

IT 5504 – Applied Machine Learning

Thushari Silva, PhD

Department of Computational Mathematics

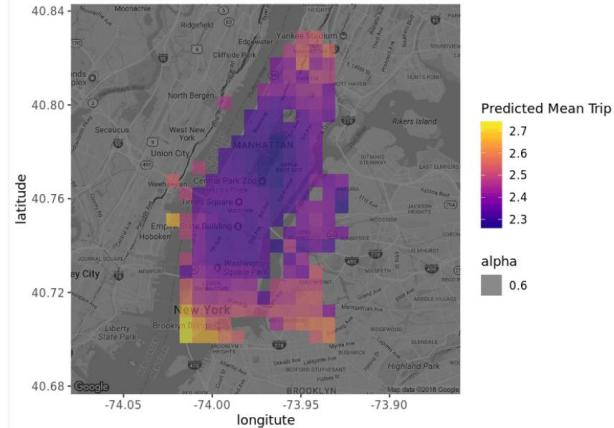
Faculty of Information Technology

University of Moratuwa



Predict Taxi Fares with Random Forests

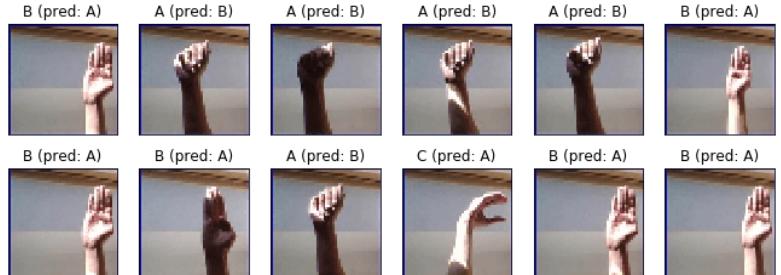
Predict the location and time to earn the biggest fare for taxi drivers – try Decision Tree and Random Forest drivers



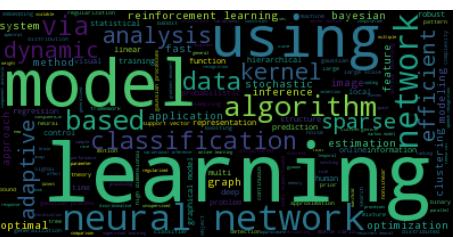
process sound files using Librosa, sound file, and sklearn for the MLPClassifier to recognize emotion from sound files.



Test Samples

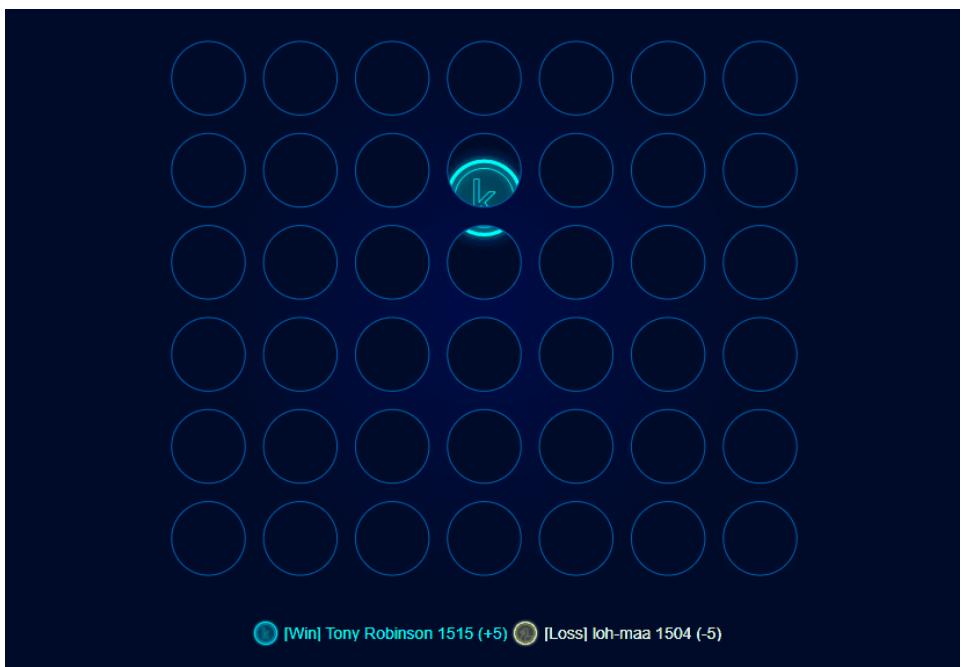


American Sign Language image classification. - CNN (Convolutional Neural Network)

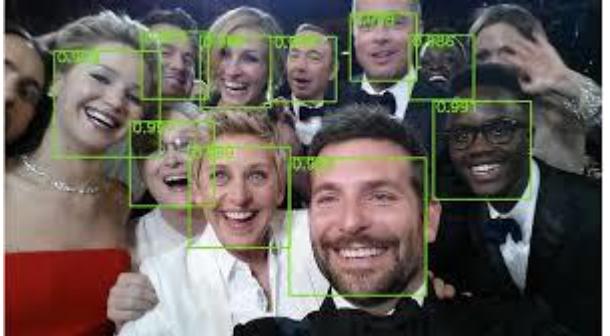
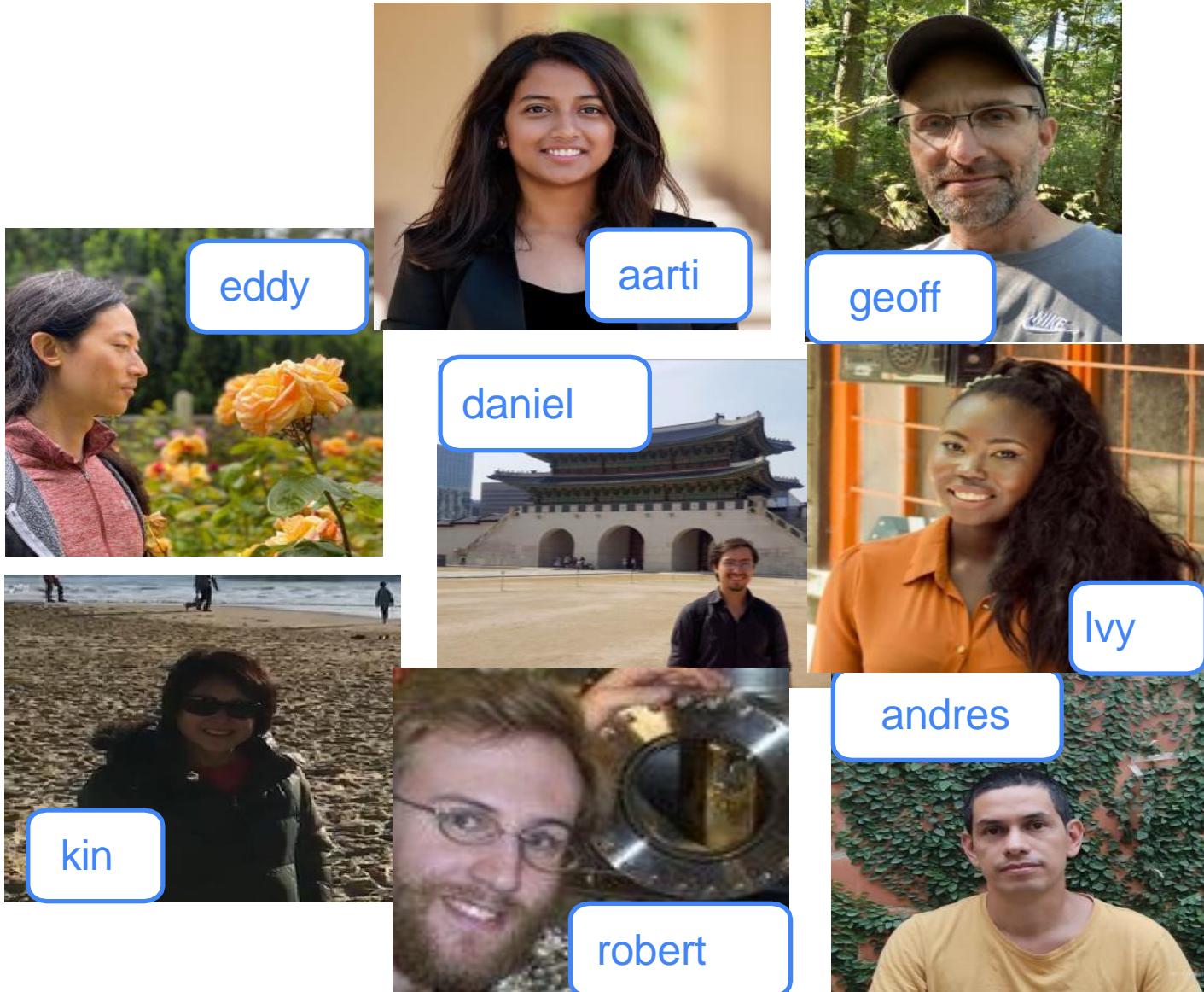


The Hottest Topics in Machine Learning

Use text processing and LDA(Linear Discriminant Analysis) to discover the latest trend in machine learning from the large collection of NIPS research papers.



RL (Reinforcement Learning) agent to compete against other Kaggle competition participants.



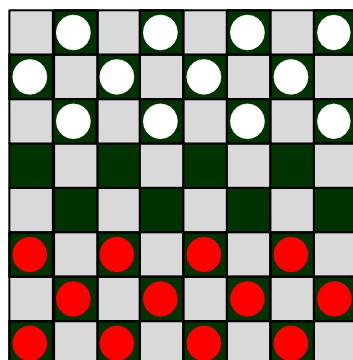
Machine Learning Overview

- What is Machine Learning ?

Machine learning

“Field of study that gives computers the ability to learn without being explicitly programmed.”

Arthur Samuel (1959)

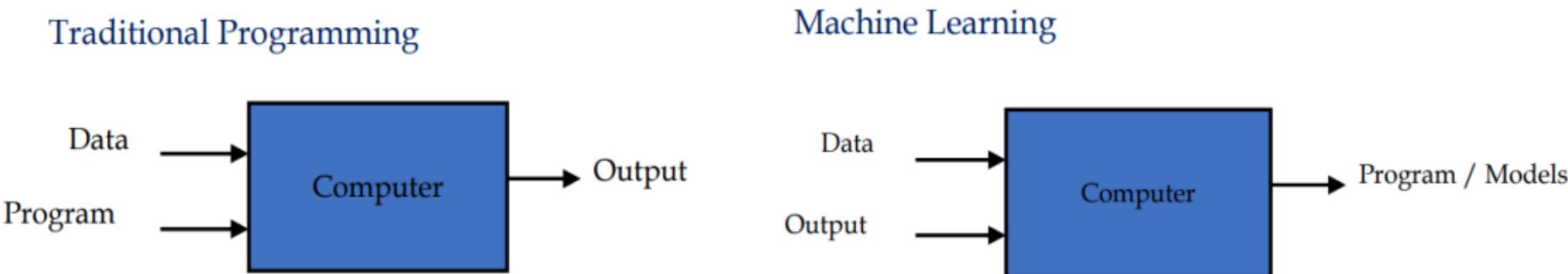


Tom Mitchell (1997): "A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P if its performance at tasks in T , as measured by P , improves with experience E ."

This definition formalizes the idea that learning involves improving performance on a task based on experience.

What is Machine Learning?

- Automating automation
- Getting computers to program themselves
- Writing software is the bottleneck
- Let the data do the work instead



ML defines methods and algorithms that can automatically detect patterns in data and then use these uncovered patterns to predict future data or to perform other kinds of decision-making under uncertainty.

Question

If the checker's program had been allowed to play only ten games (instead of tens of thousands) against itself, a much smaller number of games, how would this have affected its performance?

- Would have made it better
- Would have made it worse

Machine learning algorithms

- **Supervised learning**
 - Classification (e.g., spam detection), Regression (e.g., predicting house prices).
 - **Common Algorithms:** Linear Regression, Logistic Regression, Support Vector Machines (SVM), Decision Trees, Random Forests, Neural Networks.
- **Unsupervised learning**
 - Clustering (e.g., customer segmentation), Association (e.g., market basket analysis), Dimensionality Reduction (e.g., PCA).
 - **Common Algorithms:** K-Means, Hierarchical Clustering, DBSCAN, Principal Component Analysis (PCA), Autoencoders.
- **Recommender systems**
- **Practical advice for applying learning algorithms**

Machine learning algorithms

- **Semi-supervised learning**
 - Semi-supervised learning falls between supervised and unsupervised learning.
 - It uses a small amount of labelled data along with a large amount of unlabeled data to improve learning accuracy.
 - Common Algorithms: Self-Training, Co-Training, Generative Models (e.g., GANs used in a semi-supervised context).
- Recommender systems
- Practical advice for applying learning algorithms

Machine learning algorithms

Reinforcement Learning:

- An agent interacts with an environment and learns to make decisions by receiving feedback in the form of rewards or penalties. The goal is to maximize cumulative rewards over time.
- **Examples:** Game playing (e.g., AlphaGo), Robotics, Autonomous vehicles.
- **Common Algorithms:** Q-Learning, Deep Q-Networks (DQN), Policy Gradients, Proximal Policy Optimization (PPO).

