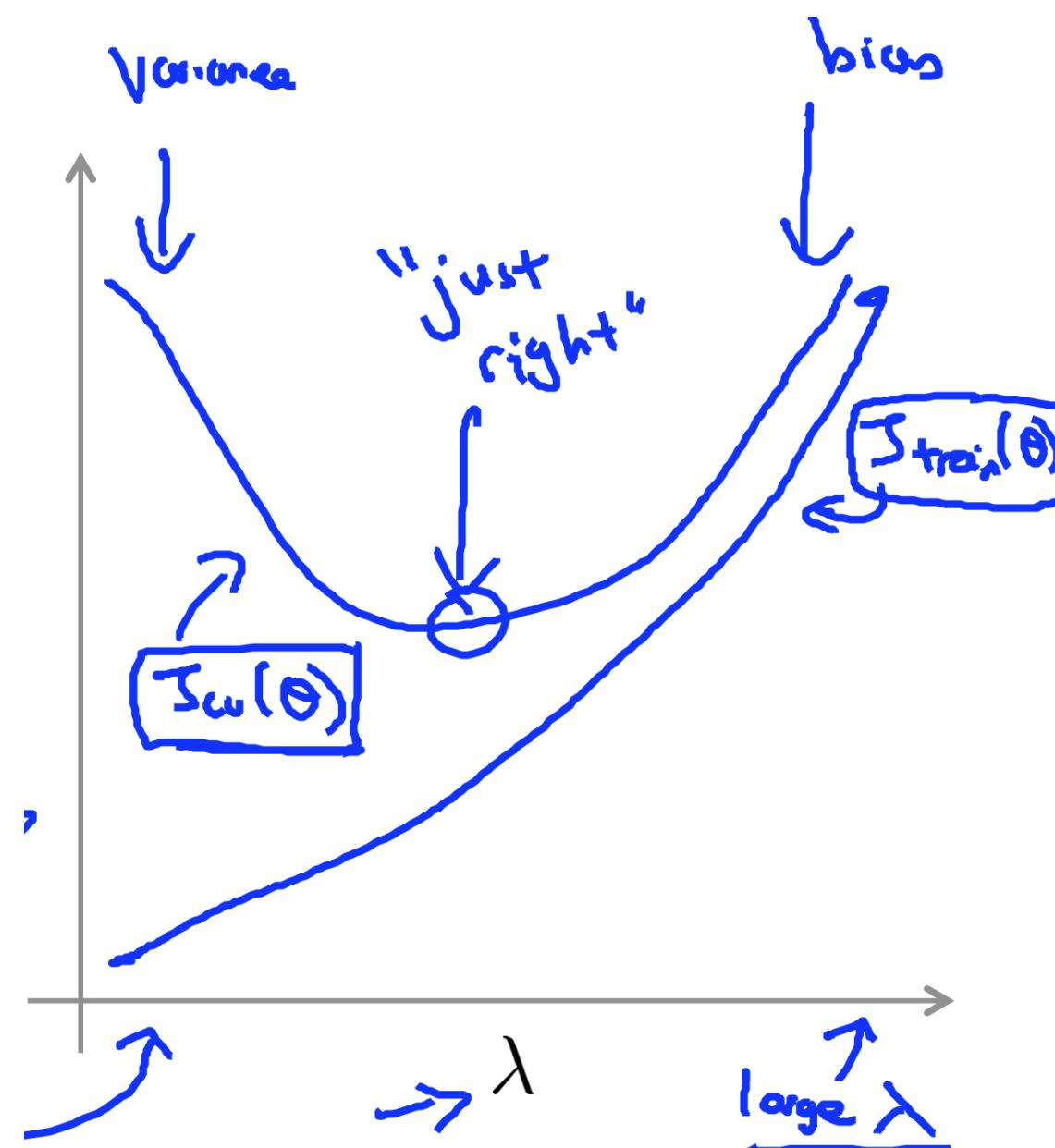
$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4 \quad \leftarrow$$

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^m \theta_j^2 \quad \leftarrow$$

$$\rightarrow J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \quad J(\theta)$$

$$J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_{\theta}(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$$

$$J_{test}(\theta) = \frac{1}{2m_{test}} \sum_{i=1}^{m_{test}} (h_{\theta}(x_{test}^{(i)}) - y_{test}^{(i)})^2$$

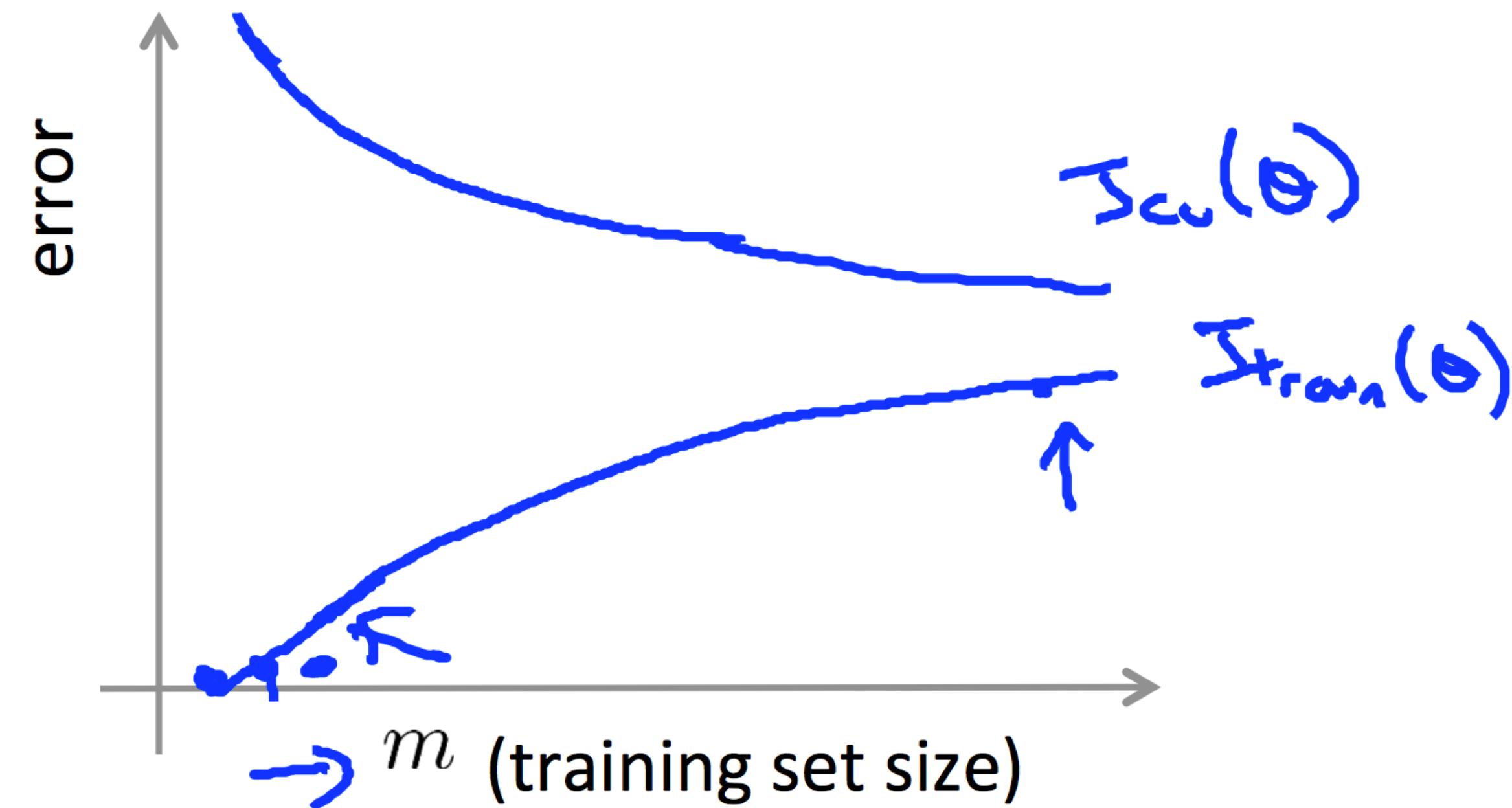


Evaluating a Learning Algorithm

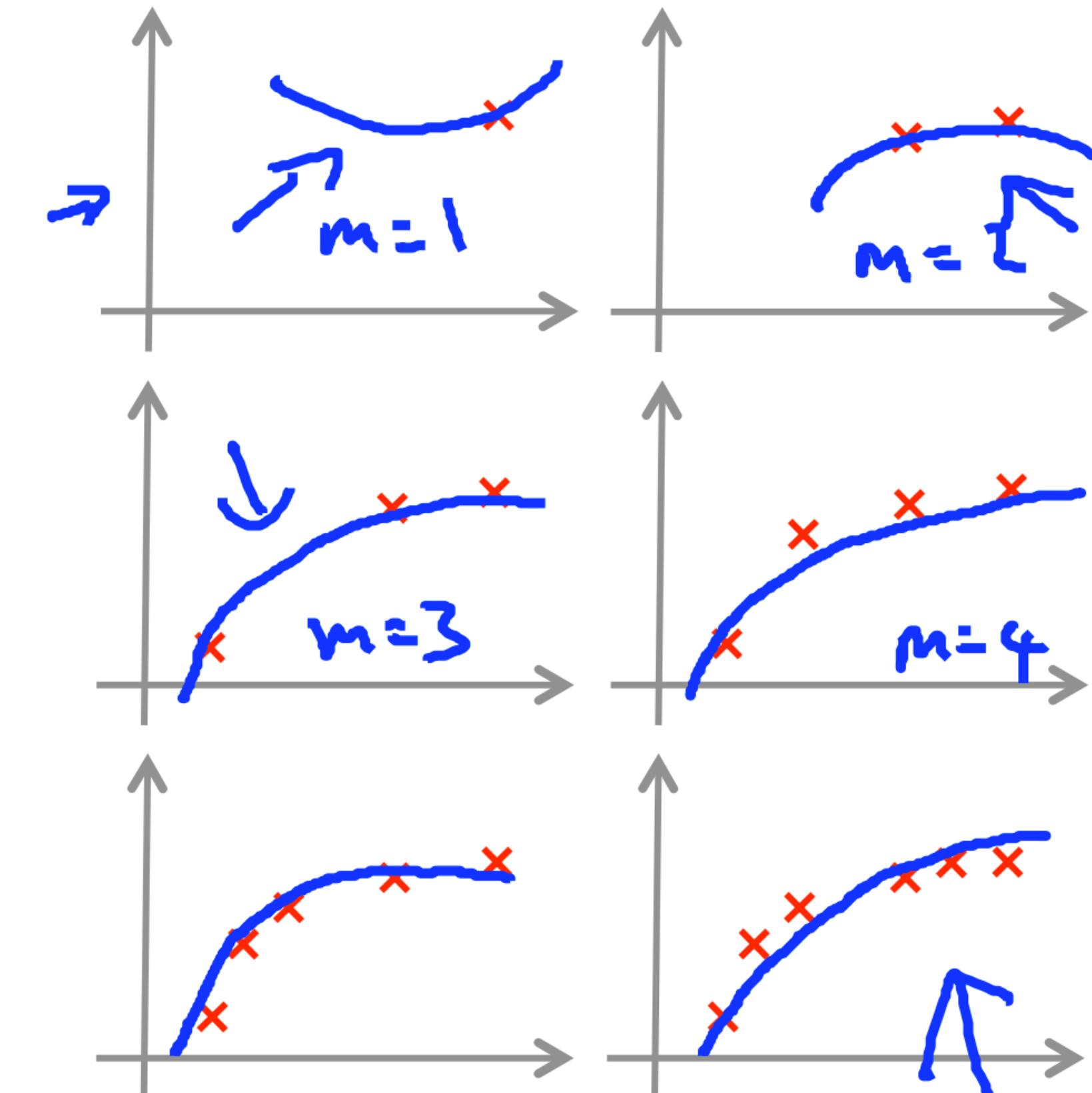
Learning curves

$$\rightarrow \underline{J_{train}(\theta)} = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$\rightarrow J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_\theta(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$$

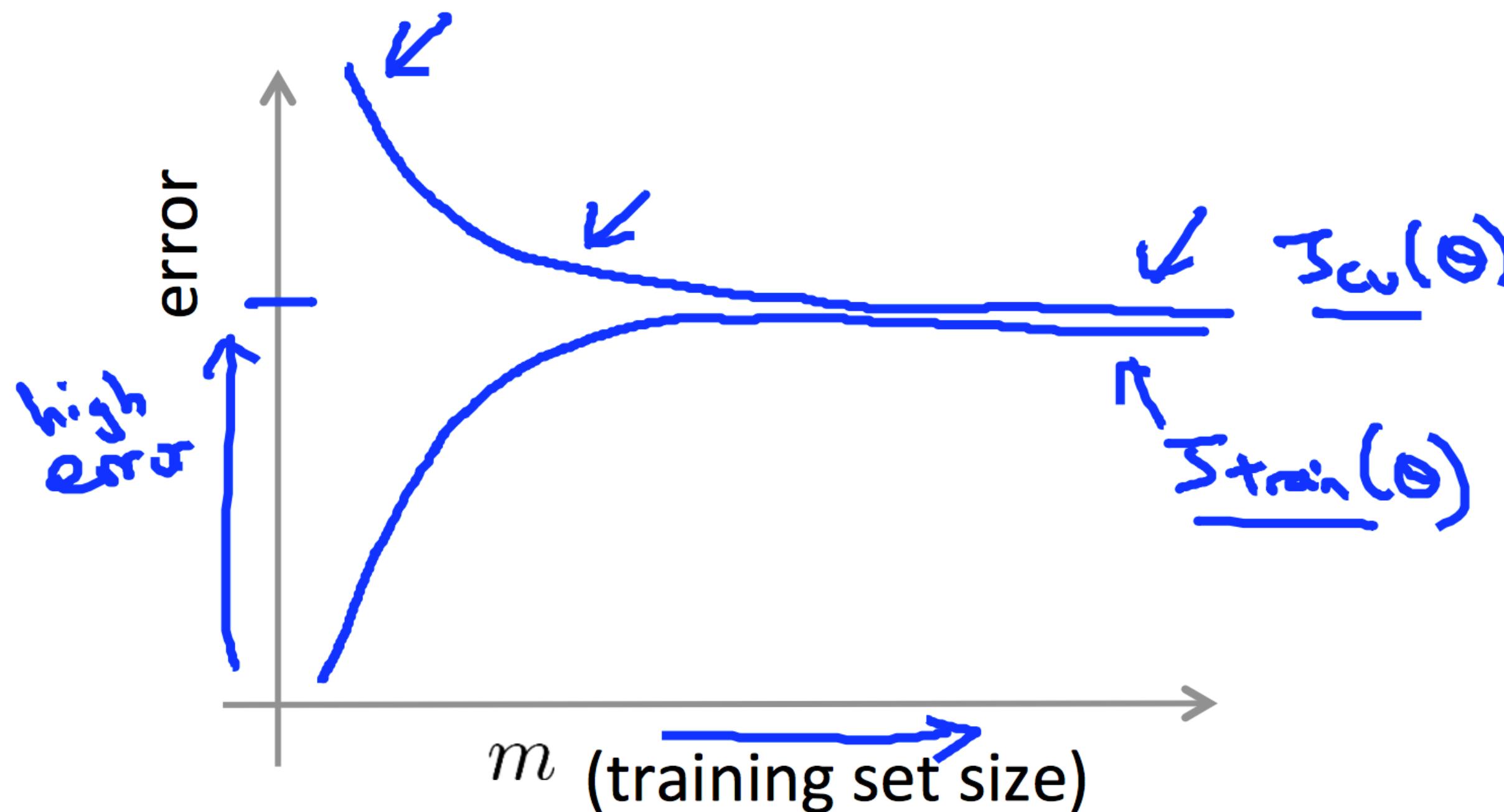


$$h_\theta(x) = \underline{\theta_0 + \theta_1 x + \theta_2 x^2}$$

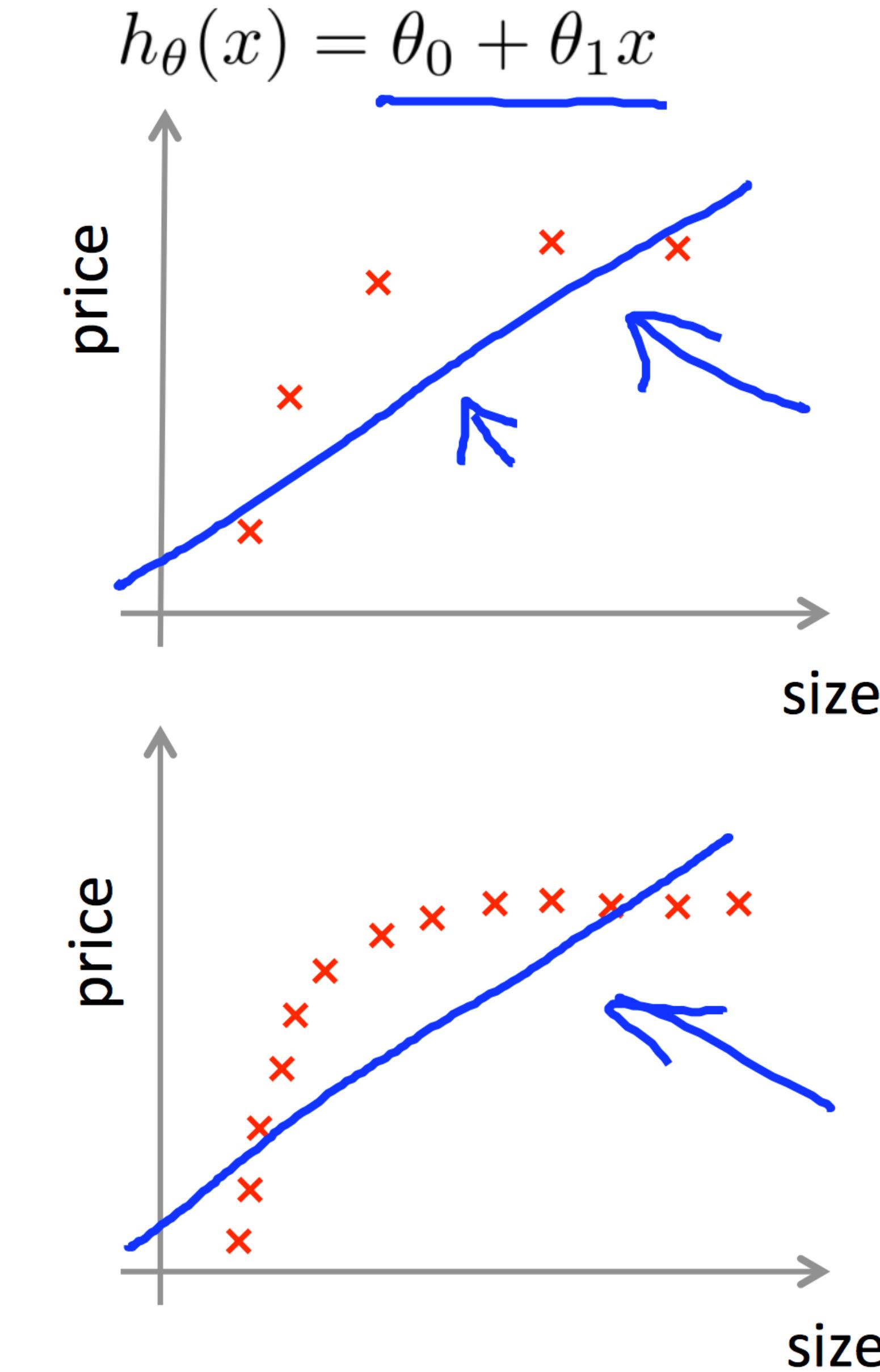


Evaluating a Learning Algorithm

High bias

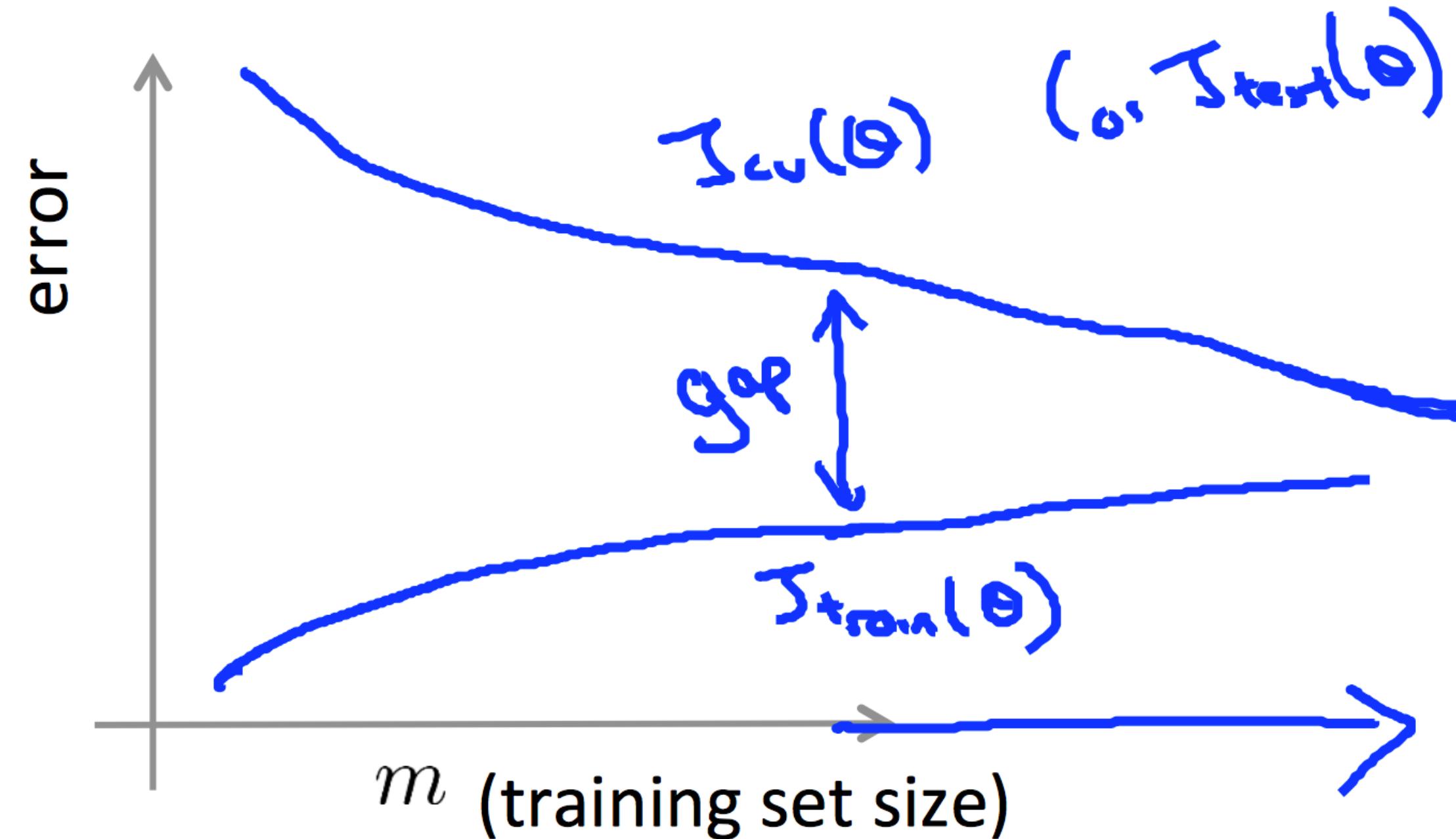


If a learning algorithm is suffering from high bias, getting more training data will not (by itself) help much.



Evaluating a Learning Algorithm

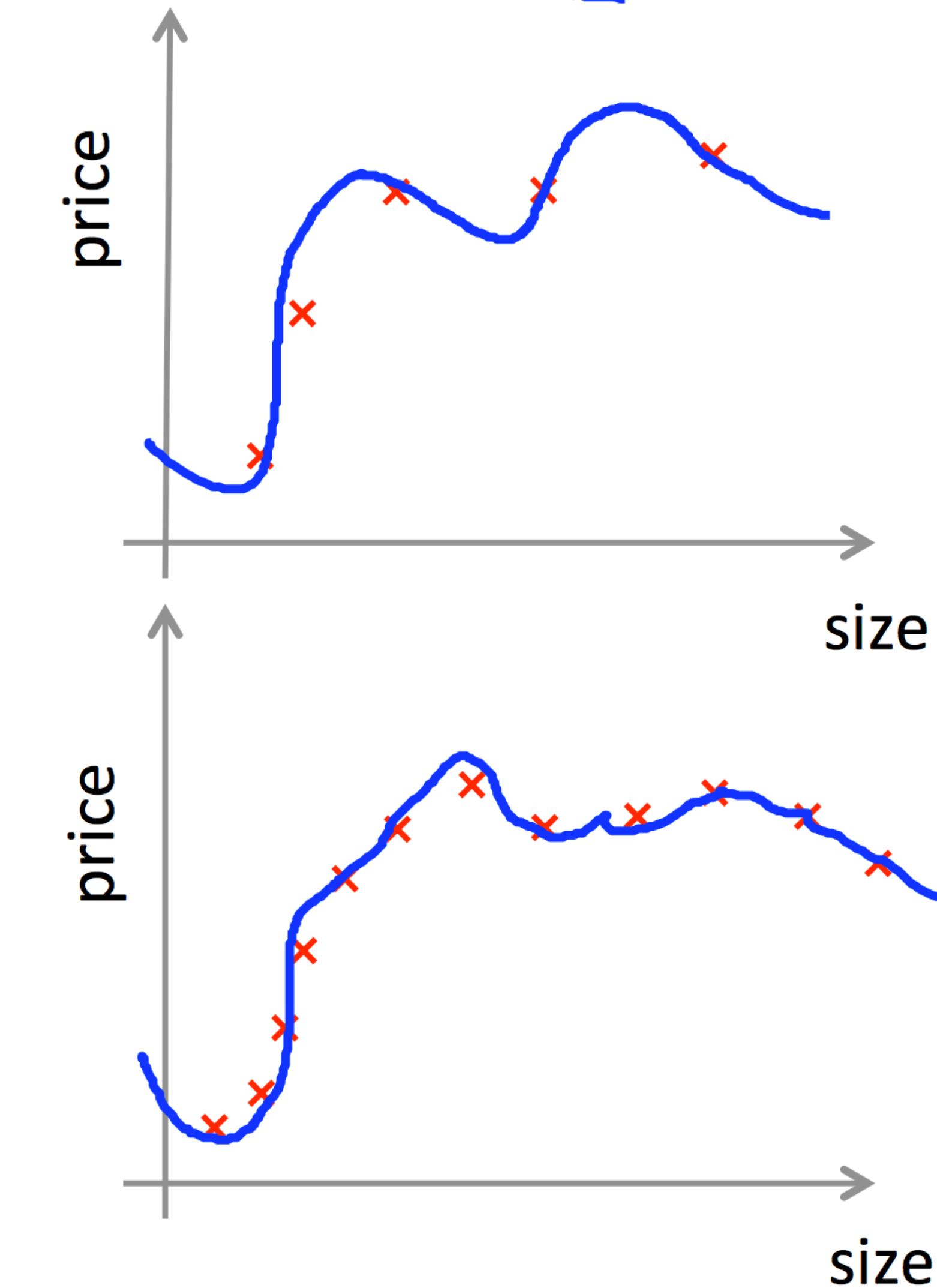
High variance



If a learning algorithm is suffering from high variance, getting more training data is likely to help. ↫

$$h_\theta(x) = \theta_0 + \theta_1 x + \dots + \theta_{100} x^{100}$$

(and small λ)



Evaluating a Learning Algorithm

Deciding what to do next revisited

Debugging a learning algorithm:

Suppose you have implemented regularized linear regression to predict housing prices. However, when you test your hypothesis in a new set of houses, you find that it makes unacceptably large errors in its prediction. What should you try next?

x

Evaluating a Learning Algorithm

Deciding what to do next revisited

Debugging a learning algorithm:

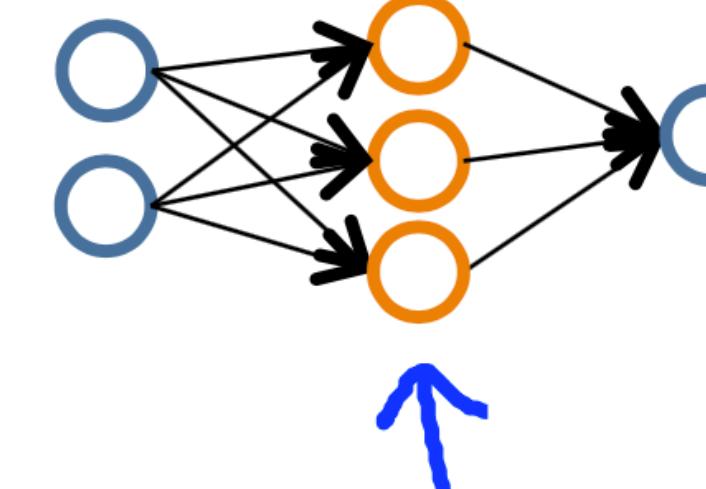
Suppose you have implemented regularized linear regression to predict housing prices. However, when you test your hypothesis in a new set of houses, you find that it makes unacceptably large errors in its prediction. What should you try next?

x

Evaluating a Learning Algorithm

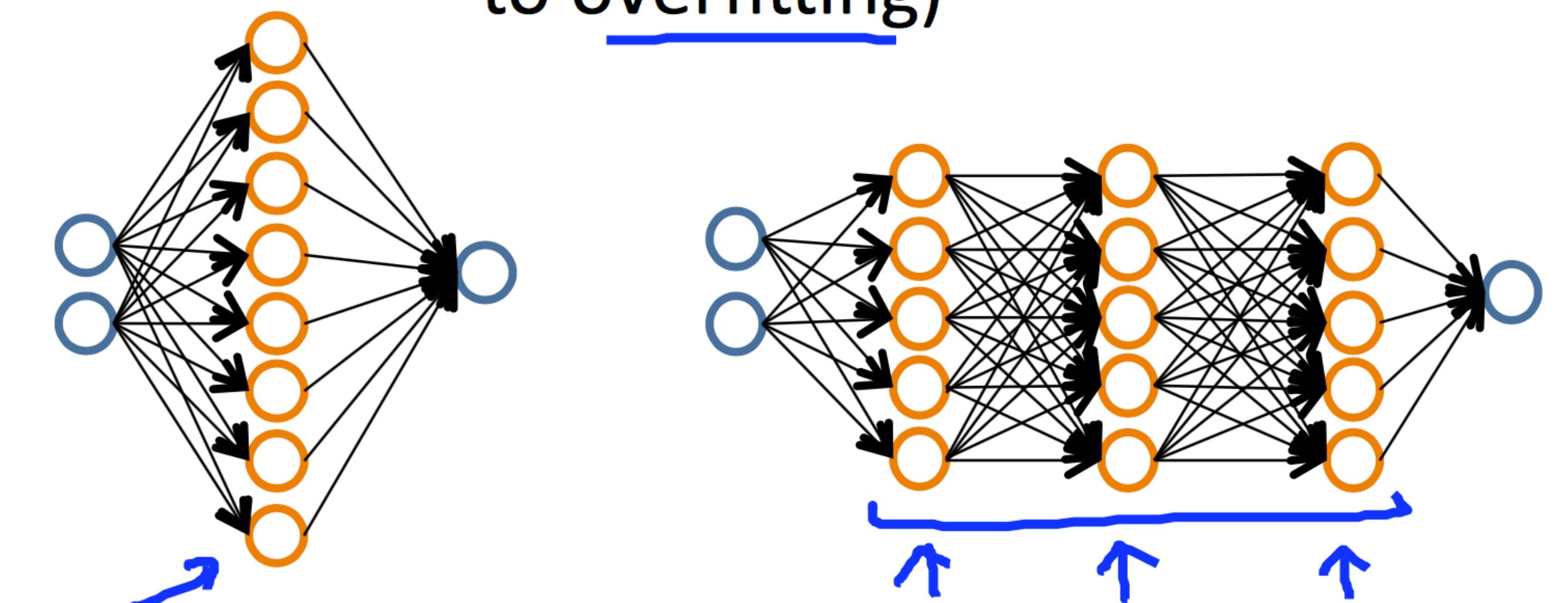
Deciding what to do next revisited

Neural networks and overfitting



Computationally cheaper

→ “Large” neural network
(more parameters; more prone
to overfitting)



$I_{c_0}(\Theta)$

↑

Machine Learning System Design

Example Building a spam classifier

Building a spam classifier

From: `cheapsales@buystufffromme.com`
To: `ang@cs.stanford.edu`
Subject: Buy now!

Deal of the week! Buy now!
Rolex w4tchs - \$100
[redacted] (any kind) - \$50
Also low cost M0rgages
available.

Spam (1)

From: Alfred Ng
To: `ang@cs.stanford.edu`
Subject: Christmas dates?

Hey Andrew,
Was talking to Mom about plans
for Xmas. When do you get off
work. Meet Dec 22?
Alf

Non-spam (0)

Machine Learning System Design

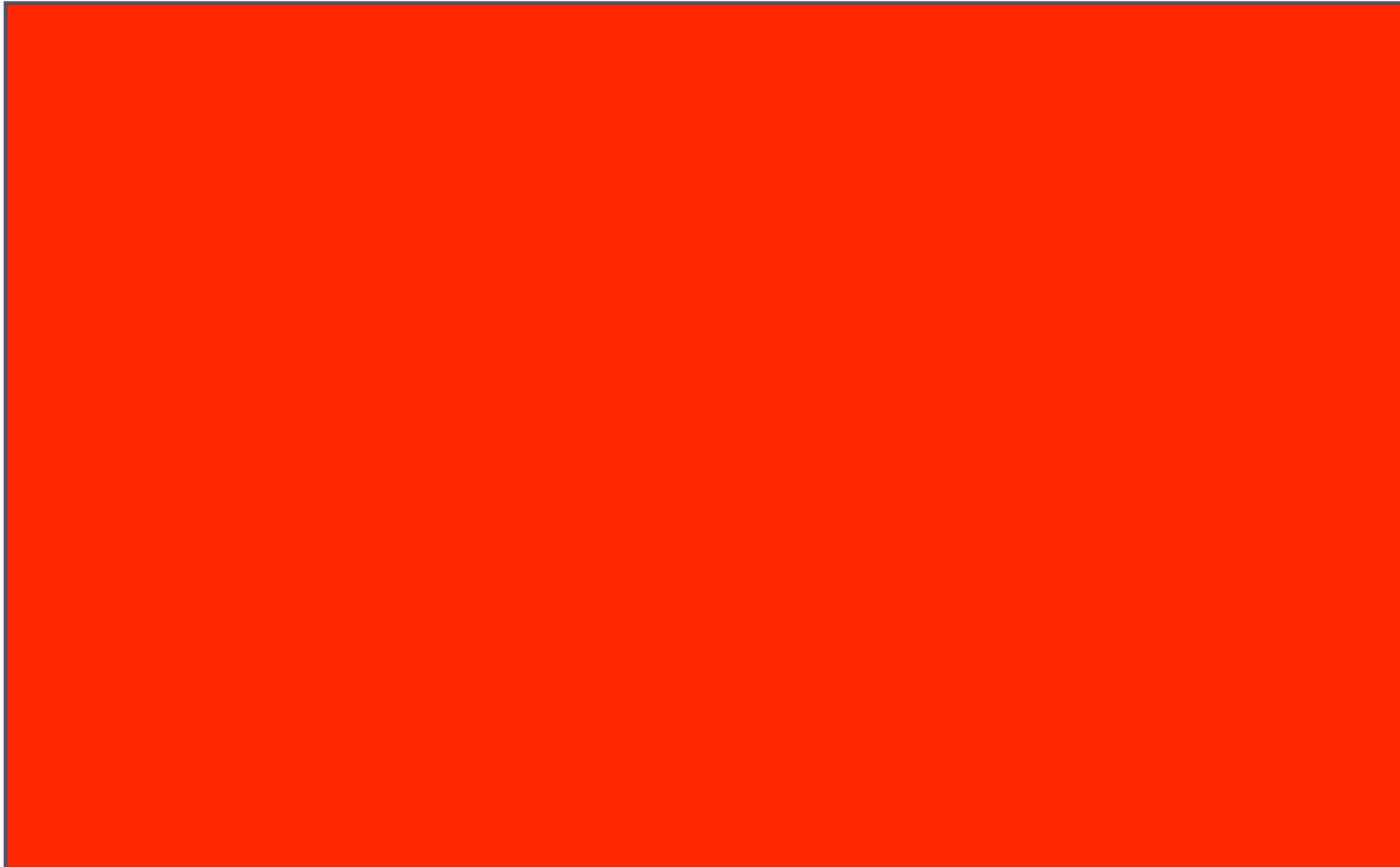
Example Building a spam classifier

Building a spam classifier

Supervised learning. $x = \text{features of email}$. $y = \text{spam (1) or not spam (0)}$.