Machine learning algorithms

- **Ensemble Learning:**

- Ensemble learning involves combining multiple models to improve overall performance.
- The idea is that a group of models working together will outperform any individual model.
- **Examples:** Random Forests (an ensemble of decision trees), Gradient Boosting Machines (e.g., XGBoost).
- **Common Techniques:** Bagging, Boosting, Stacking.

- **Online Learning:**

- Online learning is a method where the model is trained incrementally as new data arrives, rather than on a fixed dataset.
- This is useful for applications where data is continuously generated.
- **Examples:** Real-time recommendation systems and adaptive spam filters.
- **Common Algorithms:** Stochastic Gradient Descent (SGD), Perceptron Algorithm.

Supervised Learning

input → output label

Learns from being given “right answers”

Flower Classification



(a)



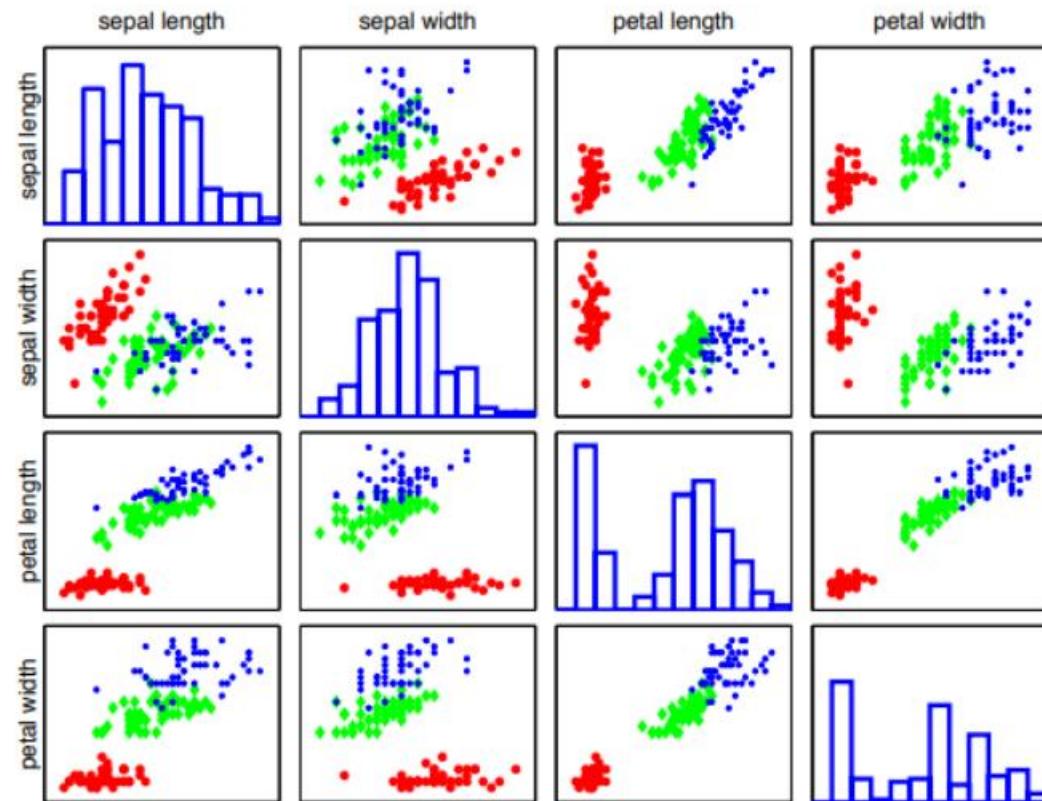
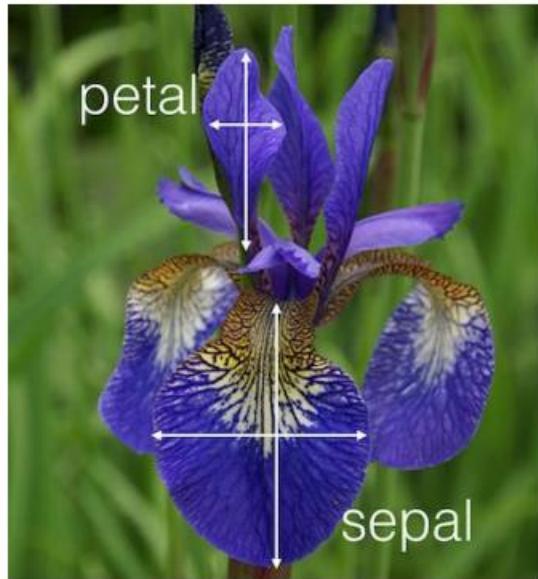
(b)



(c)

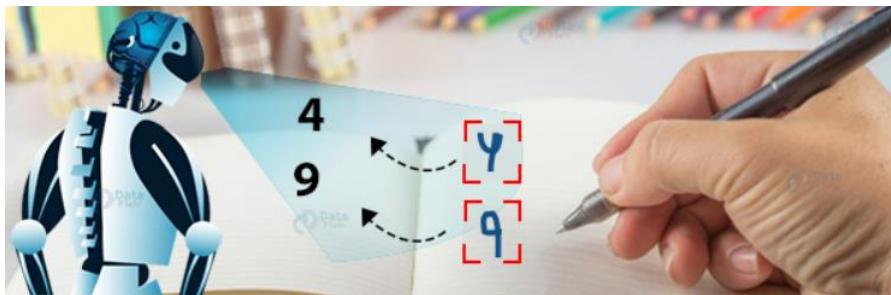
Three types of iris flowers: setosa, versicolor and virginica

Flower Classification



Red circle = setosa, green diamond = versicolor, blue star = virginica.

Image classification and handwriting recognition



true class = 7



true class = 2



true class = 1



true class = 0



true class = 4



true class = 1



true class = 4



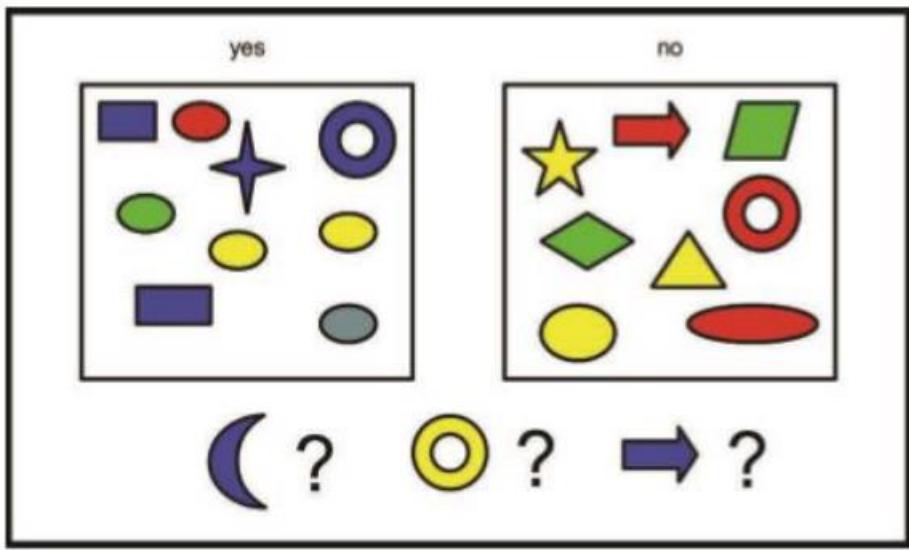
true class = 9



true class = 5



- Supervised learning e.g.



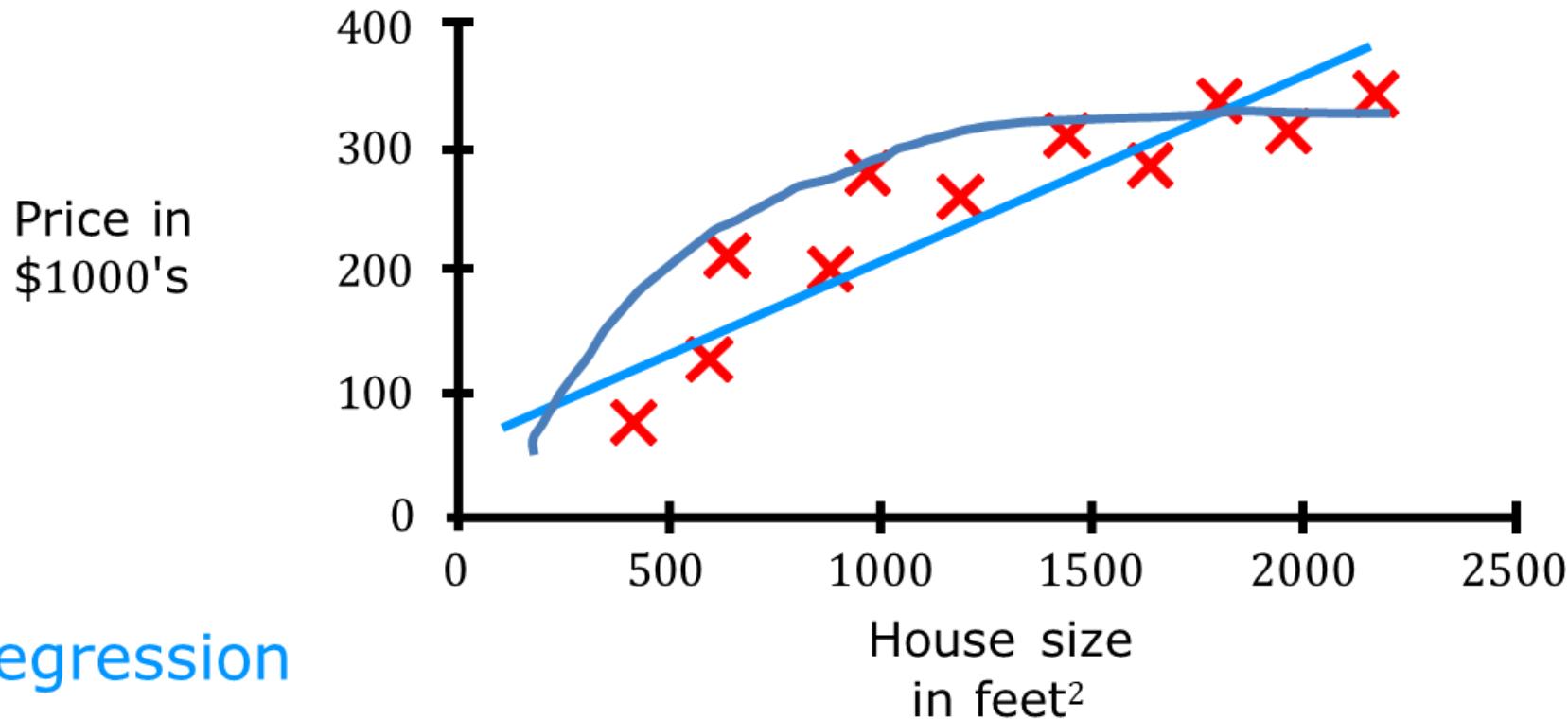
D features (attributes)

N cases

Color	Shape	Size (cm)	Label
Blue	Square	10	1
Red	Ellipse	2.4	1
Red	Ellipse	20.7	0

Input (X)	Output (Y)	Application
email	spam? (0/1)	spam filtering
audio	text transcripts	speech recognition
English	Spanish	machine translation
ad, user info	click? (0/1)	online advertising
image, radar info	position of other cars	self-driving car
image of phone	defect? (0/1)	visual inspection

Regression: Housing price prediction

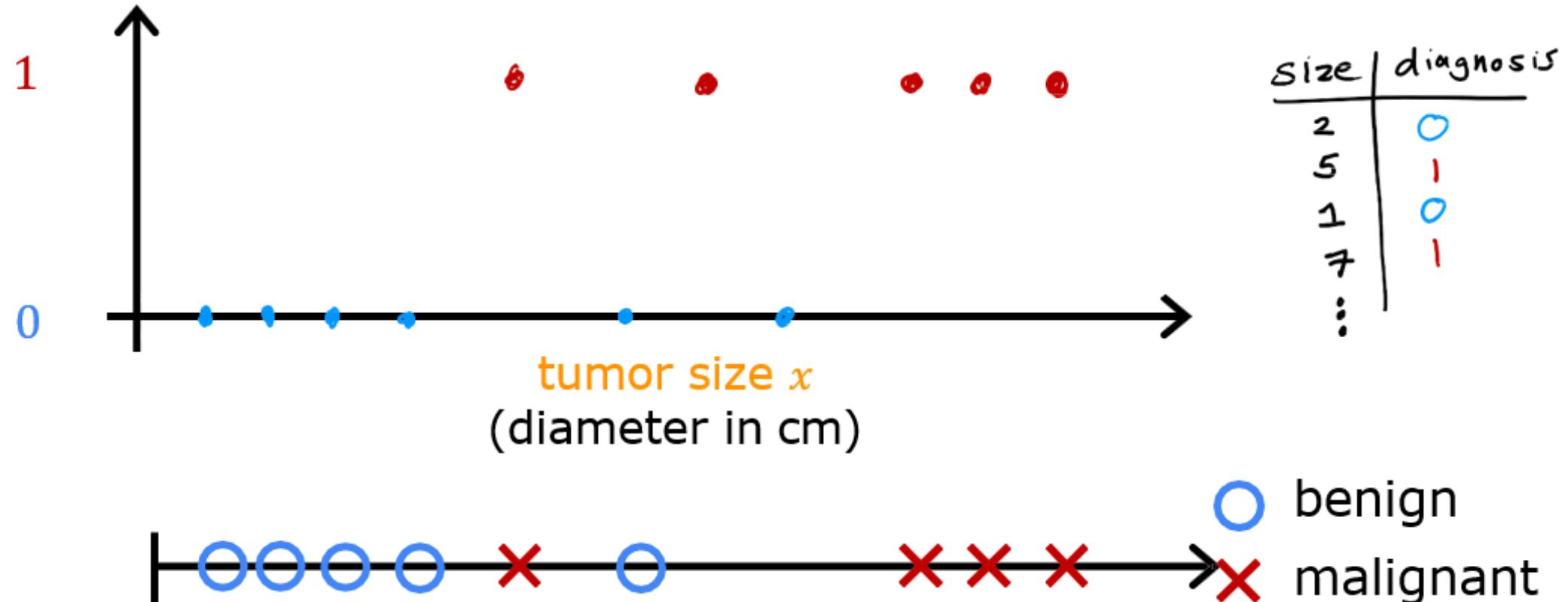


Regression
Predict a number
infinitely many possible outputs

Classification: Breast cancer detection



malignant benign



Classification: Breast cancer detection

○ benign

✗ malignant type 1

△ malignant type 2

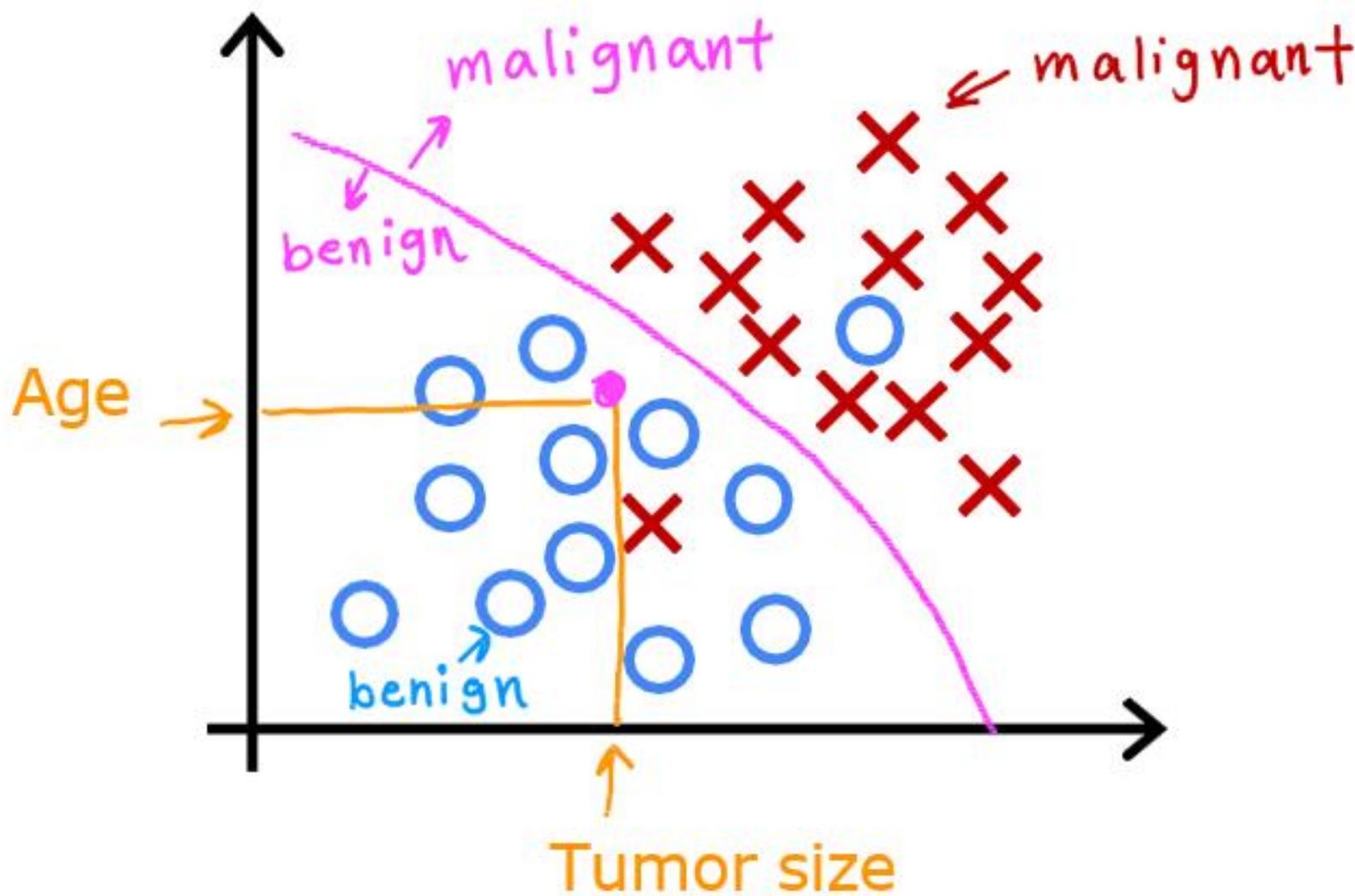


Classification

predict categories

small number of possible outputs

Two or more inputs



Supervised learning

Learns from being given “right answers”

Regression

Predict a number

infinitely many possible outputs

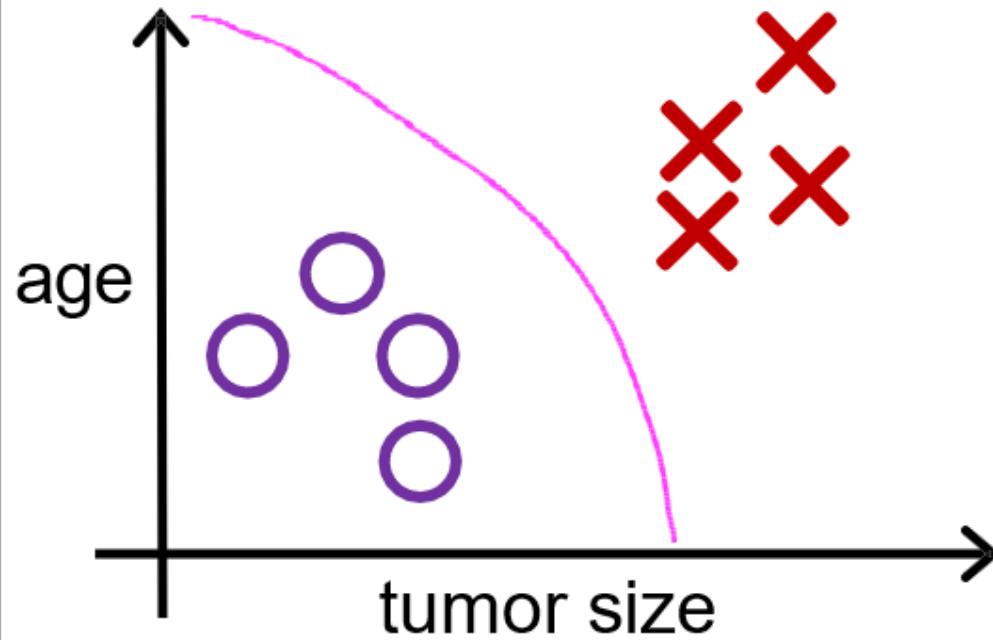
Classification

predict categories

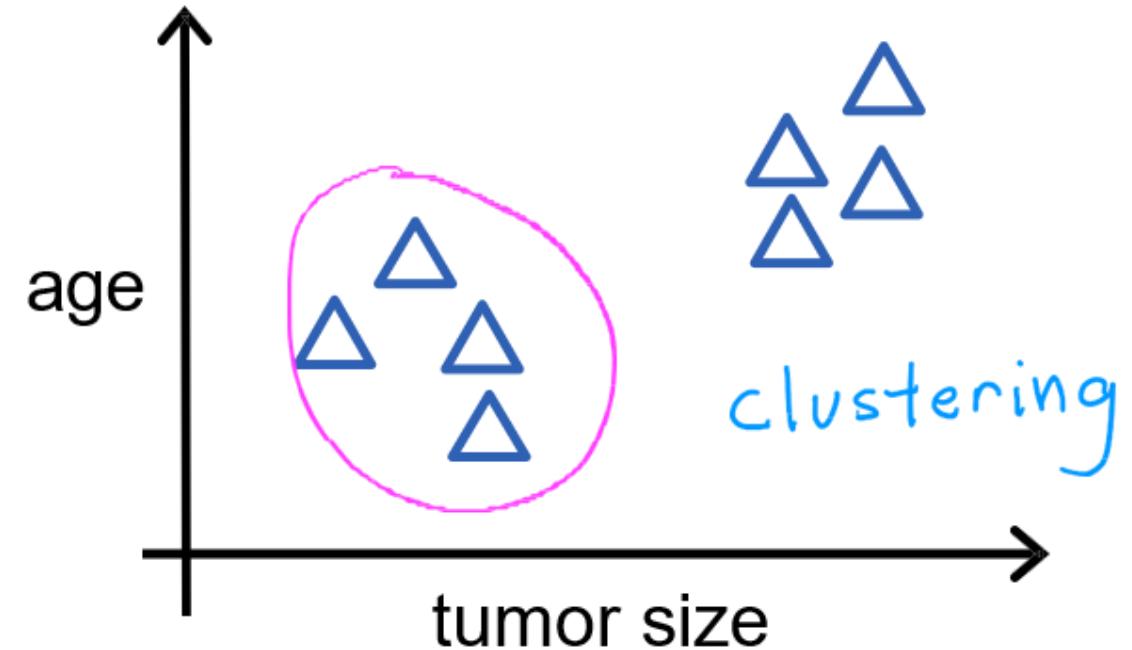
small number of possible outputs

Unsupervised Learning

Supervised learning
Learn from data **labeled**
with the “**right answers**”



Unsupervised learning
Find something interesting
in **unlabeled** data.