Features x : Choose 100 words indicative of spam/not spam.

E.g. deal, buy, discount, andrew, now, ...



$$x_j = \begin{cases} 1 & \text{if word } j \text{ appears in email} \\ 0 & \text{otherwise.} \end{cases}$$

From: `cheapsales@buystufffromme.com`
To: `ang@cs.stanford.edu`
Subject: Buy now!

Deal of the week! Buy now!

Note: In practice, take most frequently occurring n words (10,000 to 50,000) in training set, rather than manually pick 100 words.

Machine Learning System Design

Example Building a spam classifier

Building a spam classifier

- E.g. “honeypot” project.
- based on email routing information ().
- for e.g. should “discount” and “discounts” be treated as the same word? How about “deal” and “Dealer”? Features about punctuation?
- to (e.g. m0rtgage, med1cine, w4tches.)

Machine Learning System Design

Error analysis

Recommended approach

- **Implement it and test it on your cross-validation data.**
- **Use cross-validation** to decide if more data, more features, etc. are likely to help.
- **Review the errors** Manually examine the examples (in cross validation set) that your algorithm made errors on. See if you spot any systematic trend in what type of examples it is making errors on.

Machine Learning System Design

Error analysis

Error Analysis

$m_{CV} = 500$ examples in cross validation set

Algorithm misclassifies 100 emails.

[REDACTED], and categorize them based on:

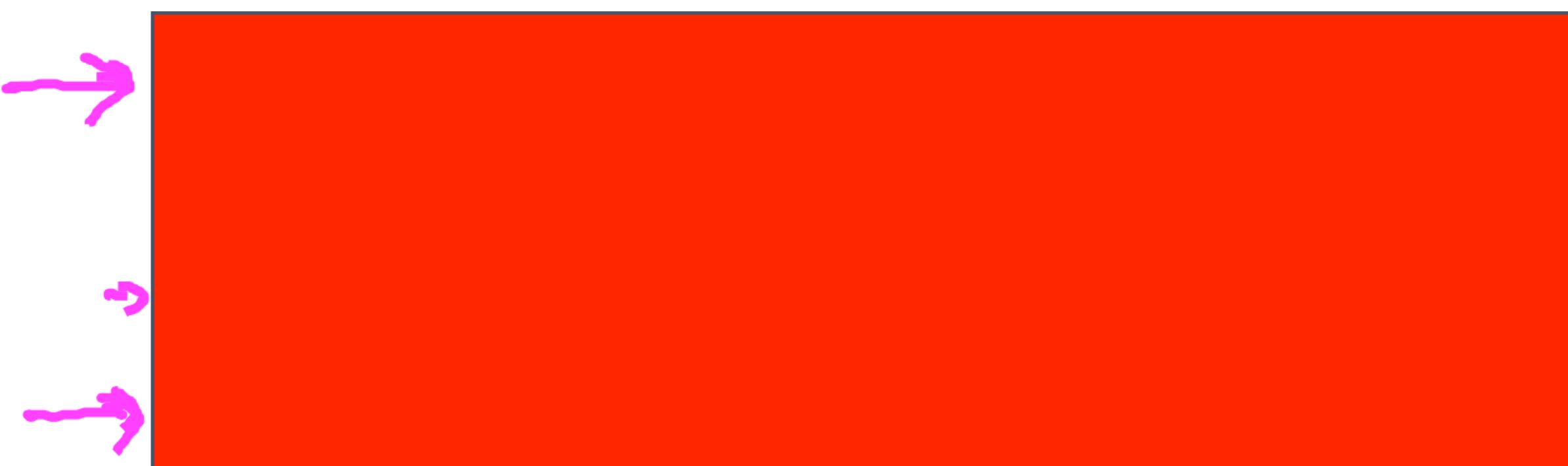
- (i) What type of email it is *pharma, replica, steal passwords, ...*
- (ii) What cues (features) you think would have helped the algorithm classify them correctly.

Pharma: 12

Replica/fake: 4

Steal passwords: 53

Other: 31



Machine Learning System Design

Handling skewed data

Cancer classification example

Train logistic regression model $h_{\theta}(x)$. ($y = 1$ if cancer, $y = 0$ otherwise)

Find that you got 1% error on test set.
(99% correct diagnoses)

Only 0.50% of patients have cancer.

skewed classes.

```
function y = predictCancer(x)
    → y = 0; %ignore x!
    return
```

0.5% error
→ 99.2% away (0.8% error)
→ 99.5% away (0.5% error)

Machine Learning System Design

Handling skewed data

Precision/Recall

$y = 1$ in presence of rare class that we want to detect

Actual class	
Predicted class	1
1	True positive
0	False positive
0	False negative
	True negative

$$y=0 \\ \text{recall} = 0$$

→ Precision

(Of all patients where we predicted $y = 1$, what fraction actually has cancer?)

$$\frac{\text{True positives}}{\#\text{predicted positive}} = \frac{\text{True positive}}{\text{True pos} + \text{False pos}}$$

→ Recall

(Of all patients that actually have cancer, what fraction did we correctly detect as having cancer?)

$$\frac{\text{True positives}}{\#\text{actual positives}} = \frac{\text{True positives}}{\text{True pos} + \text{False neg}}$$

Machine Learning System Design

Trading off precision and recall

Trading off precision and recall

→ Logistic regression: $0 \leq h_\theta(x) \leq 1$

Predict 1 if $h_\theta(x) \geq 0.5$ ~~0.7~~ ~~0.9~~ ~~0.3~~ ↗

Predict 0 if $h_\theta(x) < 0.5$ ~~0.7~~ ~~0.9~~ ~~0.3~~

→ Suppose we want to predict $y = 1$ (cancer) only if very confident.

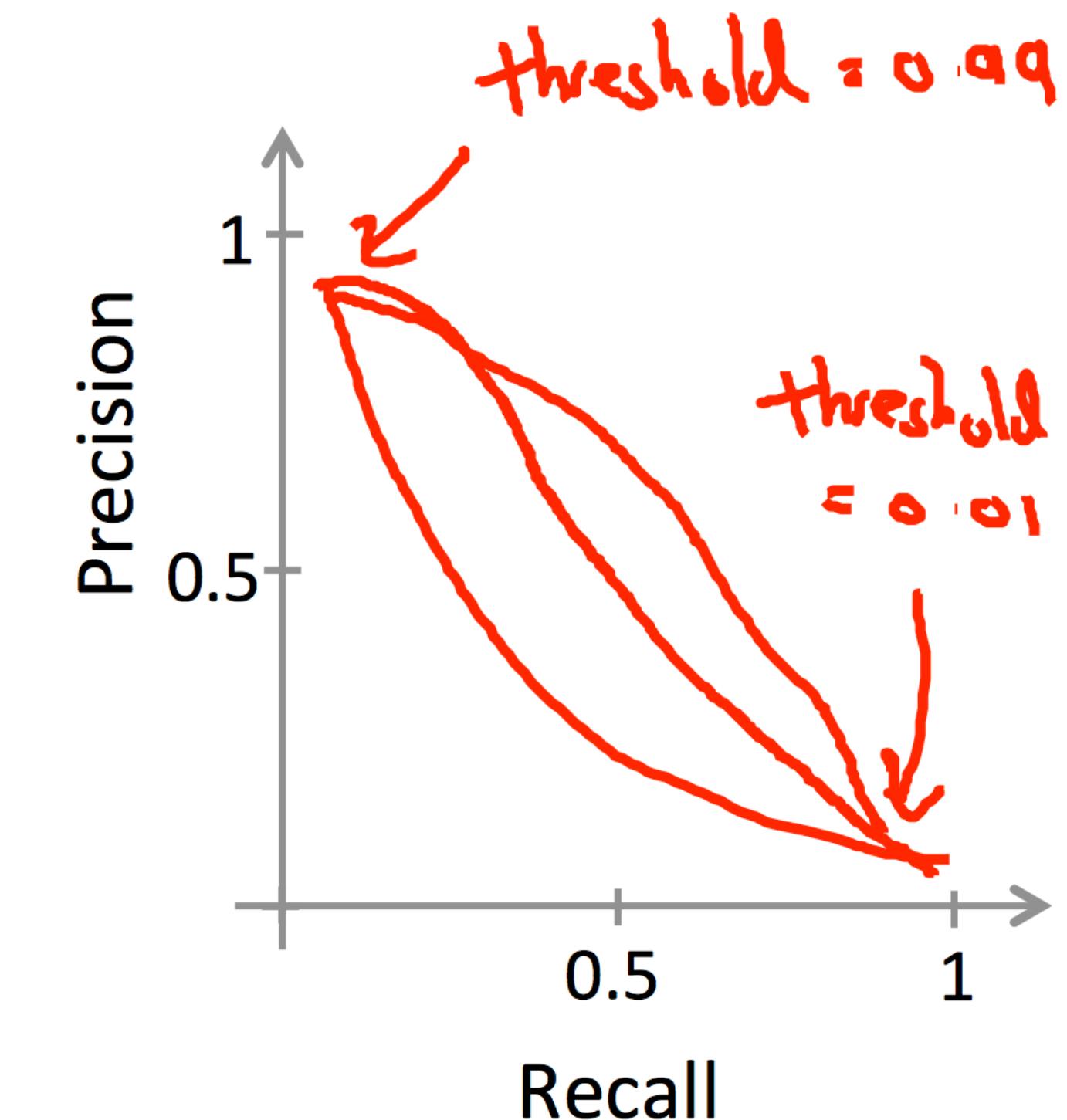
→ Higher precision, lower recall

→ Suppose we want to avoid missing too many cases of cancer (avoid false negatives).

→ Higher recall, lower precision.

More generally: Predict 1 if $h_\theta(x) \geq \text{threshold}$ ↗

$$\rightarrow \text{precision} = \frac{\text{true positives}}{\text{no. of predicted positive}}$$
$$\rightarrow \text{recall} = \frac{\text{true positives}}{\text{no. of actual positive}}$$



Machine Learning System Design

Trading off precision and recall

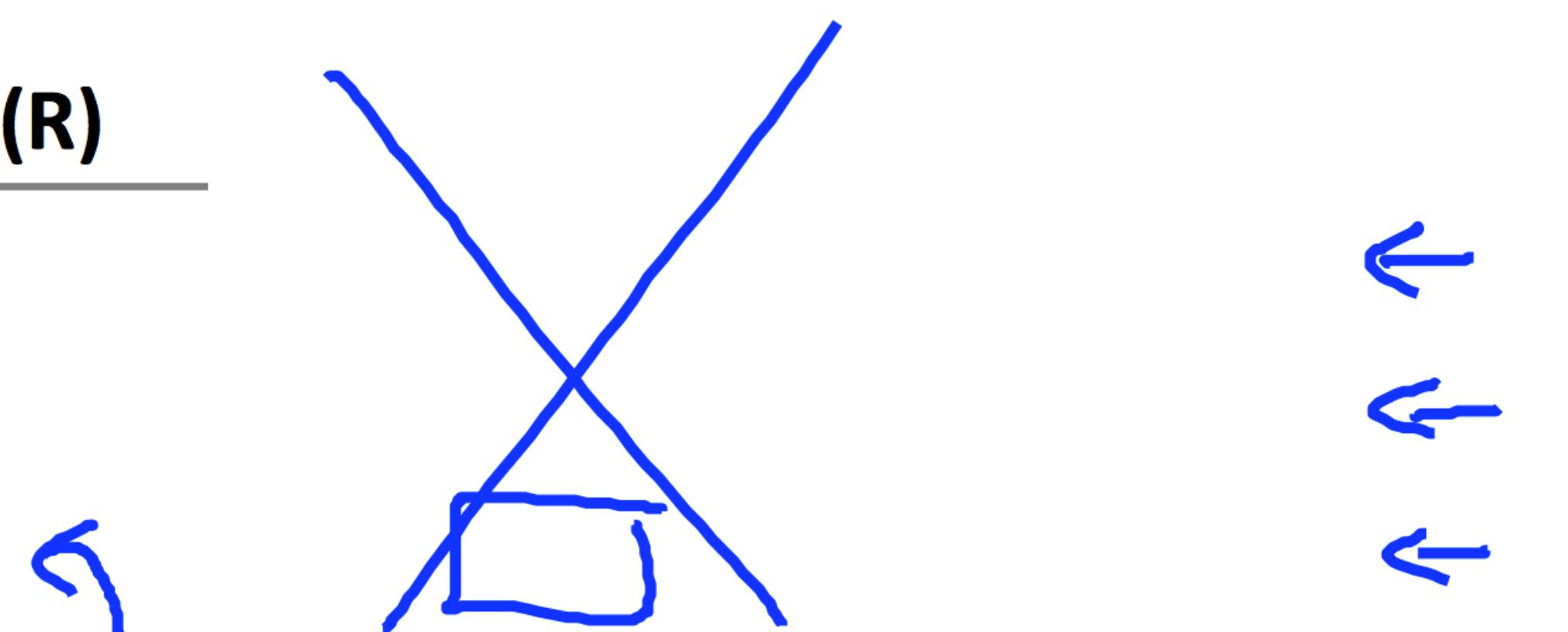
F_1 Score (F score)

How to compare precision/recall numbers?

	Precision(P)	Recall (R)
Algorithm 1	<u>0.5</u>	<u>0.4</u>
Algorithm 2	<u>0.7</u>	<u>0.1</u>
Algorithm 3	<u>0.02</u>	1.0

Average: ~~$\frac{P+R}{2}$~~

$$F_1 \text{ Score: } 2 \frac{PR}{P+R}$$



$$P=0 \quad \text{or} \quad R=0 \quad \Rightarrow F\text{-score} = 0$$

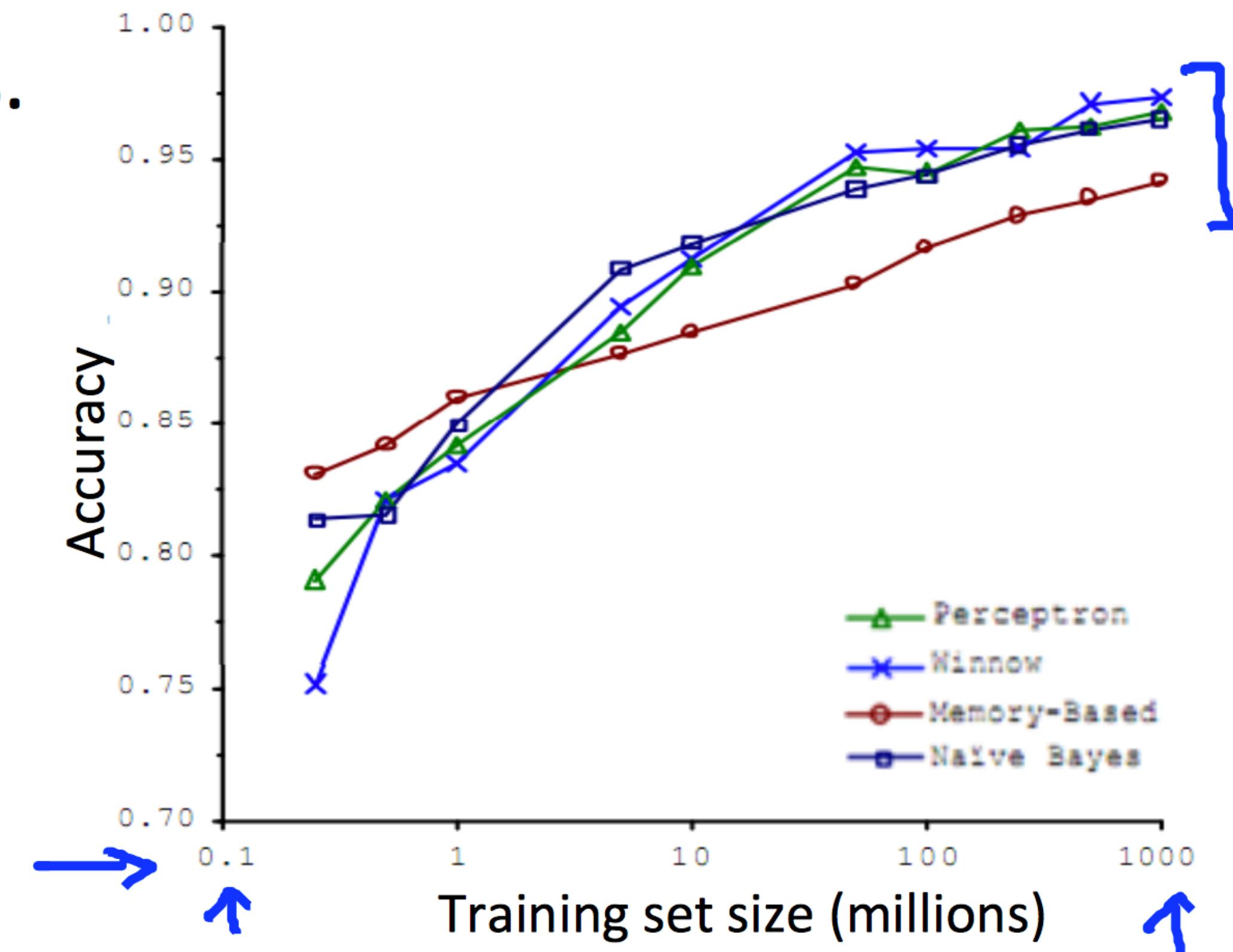
$$P=1 \quad \text{and} \quad R=1 \quad \Rightarrow F\text{-score} = 1$$

Machine Learning System Design

Designing a high accuracy learning system

E.g. Classify between confusable words.

{to, two, too}, {then, than}
For breakfast I ate two eggs.



“It’s not who has the best algorithm that wins.

It’s who has the most data.”

Using large data sets

Machine Learning System Design



$J_{train}(\theta)$ will be small.

low variance ↙

$J_{test}(\theta)$ will be small

Unsupervised Learning

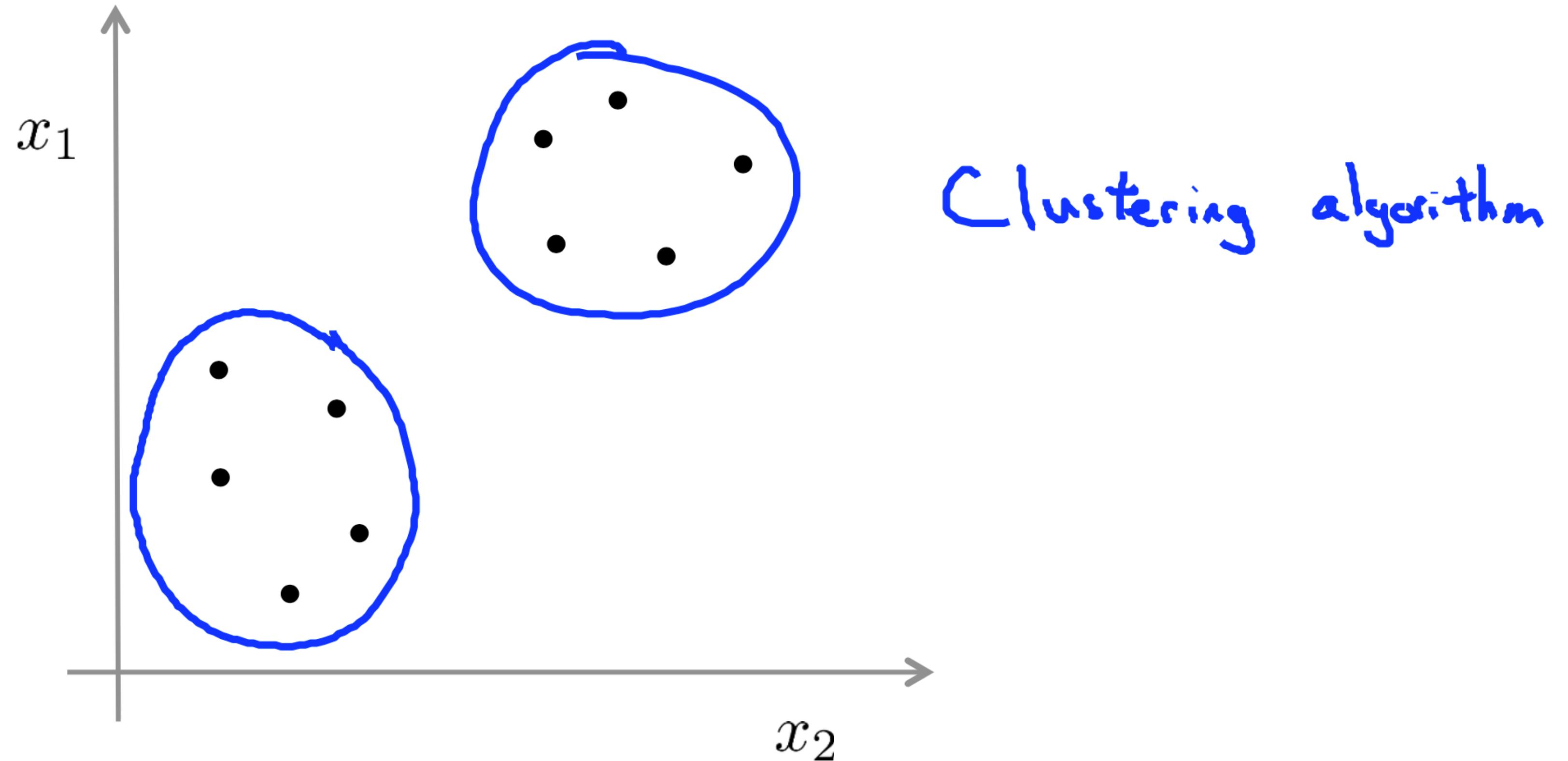
Unsupervised learning, on the other hand, allows us to approach problems with little or no idea what our results should look like. We can derive structure from data where we don't necessarily know the effect of the variables.

We can derive this structure by clustering the data based on relationships among the variables in the data.

With unsupervised learning there is no feedback based on the prediction results, i.e., there is no teacher to correct you.

Clustering

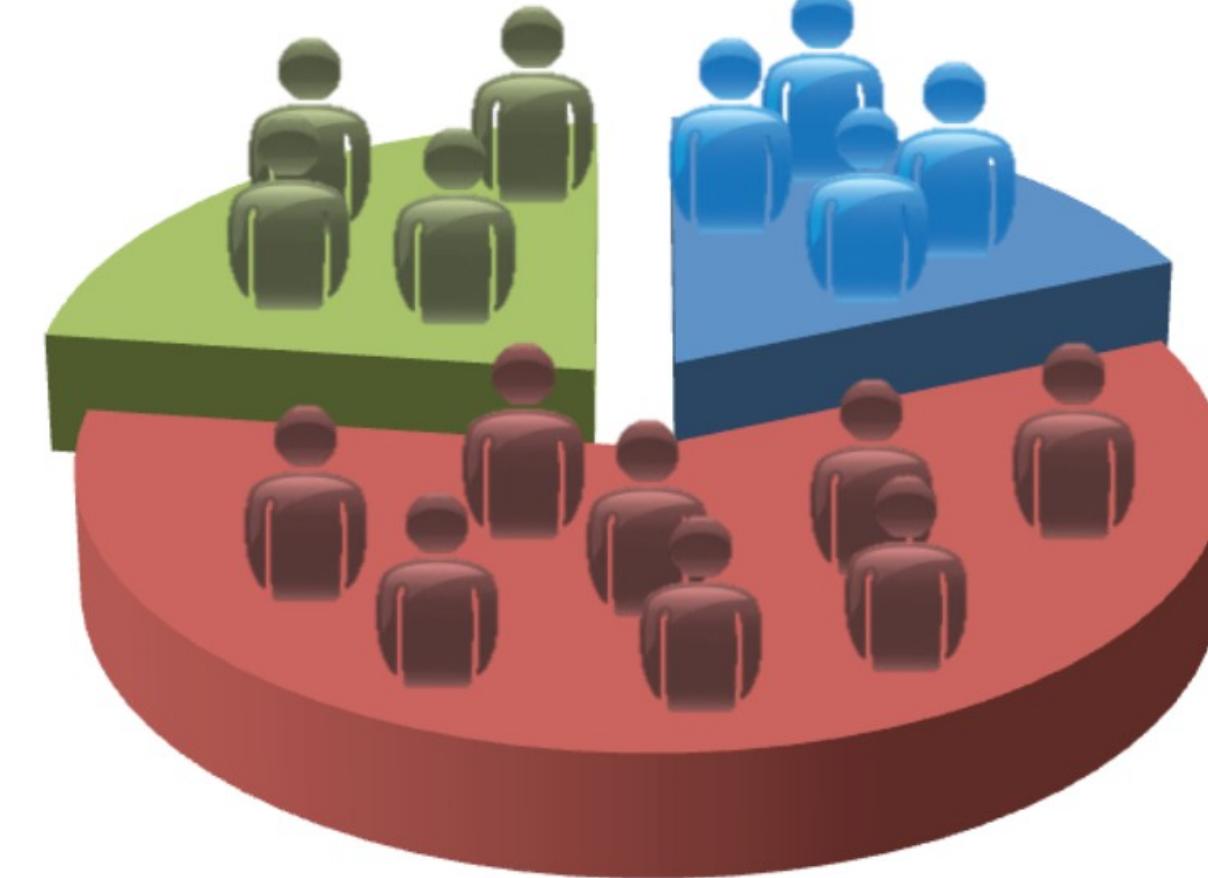
Unsupervised learning



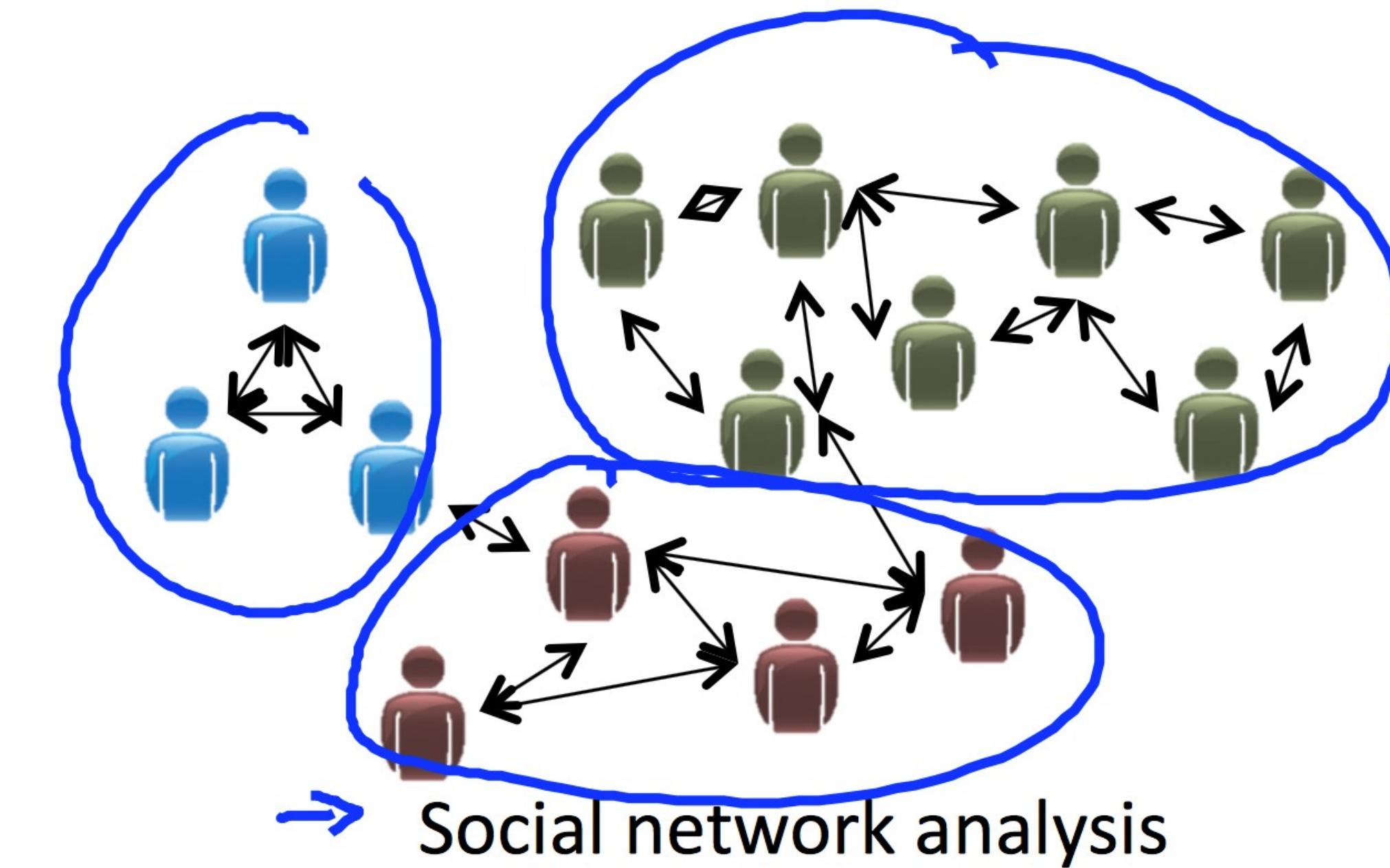
Training set: $\{\underline{x}^{(1)}, \underline{x}^{(2)}, \underline{x}^{(3)}, \dots, \underline{x}^{(m)}\}$ ←