Clustering: Google news

Giant panda gives birth to rare twin cubs at Japan's oldest zoo

USA TODAY · 6 hours ago

- **Giant panda gives birth to twin cubs at Japan's oldest zoo**

CBS News · 7 hours ago

- **Giant panda gives birth to twin cubs at Tokyo's Ueno Zoo**

WHBL News · 16 hours ago

- **A Joyful Surprise at Japan's Oldest Zoo: The Birth of Twin Pandas**

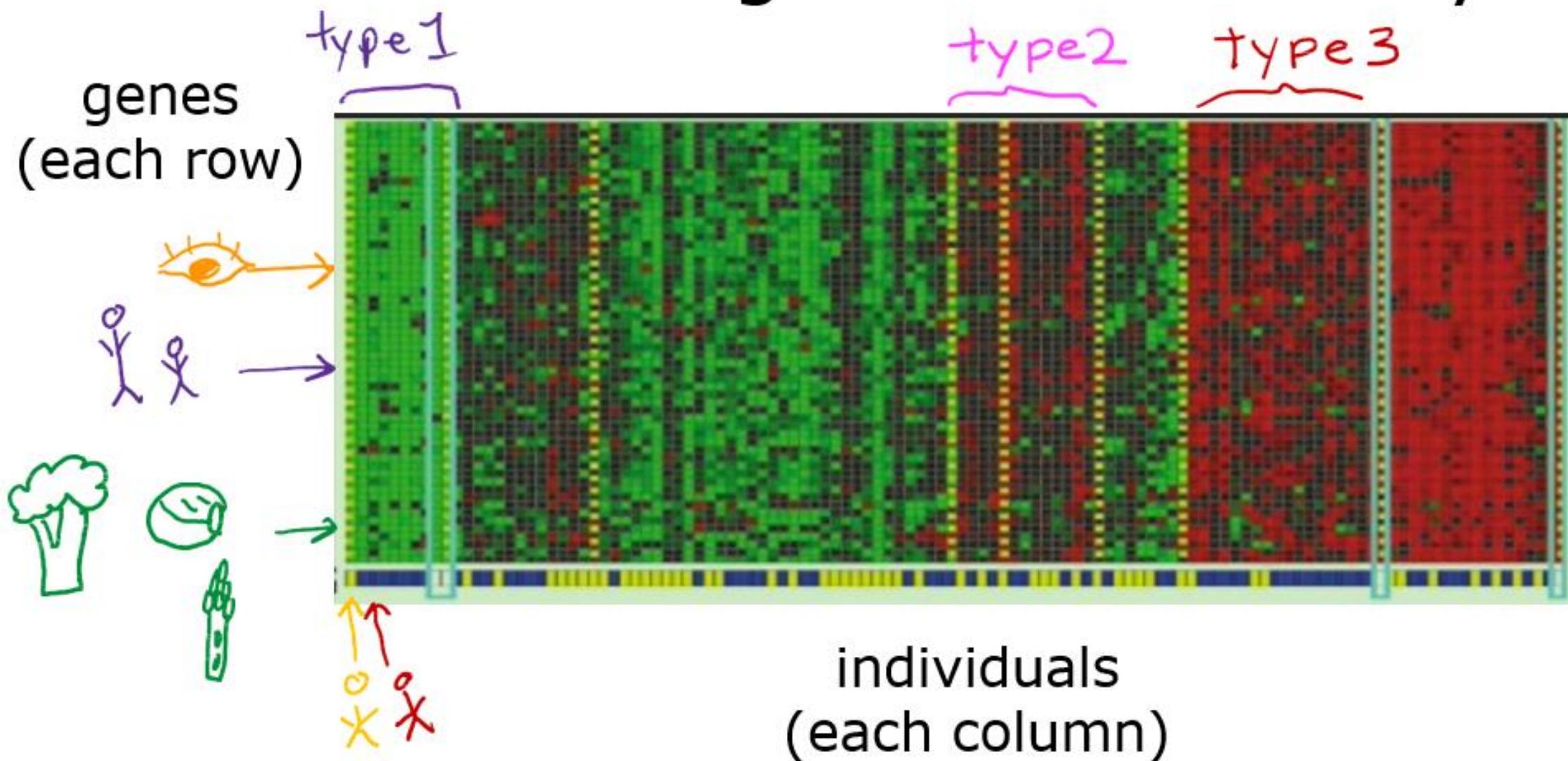
The New York Times · 1 hour ago

- **Twin Panda Cubs Born at Tokyo's Ueno Zoo**

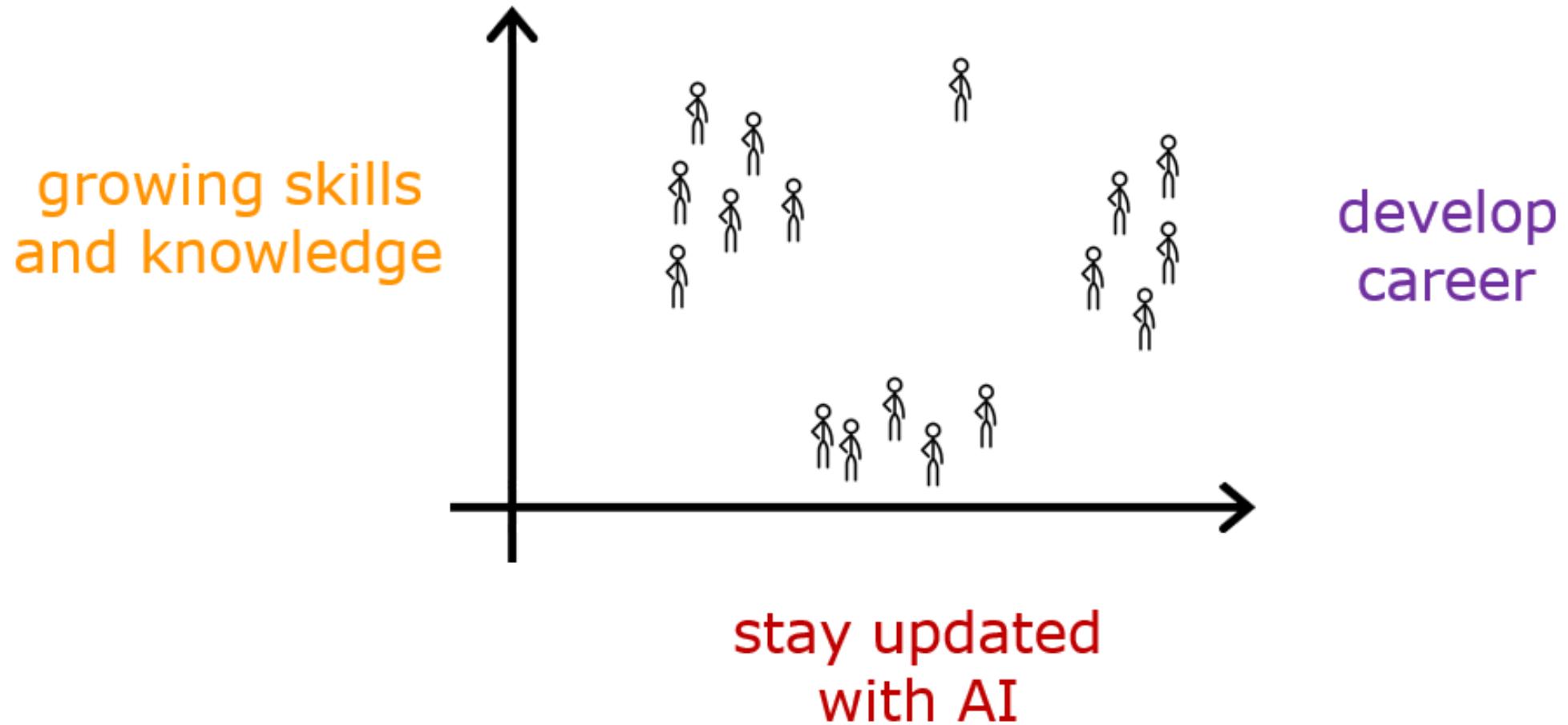
PEOPLE · 6 hours ago

 View Full Coverage

Clustering: DNA microarray



Clustering: Grouping customers



Unsupervised learning

Data only comes with inputs x , but not output labels y .
Algorithm has to find **structure** in the data.

Clustering

Group similar data points together.

Dimensionality reduction

Compress data using fewer numbers.

Anomaly detection

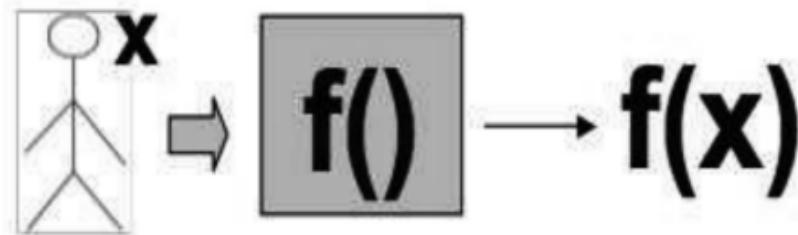
Find unusual data points.

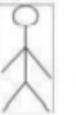
Question

Of the following examples, which would you address using an **unsupervised** learning algorithm?

- Given email labeled as spam/not spam, learn a spam filter.
- Given a set of news articles found on the web, group them into sets of articles about the same story.
- Given a database of customer data, automatically discover market segments and group customers into different market segments.
- Given a dataset of patients diagnosed as either having diabetes or not, learn to classify new patients as having diabetes or not

Representing Data



- How do we represent  mathematically?
- Depends on what we're trying to do:
 - deciding to loan money?
 - predicting gender?
- Represent  as a set of attribute-value pairs
 - example: $x = \{\text{height}=\text{180cm}, \text{eyes}=\text{"blue"}, \text{job}=\text{"student"}\}$

Attribute-value pairs

- $x = \{\text{height}=\text{180cm}, \text{eyes}=\text{"blue"}, \text{job}=\text{"student"}\}$
- un-ordered “bag-of-features”
 - if structure is essential – embed it in the attributes
- Have to convert any dataset to this form
- Generally three types of attributes:
 - categorical: *red, blue, brown, yellow*
 - ordinal: *poor, satisfactory, good, excellent*
 - numeric: *-3.14, 6E23, 0, 1*

Categorical Attributes

- Each instance falls into one of a set of categories
 - **genre**: *{classical, jazz, rock, techno}*
 - categories are mutually exclusive
- Categories usually encoded as numbers
 - no natural ordering to categories
 - only equality testing ($=, \neq$) is meaningful
- Synonymy a major challenge for real datasets:
 - e.g. social tags: *country == folk?* *house == techno?*

Ordinal Attribute

- Instance falls into one of a set of categories
- There is a natural ordering to categories
 - **education level:** $\{none, school, university, post-graduate\}$
 - **Likert scale:** $\{disagree, neutral, agree, strongly\ agree\}$
- Encoded as numbers to preserve ordering
 - meaningful to compare values: $(<, =, >)$
 - should not add / multiply / measure “distance”
- Sometimes hard to differentiate from categorical:
 - does $\{single, married, divorced\}$ have a natural ordering?

Numeric Attribute

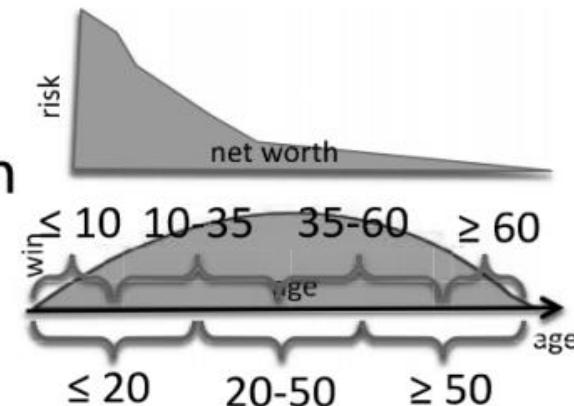
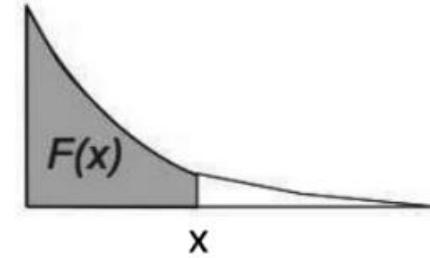
- Integers or real numbers
 - meaningful to add, multiply, compute mean / variance
 - integers not always the same as real numbers
- Usually want to normalize values (why?)
 - zero mean, unit variance: $x' = (x - \text{mean}) / \text{st.dev}$
 - sometimes want [0,1]: $x' = (x - \text{min}) / (\text{max} - \text{min})$
- Sensitive to extreme (unusually large/small) values
 - e.g. height: {165,167,171,175,176,181,183,1820}[cm]



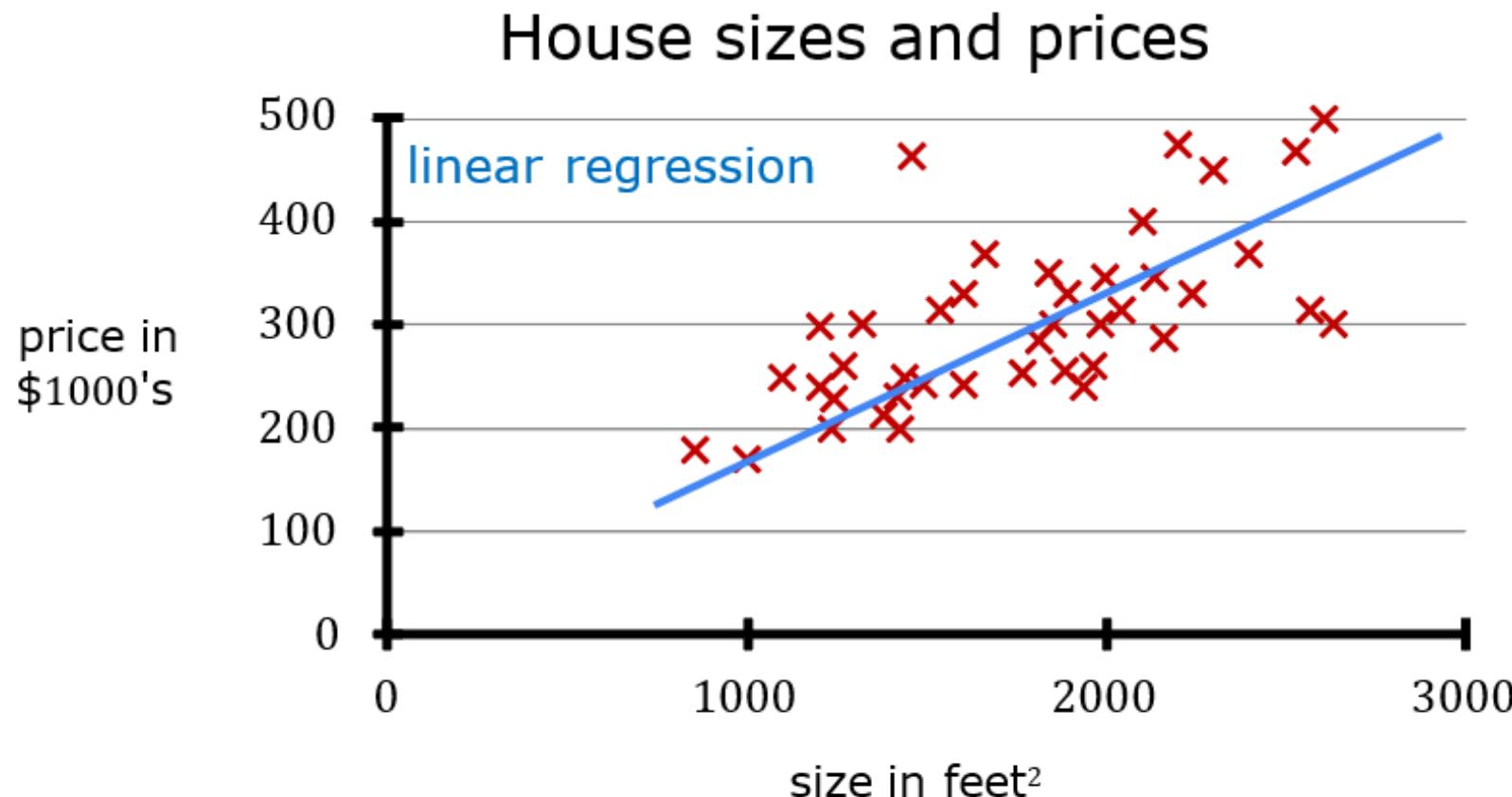
- must handle this before normalization

Numeric Attribute Issues

- Skewed distributions
 - systematic extreme values
 - affects regression, kNN, NB; but not DTs
 - simple fix: $\log(x)$ or $\text{atan}(x)$, then normalize
 - cumulative distribution function: $x' = F(x) = P(X \leq x)$
- Non-monotonic effect of attributes
 - affects regression, NB, DTs(gain); less important for kNN
 - monotonic: net worth and lending risk
 - higher net worth \rightarrow lower lending risk
 - non-monotonic: age \rightarrow win a marathon
 - sweet spot: not too young, not too old
 - simple fix: quantization
 - can be unsupervised, overlapping



Linear regression model with one variable

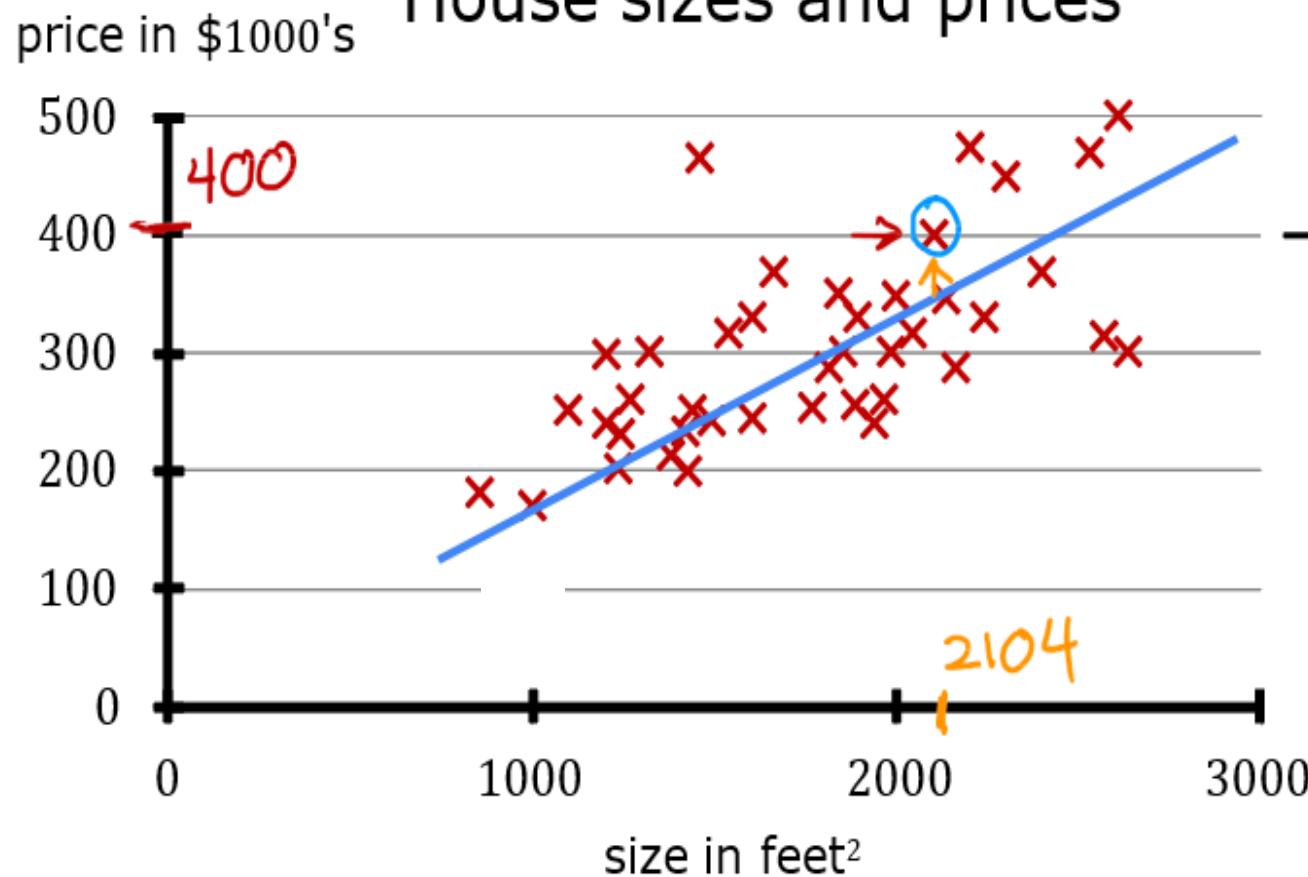


Regression model
Predicts numbers
Infinitely many possible outputs

Supervised learning model
Data has "right answers"

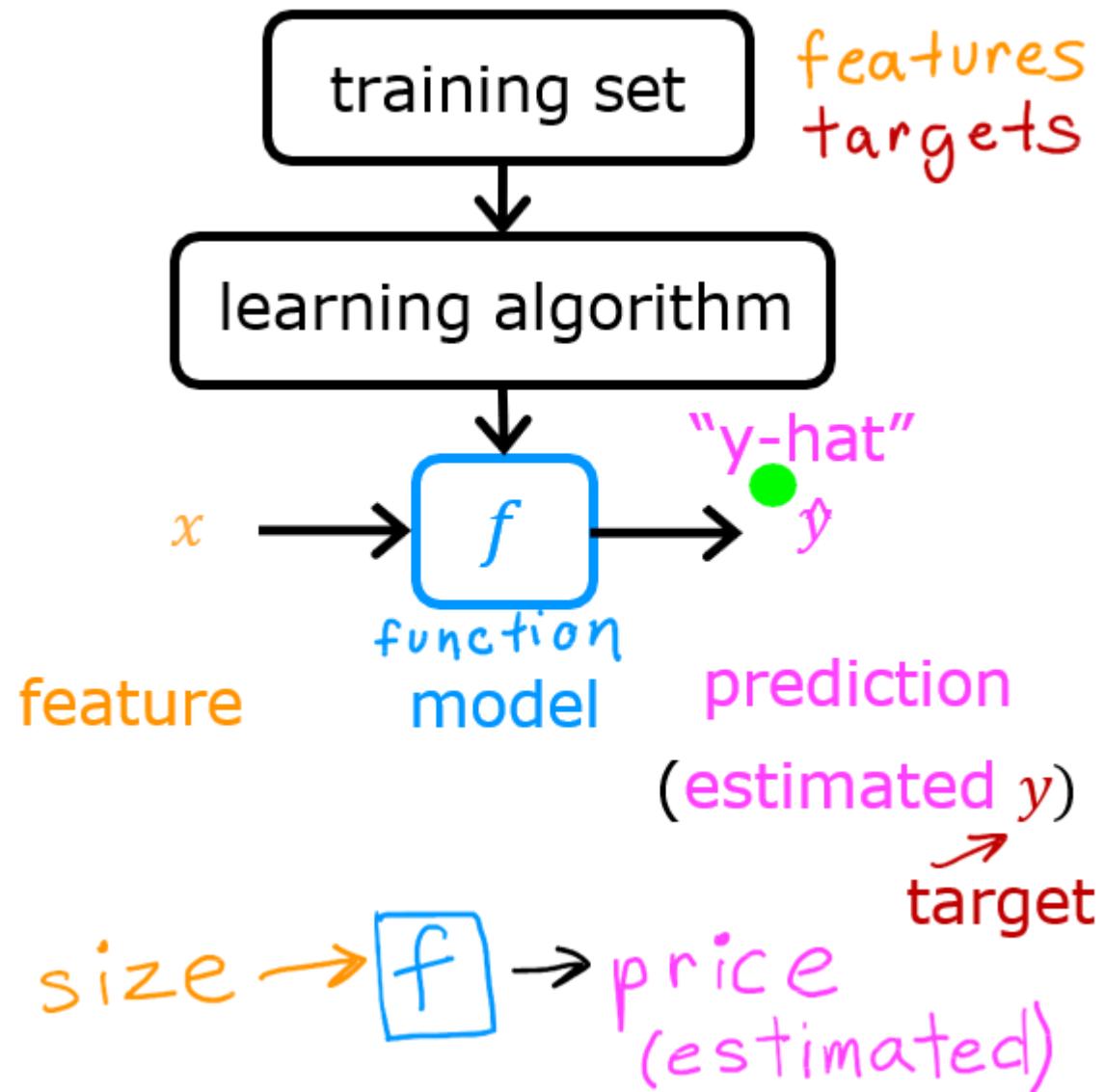
Classification model
Predicts categories
Small number of possible outputs

House sizes and prices



Data table

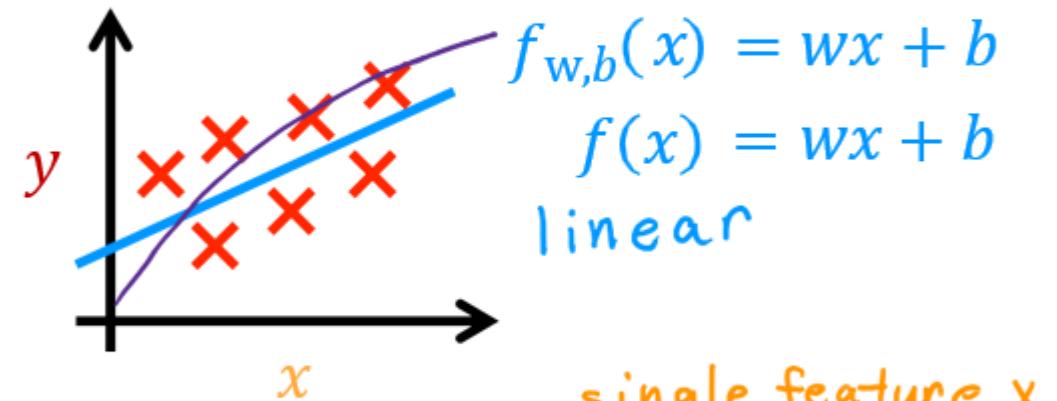
size in feet ²	price in \$1000's
2104	400
1416	232
1534	315
852	178
...	...
3210	870



How to represent f ?

$$f_{w,b}(x) = wx + b$$

$$f(x)$$



Linear regression with one variable.
size
Univariate linear regression.
one variable

Linear Regression with One Variable

Cost Function

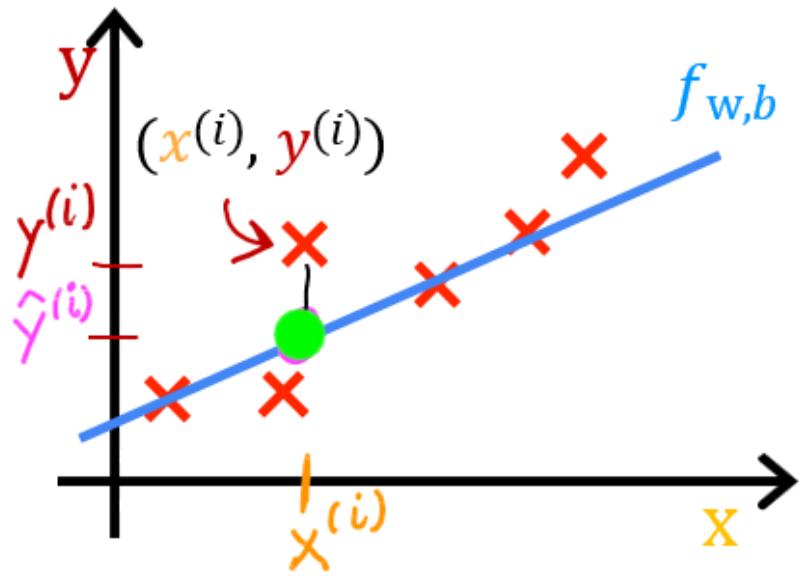
Training set

features	targets
size in feet ² (x)	price \$1000's (y)
2104	460
1416	232
1534	315
852	178
...	...

$$\text{Model: } f_{w,b}(x) = wx + b$$

w, b : parameters
coefficients
weights

What do w, b do?



$$\hat{y}^{(i)} = f_{w,b}(x^{(i)}) \leftarrow$$

$$f_{w,b}(x^{(i)}) = w x^{(i)} + b$$

Cost function: Squared error cost function

$$\bar{J}(w, b) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2$$

m = number of training examples

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2$$

intuition (next!)