Clustering

Applications of clustering



→ Market segmentation



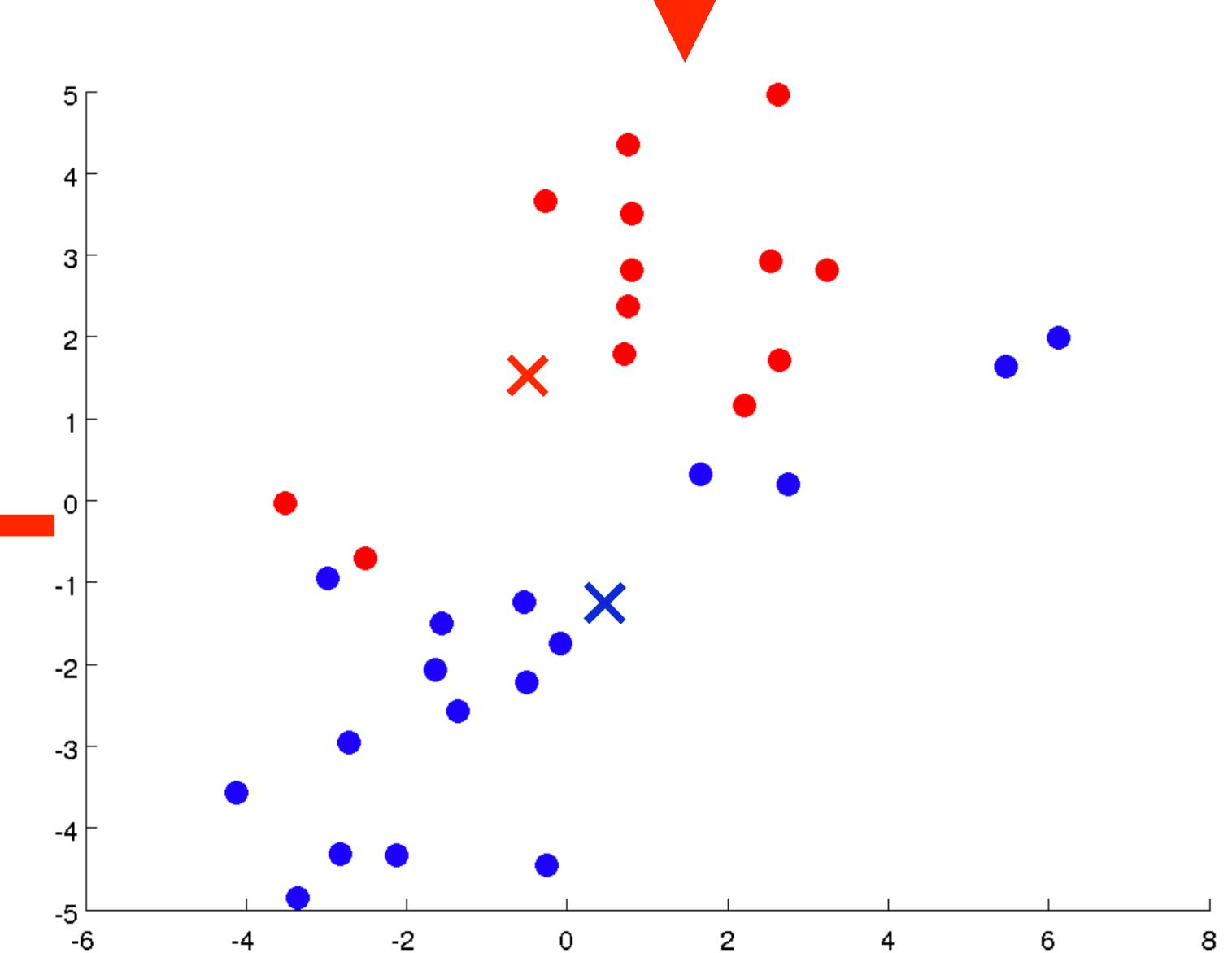
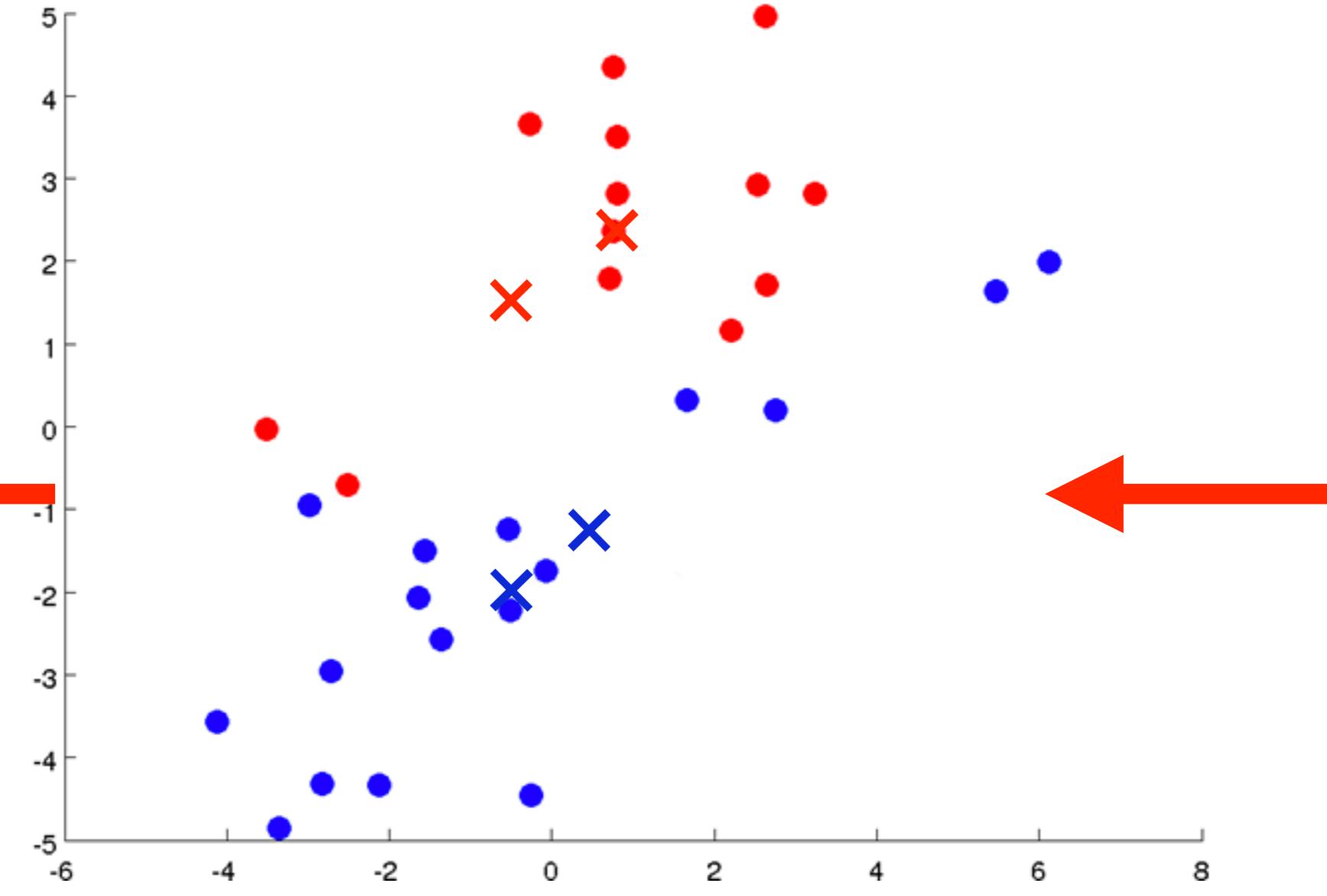
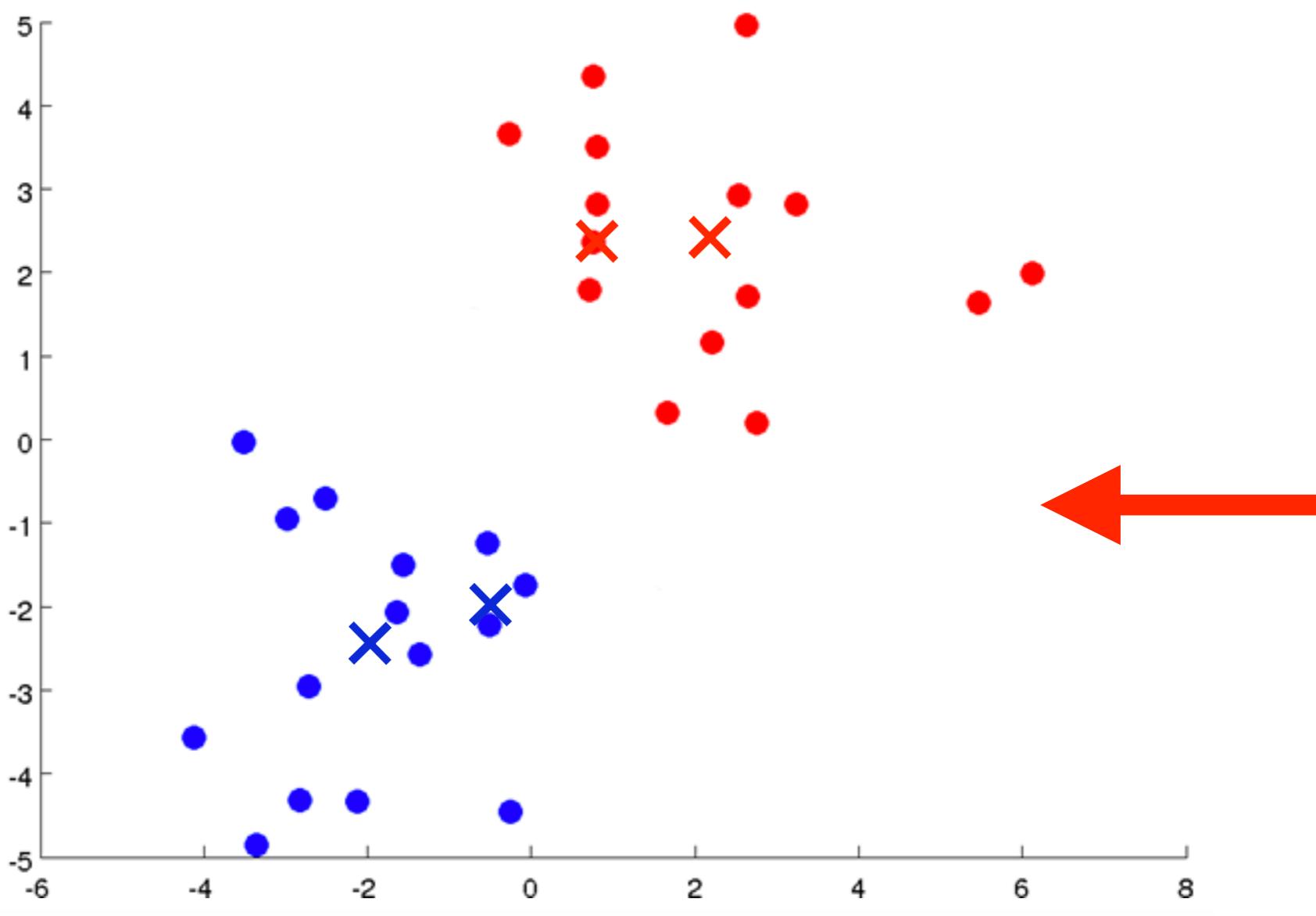
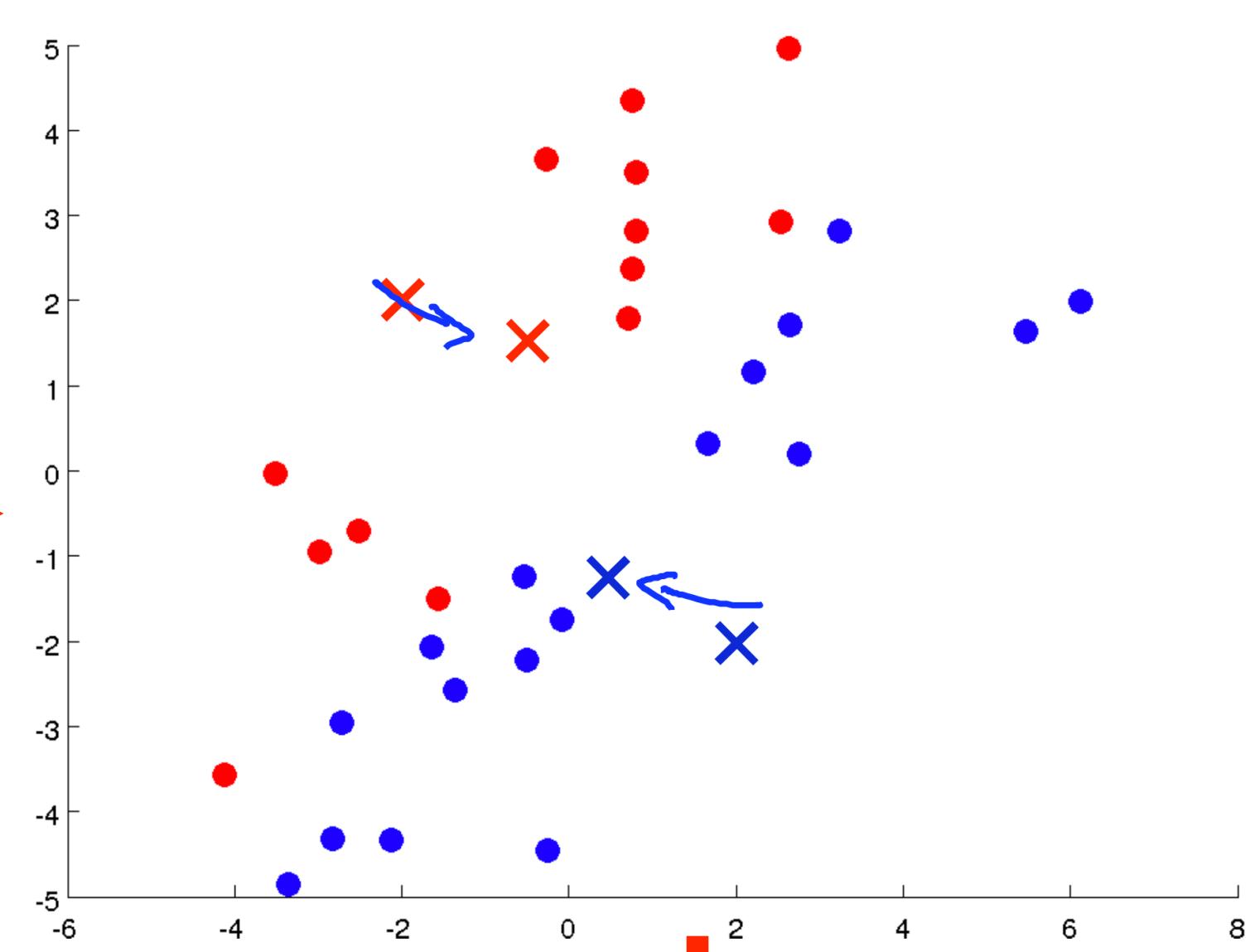
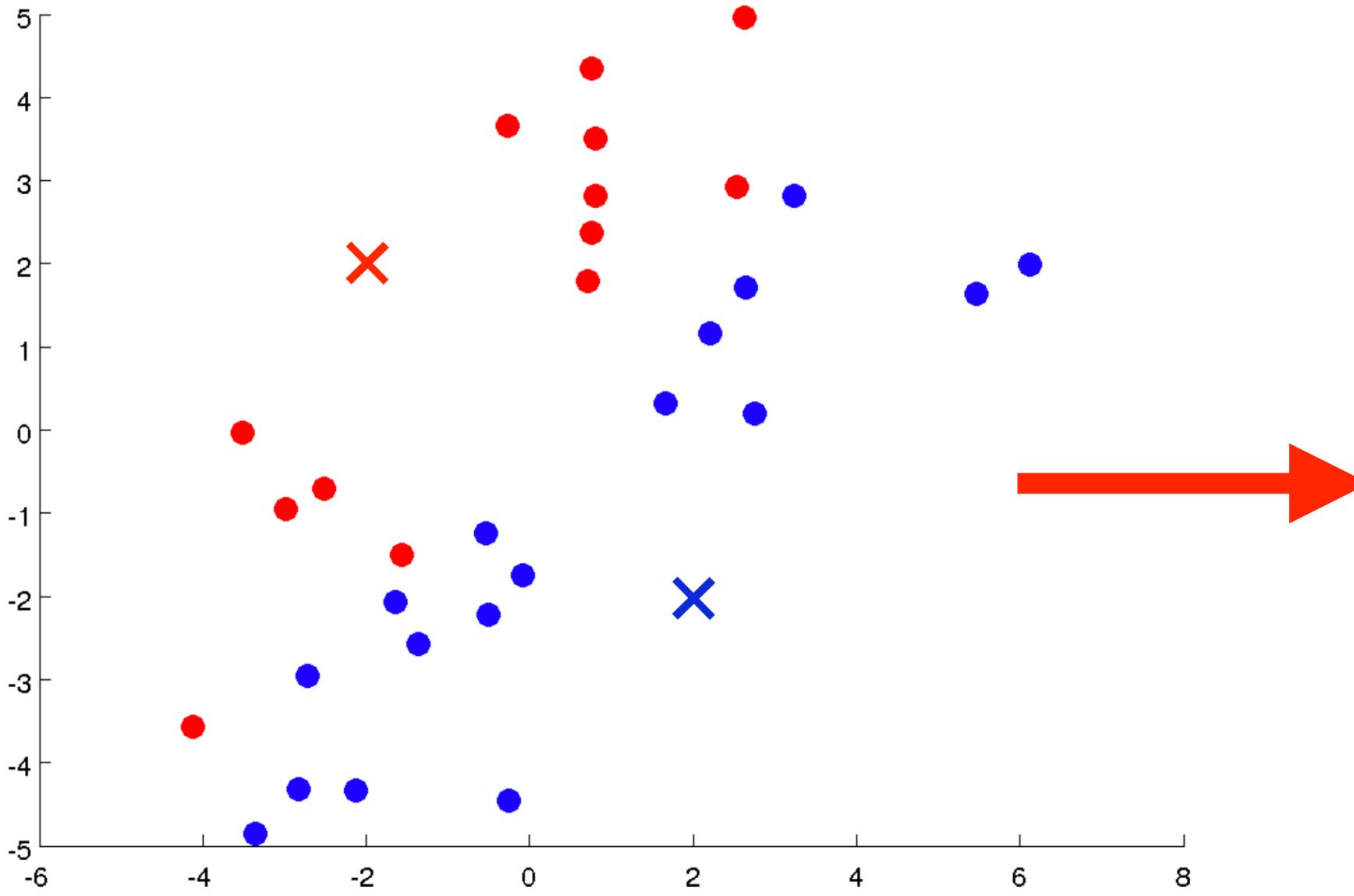
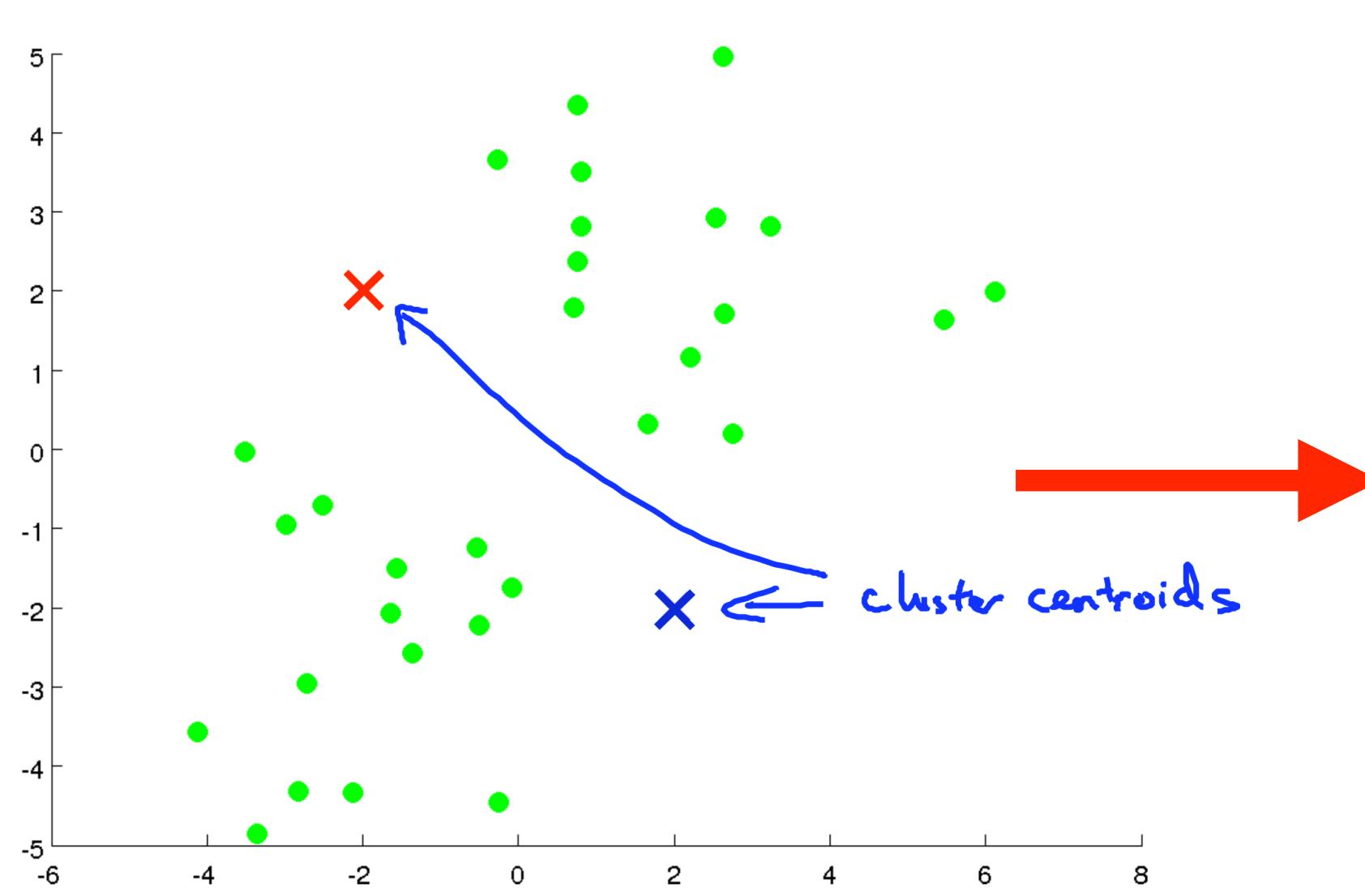
→ Organize computing clusters



→ Astronomical data analysis

K-Means Algorithm

Clustering



Clustering

K-means algorithm

$$\begin{array}{c} \mu_1 \\ \times \\ \mu_2 \\ \times \end{array}$$

Randomly initialize K cluster centroids $\underline{\mu}_1, \underline{\mu}_2, \dots, \underline{\mu}_K \in \mathbb{R}^n$

Repeat {

Cluster assignment step

for $i = 1$ to m
 $\underline{c}^{(i)} :=$ index (from 1 to K) of cluster centroid
 closest to $x^{(i)}$

$$\min_{\underline{c}^{(i)}} \|\underline{x}^{(i)} - \underline{\mu}_k\|^2$$

for $k = 1$ to K

$\rightarrow \underline{\mu}_k :=$ average (mean) of points assigned to cluster k

$$\underline{x}^{(1)}, \underline{x}^{(2)}, \underline{x}^{(3)}, \underline{x}^{(4)} \rightarrow \underline{c}^{(1)}=2, \underline{c}^{(2)}=2, \underline{c}^{(3)}=2, \underline{c}^{(4)}=2$$

$$\underline{\mu}_2 = \frac{1}{4} \left[\underline{x}^{(1)} + \underline{x}^{(2)} + \underline{x}^{(3)} + \underline{x}^{(4)} \right] \in \mathbb{R}^n$$

Clustering

Optimization objective

K-means optimization objective

- $c^{(i)}$ = index of cluster (1,2,...,K) to which example $x^{(i)}$ is currently assigned
- μ_k = cluster centroid \underline{k} ($\mu_k \in \mathbb{R}^n$) K $k \in \{1, 2, \dots, K\}$
- $\mu_{c^{(i)}}$ = cluster centroid of cluster to which example $x^{(i)}$ has been assigned $x^{(i)} \rightarrow \underline{5}$ $\underline{c^{(i)}} = \underline{5}$ $\underline{\mu_{c^{(i)}}} = \underline{\mu_5}$

Clustering

K-means algorithm

Randomly initialize K cluster centroids $\mu_1, \mu_2, \dots, \mu_K \in \mathbb{R}^n$

```
Repeat {  
    for  $i = 1$  to  $m$   
         $c^{(i)} :=$  index (from 1 to  $K$ ) of cluster centroid  
        closest to  $x^{(i)}$   
    for  $k = 1$  to  $K$   
         $\mu_k :=$  average (mean) of points assigned to cluster  $k$   
}
```

Random initialization

Clustering

Random initialization

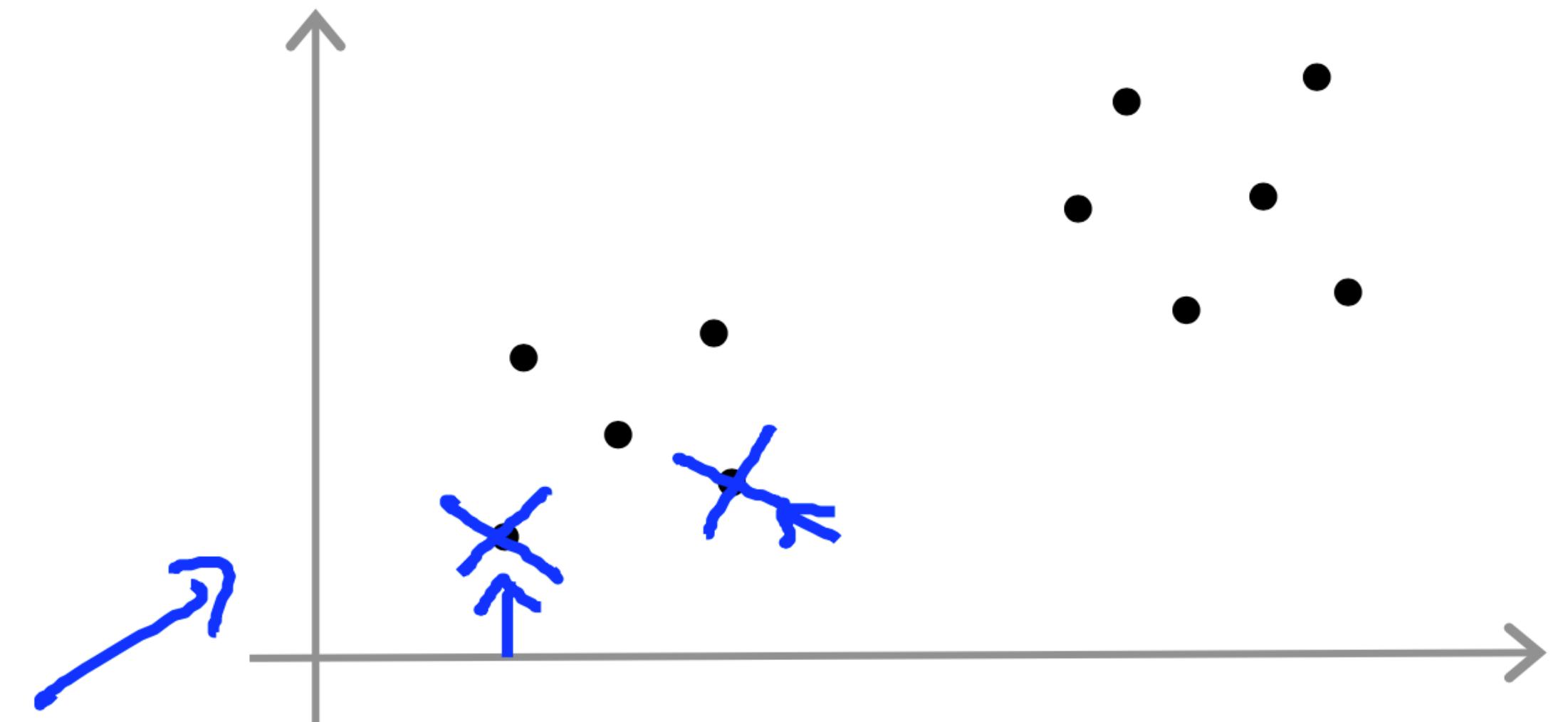
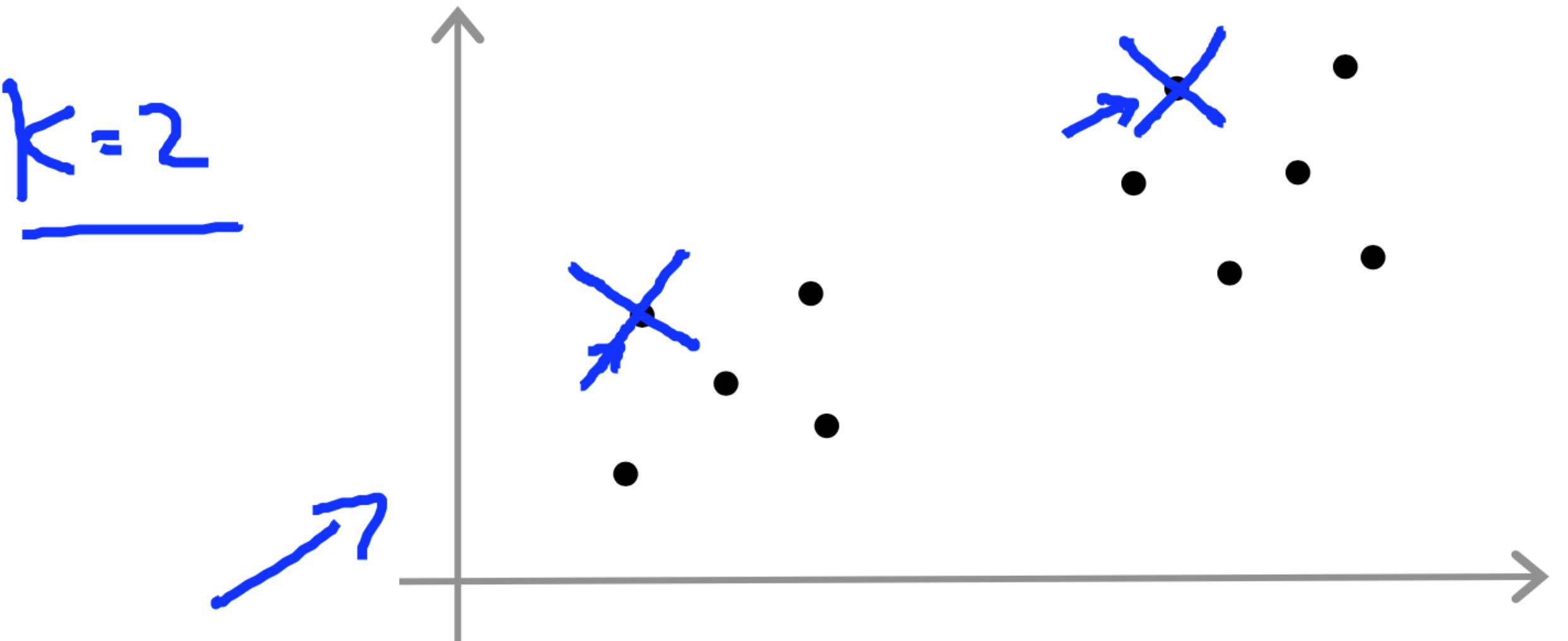
Should have $K < m$

Randomly pick K training examples.

Set μ_1, \dots, μ_K equal to these K examples.

$$\mu_1 = x^{(i)}$$

$$\mu_2 = x^{(j)}$$



Clustering

Random initialization

For i = 1 to 100 {

 Randomly initialize K-means.

 Run K-means. Get $c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K$.

 Compute cost function (distortion)

$J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$

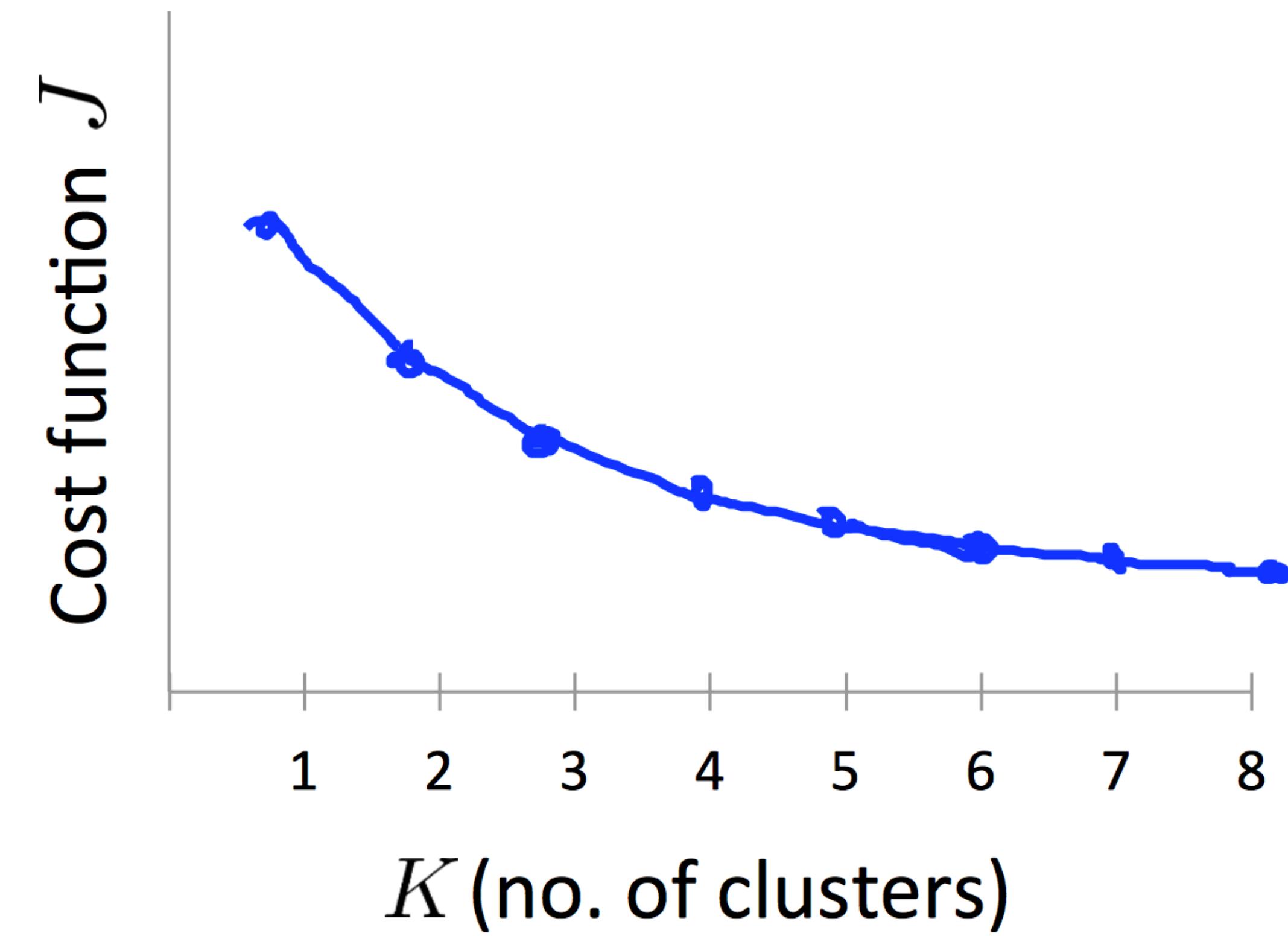
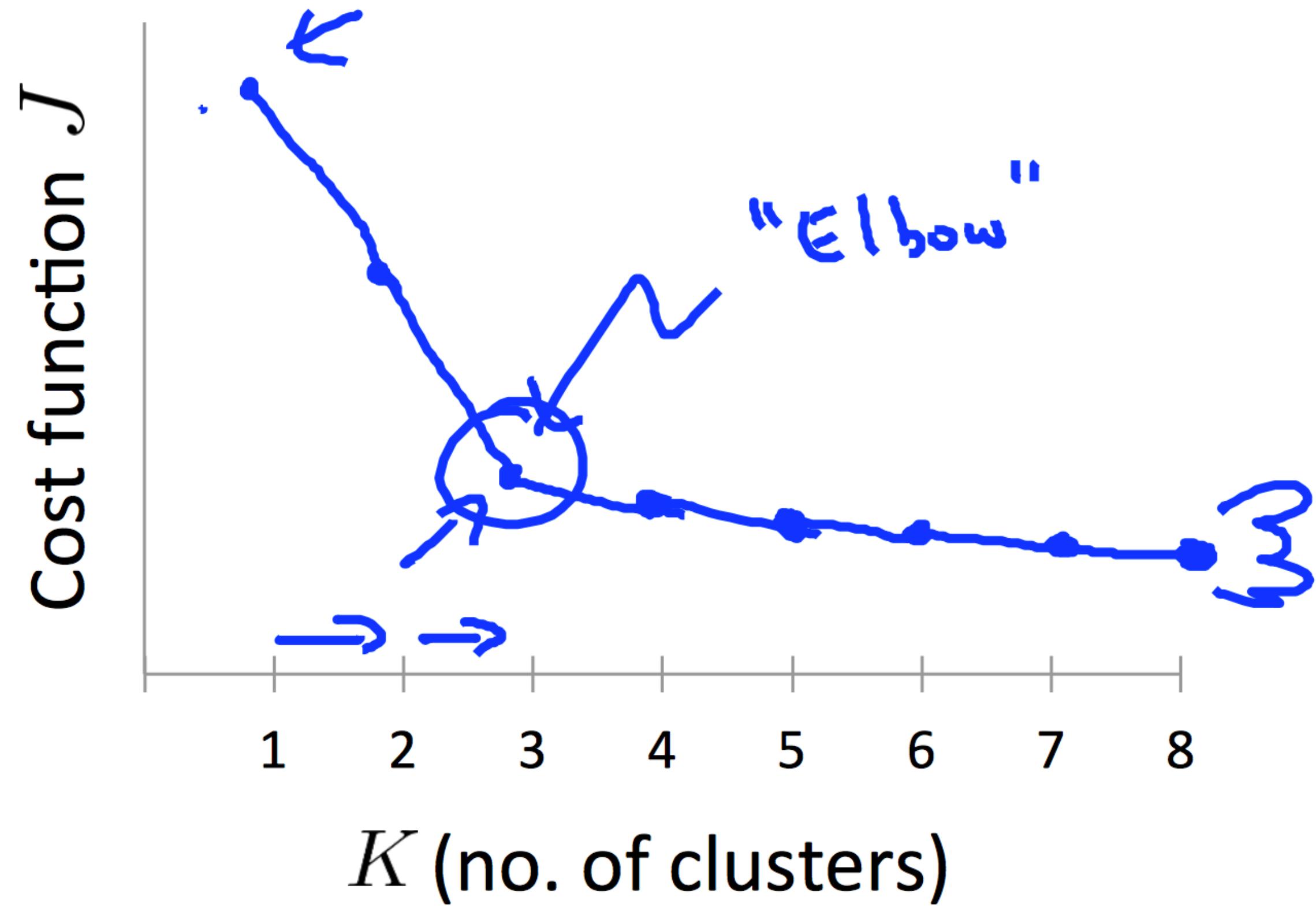
}

Pick clustering that gave lowest cost $J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$

Clustering

Choosing the value of K

Elbow method:

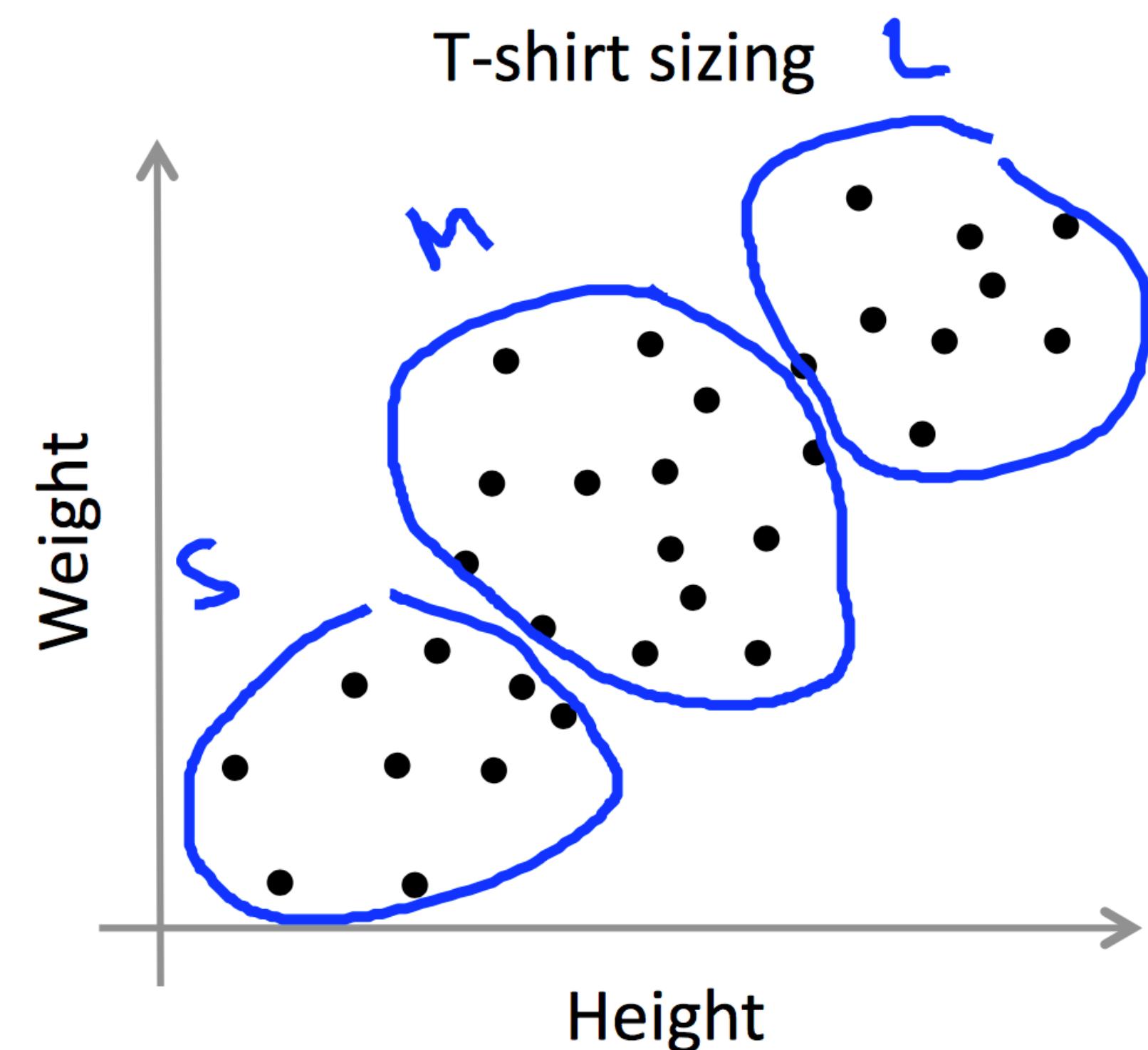


Choosing the number of clusters

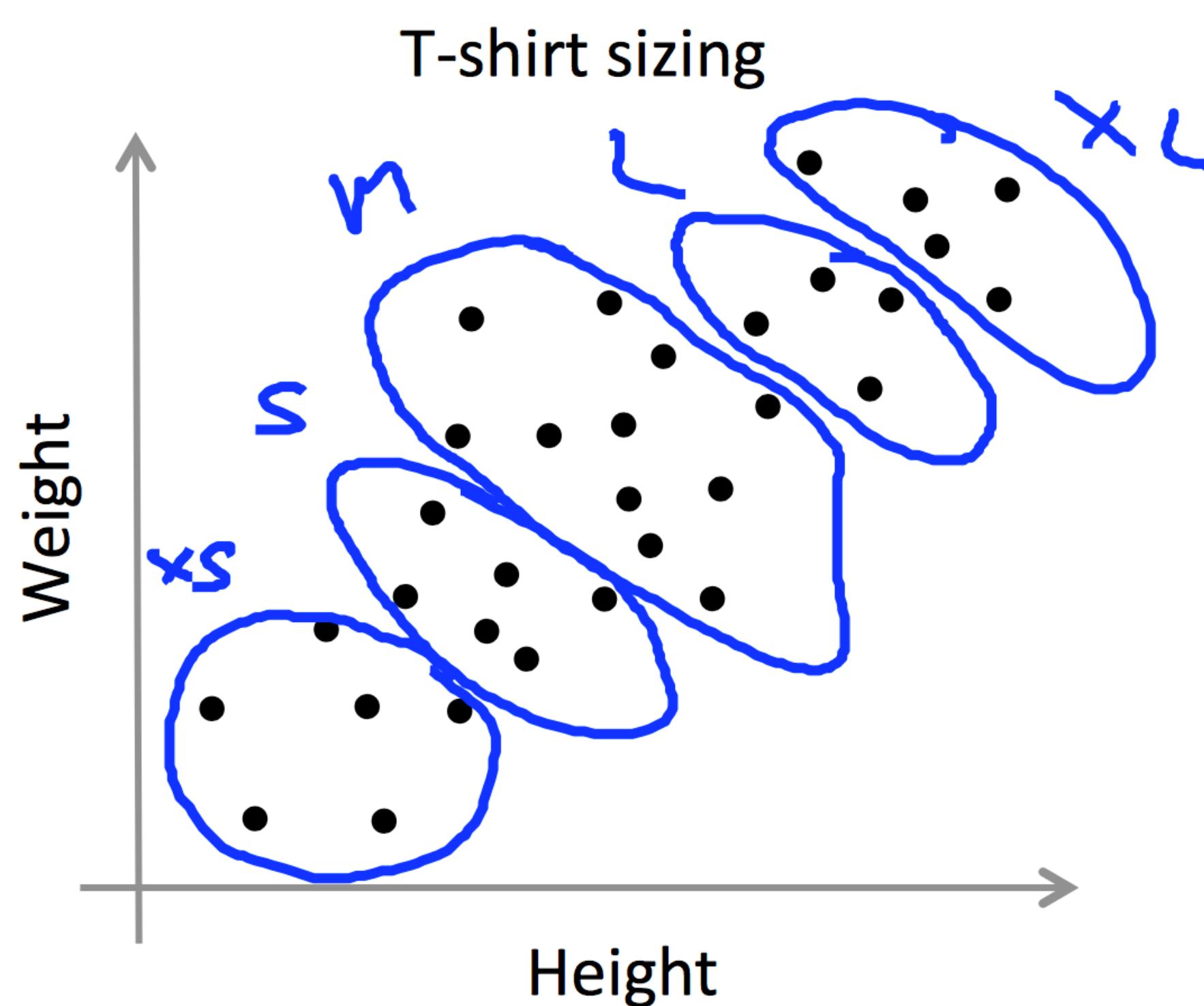
Clustering

$k=3$ S, M, L

E.g.



$k=5$ XS, S, M, L, XL



Dimensionality Reduction

In this module, we introduce Principal Components Analysis, and show how it can be used for data compression to speed up learning algorithms as well as for visualizations of complex datasets.

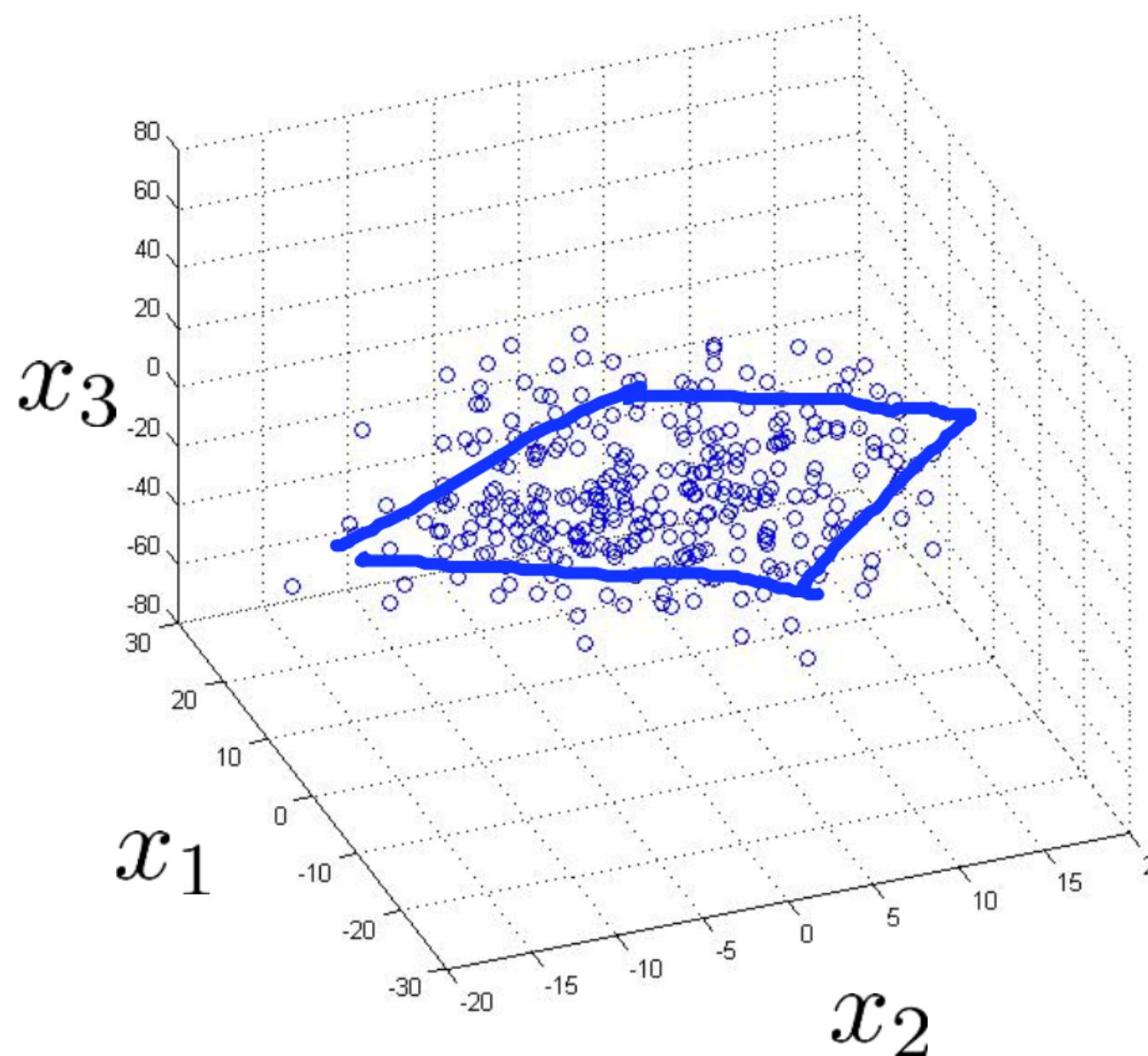
Dimensionality Reduction

Motivation I: Data Compression

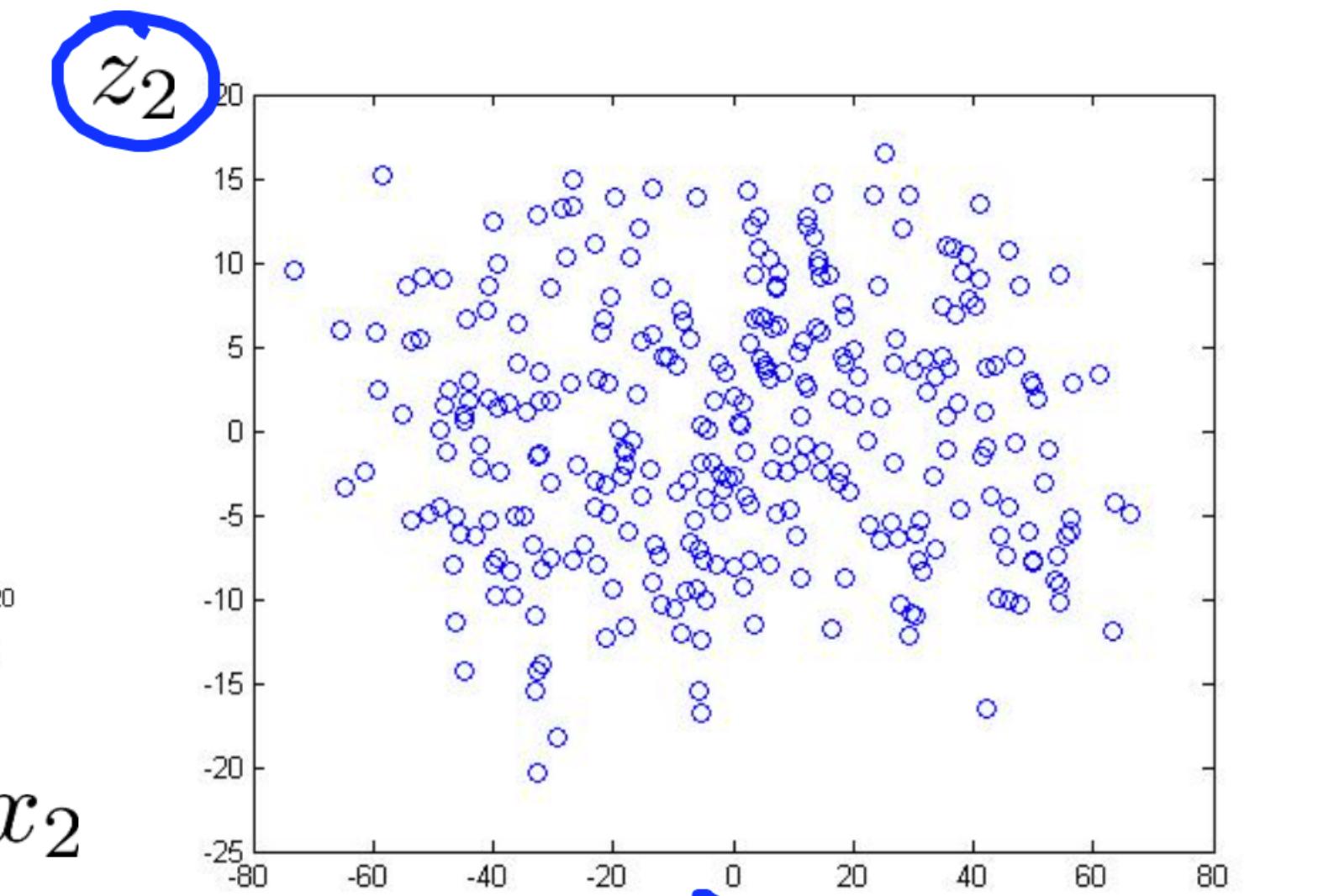
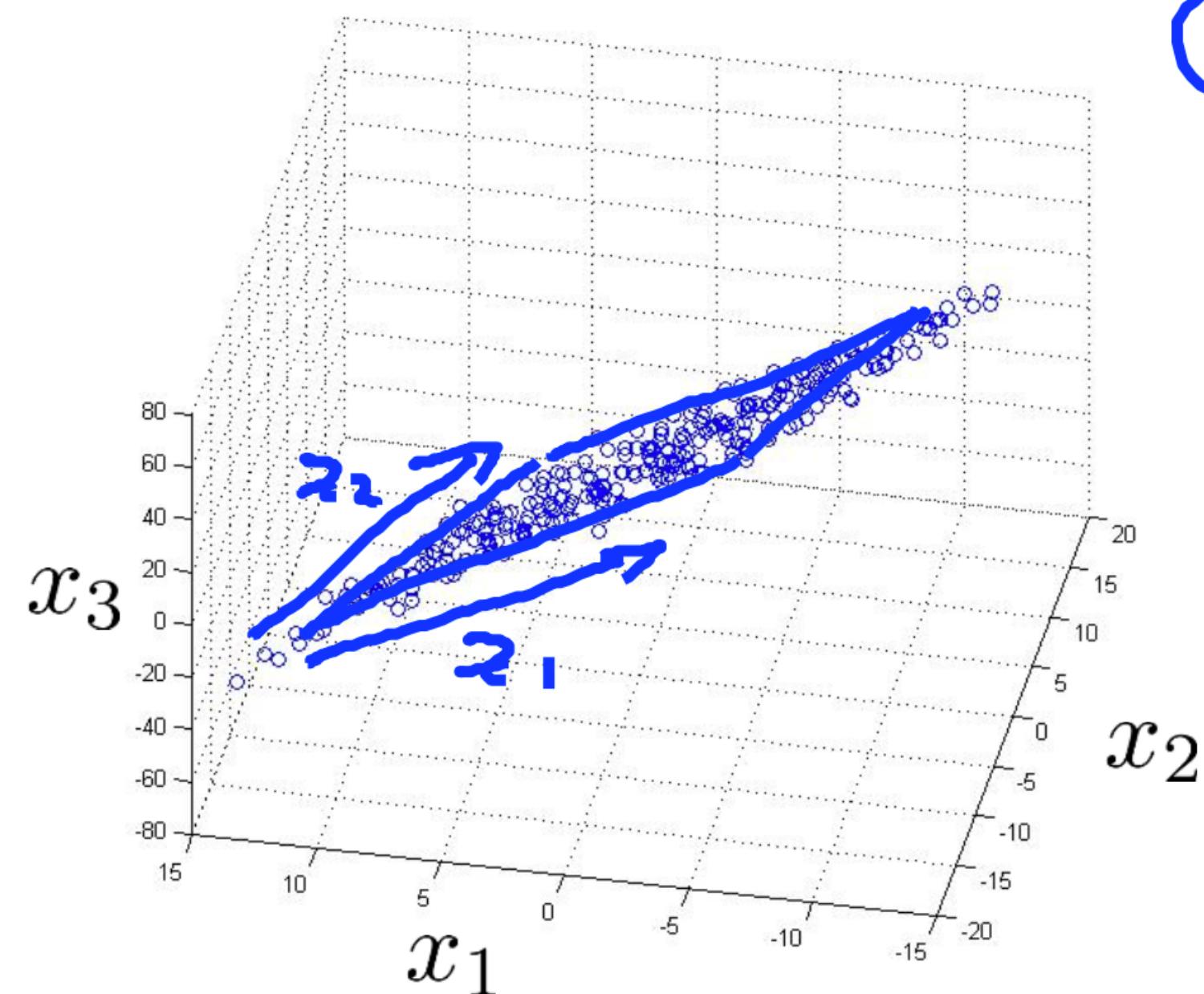
Data Compression

$$10000 \rightarrow 1000$$

Reduce data from 3D to 2D



$$x^{(i)} \in \mathbb{R}^3$$



$$z^{(i)} \in \mathbb{R}^2$$

$$\underline{z} = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$$

$$\underline{\underline{z}}^{(i)} = \begin{bmatrix} z_1^{(i)} \\ z_2^{(i)} \end{bmatrix}$$

Dimensionality Reduction

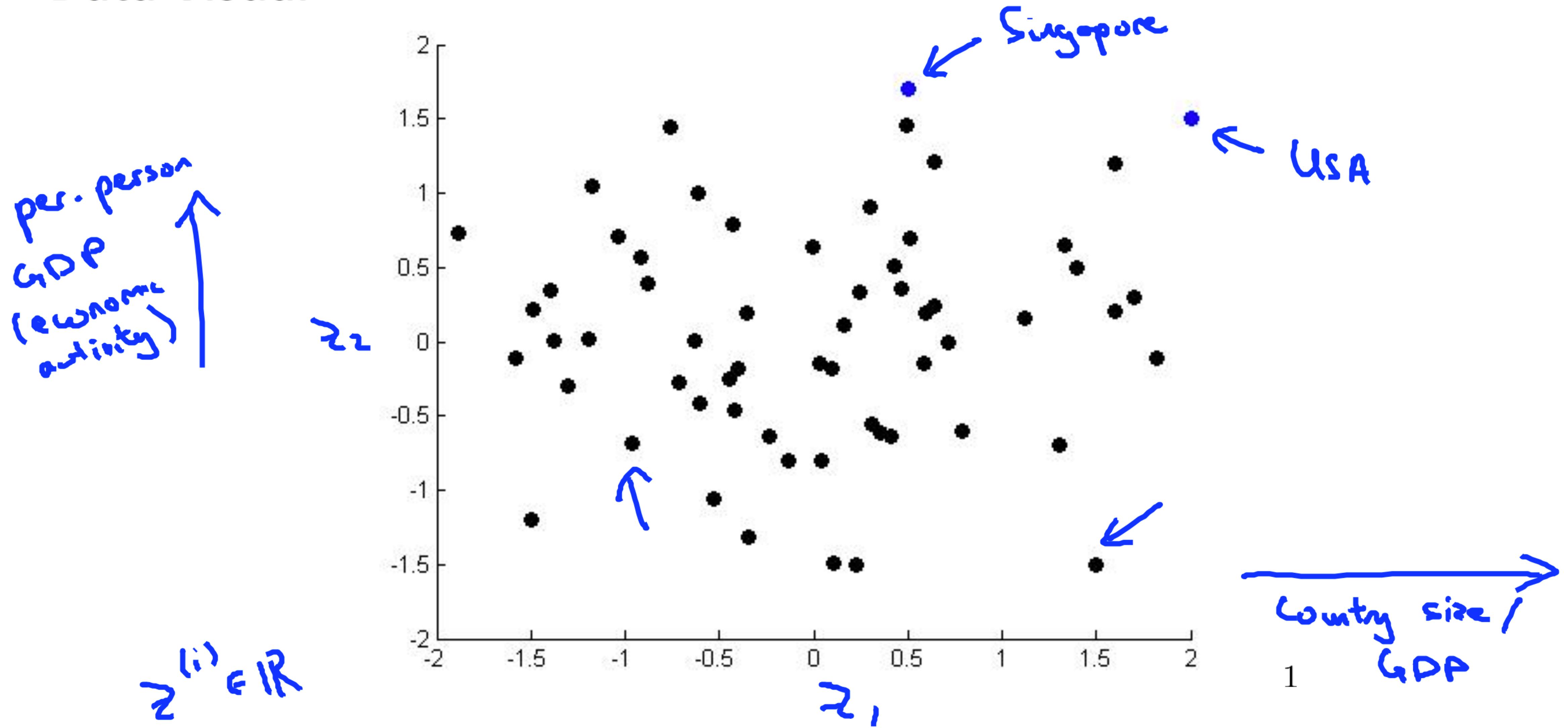
Motivation II: Visualization

Data Visualization

Dimensionality Reduction

Motivation II: Visualization

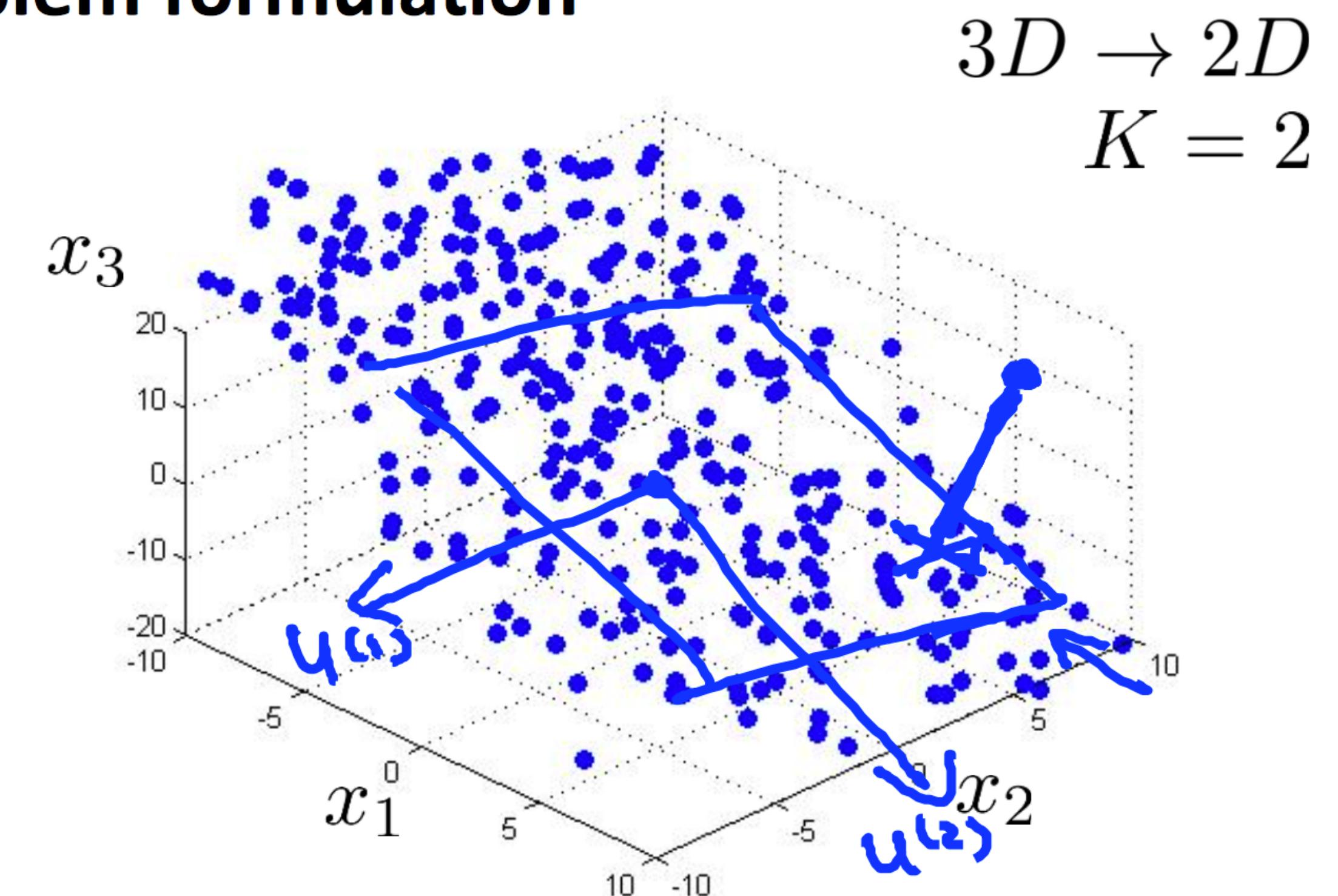
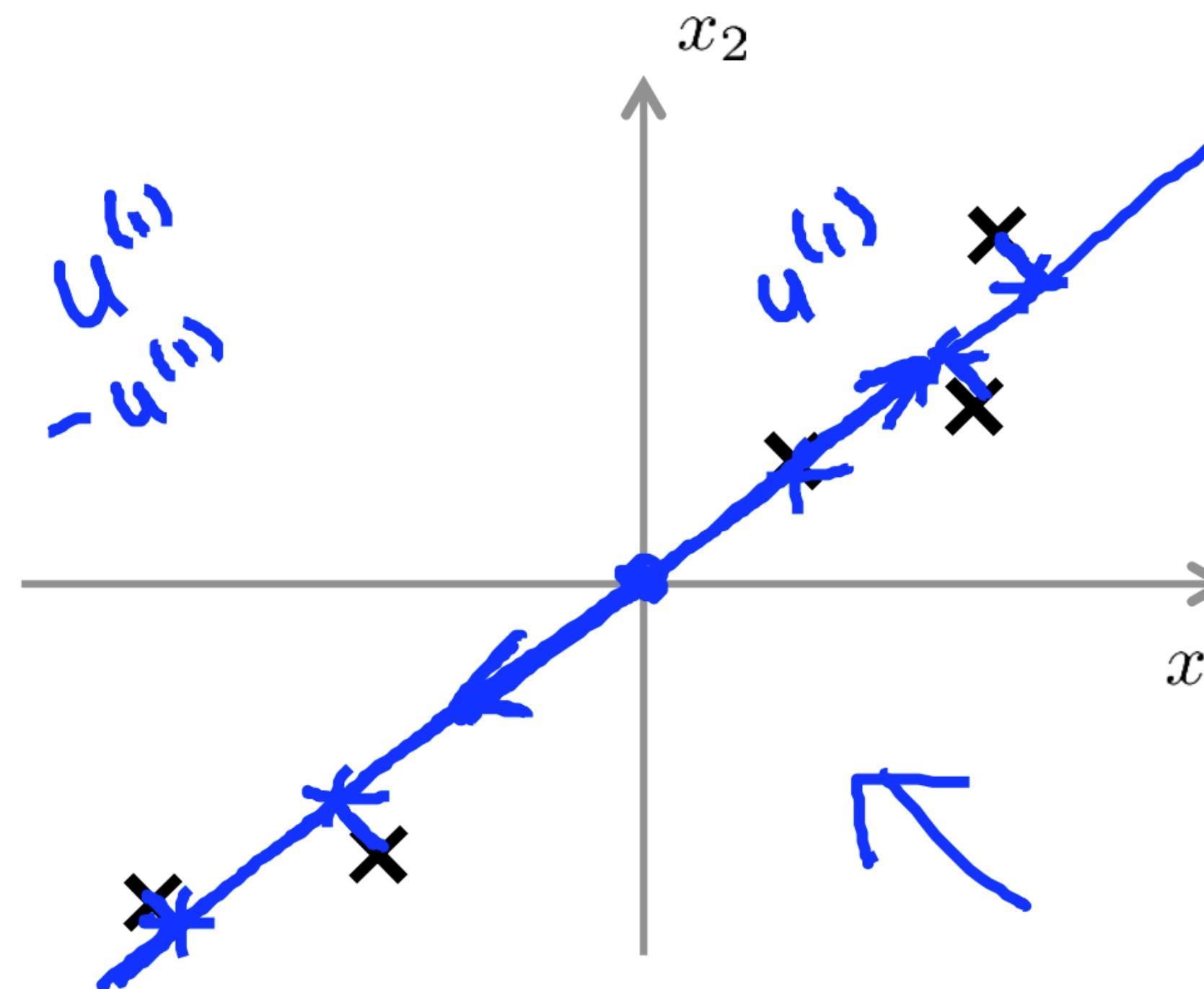
Data Visualization



Dimensionality Reduction

Principle component analysis

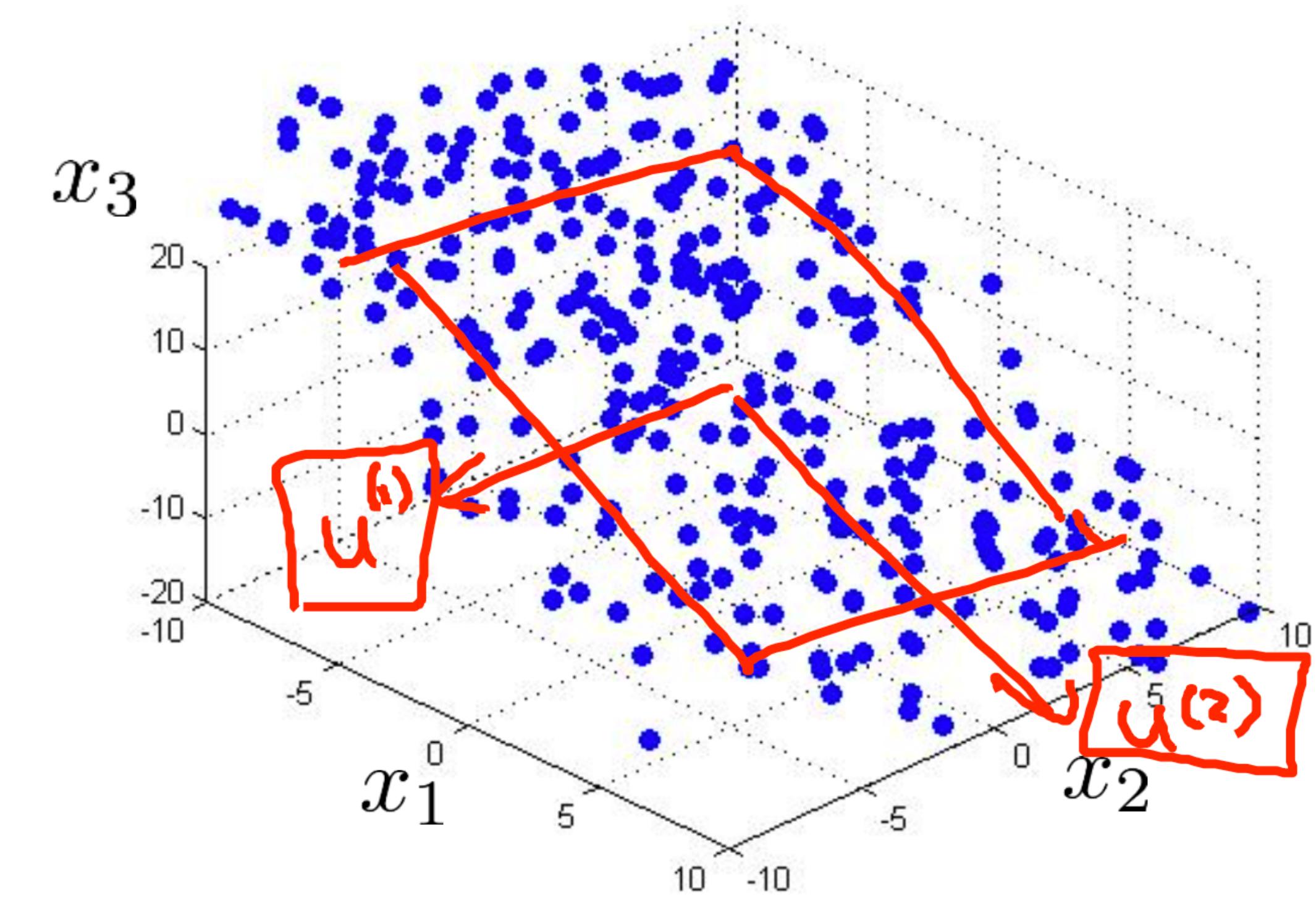
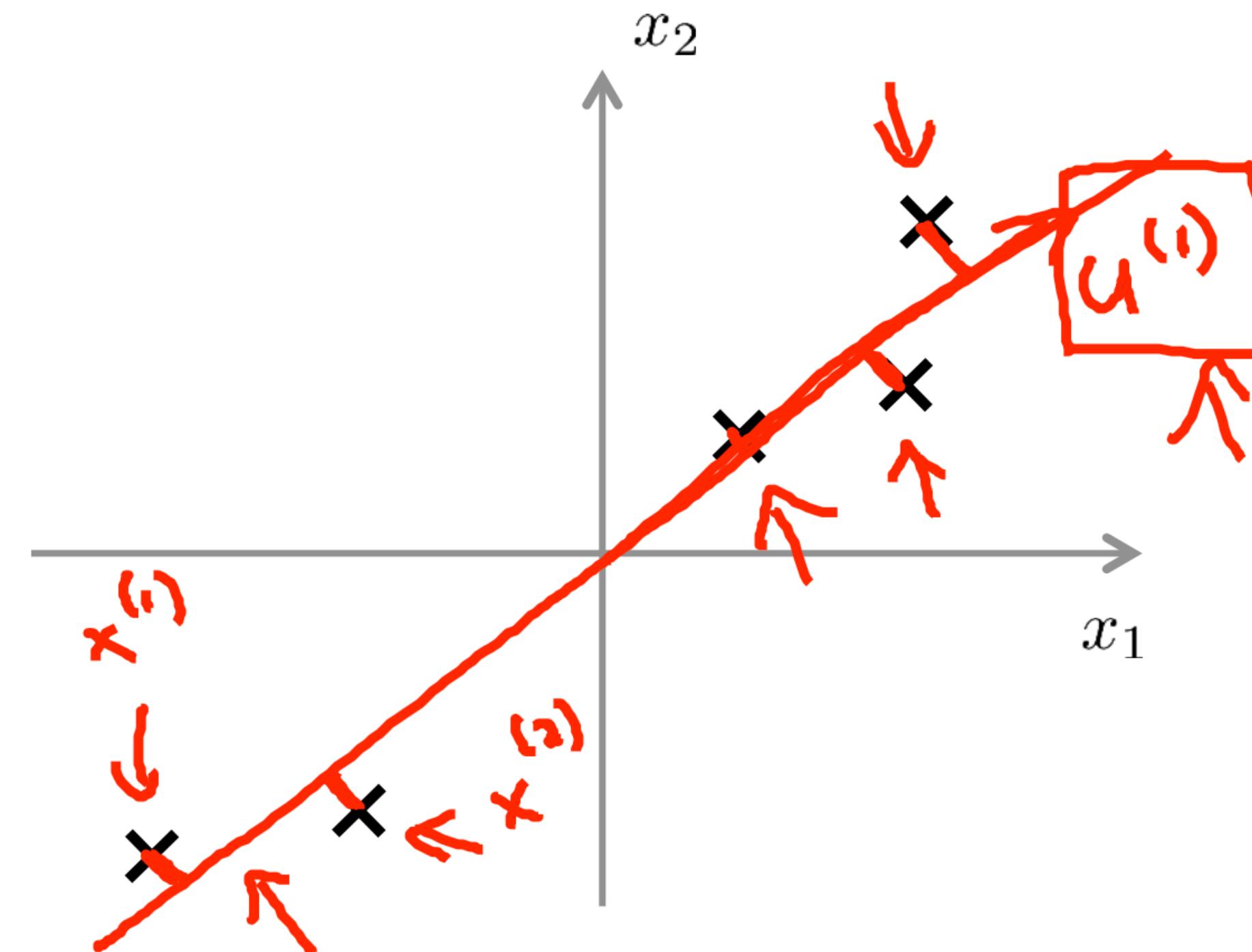
Principal Component Analysis (PCA) problem formulation



Dimensionality Reduction

Principle component analysis

Principal Component Analysis (PCA) algorithm



Reduce data from 2D to 1D

$$x^{(i)} \in \mathbb{R}^2 \rightarrow z^{(i)} \in \mathbb{R}$$

Diagram: A horizontal line with data points $x^{(i)}$ and a red arrow pointing to a box labeled E_1 .