Find w, b :

$\hat{y}^{(i)}$ is close to $y^{(i)}$ for all $(x^{(i)}, y^{(i)})$.

model:

$$f_{w,b}(x) = wx + b$$

parameters:

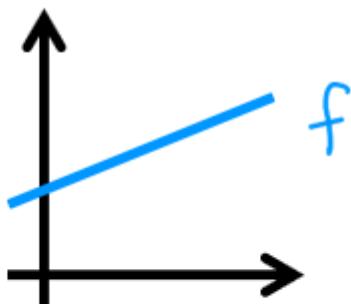
$$w, b$$

cost function:

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2$$

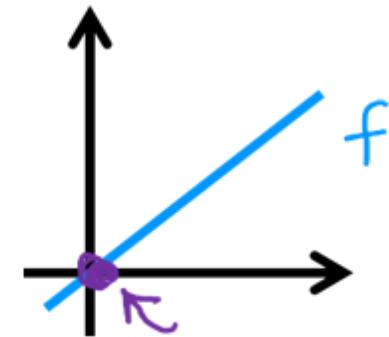
goal:

$$\underset{w,b}{\text{minimize}} J(w, b)$$



simplified

$$f_w(x) = \underbrace{wx}_{w} + \underbrace{b}_{\emptyset}$$



$$J(w) = \frac{1}{2m} \sum_{i=1}^m (\underbrace{f_w(x^{(i)})}_{wx^{(i)}} - y^{(i)})^2$$

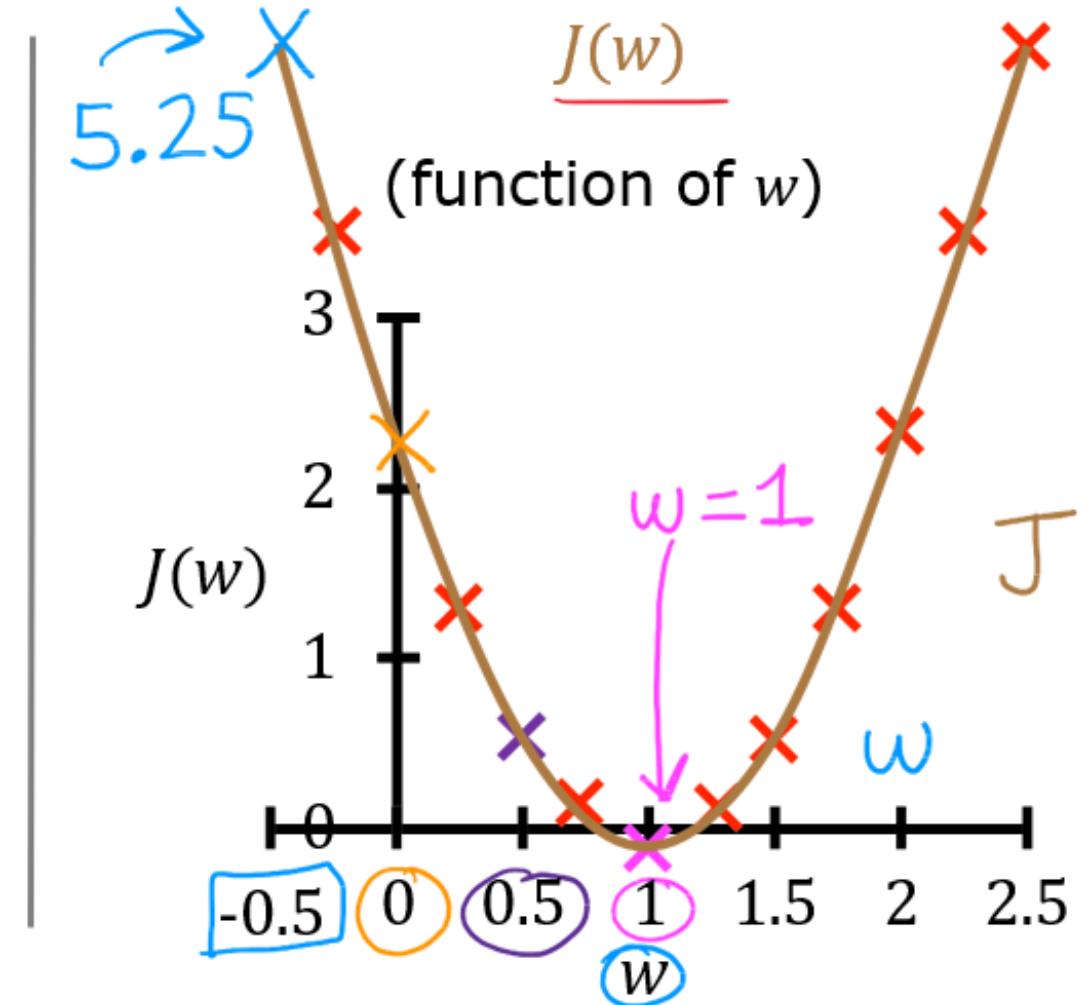
$$\underset{w}{\text{minimize}} J(w)$$

goal of linear regression:

$$\underset{w}{\text{minimize}} J(w)$$

general case:

$$\underset{w,b}{\text{minimize}} J(w, b)$$



choose w to minimize $J(w)$

Model

$$f_{w,b}(x) = wx + b$$

Parameters

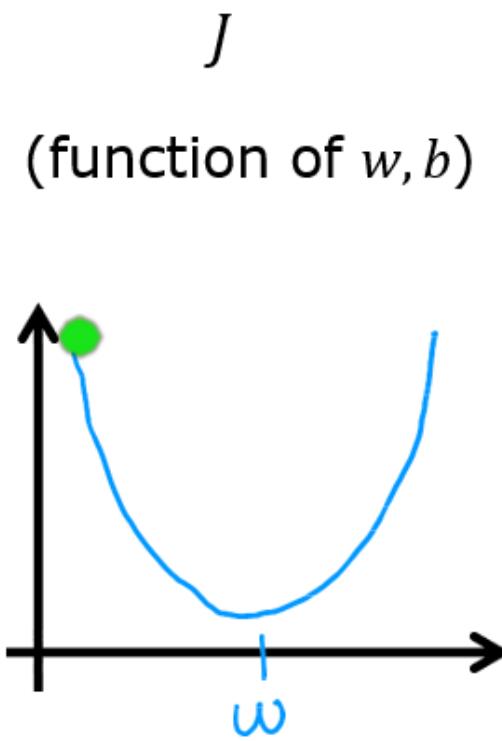
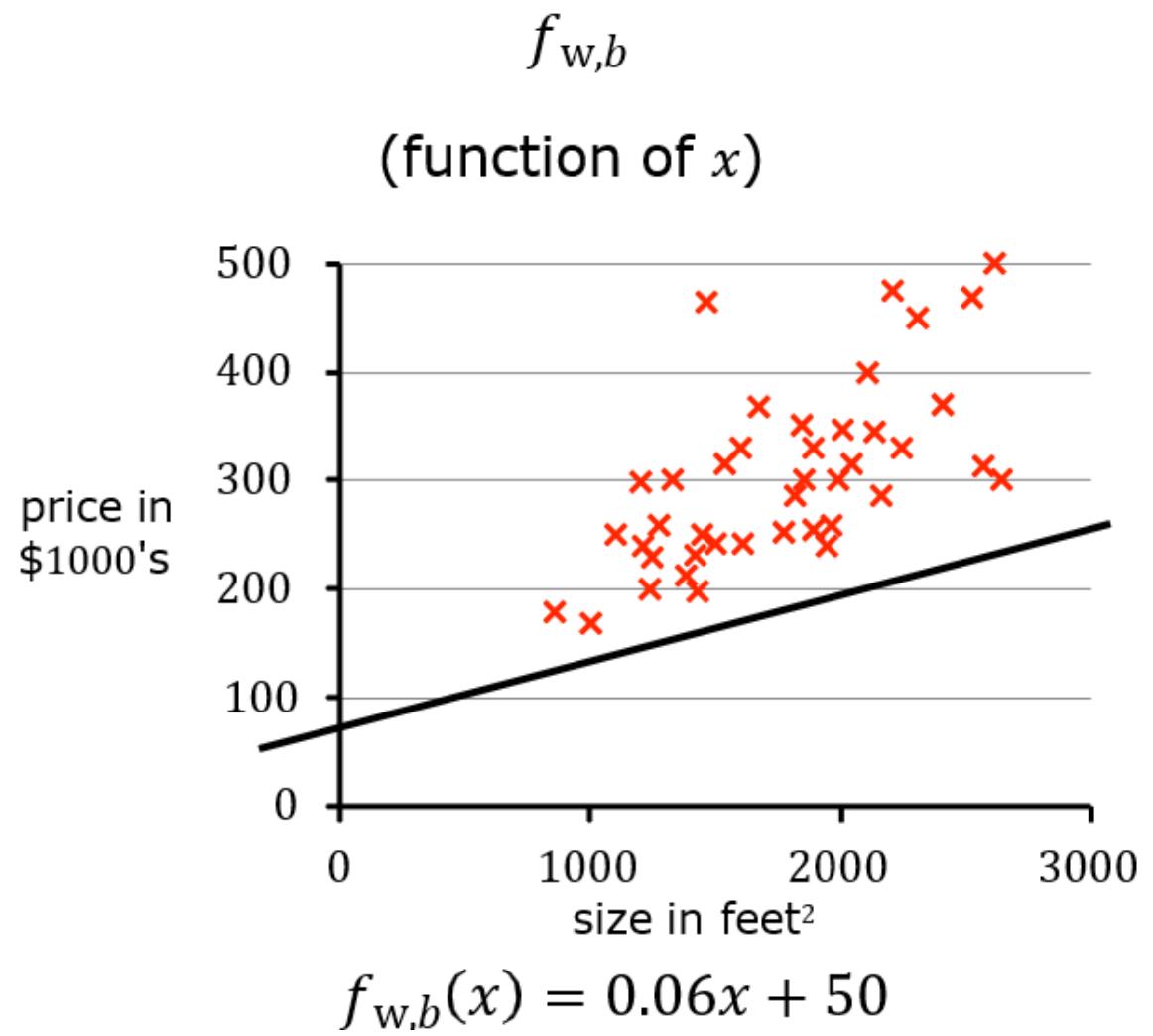
$$w, b$$

Cost Function

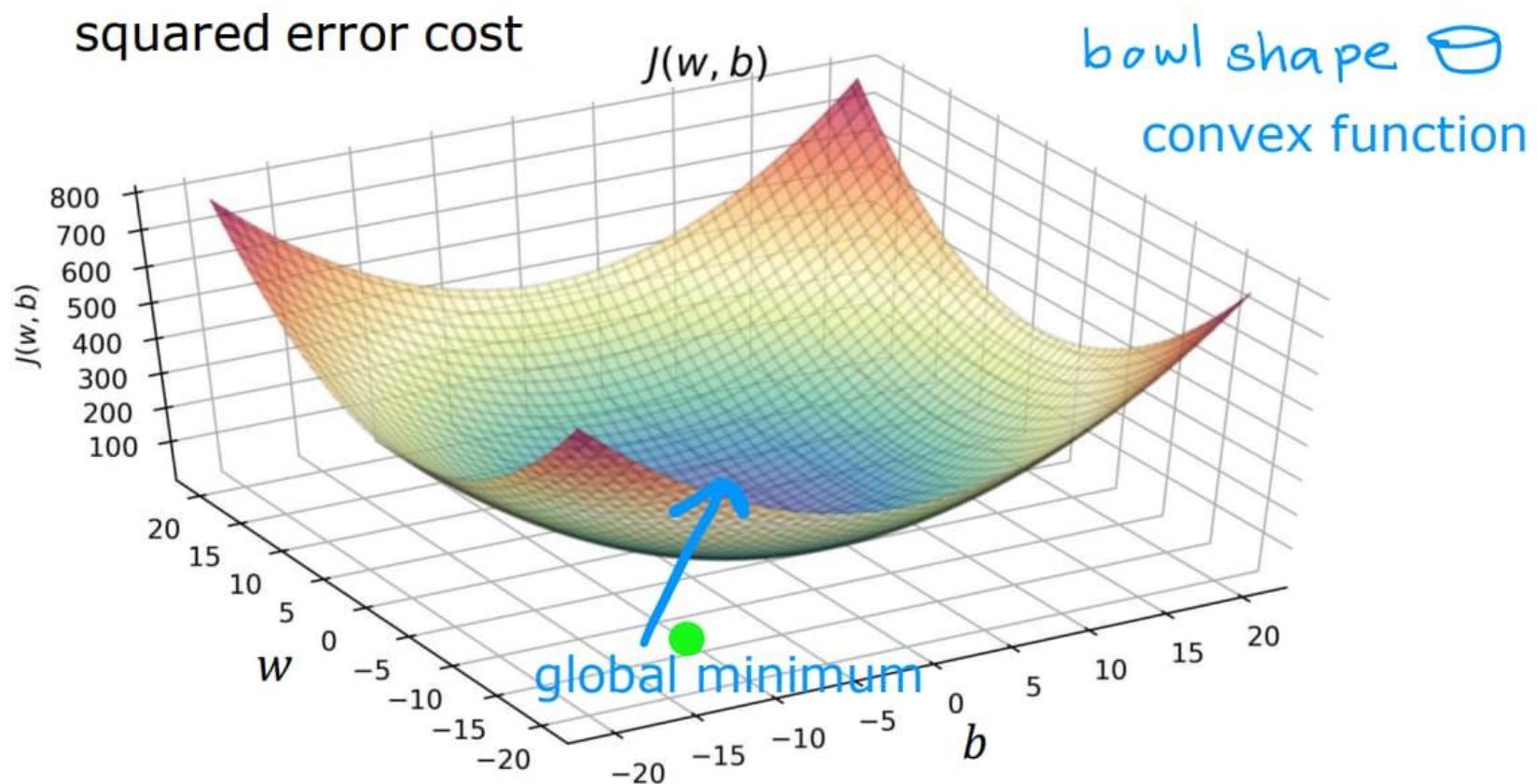
$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2$$

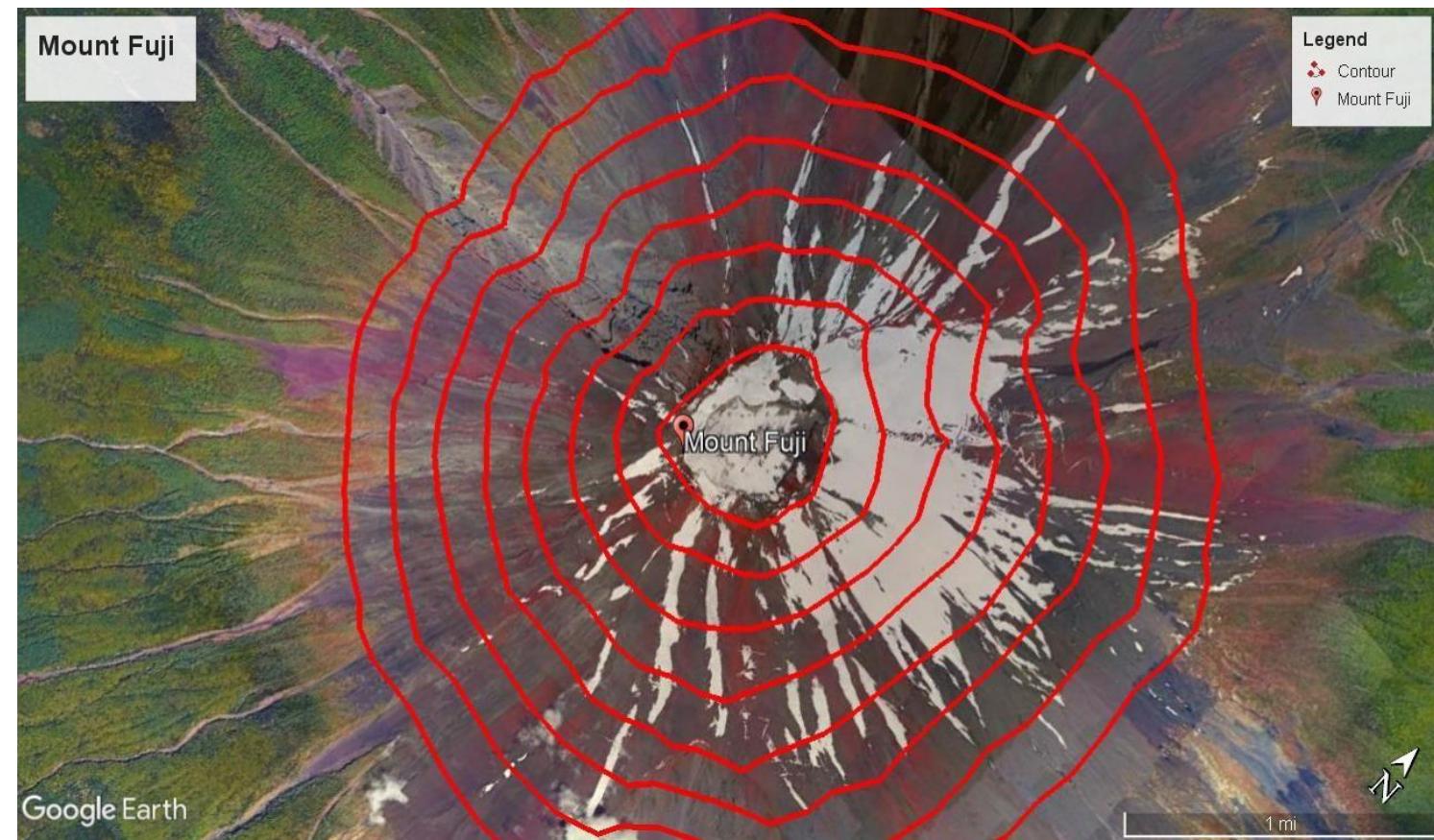
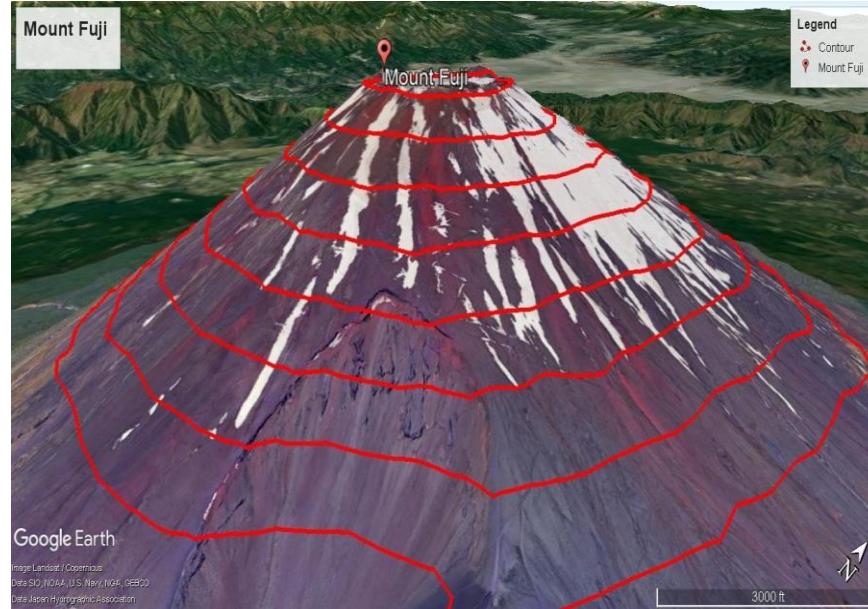
Objective

$$\underset{w,b}{\text{minimize}} J(w, b)$$



Local/Global minimum in cost Function





Have some function $J(w, b)$

for linear regression
or any function

Want $\min_{w, b} J(w, b)$

$\min_{w_1, \dots, w_n, b} J(w_1, w_2, \dots, w_n, b)$

Outline:

Start with some w, b (set $w=0, b=0$)

Keep changing w, b to reduce $J(w, b)$

Until we settle at or near a minimum

may have >1 minimum



Gradient descent algorithm

Repeat until convergence

$$w = w - \alpha \frac{\partial}{\partial w} J(w, b)$$

$$b = b - \alpha \frac{\partial}{\partial b} J(w, b)$$

Learning rate
Derivative

Simultaneously
update w and b

Assignment

$$a = c$$

\swarrow

$$a = a + 1$$

Code

Truth assertio

$$a = c$$

\times

$$a = a + 1$$

Math
 $a == c$

Correct: Simultaneous update

$$tmp_w = w - \alpha \frac{\partial}{\partial w} J(w, b)$$

$$tmp_b = b - \alpha \frac{\partial}{\partial b} J(w, b)$$

$$w = tmp_w$$

$$b = tmp_b$$

Incorrect

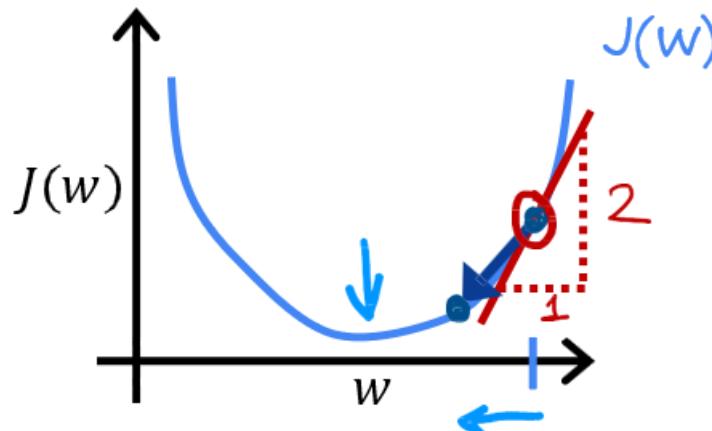
$$tmp_w = w - \alpha \frac{\partial}{\partial w} J(w, b)$$

$$w = tmp_w$$

$$tmp_b = b - \alpha \frac{\partial}{\partial b} J(w, b)$$

$$b = tmp_b$$

Gradient Descent Intuition

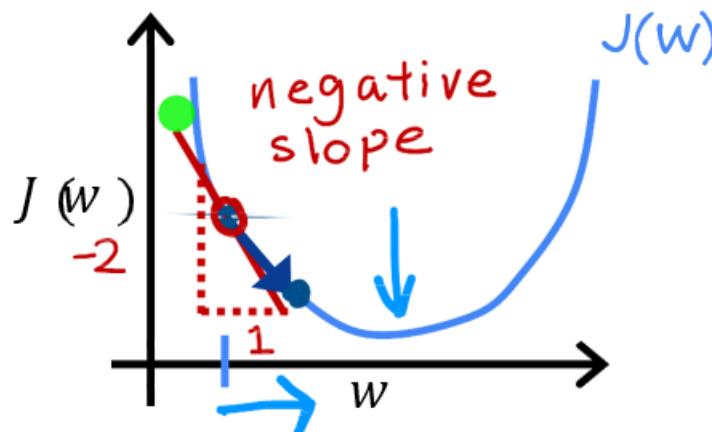


$$w = w - \alpha \frac{\frac{d}{dw} J(w)}{> 0}$$

$w = w - \underline{\alpha} \cdot (\text{positive number})$

$$\frac{d}{dw} J(w) < 0$$

$w = w - \alpha \cdot (\text{negative number})$



Learning Rate

$$w = w - \alpha \frac{d}{dw} J(w)$$

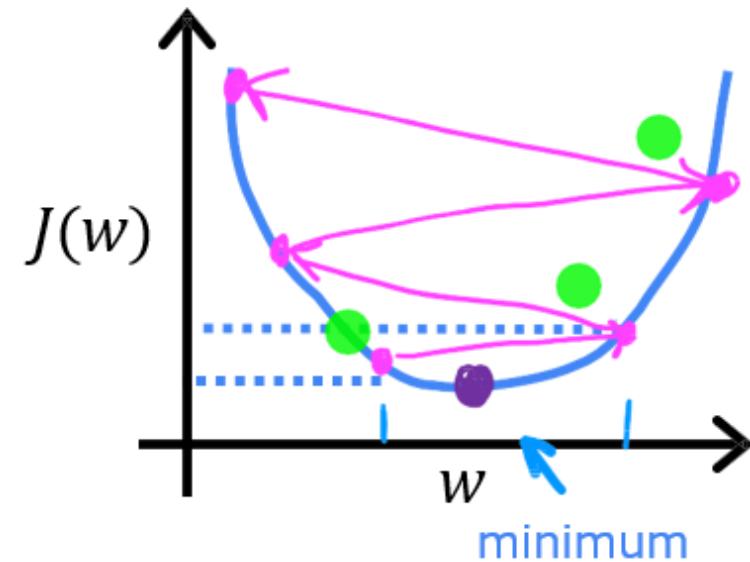
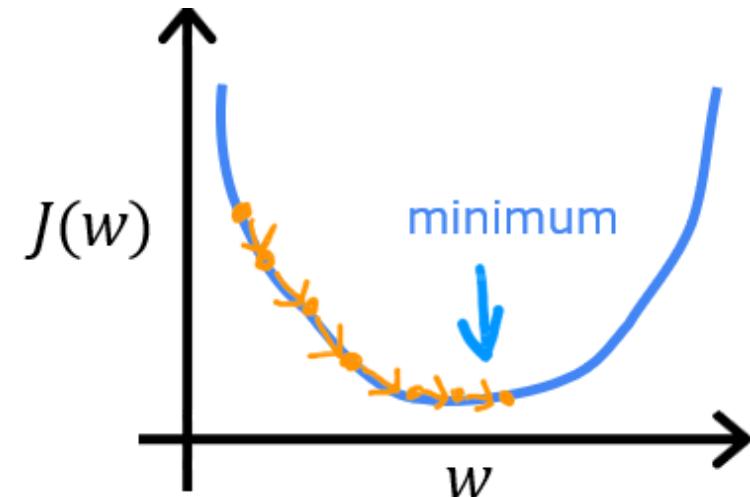
If α is too small...