Reduce data from 3D to 2D

$$x^{(i)} \in \mathbb{R}^3 \rightarrow z^{(i)} \in \mathbb{R}^2$$
$$z = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$$

Dimensionality Reduction

Principle component analysis: Algorithm

Principal Component Analysis (PCA) algorithm

Reduce data from n -dimensions to k -dimensions

Dimensionality Reduction

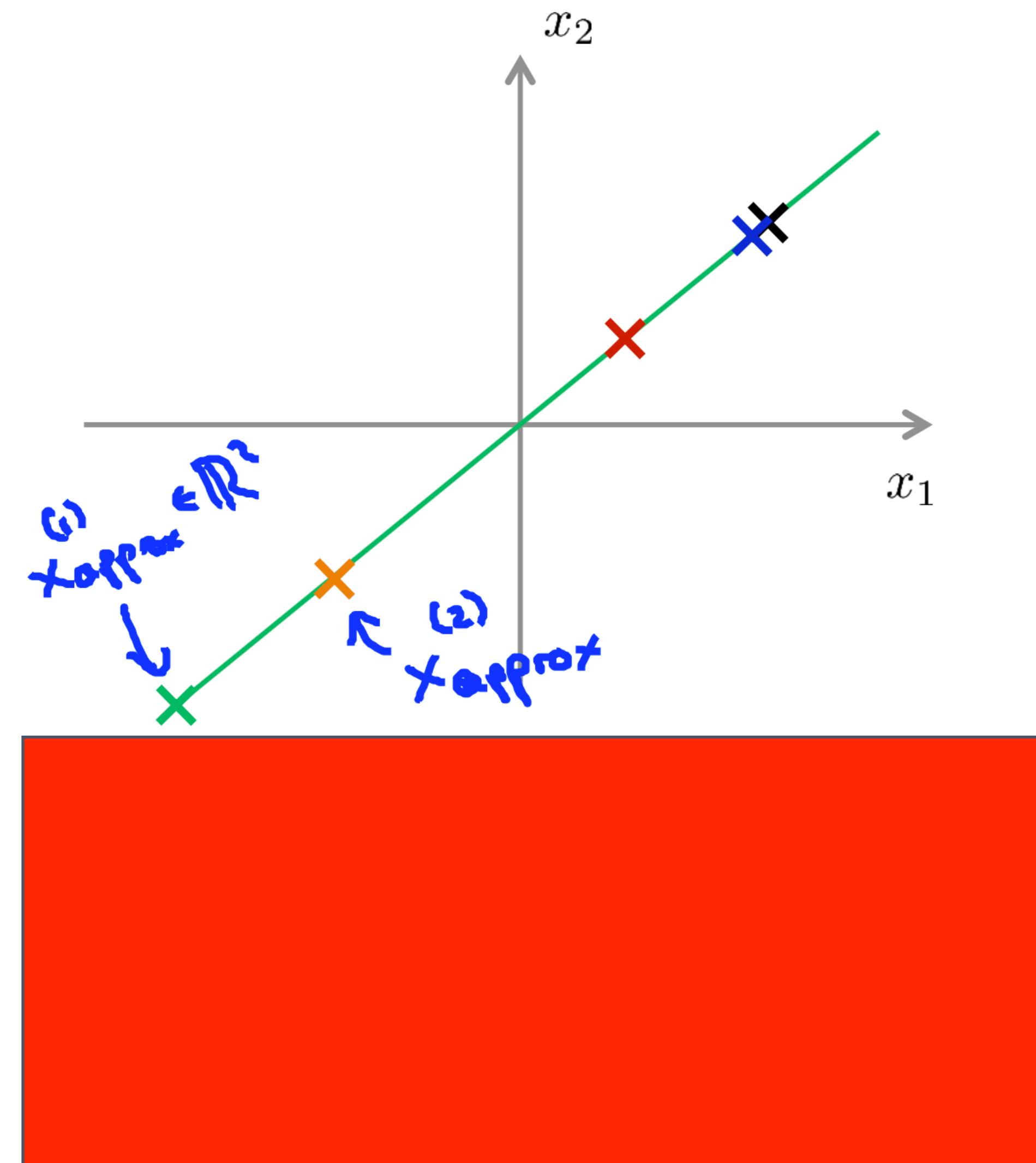
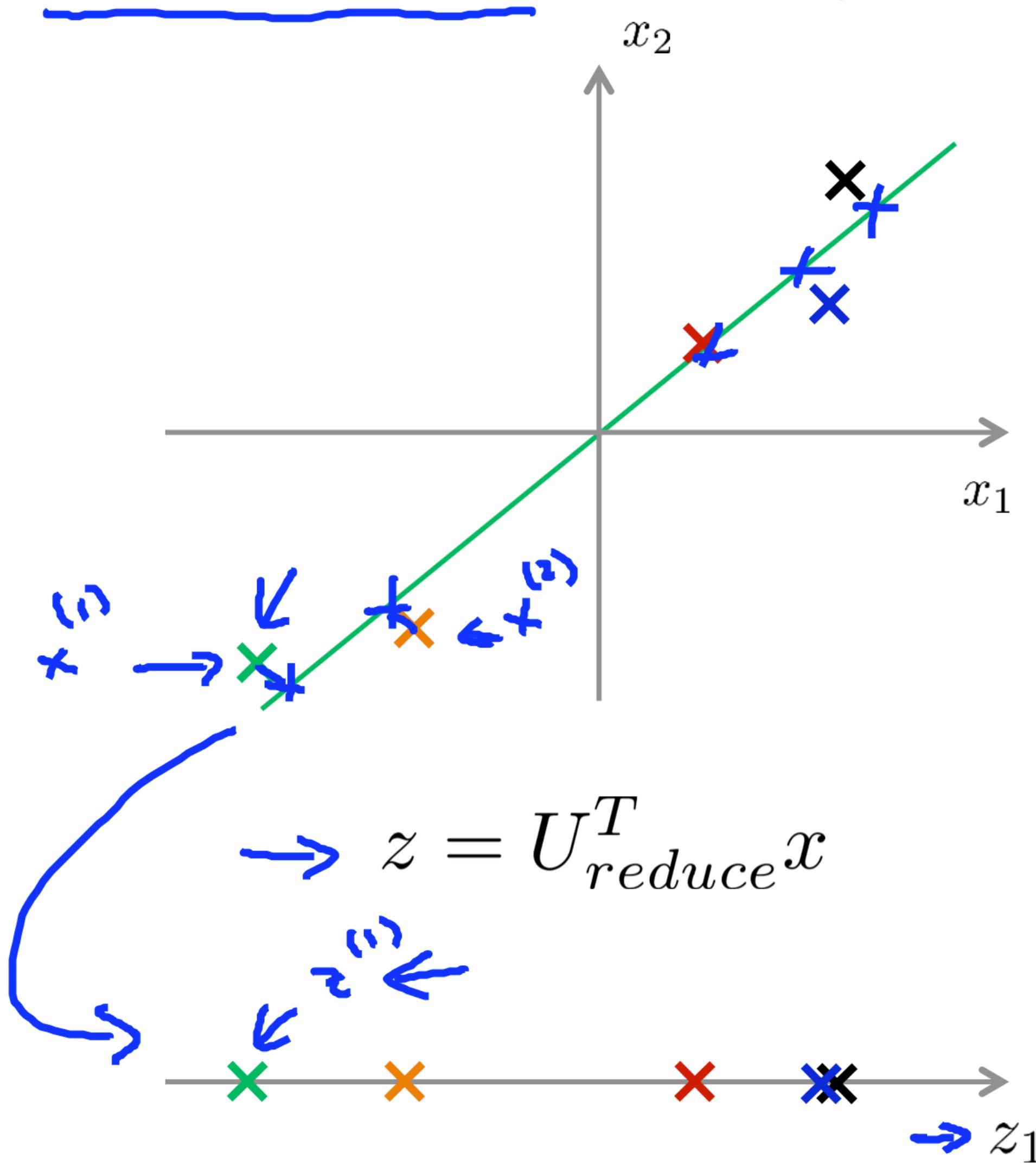
Principle component analysis: Algorithm

Principal Component Analysis (PCA) algorithm summary

- After mean normalization (ensure every feature has zero mean) and optionally feature scaling:

Dimensionality Reduction

Reconstruction from compressed representation



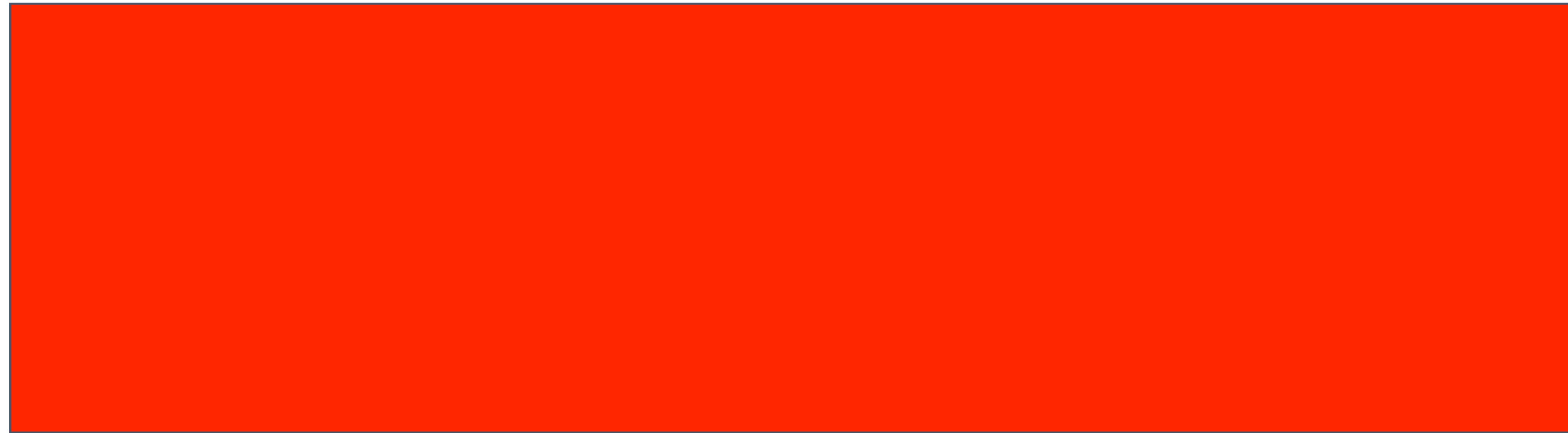
Dimensionality Reduction

Choosing the number of principal components

Choosing k (number of principal components)

Average squared projection error: $\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \hat{x}^{(i)}_{\text{approx}}\|^2$

Total variation in the data: $\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2$



• “99% of variance is retained”
~~95 to 96%~~

Dimensionality Reduction

Choosing the number of principal components

Choosing k (number of principal components)

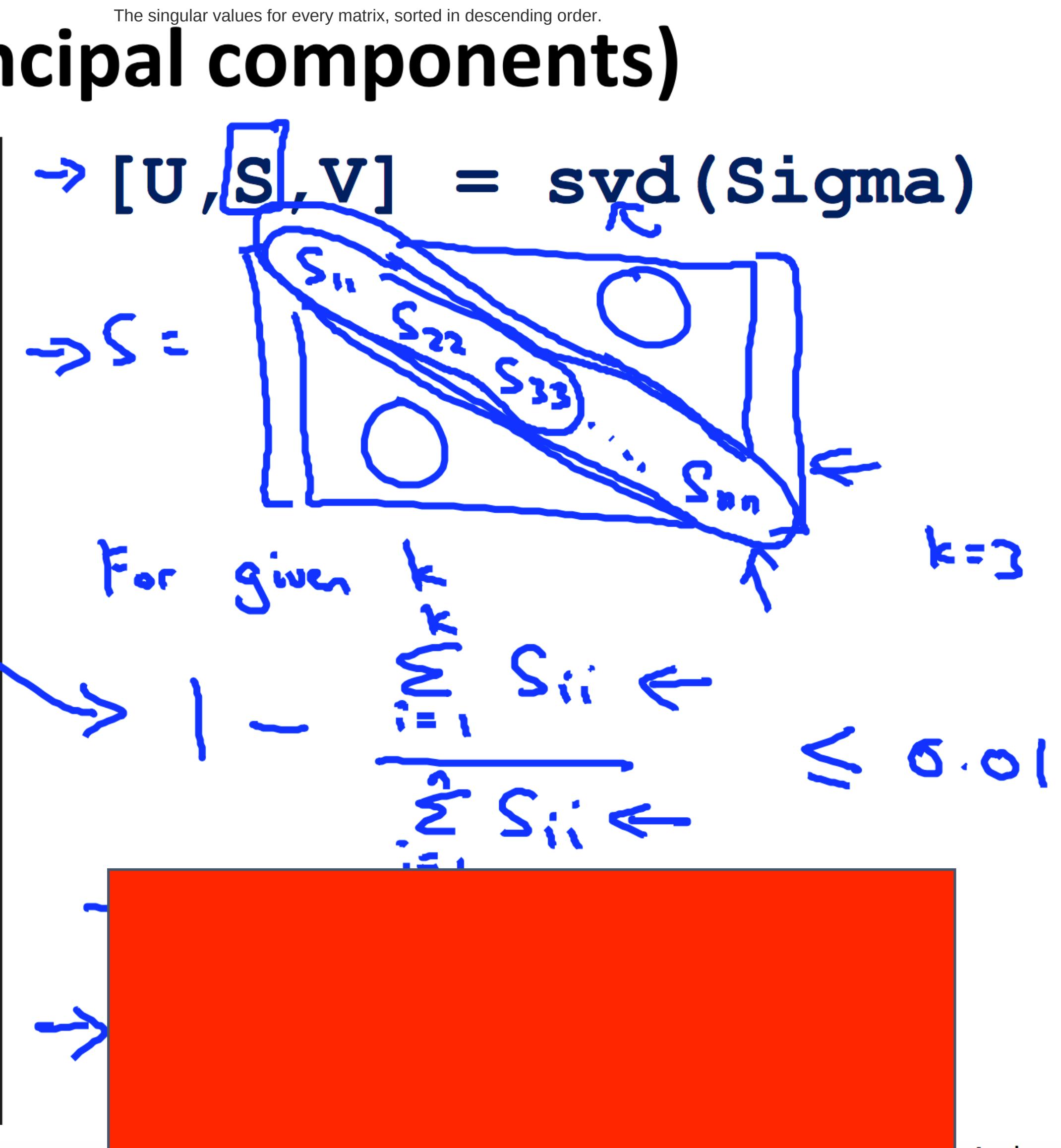
Algorithm:

Try PCA with $k = 1$ ~~$k=2$~~ ~~$k=3$~~ ~~$k=4$~~

Compute $U_{reduce}, z^{(1)}, z^{(2)}, \dots, z^{(m)}, x_{approx}^{(1)}, \dots, x_{approx}^{(m)}$



$k = 17$



Dimensionality Reduction

Supervised learning speedup

- $(\underline{x}^{(1)}, y^{(1)}), (\underline{x}^{(2)}, y^{(2)}), \dots, (\underline{x}^{(m)}, y^{(m)})$

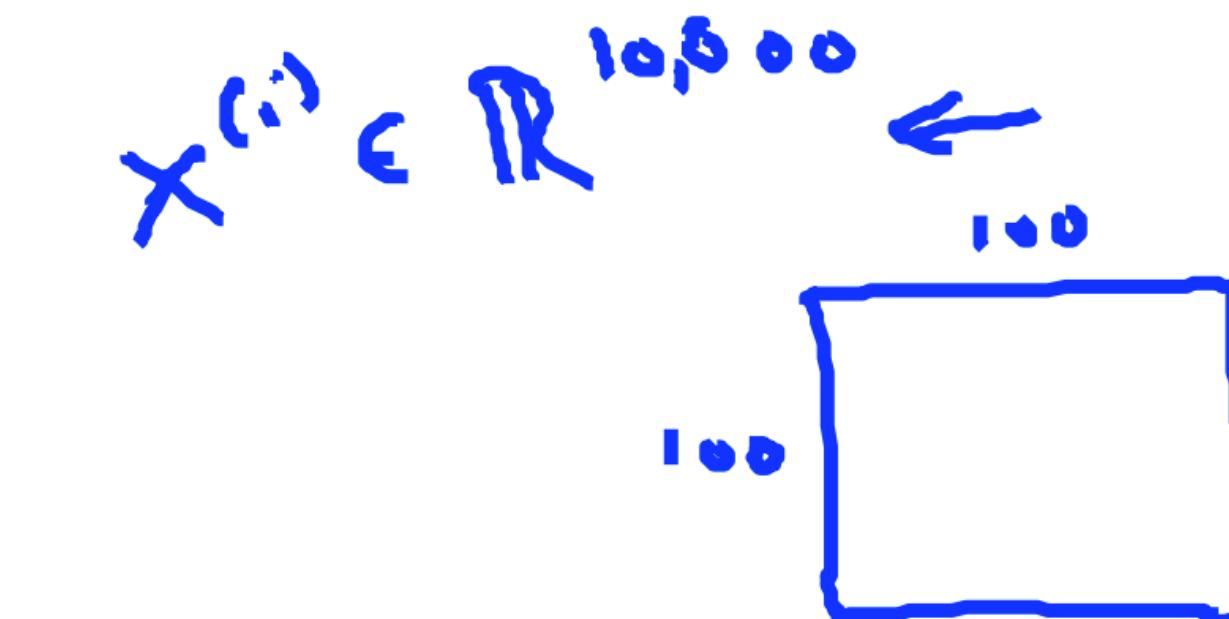
Extract inputs:

Unlabeled dataset:

U reduce

New training set:

- $(z^{(1)}, y^{(1)}), (z^{(2)}, y^{(2)}), \dots, (z^{(m)}, y^{(m)})$



$$h_{\theta}(z) = \frac{1}{1 + e^{-\theta^T z}}$$

Note: Mapping $x^{(i)} \rightarrow z^{(i)}$ should be defined by running PCA
only on the training set. This mapping can be applied as well to
the examples $x_{cv}^{(i)}$ and $x_{test}^{(i)}$ in the cross validation and test

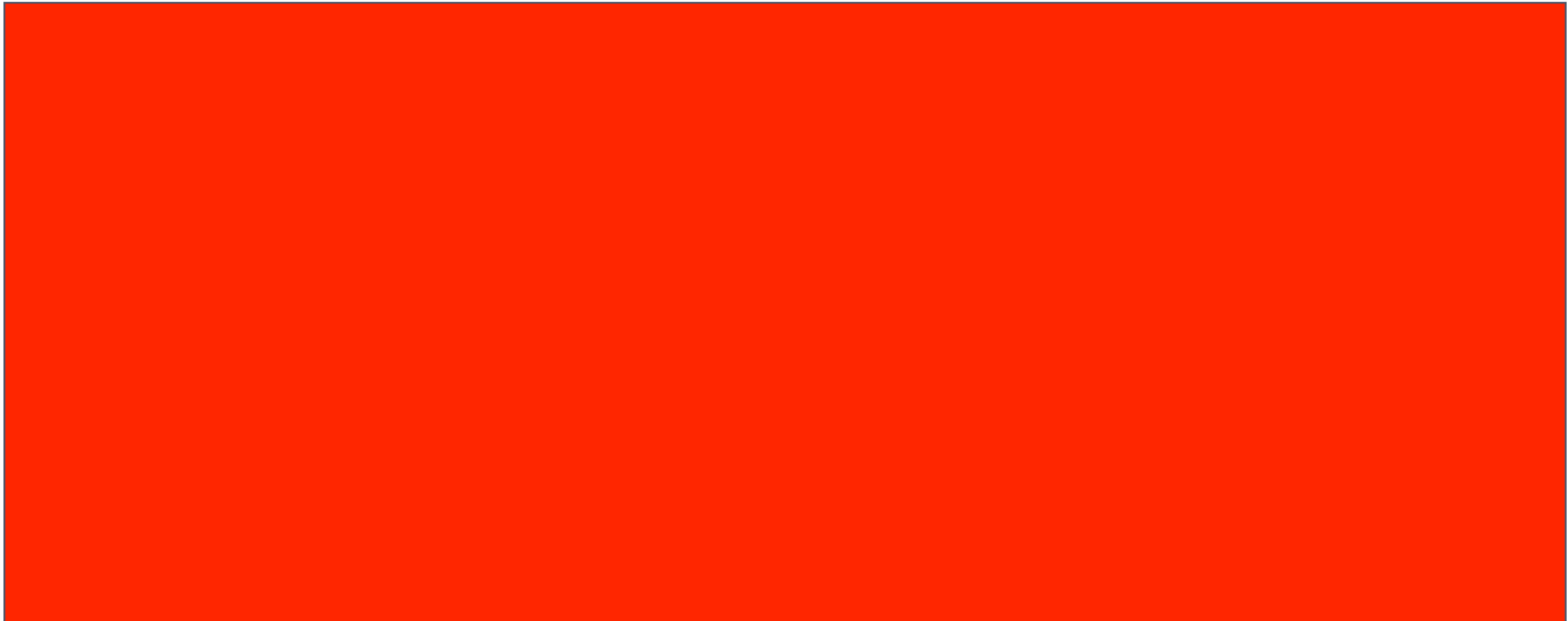
x → z

Applying PCA

Dimensionality Reduction

Advice for applying PCA

Application of PCA



Dimensionality Reduction

Use $\underline{z^{(i)}}$ instead of $\underline{x^{(i)}}$ to reduce the number of features to $\underline{k} < \underline{n}$. — 10000

Thus, fewer features, less likely to overfit.

Bad!

Dimensionality Reduction

PCA is sometimes used where it shouldn't be

Design of ML system:

- - Get training set $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$
- - ~~Run PCA to reduce $x^{(i)}$ in dimension to get $z^{(i)}$~~
- - Train logistic regression on $\{(z^{(1)}, y^{(1)}), \dots, (z^{(n)}, y^{(m)})\}$
- - Test on test set: Map $x_{test}^{(i)}$ to $z_{test}^{(i)}$. Run $h_\theta(z)$ on $\{(z_{test}^{(1)}, y_{test}^{(1)}), \dots, (z_{test}^{(m)}, y_{test}^{(m)})\}$

Recommender Systems

Recommender systems, which are used by companies like Amazon, Netflix and Apple to recommend products to their users. Recommender systems look at patterns of activities between different users and different products to produce these recommendations. In these lessons, we introduce recommender algorithms such as the collaborative filtering algorithm and low-rank matrix factorization.



Introduction to Python Recommendation System

1. Popularity-Based Recommender Systems
2. Correlation-Based Recommender Systems
3. Classification-Based Collaborative Filtering
4. Model-Based Collaborative Filtering
5. Content-Based Recommender Systems
6. Evaluating of Recommender Systems

Large Scale Machine Learning

Machine learning works best when there is an abundance of data to leverage for training. With the amount data that many websites/companies are gathering today, knowing how to handle ‘big data’ is one of the most sought after skills in Silicon Valley.

Gradient Descent with Large Datasets

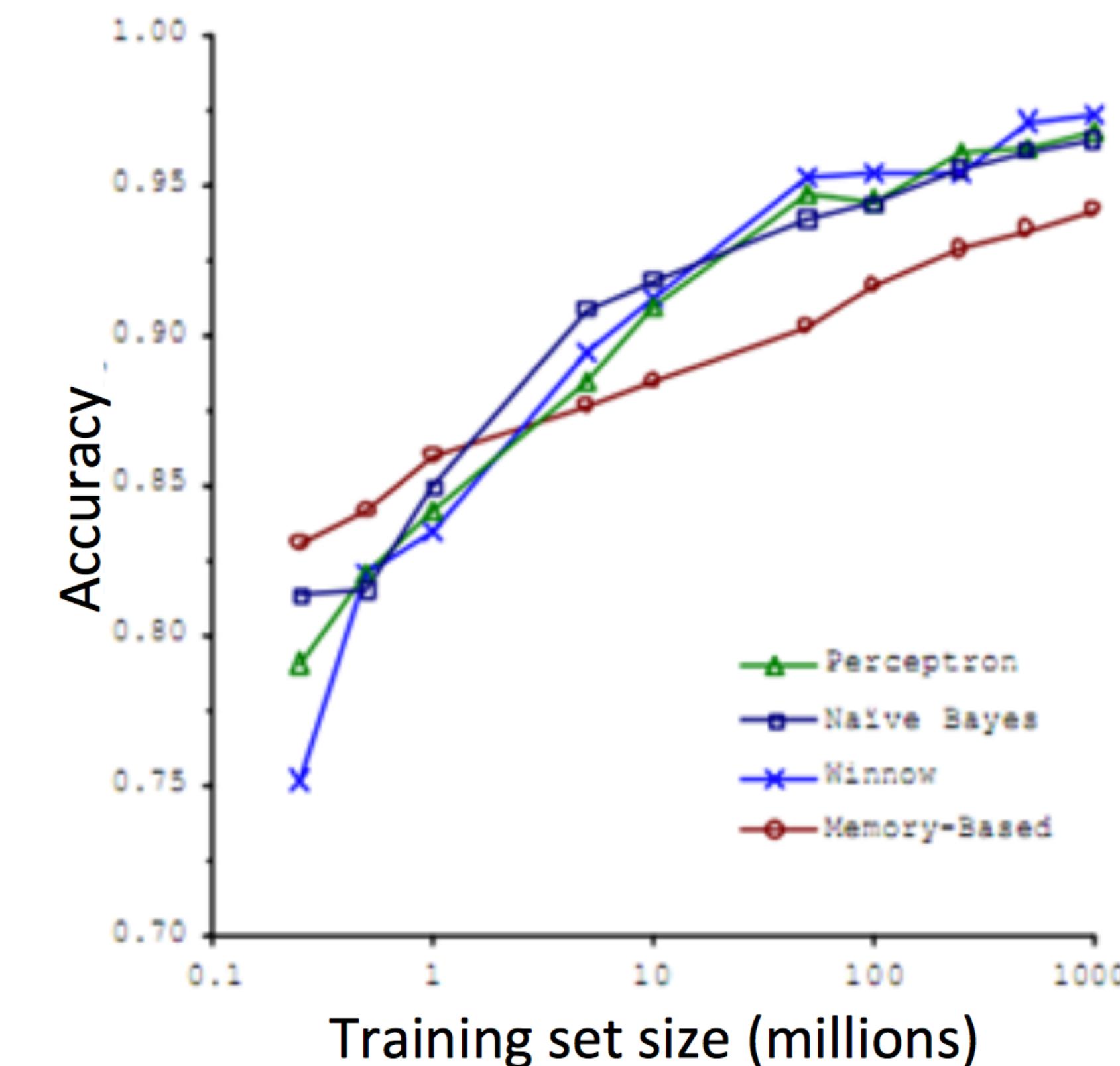
Learning with Large Datasets

Machine learning and data

Classify between confusable words.

E.g., {to, two, too}, {then, than}.

For breakfast I ate two eggs.



“It’s not who has the best algorithm that wins.

It’s who has the most data.”

Gradient Descent with Large Datasets

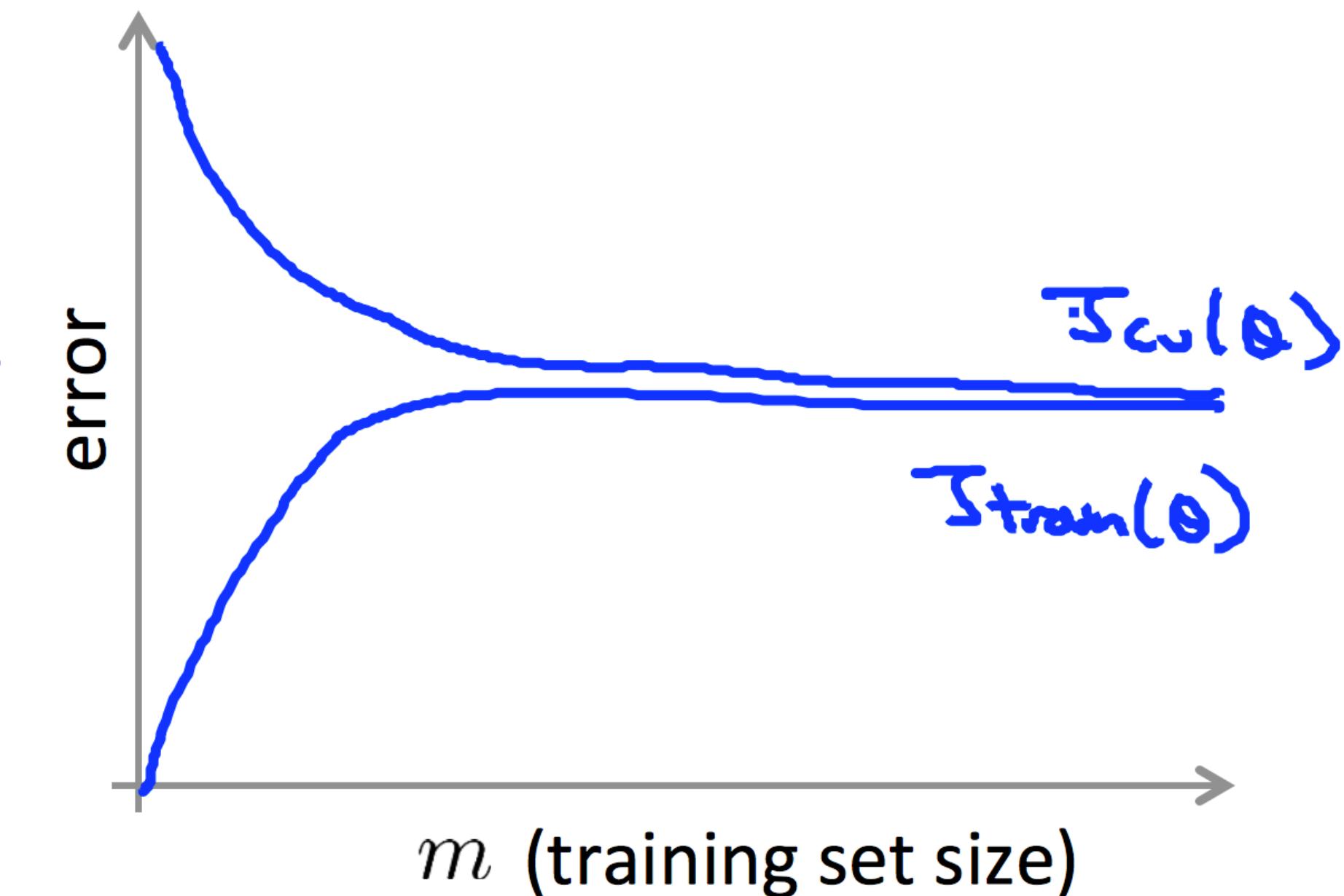
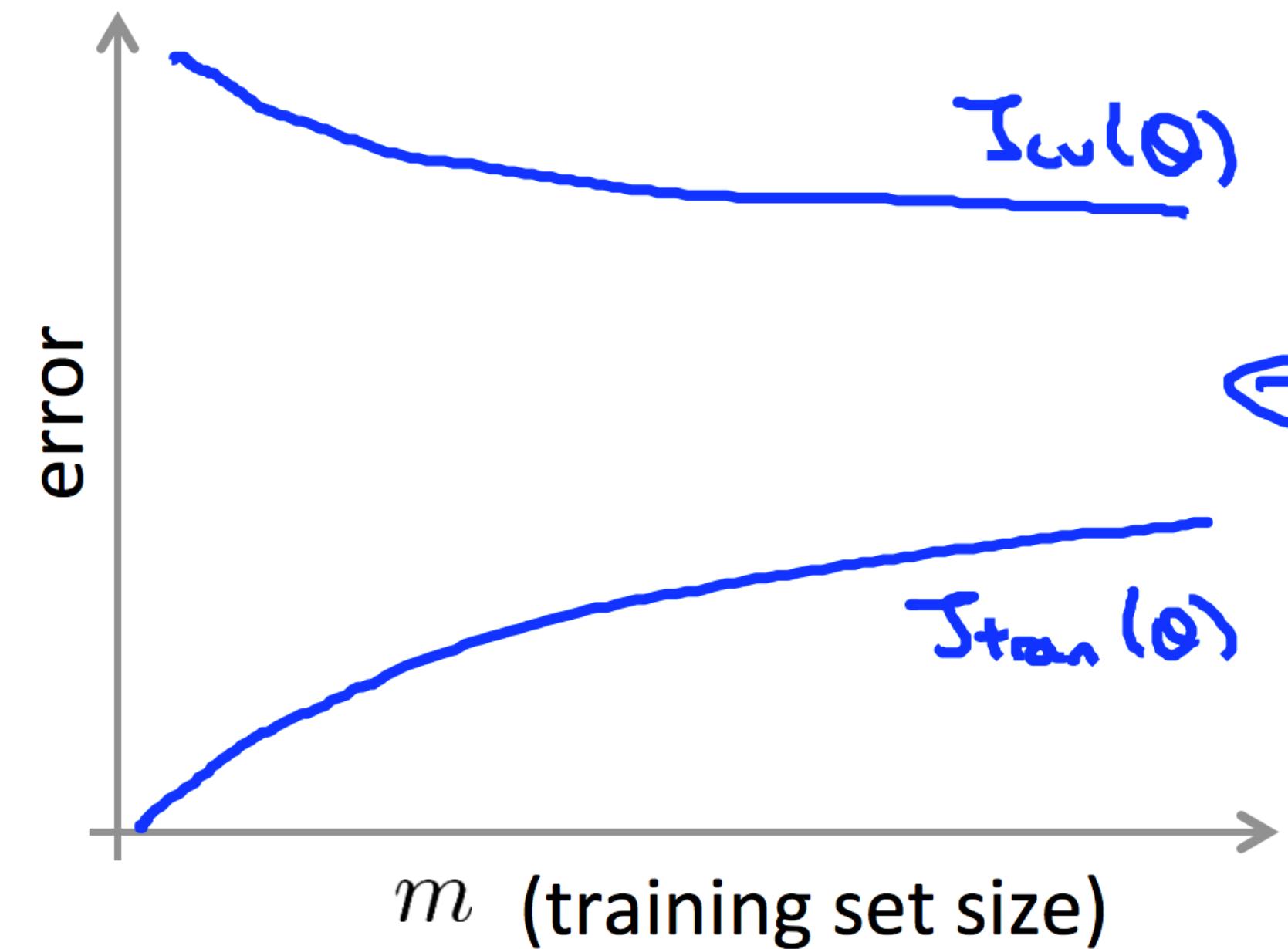
Learning with Large Datasets

Learning with large datasets

$$m = 100,000,000$$

$$m = 1,000?$$

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

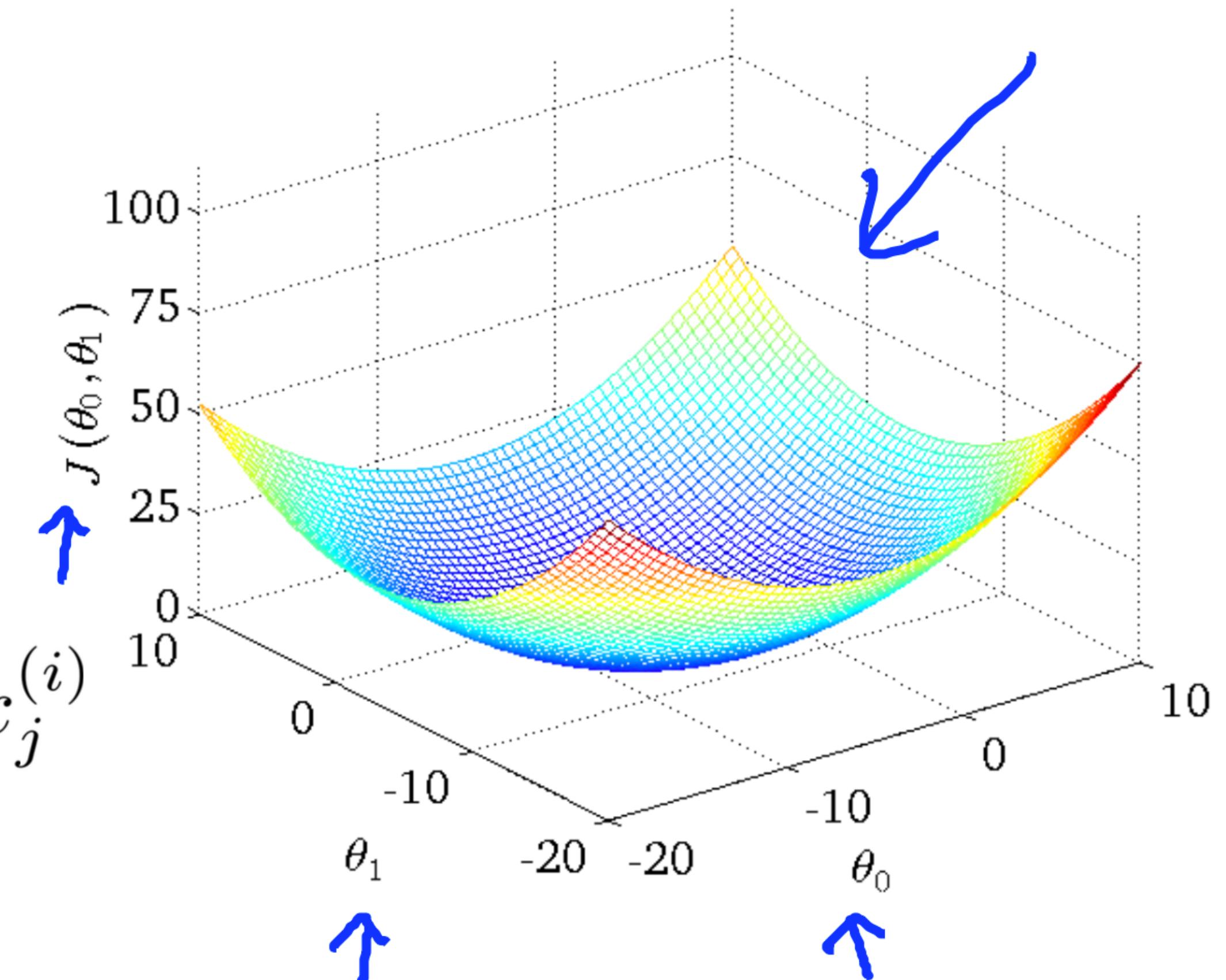


Gradient Descent with Large Datasets

Stochastic Gradient Descent

Linear regression with gradient descent Recap

- $$h_{\theta}(x) = \sum_{j=0}^n \theta_j x_j$$
- $$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$
- Repeat {
 - $\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$
(for every $j = 0, \dots, n$)



Gradient Descent with Large Datasets

Stochastic Gradient Descent

Linear regression with gradient descent

$$h_{\theta}(x) = \sum_{j=0}^n \theta_j x_j$$

$$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Repeat {

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

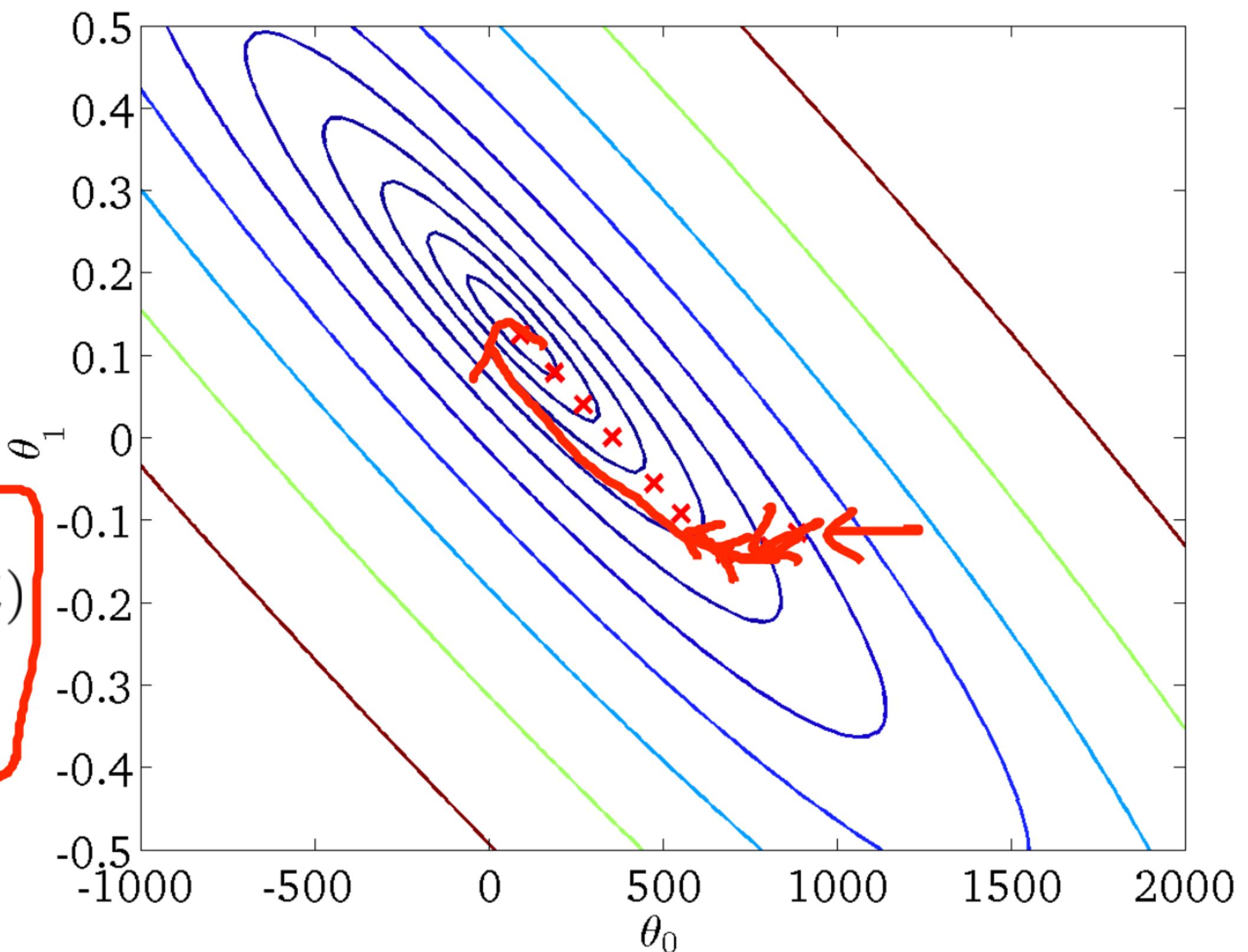
(for every $j = 0, \dots, n$)

}

$$M = \underline{\underline{300,000,000}}$$

Batch gradient descent

Recap



Gradient Descent with Large Datasets

Stochastic Gradient Descent

Batch gradient descent

$$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

Repeat {

$$\Rightarrow \theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

$$\frac{\partial}{\partial \theta_j} J_{train}(\theta)$$

(for every $j = 0, \dots, n$)

}

$$m = 300,000,000$$

Stochastic gradient descent

$$\Rightarrow \underbrace{cost(\theta, (x^{(i)}, y^{(i)}))}_{\text{one cost}} = \frac{1}{2} (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$J_{train}(\theta) = \frac{1}{m} \sum_{i=1}^m cost(\theta, (x^{(i)}, y^{(i)}))$$



Gradient Descent with Large Datasets

Stochastic Gradient Descent

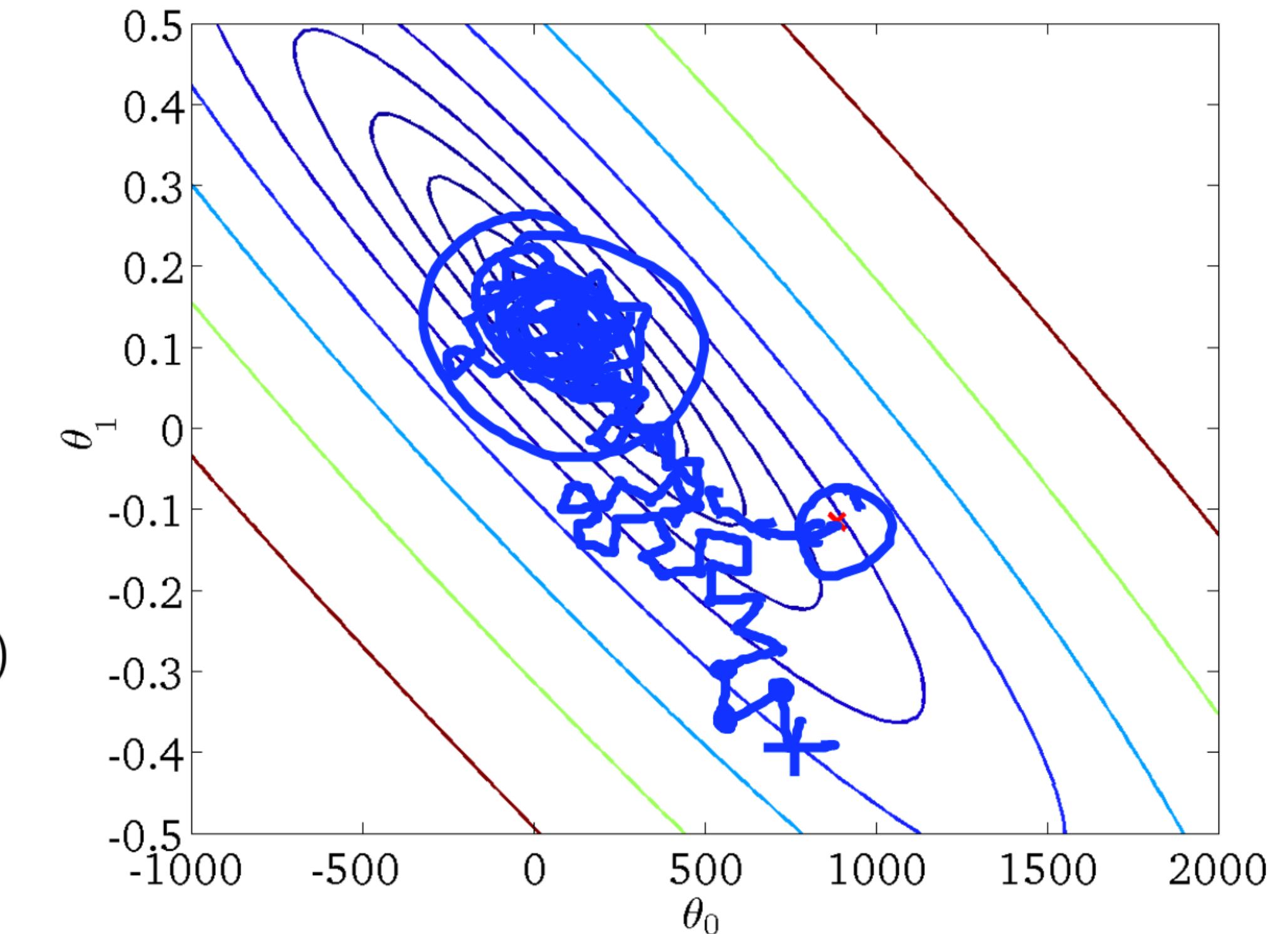
Stochastic gradient descent

1. Randomly shuffle (reorder) training examples

```

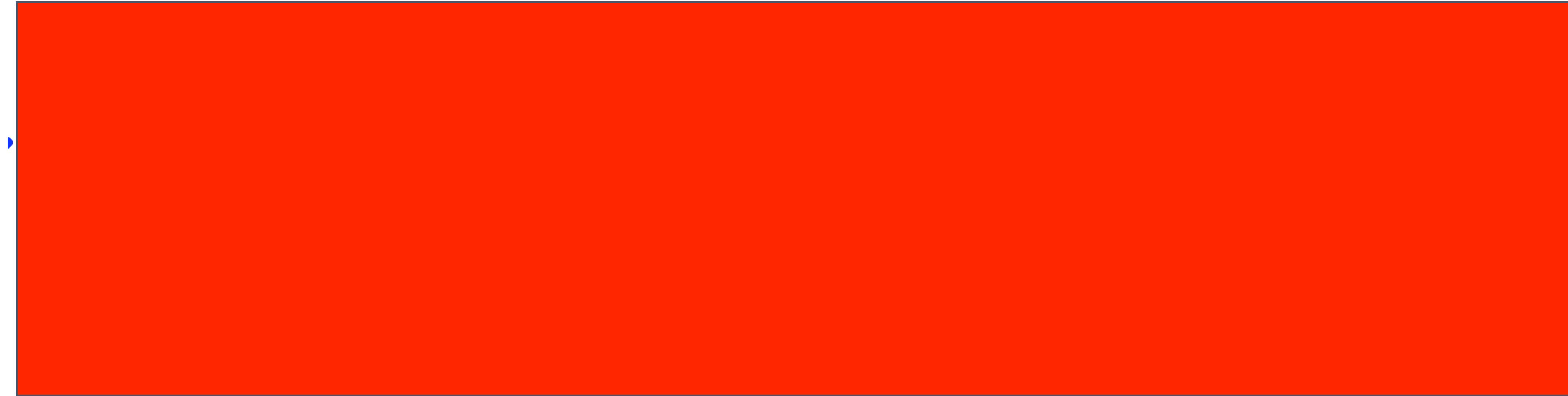
2. ↶ Repeat { 1-10x
    {
        for i := 1, ..., m {
            →  $\theta_j := \theta_j - \alpha(h_\theta(x^{(i)}) - y^{(i)})x_j^{(i)}$ 
                (for  $j = 0, \dots, n$ 
every ) →  $m = 300,000,000$ 
        }
    }
}

```



Gradient Descent with Large Datasets

Mini-Batch Gradient Descent



$b = \text{mini-batch size}$.
Get $b = 10$ examples
 $\rightarrow \theta_j := \theta_j - \alpha \frac{1}{10} \sum_{k=i}^{i+9} (\hat{y}_k - y_k) \cdot x_j^{(k)}$
 $i := i + 10$

Gradient Descent with Large Datasets

Mini-Batch Gradient Descent

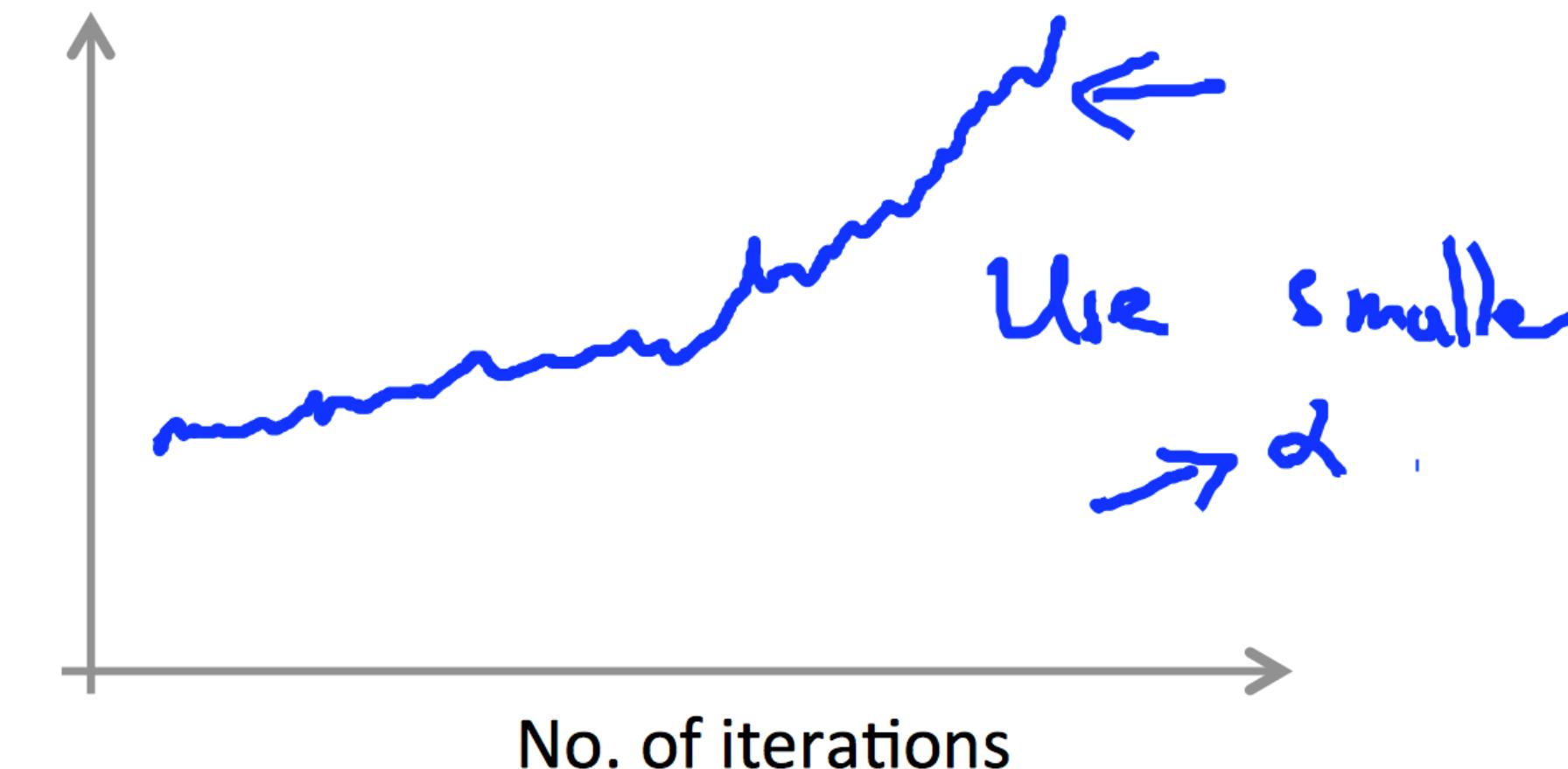
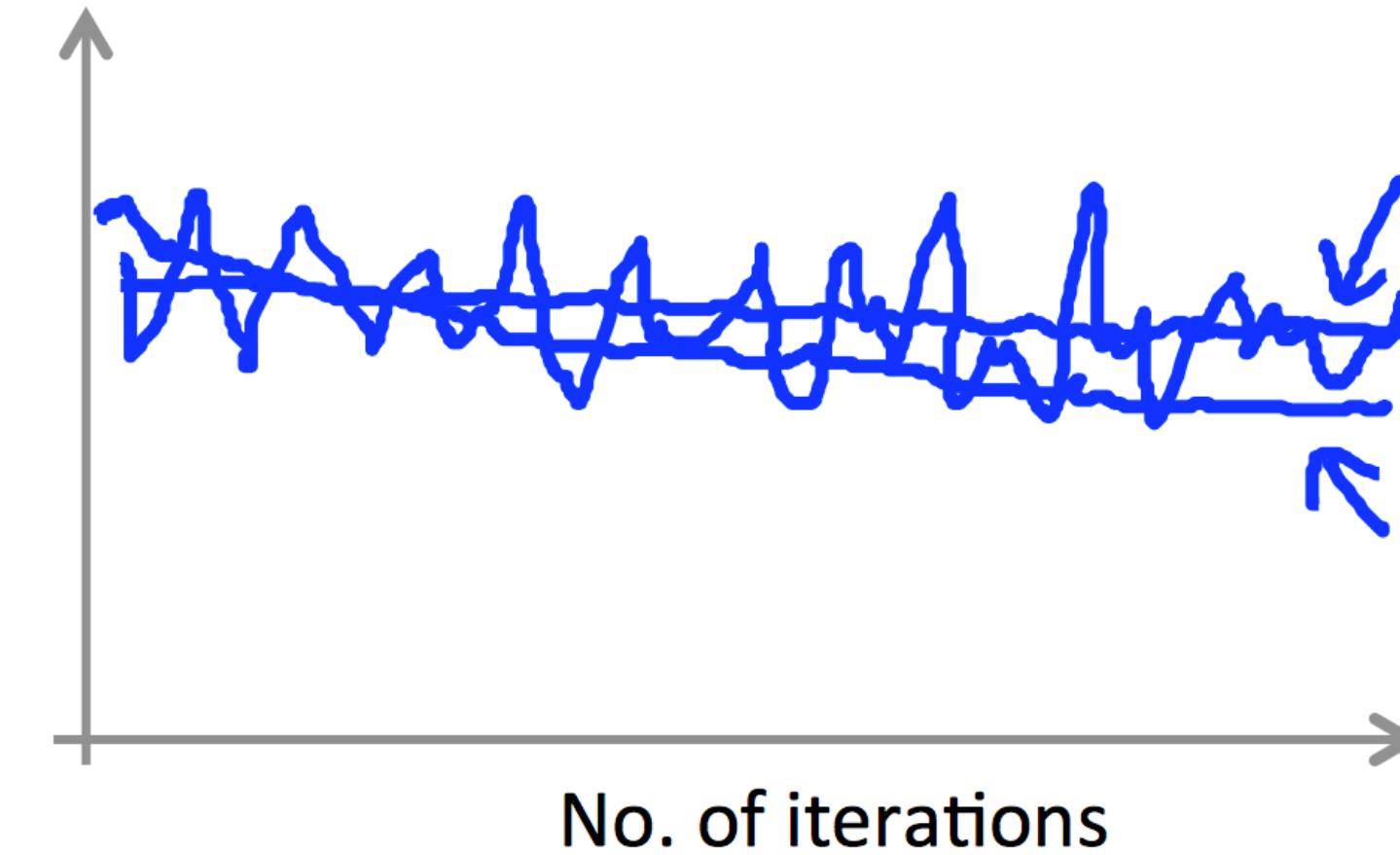
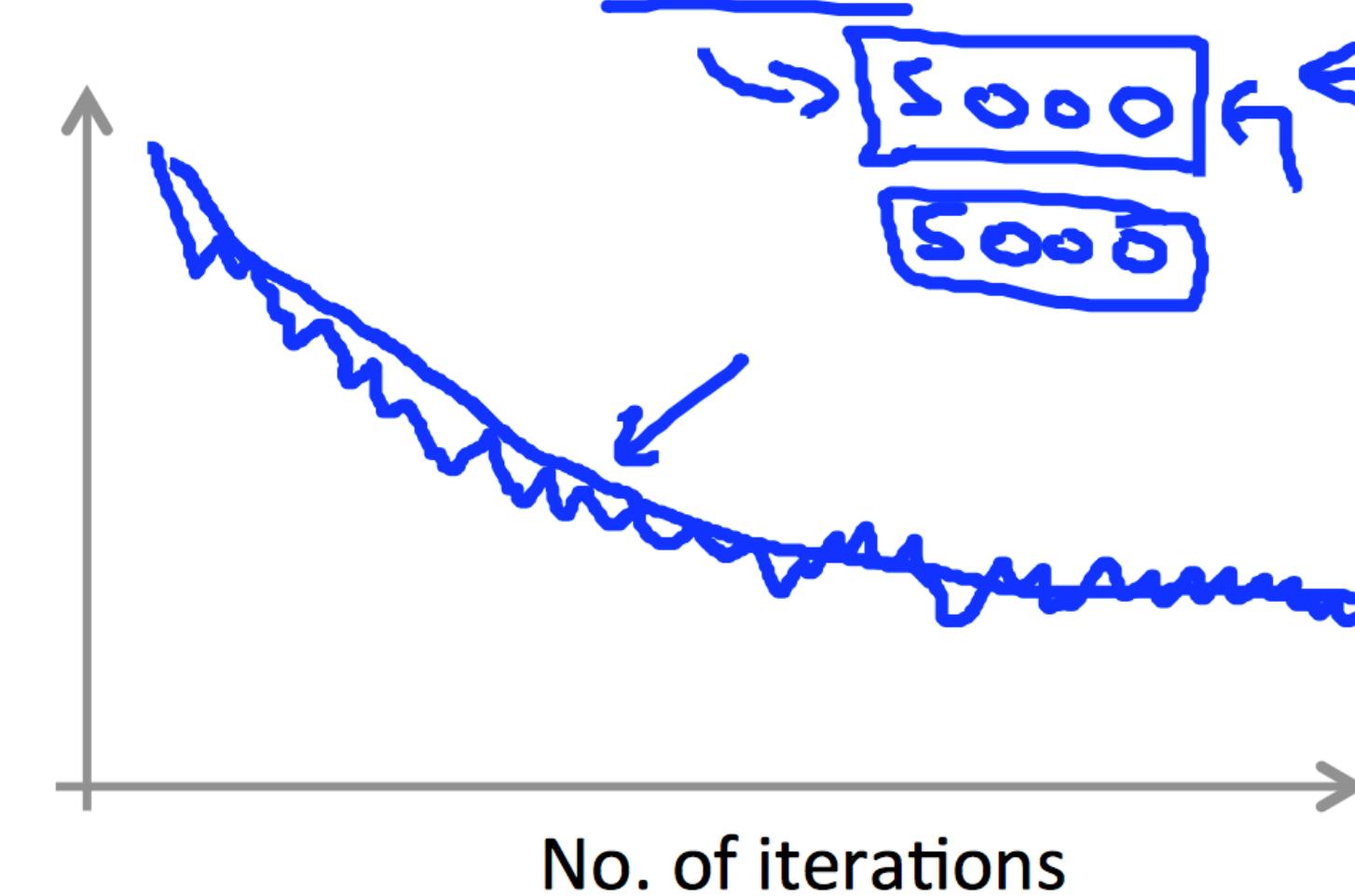
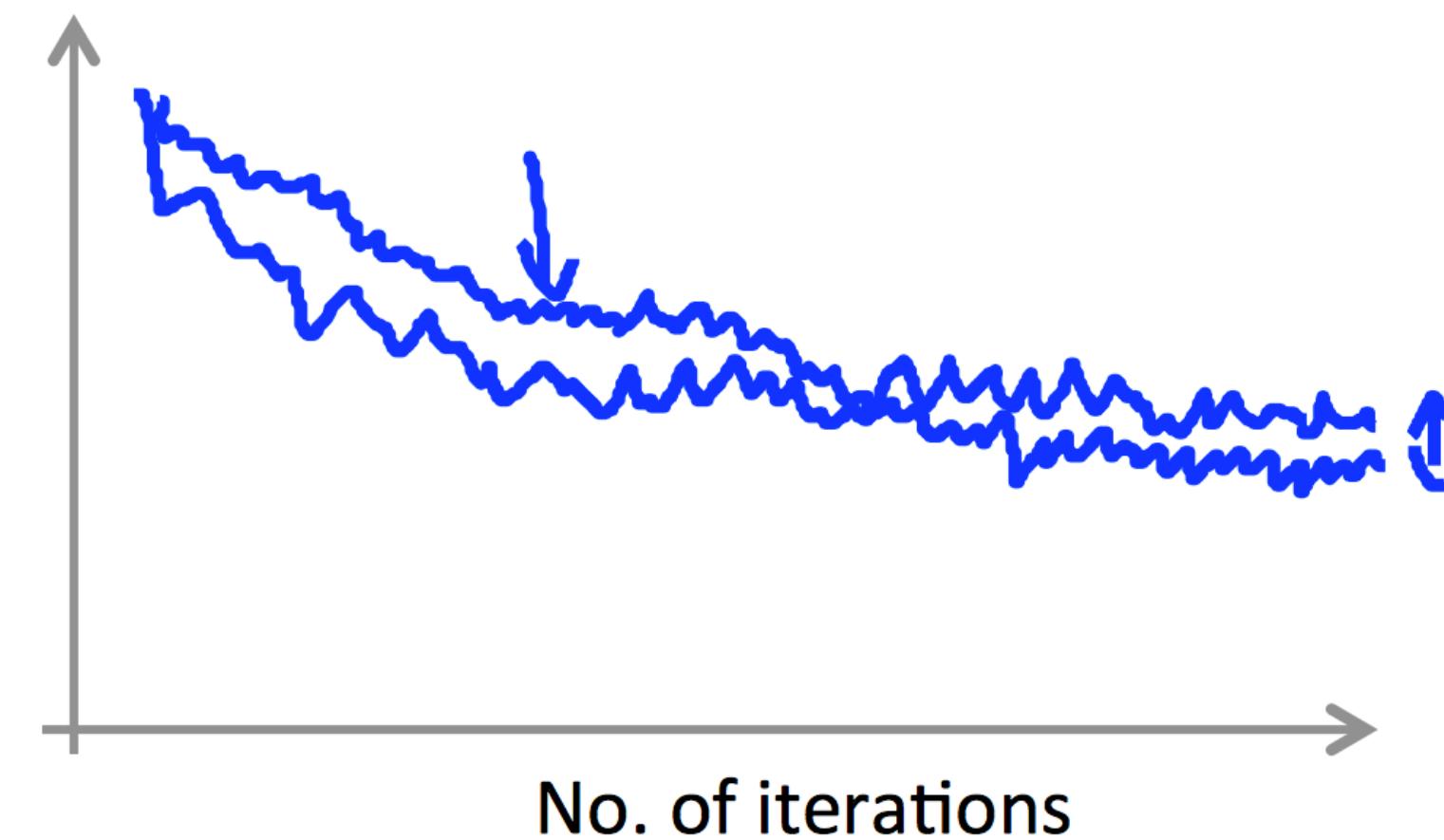


Gradient Descent with Large Datasets

Stochastic Gradient Descent Convergence

Checking for convergence

Plot $\text{cost}(\theta, (x^{(i)}, y^{(i)}))$, averaged over the last 1000 (say) example:



Gradient Descent with Large Datasets

Stochastic Gradient Descent Convergence

Checking for convergence

Batch gradient descent:

→ Plot $J_{train}(\theta)$ as a function of the number of iterations of gradient descent.

$$\rightarrow J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

M = 300, 500, 800

Gradient Descent with Large Datasets

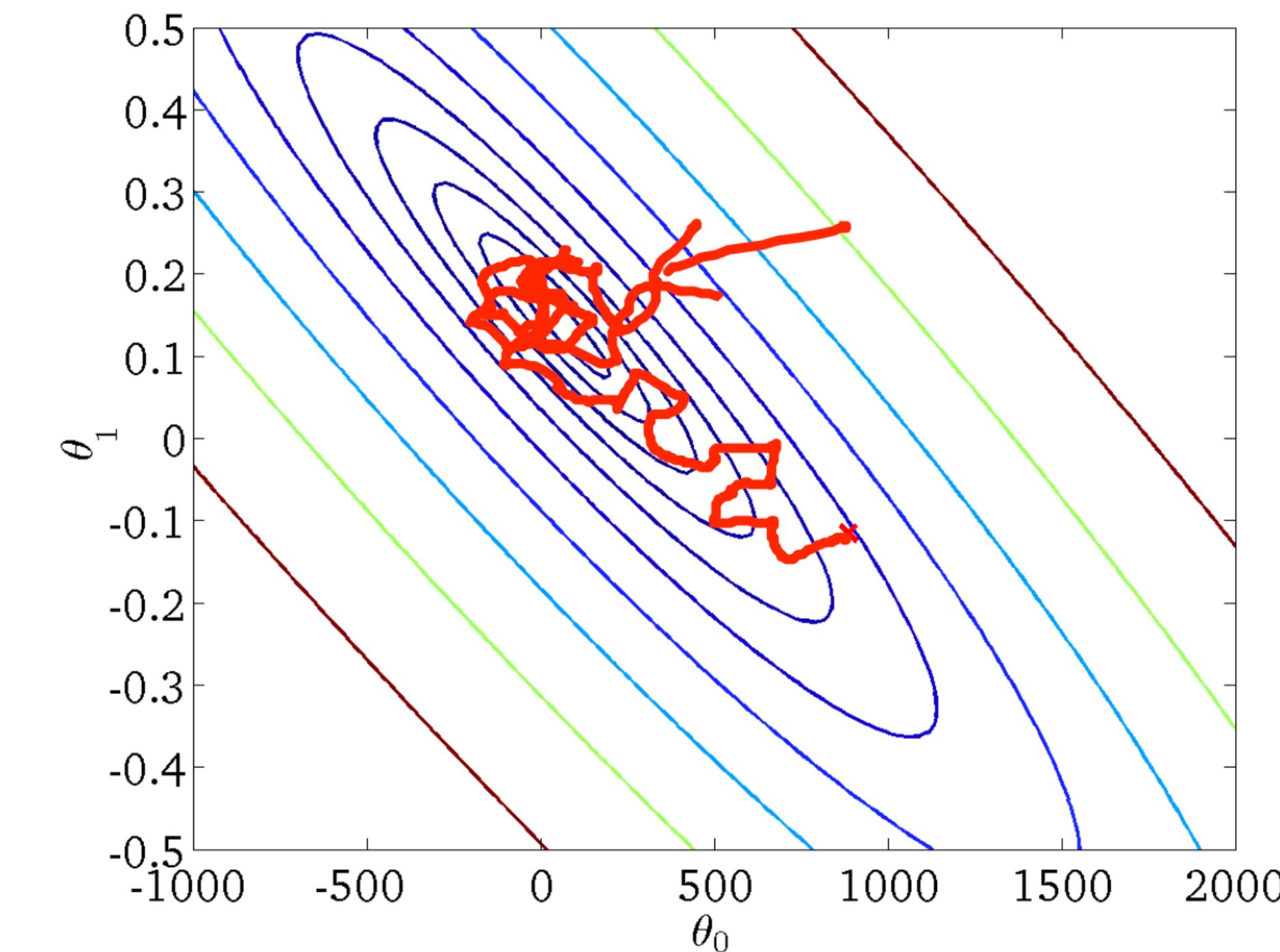
Stochastic Gradient Descent Convergence

Stochastic gradient descent

$$cost(\theta, (x^{(i)}, y^{(i)})) = \frac{1}{2}(h_\theta(x^{(i)}) - y^{(i)})^2$$

$$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m cost(\theta, (x^{(i)}, y^{(i)}))$$

1. Randomly shuffle dataset.
2. Repeat {
 - for $i := 1, \dots, m$ {
 - $\theta_j := \theta_j - \alpha(h_\theta(x^{(i)}) - y^{(i)})x_j^{(i)}$
 - (for $j = 0, \dots, n$)



Gradient Descent with Large Datasets

Stochastic Gradient Descent Convergence

Stochastic gradient descent

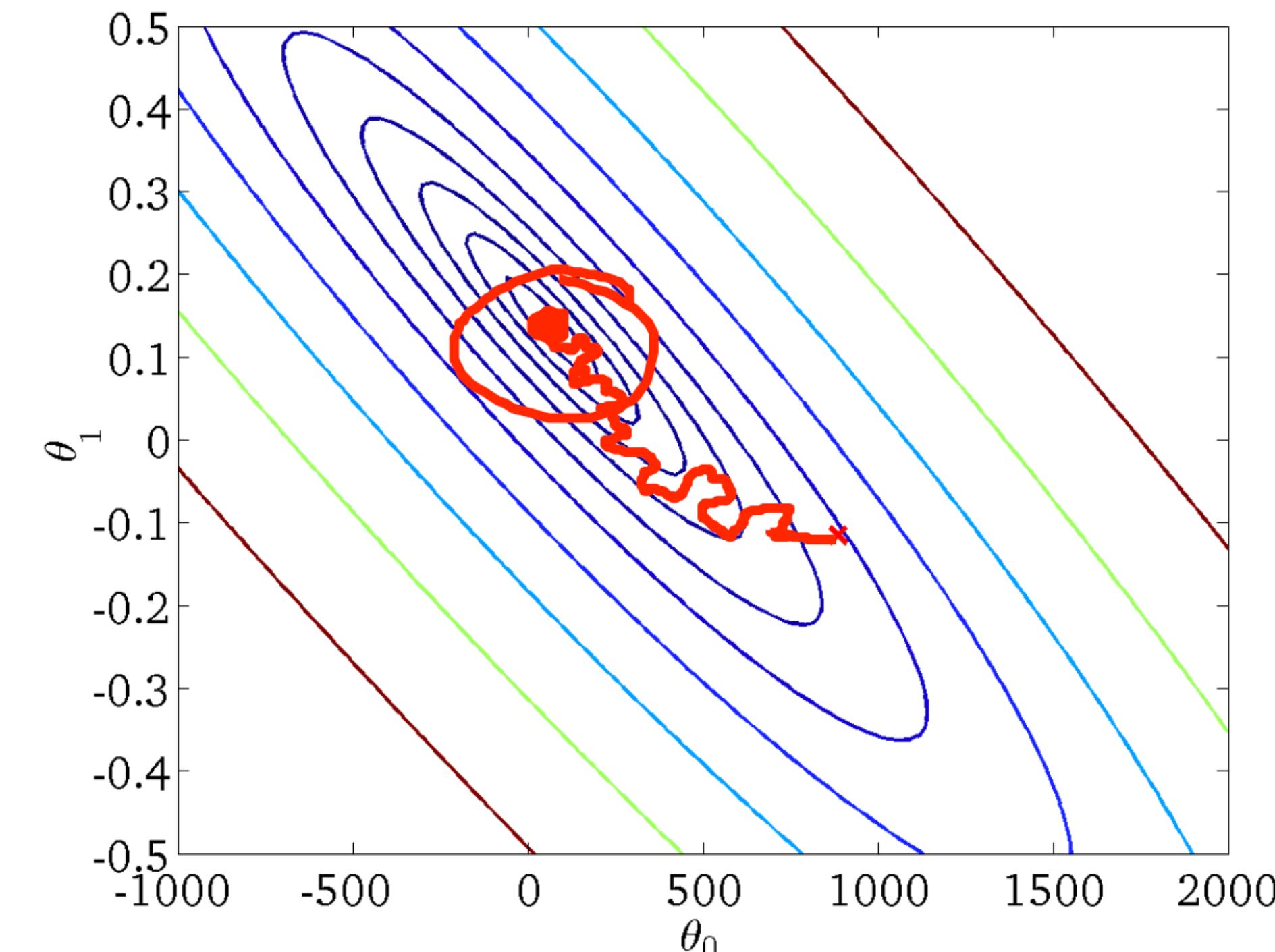
$$cost(\theta, (x^{(i)}, y^{(i)})) = \frac{1}{2}(h_\theta(x^{(i)}) - y^{(i)})^2$$

$$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m cost(\theta, (x^{(i)}, y^{(i)}))$$

1. Randomly shuffle dataset.
2. Repeat {

```
    for i := 1, ..., m      {
        theta_j := theta_j - alpha(h_theta(x^{(i)}) - y^{(i)})x_j^{(i)}
        (for j = 0, ..., n)
    }
```

```
}
```



Advanced Topics

Online learning

Shipping service website where user comes, specifies origin and destination, you offer to ship their package for some asking price, and users sometimes choose to use your shipping service ($y = 1$), sometimes not ($y = 0$).

Features x capture properties of user, of origin/destination and asking price. We want to learn $p(y = 1|x; \theta)$ to optimize price.

Repeat forever {
Get (x, y) corresponding to user.
Update θ using (x, y) :
 $\rightarrow \theta_j := \theta_j - \alpha (h_\theta(x) - y) \cdot x_j$ $(j=0, \dots, n)$
 \rightarrow Can adapt to changing user preference.

Advanced Topics

Other online learning example:

Product search (learning to search)

User searches for “Android phone 1080p camera” 

Have 100 phones in store. Will return 10 results.

→ $x = \text{features of phone, how many words in user query match name of phone, how many words in query match description of phone, etc.}$

→ $y = 1$ if user clicks on link. $y = 0$

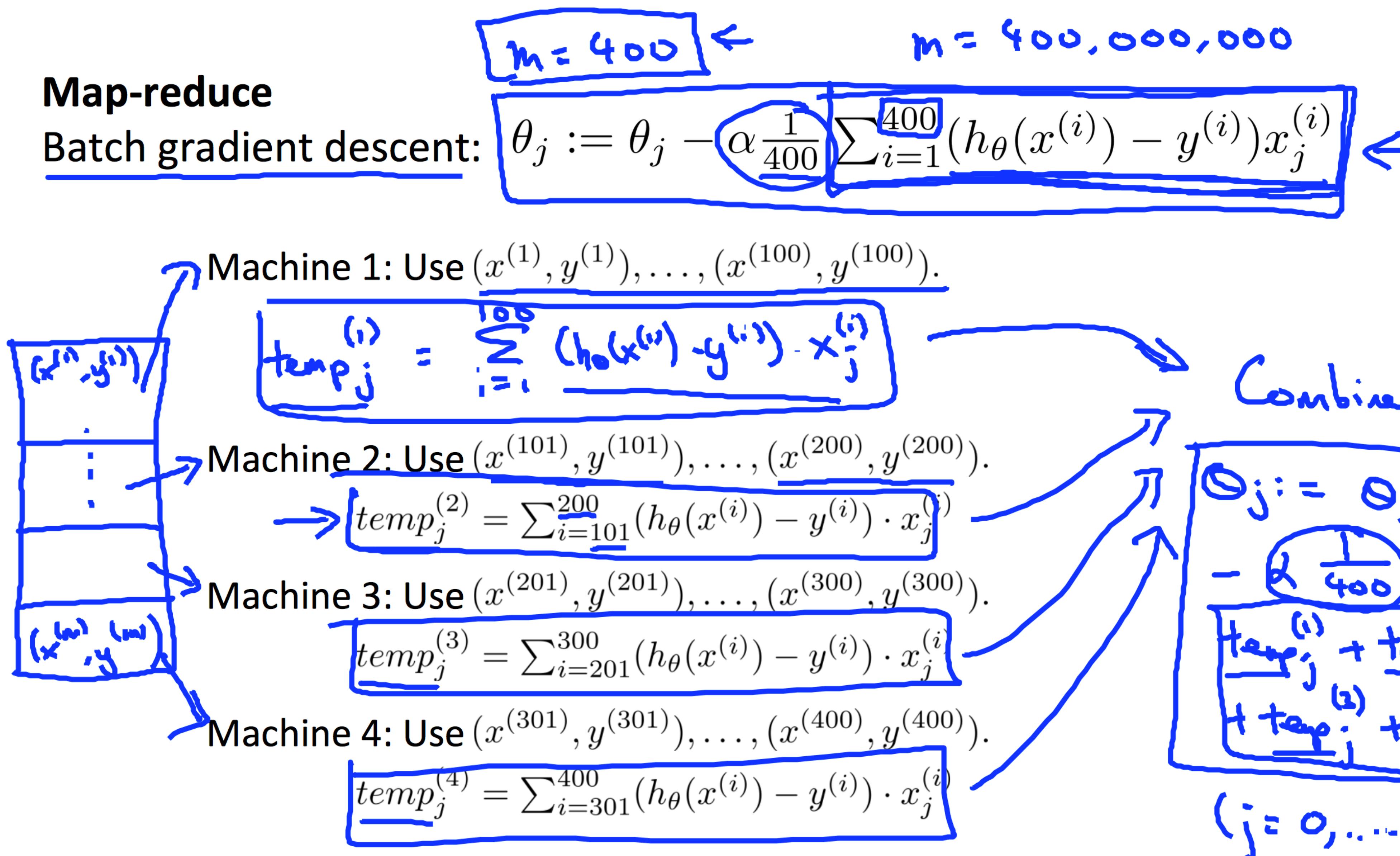
(x, y) 
otherwise 

→ Learn $p(y = 1|x; \theta)$.  predicted CTR

→ Use to show user the 10 phones they’re most likely to click on.

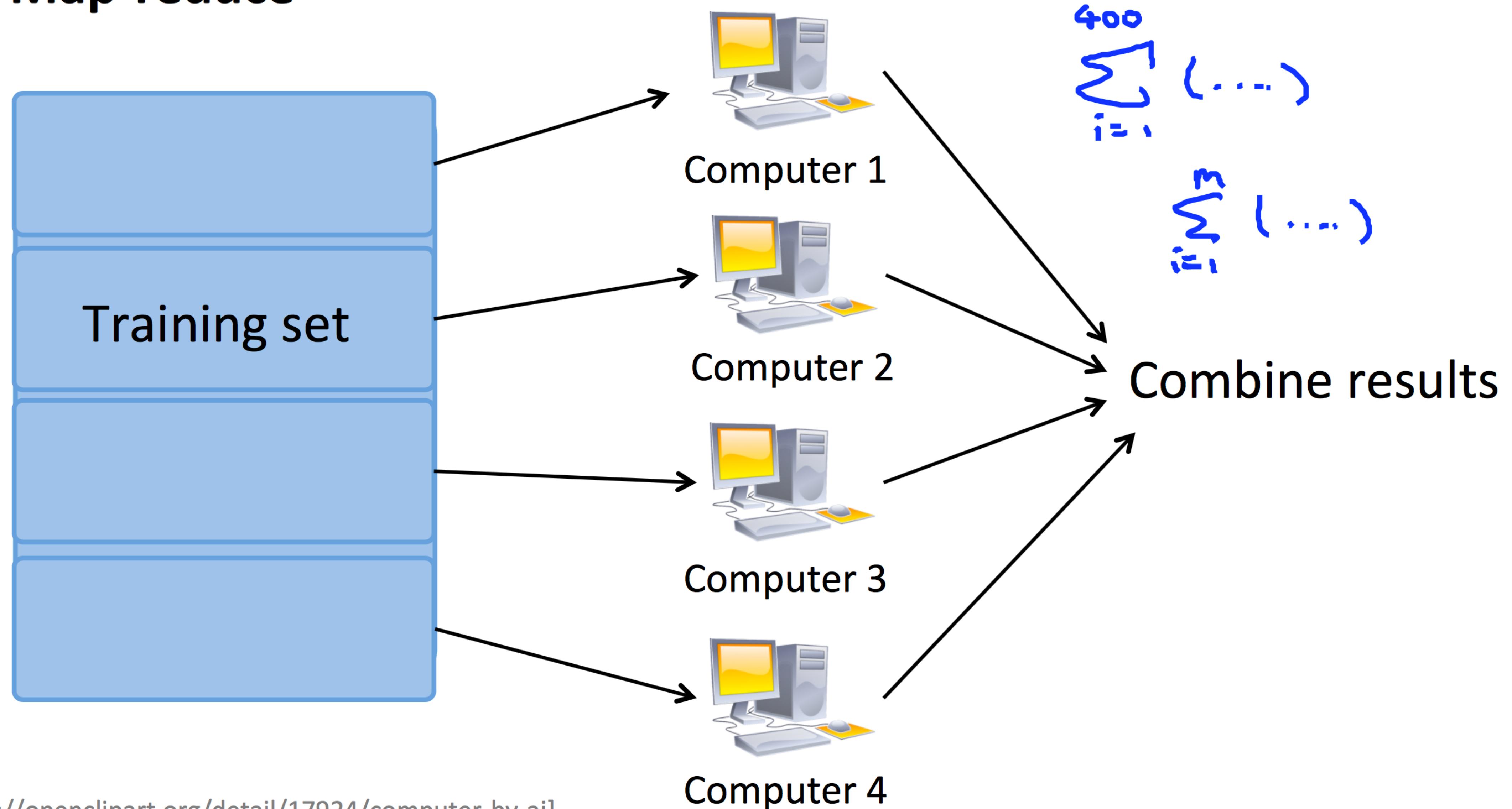
Other examples: Choosing special offers to show user; customized selection of news articles; product recommendation; ...

Advanced Topics



Advanced Topics

Map-reduce



Map-reduce and summation over the training set

Many learning algorithms can be expressed as computing sums of functions over the training set.

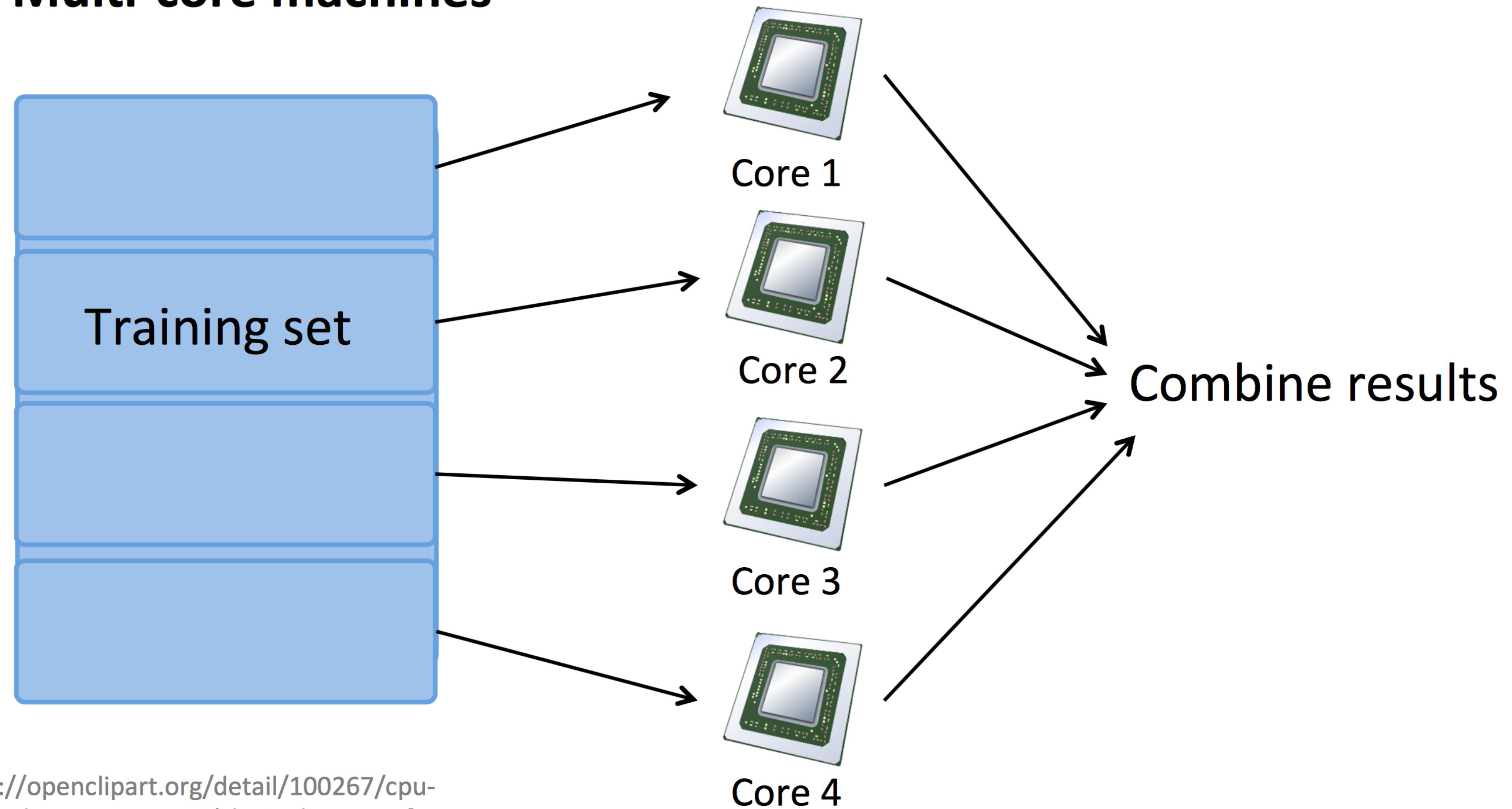
E.g. for advanced optimization, with logistic regression, need:

$$\Rightarrow \underline{J_{train}(\theta)} = -\frac{1}{m} \sum_{i=1}^m \underbrace{y^{(i)} \log h_\theta(x^{(i)}) - (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)}))}_{\textcircled{C}}$$

$$\Rightarrow \underline{\frac{\partial}{\partial \theta_j} J_{train}(\theta)} = \frac{1}{m} \sum_{i=1}^m \underbrace{(h_\theta(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}}_{\textcircled{C}}$$

$\text{temp}^{(i)}$ $\text{temp}_j^{(i)} \leftarrow$

Multi-core machines



Recursive Feature Elimination

The Recursive Feature Elimination (RFE) method is a feature selection approach. It works by recursively removing attributes and building a model on those attributes that remain. It uses the model accuracy to identify which attributes (and combination of attributes) contribute the most to predicting the target attribute.

This recipe shows the use of RFE on the iris dataset to select 3 attributes.

```
# Recursive Feature Elimination
from sklearn import datasets
from sklearn.feature_selection import RFE
from sklearn.linear_model import LogisticRegression
# load the iris datasets
dataset = datasets.load_iris()
# create a base classifier used to evaluate a subset of attributes
model = LogisticRegression()
# create the RFE model and select 3 attributes
rfe = RFE(model, 3)
rfe = rfe.fit(dataset.data, dataset.target)
# summarize the selection of the attributes
print(rfe.support_)
print(rfe.ranking_)
```

Feature Importance

Methods that use ensembles of decision trees (like Random Forest and Extra Trees) can also compute the relative importance of each attribute. These importance values can be used to inform a feature selection process.

This recipe shows the construction of an Extra Trees ensemble of the iris dataset and the display of the relative feature importance.

```
# Feature Importance
from sklearn import datasets
from sklearn import metrics
from sklearn.ensemble import ExtraTreesClassifier
# load the iris datasets
dataset = datasets.load_iris()
# fit an Extra Trees model to the data
model = ExtraTreesClassifier()
model.fit(dataset.data, dataset.target)
# display the relative importance of each attribute
print(model.feature_importances_)
```

Advanced Topics

Cross Validation

Cross validation is a way of evaluating a model on a dataset. It provides an estimation of the accuracy of the model if it were to make predictions on previously unseen data. Cross validation estimations are used to aid in the selection of a robust model that is fit for purpose.

This recipe estimates the performance of logistic regression on the iris dataset using k-fold cross validation with 10 folds.

```
# Cross Validation Classification
import numpy as np
from sklearn import datasets
from sklearn.linear_model import LogisticRegression
from sklearn import cross_validation
# load the iris datasets
dataset = datasets.load_iris()
# prepare cross validation folds
num_folds = 10
num_instances = len(dataset.data)
kfold = cross_validation.KFold(n=num_instances, n_folds=num_folds)
# prepare a Logistic Regression model
model = LogisticRegression()
# evaluate the model k-fold cross validation
results = cross_validation.cross_val_score(model, dataset.data, dataset.target, cv=kfold)
# display the mean classification accuracy on each fold
print(results)
# display the mean and stdev of the classification accuracy
print(results.mean())
print(results.std())
```

Grid Search Parameter Tuning

Machine learning models are parameterized so that their behavior can be tuned for a given problem. Models can have many parameters and finding the best combination of parameters can be treated as a search problem.

Grid search is an approach to parameter tuning that will methodically build and evaluate a model for each combination of algorithm parameters specified in a grid.

This recipe evaluates different alpha values for the Ridge Regression algorithm on the diabetes dataset.

```
# Grid Search for Algorithm Tuning
import numpy as np
from sklearn import datasets
from sklearn.linear_model import Ridge
from sklearn.grid_search import GridSearchCV
# load the diabetes datasets
dataset = datasets.load_diabetes()
# prepare a range of alpha values to test
alphas = np.array([1,0.1,0.01,0.001,0.0001,0])
# create and fit a ridge regression model, testing each alpha
model = Ridge()
grid = GridSearchCV(estimator=model, param_grid=dict(alpha=alphas))
grid.fit(dataset.data, dataset.target)
print(grid)
# summarize the results of the grid search
print(grid.best_score_)
print(grid.best_estimator_.alpha)
```

Random Search Parameter Tuning

Random search is an approach to parameter tuning that will sample algorithm parameters from a random distribution (i.e. uniform) for a fixed number of iterations. A model is constructed and evaluated for each combination of parameters chosen.

This recipe evaluates different alpha random values between 0 and 1 for the Ridge Regression algorithm on the diabetes dataset.

```
# Randomized Search for Algorithm Tuning
import numpy as np
from scipy.stats import uniform as sp_rand
from sklearn import datasets
from sklearn.linear_model import Ridge
from sklearn.grid_search import RandomizedSearchCV
# load the diabetes datasets
dataset = datasets.load_diabetes()
# prepare a uniform distribution to sample for the alpha parameter
param_grid = {'alpha': sp_rand()}
# create and fit a ridge regression model, testing random alpha values
model = Ridge()
rsearch = RandomizedSearchCV(estimator=model, param_distributions=param_grid, n_iter=100)
rsearch.fit(dataset.data, dataset.target)
print(rsearch)
# summarize the results of the random parameter search
print(rsearch.best_score_)
print(rsearch.best_estimator_.alpha)
```

Random Search Parameter Tuning

Random search is an approach to parameter tuning that will sample algorithm parameters from a random distribution (i.e. uniform) for a fixed number of iterations. A model is constructed and evaluated for each combination of parameters chosen.

This recipe evaluates different alpha random values between 0 and 1 for the Ridge Regression algorithm on the diabetes dataset.

```
# Randomized Search for Algorithm Tuning
import numpy as np
from scipy.stats import uniform as sp_rand
from sklearn import datasets
from sklearn.linear_model import Ridge
from sklearn.grid_search import RandomizedSearchCV
# load the diabetes datasets
dataset = datasets.load_diabetes()
# prepare a uniform distribution to sample for the alpha parameter
param_grid = {'alpha': sp_rand()}
# create and fit a ridge regression model, testing random alpha values
model = Ridge()
rsearch = RandomizedSearchCV(estimator=model, param_distributions=param_grid, n_iter=100)
rsearch.fit(dataset.data, dataset.target)
print(rsearch)
# summarize the results of the random parameter search
print(rsearch.best_score_)
print(rsearch.best_estimator_.alpha)
```

Advanced Topics



Cornell University
Library

We

arXiv.org > cs > arXiv:1309.0238

Search or Article ID inside arXiv

All papers



Broaden your

(Help | Advanced search)

Computer Science > Learning

API design for machine learning software: experiences from the scikit-learn project

Lars Buitinck (ILPS), Gilles Louppe, Mathieu Blondel, Fabian Pedregosa (INRIA Saclay – Ile de France), Andreas Mueller, Olivier Grisel, Vlad Niculae, Peter Prettenhofer, Alexandre Gramfort (INRIA Saclay – Ile de France, LTCI), Jaques Grobler (INRIA Saclay – Ile de France), Robert Layton, Jake Vanderplas, Arnaud Joly, Brian Holt, Gaël Varoquaux (INRIA Saclay – Ile de France)

(Submitted on 1 Sep 2013)

Scikit-learn is an increasingly popular machine learning library. Written in Python, it is designed to be simple and efficient, accessible to non-experts, and reusable in various contexts. In this paper, we present and discuss our design choices for the application programming interface (API) of the project. In particular, we describe the simple and elegant interface shared by all learning and processing units in the library and then discuss its advantages in terms of composition and reusability. The paper also comments on implementation details specific to the Python ecosystem and analyzes obstacles faced by users and developers of the library.

Subjects: [Learning \(cs.LG\)](#); Mathematical Software (cs.MS)

Journal reference: European Conference on Machine Learning and Principles and Practices of Knowledge Discovery in Databases (2013)

Cite as: [arXiv:1309.0238 \[cs.LG\]](#)

(or [arXiv:1309.0238v1 \[cs.LG\]](#) for this version)

Submission history

From: Gael Varoquaux [[view email](#)]

[v1] Sun, 1 Sep 2013 16:22:48 GMT (28kb)

[Which authors of this paper are endorsers?](#) | [Disable MathJax](#) ([What is MathJax?](#))

<https://arxiv.org/abs/1309.0238>

Deep Learning



Upload your happiest moments in FY18 Asia GBB Ready - 🐦 #AsiaICGBB

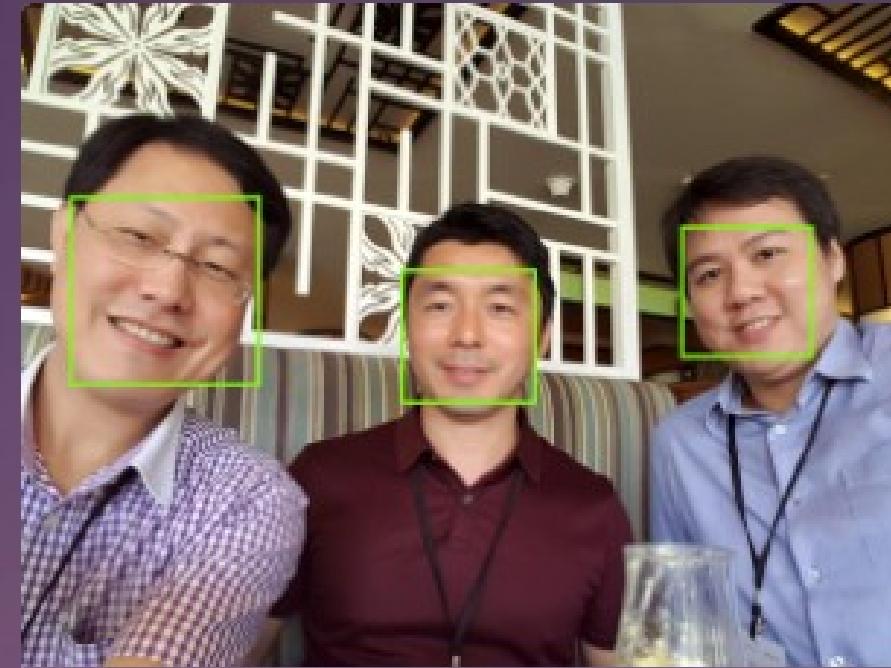
Dont' have 🐦? Attach the photos and 📧 wely.lau@microsoft.com with subject containing #AsiaICGBB

[Top 20 the happiest photos](#)

[Top 20 the unhappiest photos](#)



👤 Christina Peh (Adecco) (OneDrive)
★ 0.99999765 😊 3



👤 TakayukiHoshino (Twitter)
★ 0.99994497 😊 3



👤 Christina Peh (Adecco) (OneDrive)
★ 0.99815755 😊 3



👤 Margaret Synan (OneDrive)
★ 0.97417196 😊 8



👤 Genevieve Ang (Adecco) (OneDrive)
★ 0.89848932 😊 10



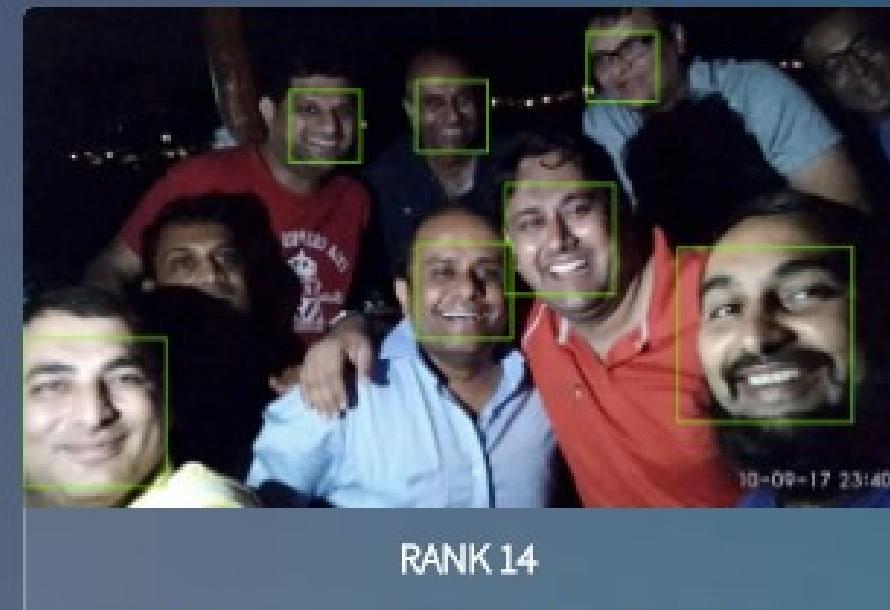
👤 Genevieve Ang (Adecco) (OneDrive)
★ 0.99999755 😊 2



👤 Genevieve Ang (Adecco) (OneDrive)
★ 0.999992192 😊 2



👤 Genevieve Ang (Adecco) (OneDrive)
★ 0.9999922286 😊 5



👤 Yusuf Rangwala (OneDrive)
★ 0.94711960 😊 7



👤 Danette Seward (OneDrive)
★ 0.86474696 😊 12



OPEN-SOURCE ML TOOLS

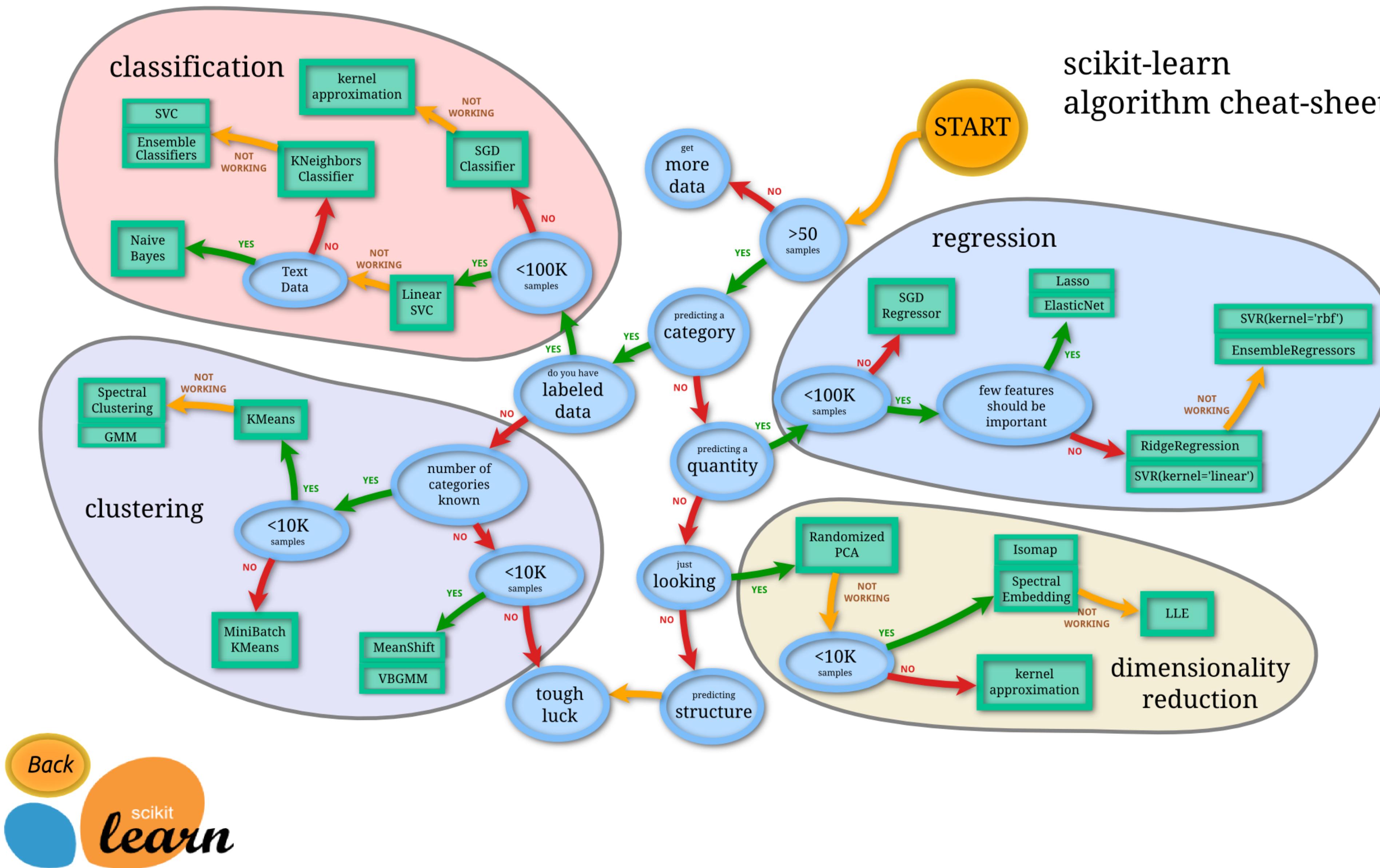
What is scikit-learn?

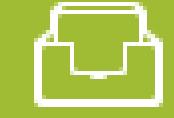
Scikit-learn provides a range of supervised and unsupervised learning algorithms via a consistent interface in Python. It is licensed under a permissive simplified BSD license and is distributed under many Linux distributions, encouraging academic and commercial use. This is stack that includes:

- \$ **NumPy**: Base n-dimensional array package
- \$ **SciPy**: Fundamental library for scientific computing
- \$ **Matplotlib**: Comprehensive 2D/3D plotting
- \$ **IPython**: Enhanced interactive console
- \$ **Sympy**: Symbolic mathematics
- \$ **Pandas**: Data structures and analysis



OPEN-SOURCE ML TOOLS





OPEN-SOURCE ML TOOLS

Regression

Models

- Linear regression
- Regularization:
Ridge (L2), Lasso (L1)

Algorithms

- Gradient descent
- Coordinate descent

Concepts

- Loss functions, bias-variance tradeoff, cross-validation, sparsity, overfitting, model selection



OPEN-SOURCE ML TOOLS

Classification

Models

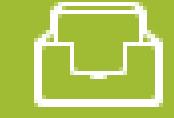
- Linear classifiers
(logistic regression, SVMs, perceptron)
- Kernels
- Decision trees

Algorithms

- Stochastic gradient descent
- Boosting

Concepts

- Decision boundaries, MLE, ensemble methods, random forests, CART, online learning



OPEN-SOURCE ML TOOLS

Clustering and Retrieval

Models

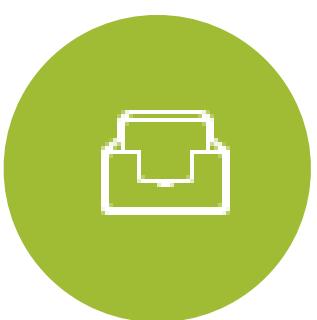
- Nearest neighbors
- Clustering, mixtures of Gaussians
- Latent Dirichlet allocation (LDA)

Algorithms

- KD-trees, locality-sensitive hashing (LSH)
- K-means
- Expectation-maximization (EM)

Concepts

- Distance metrics, approximation algorithms, hashing, sampling algorithms, scaling up with map-reduce



GETTING STARTED WITH SCIKIT-LEARN



ANACONDA®

ACCELERATE

Time-to-Value



INNOVATE faster through managed agile experimentation

MOVE from analysis to deployment immediately

DELIVER high performance analytics processing

CONNECT

Data, Analytics & Compute



LEVERAGE innovative open source analytics to extract value from data

MAXIMIZE your computational power to easily analyze all your data

CONNECT and integrate all your data sources for predictive models

EMPOWER

Data Science Teams



ITERATE quickly to create powerful analysis and predictive models

COLLABORATE and share with your data science team

PUBLISH interactive results to the business



CLOUD-BASED ML TOOLS

Amazon Machine Learning

Amazon Machine Learning makes it easy for developers of all skill levels to use machine learning (ML) technology. Amazon Machine Learning is a managed service for building ML models and generating predictions that enable the development of robust, scalable smart applications.

[Get started](#)

AWS ML

Machine Learning
Powerful cloud based analytics, now part of Cortana Intelligence Suite

[Get started now >](#)

Learn more about Cortana Intelligence Suite >

bigml PRIVATE DEPLOYMENTS EVENTS GALLERY DOCUMENTATION HELP & SUPPORT Login AU

FREE PRO Subscription for students and educators Sign up NOW with your .edu email

Single Platform for All Predictive Use Cases Machine Learning built for Developers

sign up here In Australia or New Zealand? sign up here Instant access. No credit card required.

BIG ML

turi create intelligence™

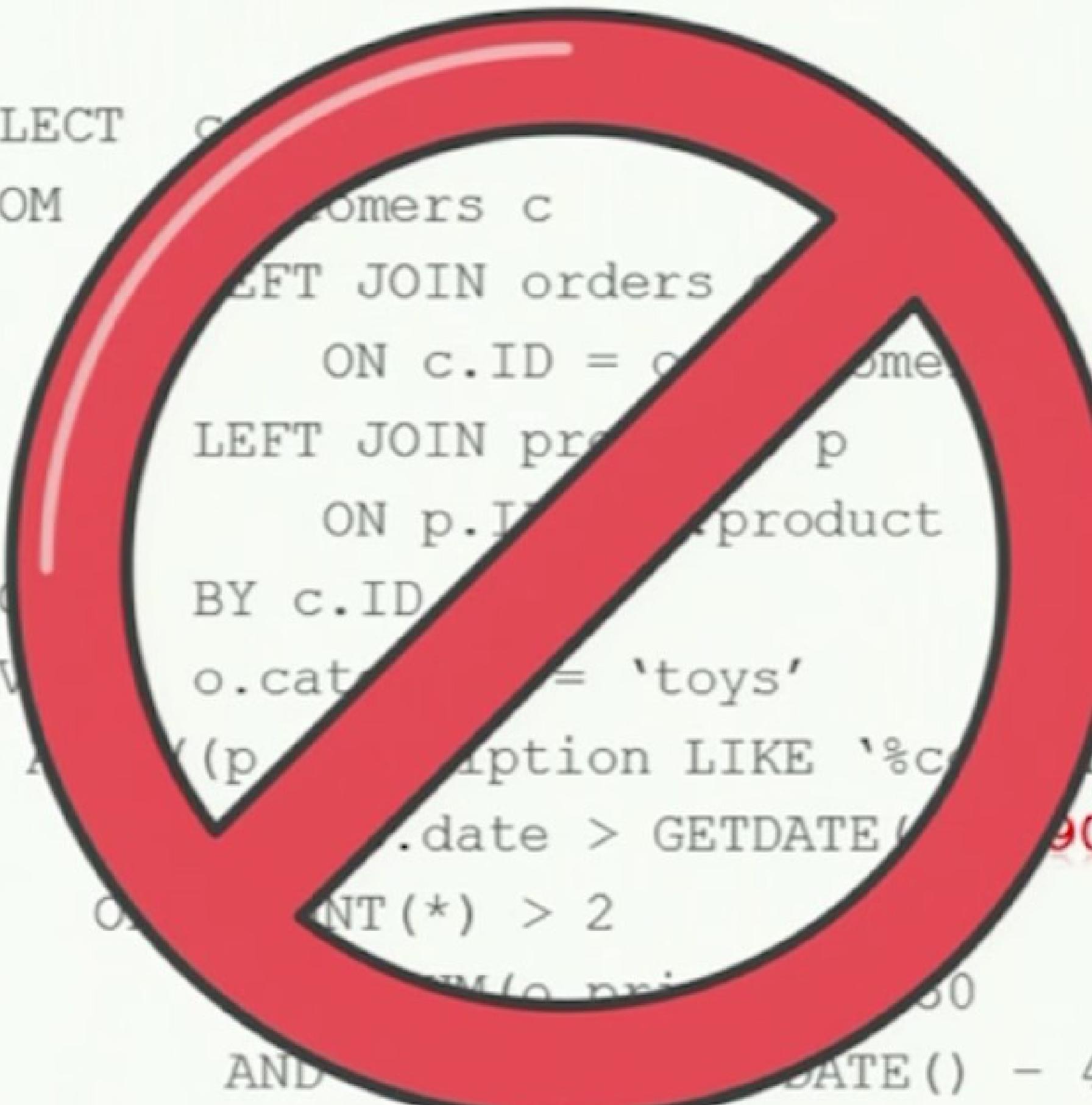
RECOMMENDER SENTIMENT ANALYSIS CHURN PREDICTOR

DATA MATCHING CLICKTHROUGH PREDICTOR LEAD SCORING

TURI ML

CLOUD-BASED ML TOOLS

Smart applications by *counterexample*



```
SELECT
  ...
FROM
  ...
  LEFT JOIN orders
    ON c.ID = o.customer_id
  LEFT JOIN products p
    ON p.ID = o.product_id
GROUP BY
  ...
HAVING
  ...
  AND p.description LIKE '%color%' OR
  ...
  AND o.date > GETDATE() - 90
  ...
  COUNT(*) > 2
  ...
  SUM(o.price) < 50
  ...
  AND DATEDIFF(o.order_date, GETDATE()) - 40
)
```

Use machine learning technology to **learn** your business rules from data!

CLOUD-BASED ML TOOLS

Three types of data-driven development



**Retrospective
analysis and
reporting**

Amazon Redshift,
Amazon RDS
Amazon S3
Amazon EMR



**Here-and-now
real-time processing
and dashboards**

Amazon Kinesis
Amazon EC2
AWS Lambda



**Predictions
to enable smart
applications**

CLOUD-BASED ML TOOLS

And a few more examples...

Fraud detection

Detecting fraudulent transactions, filtering spam emails, flagging suspicious reviews,...

Personalization

Recommending content, predictive content loading, improving user experience,...

Targeted marketing

Matching customers and offers, choosing marketing campaigns, cross-selling and up-selling,...

Content classification

Categorizing documents, matching hiring managers and resumes,...

Churn prediction

Finding customers who are likely to stop using the service, upgrade targeting,...

Customer support

Predictive routing of customer emails, social media listening,...

Three supported types of predictions



Binary classification

Predict the answer to a Yes/No question

Multiclass classification

Predict the correct category from a list

Regression

Predict the value of a numeric variable

Building smart applications with Amazon ML

1

Train
model

2

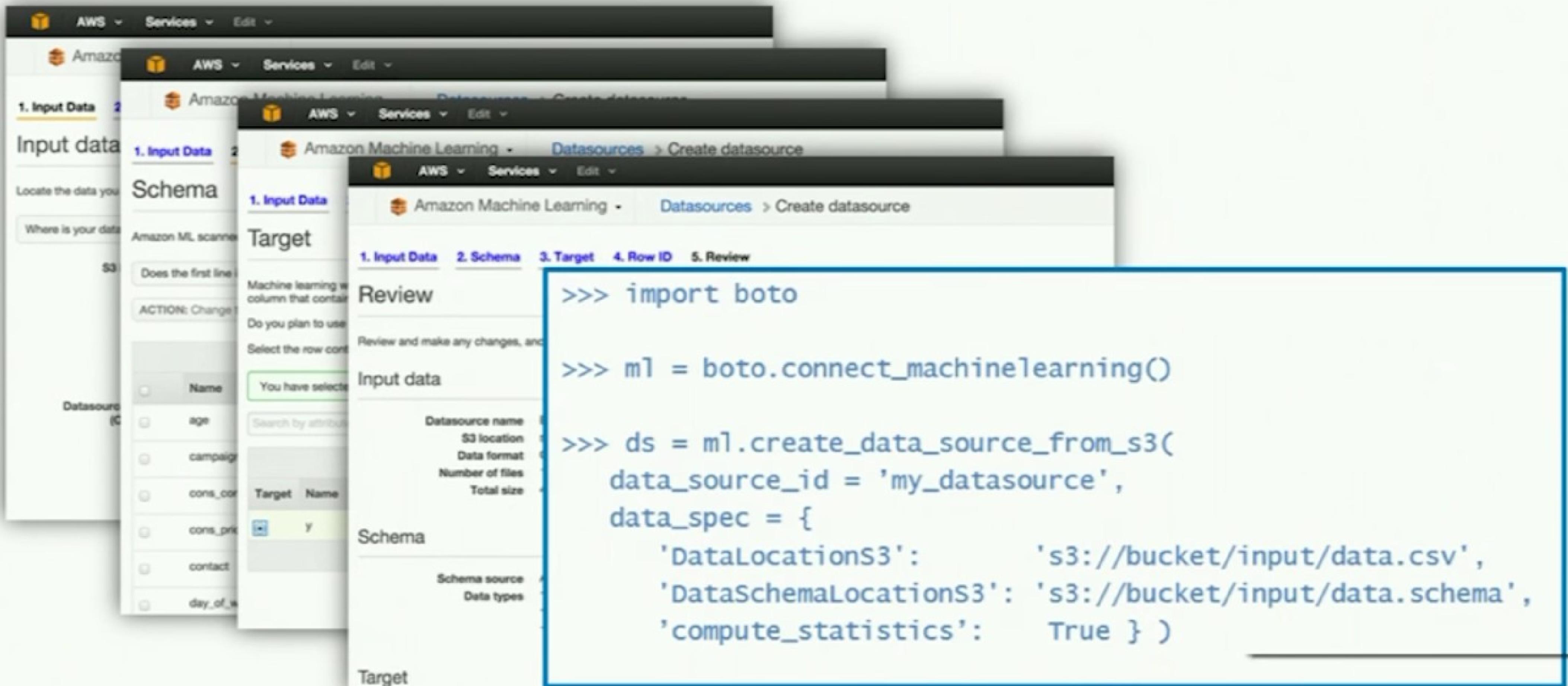
Evaluate and
optimize

3

Retrieve
predictions

CLOUD-BASED ML TOOLS

Create a datasource object

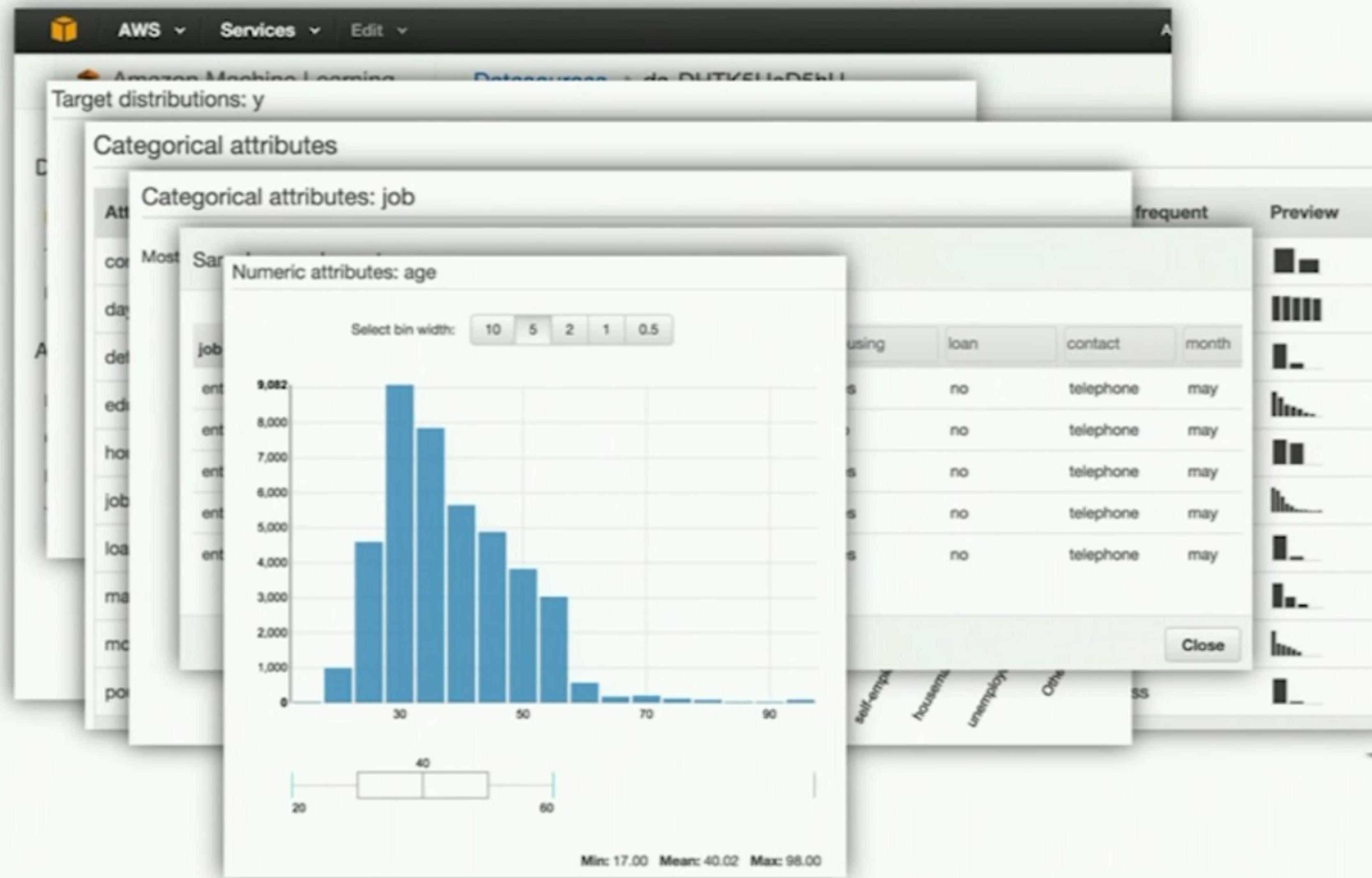


The screenshot shows the AWS Amazon Machine Learning 'Create datasource' wizard. It consists of five steps: 1. Input Data, 2. Schema, 3. Target, 4. Row ID, and 5. Review. The 'Review' step is highlighted with a blue box. Inside this box, the following Python code is displayed:

```
>>> import boto  
  
>>> ml = boto.connect_machinelearning()  
  
>>> ds = ml.create_data_source_from_s3(  
    data_source_id = 'my_datasource',  
    data_spec = {  
        'DataLocationS3':      's3://bucket/input/data.csv',  
        'DataSchemaLocationS3': 's3://bucket/input/data.schema',  
        'compute_statistics': True } )
```

CLOUD-BASED ML TOOLS

Explore and understand your data



Machine
Management,

CLOUD-BASED ML TOOLS

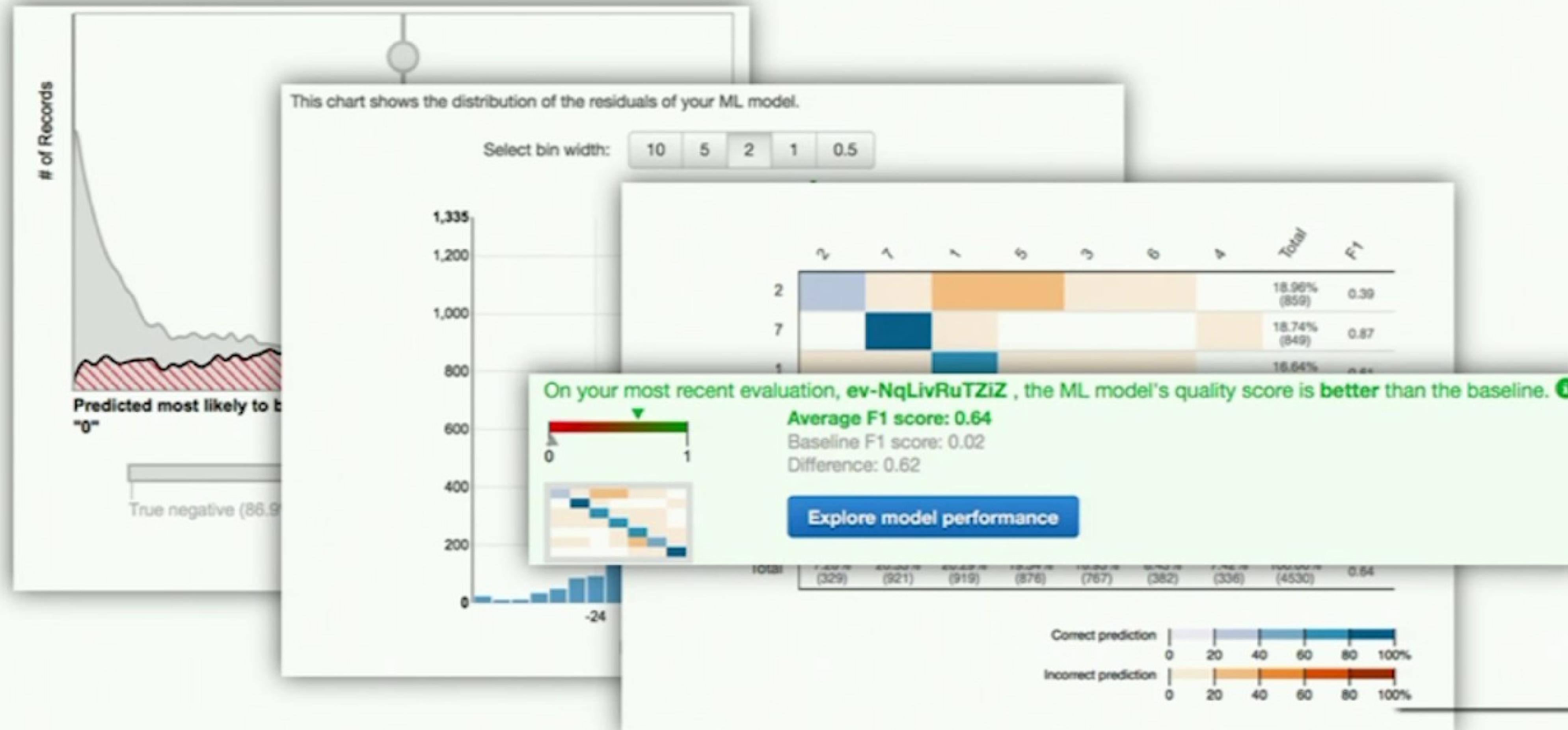
Train your model

The screenshot shows the AWS Amazon Machine Learning 'Create ML model' wizard. The current step is 'Review'. A code editor window is overlaid on the right side of the interface, displaying Python code for creating an ML model:

```
>>> import boto  
  
>>> m1 = boto.connect_machinelearning()  
  
>>> model = m1.create_ml_model(  
    ml_model_id = 'my_model',  
    ml_model_type = 'REGRESSION',  
    training_data_source_id = 'my_datasource')
```

CLOUD-BASED ML TOOLS

Explore model quality





CLOUD-BASED ML TOOLS

Accuracy

Accuracy (ACC) measures the fraction of correct predictions. The range is 0 to 1. A larger value indicates better predictive accuracy:

$$Accuracy = \frac{TP + TN}{TP + FP + FN + TN}$$

Precision

Precision measures the fraction of actual positives among those examples that are predicted as positive. The range is 0 to 1. A larger value indicates better predictive accuracy:

$$Precision = \frac{TP}{TP + FP}$$

Recall

Recall measures the fraction of actual positives that are predicted as positive. The range is 0 to 1. A larger value indicates better predictive accuracy:

$$Recall = \frac{TP}{TP + FN}$$

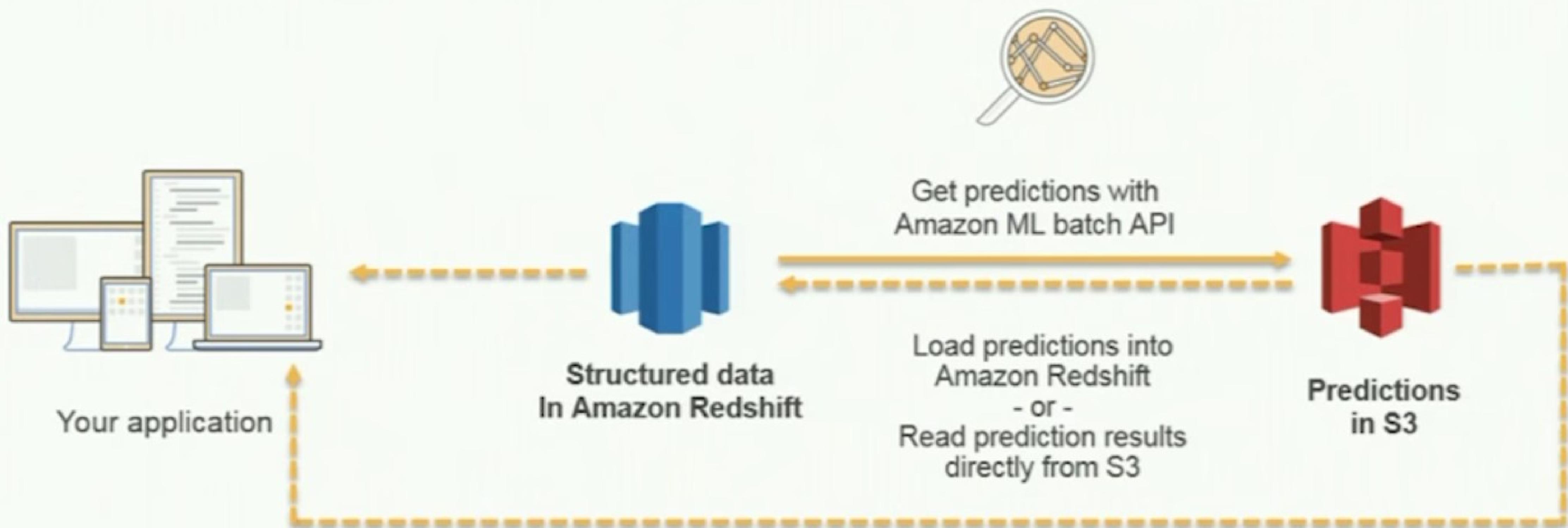
False Positive Rate

The *false positive rate* (FPR) measures the false alarm rate or the fraction of actual negatives that are predicted as positive. The range is 0 to 1. A smaller value indicates better predictive accuracy:

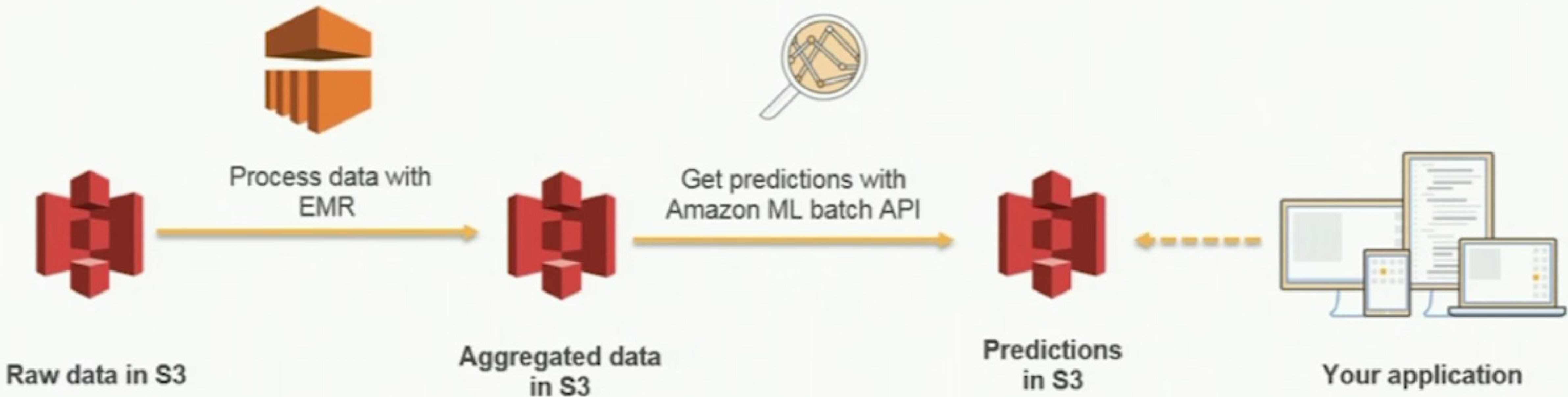
$$FPR = \frac{FP}{FP + TN}$$

		Prediction	
		Positive	Negative
Actual	Positive	TP	FN
	Negative	FP	TN

Batch predictions with Amazon Redshift



Batch predictions with EMR

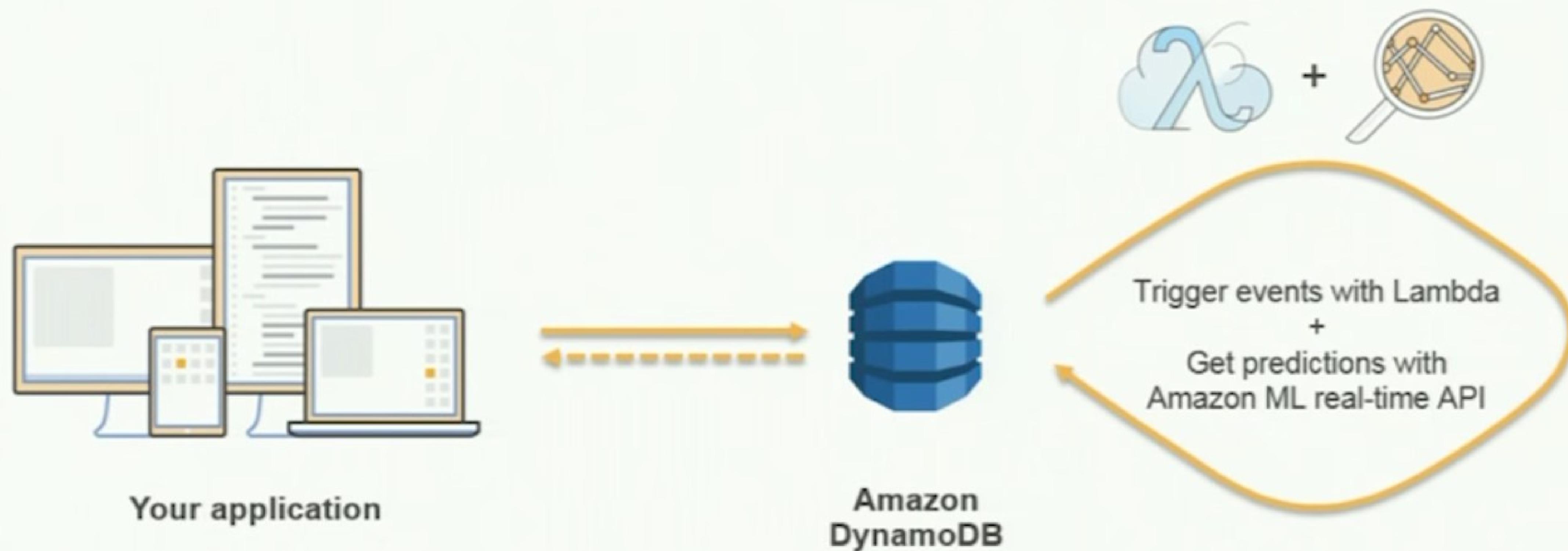


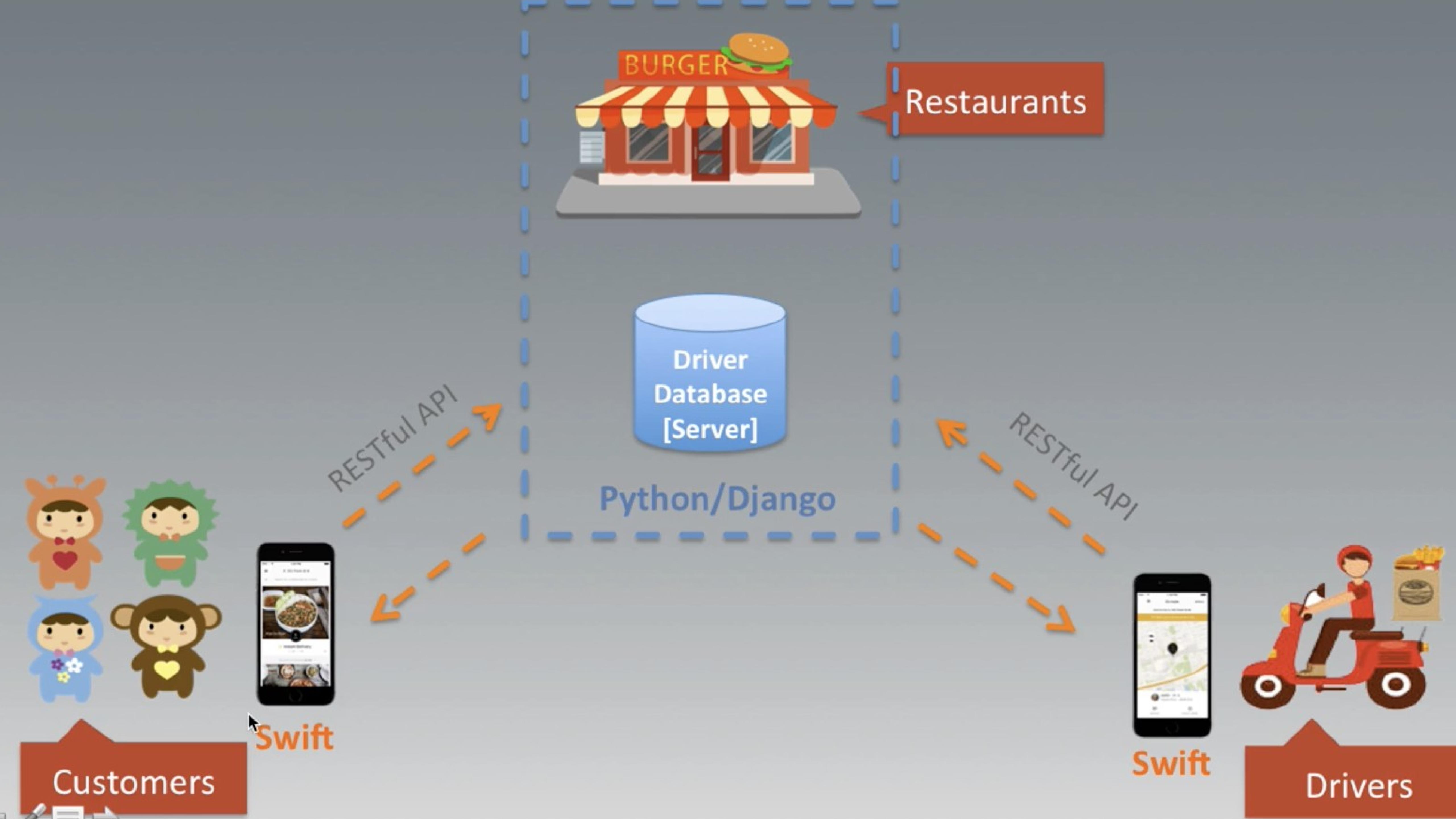
Real-time predictions for interactive applications



CLOUD-BASED ML TOOLS

Adding predictions to an existing data flow





SO, WHAT IF I WANT TO LEARN HOW TO BUILD DATA

Code4Startup

One Month



Udemy

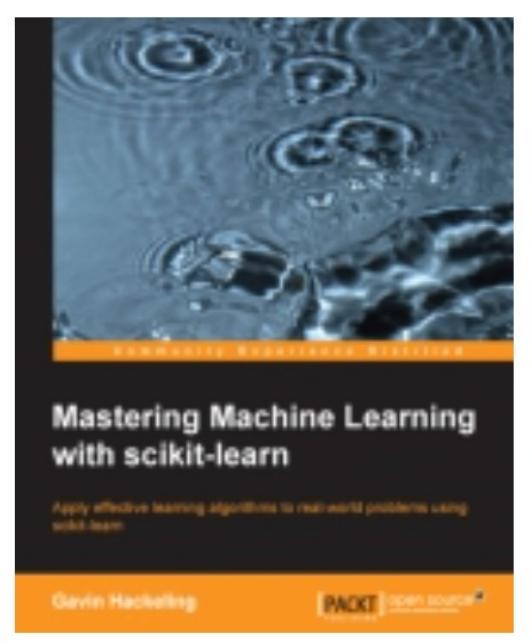
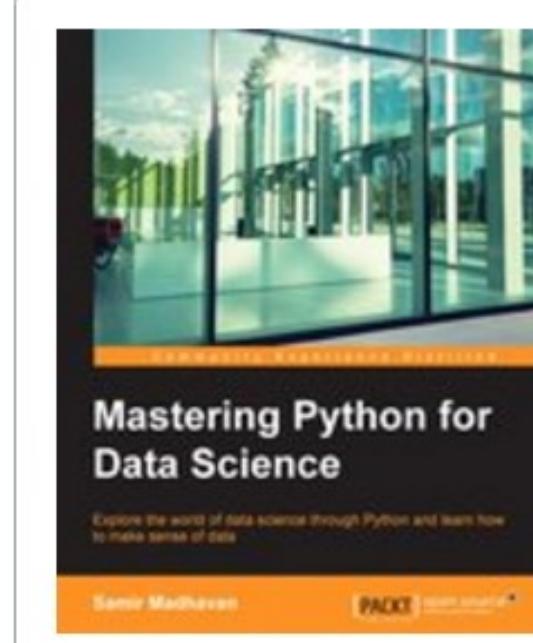
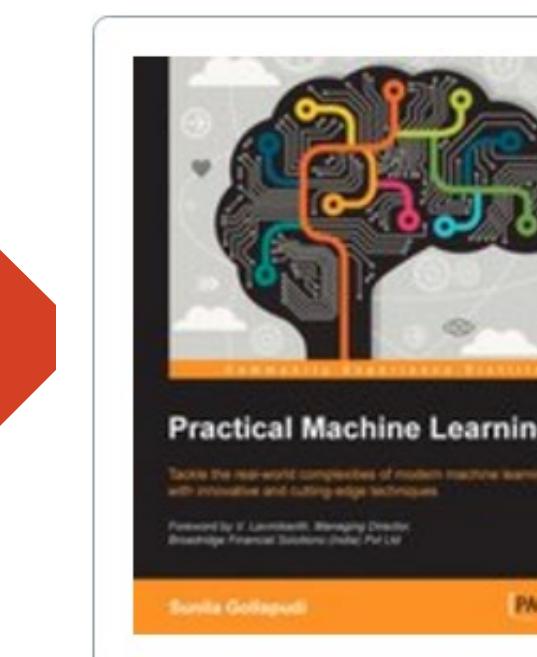
treehouse™



coursera

Code School
a Pluralsight company

codecademy



< P

Ma
By: C
Publ
Pub.
Print
Web
Page

Sub

CHALLENGE YOURSELF WITH REAL-WORLD ML

kaggle

Search kaggle

Competitions

Datasets

Kernels

Discussion

Jobs

Sign Up

Log In

PROBLEMS

Welcome to Kaggle Competitions

Challenge yourself with real-world machine learning problems



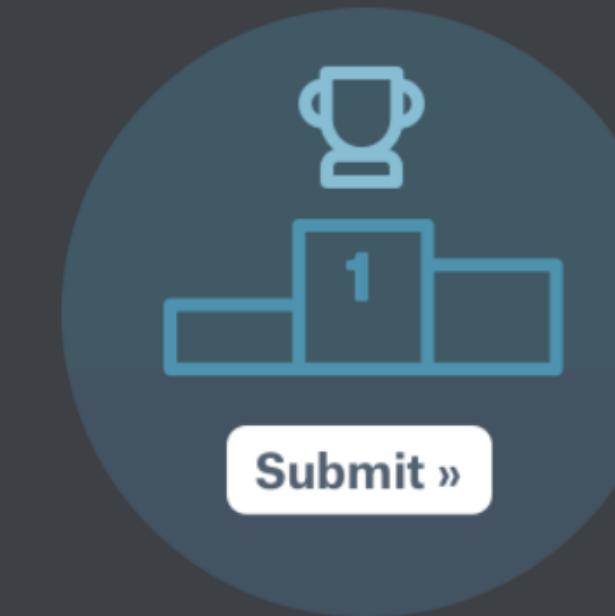
New to Data Science?

Get started with a tutorial on our most popular competition for beginners, [Titanic: Machine Learning from Disaster](#).



Build a Model

Get the data & use whatever tools or methods you prefer to make predictions.



Make a Submission

Upload your prediction file for real-time scoring & a spot on the leaderboard.