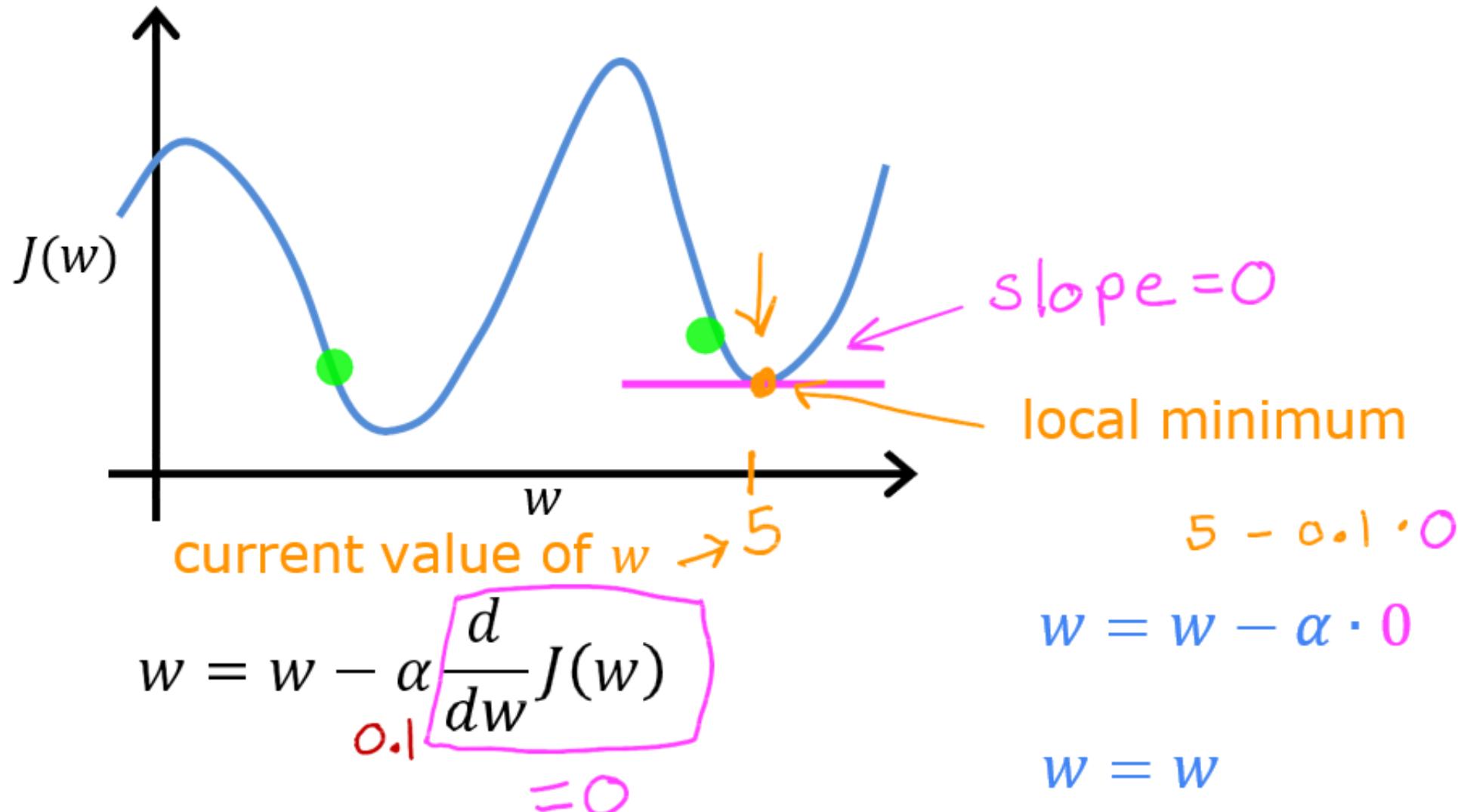
Gradient descent may be slow.

If α is too large...

Gradient descent may:

- Overshoot, never reach minimum
- Fail to converge, diverge





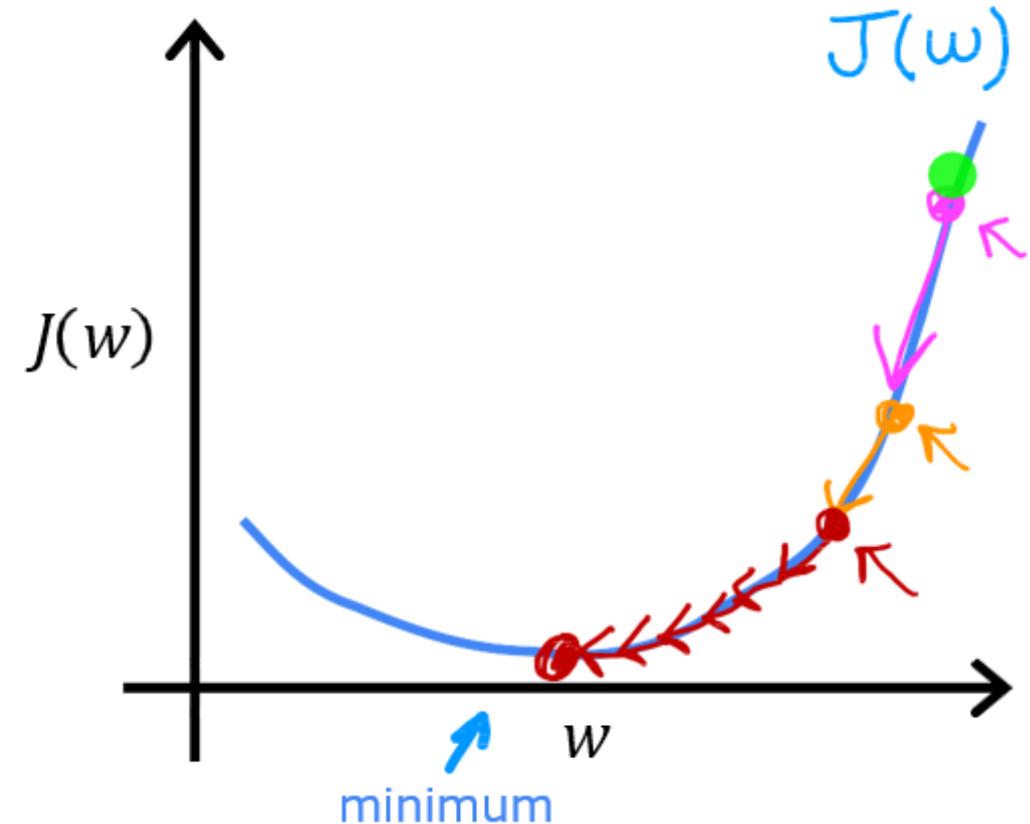
Can reach local minimum with fixed learning rate

$$w = w - \alpha \frac{d}{dw} J(w)$$

smaller
not as large
large

Near a local minimum,
- Derivative becomes smaller
- Update steps become smaller

Can reach minimum without
decreasing learning rate α



Linear regression model

$$f_{w,b}(x) = wx + b$$

Cost function

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2$$

Gradient descent algorithm

repeat until ~~converge~~ {

$$w = w - \alpha \frac{\partial}{\partial w} J(w, b)$$

$$b = b - \alpha \frac{\partial}{\partial b} J(w, b)$$

}

$$\frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) x^{(i)}$$

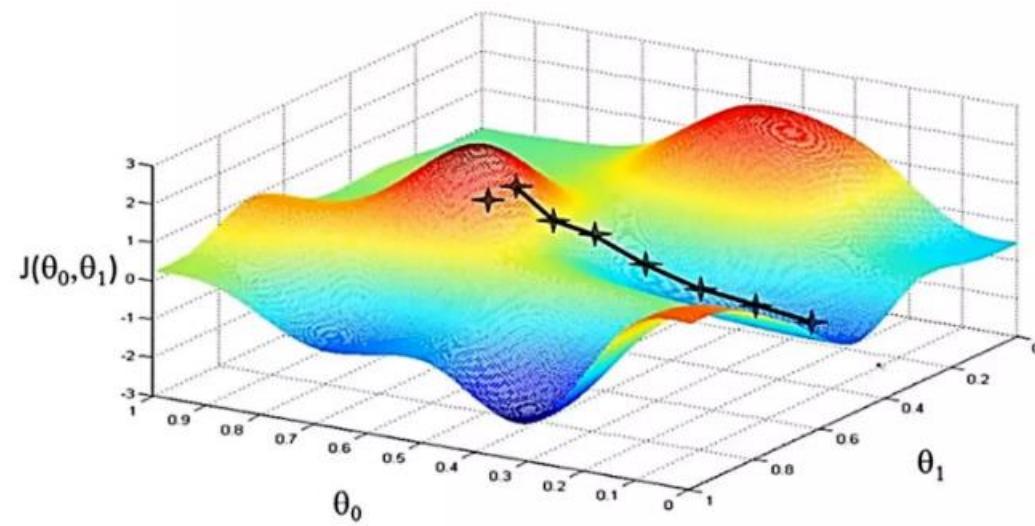
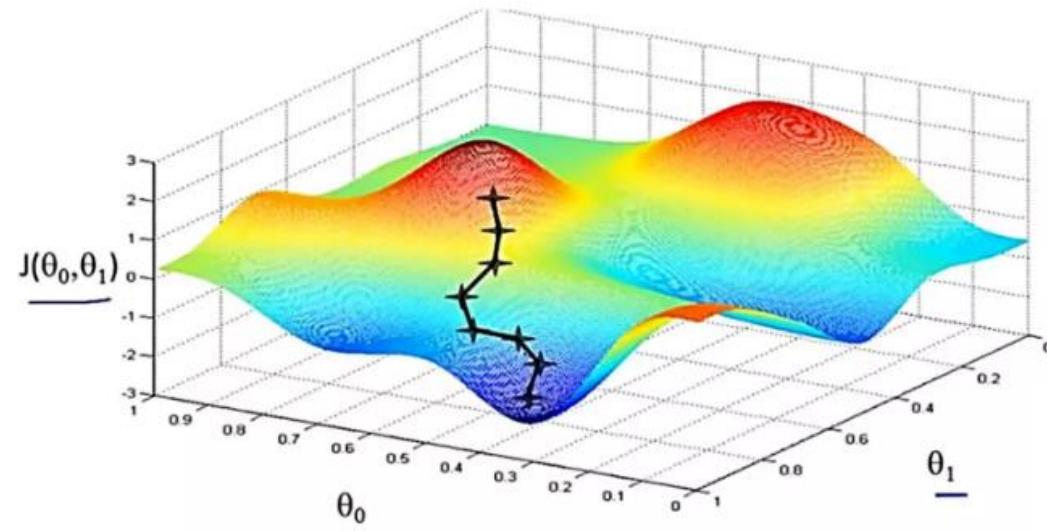
$$\frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})$$

(Optional)

$$\frac{\partial}{\partial w} J(w, b) = \frac{\partial}{\partial w} \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2 = \frac{\partial}{\partial w} \frac{1}{2m} \sum_{i=1}^m (wx^{(i)} + b - y^{(i)})^2$$
$$= \cancel{\frac{1}{2m}} \sum_{i=1}^m (wx^{(i)} + b - y^{(i)}) \cancel{2x^{(i)}} = \boxed{\frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})x^{(i)}}$$

$$\frac{\partial}{\partial b} J(w, b) = \frac{\partial}{\partial b} \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2 = \frac{\partial}{\partial b} \frac{1}{2m} \sum_{i=1}^m (wx^{(i)} + b - y^{(i)})^2$$
$$= \cancel{\frac{1}{2m}} \sum_{i=1}^m (wx^{(i)} + b - y^{(i)}) \cancel{2} = \boxed{\frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})}$$

no $x^{(i)}$



Gradient descent algorithm

repeat until convergence {

$$w = w - \alpha \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) x^{(i)}$$

$$b = b - \alpha \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})$$

}

Update
 w and b
simultaneously

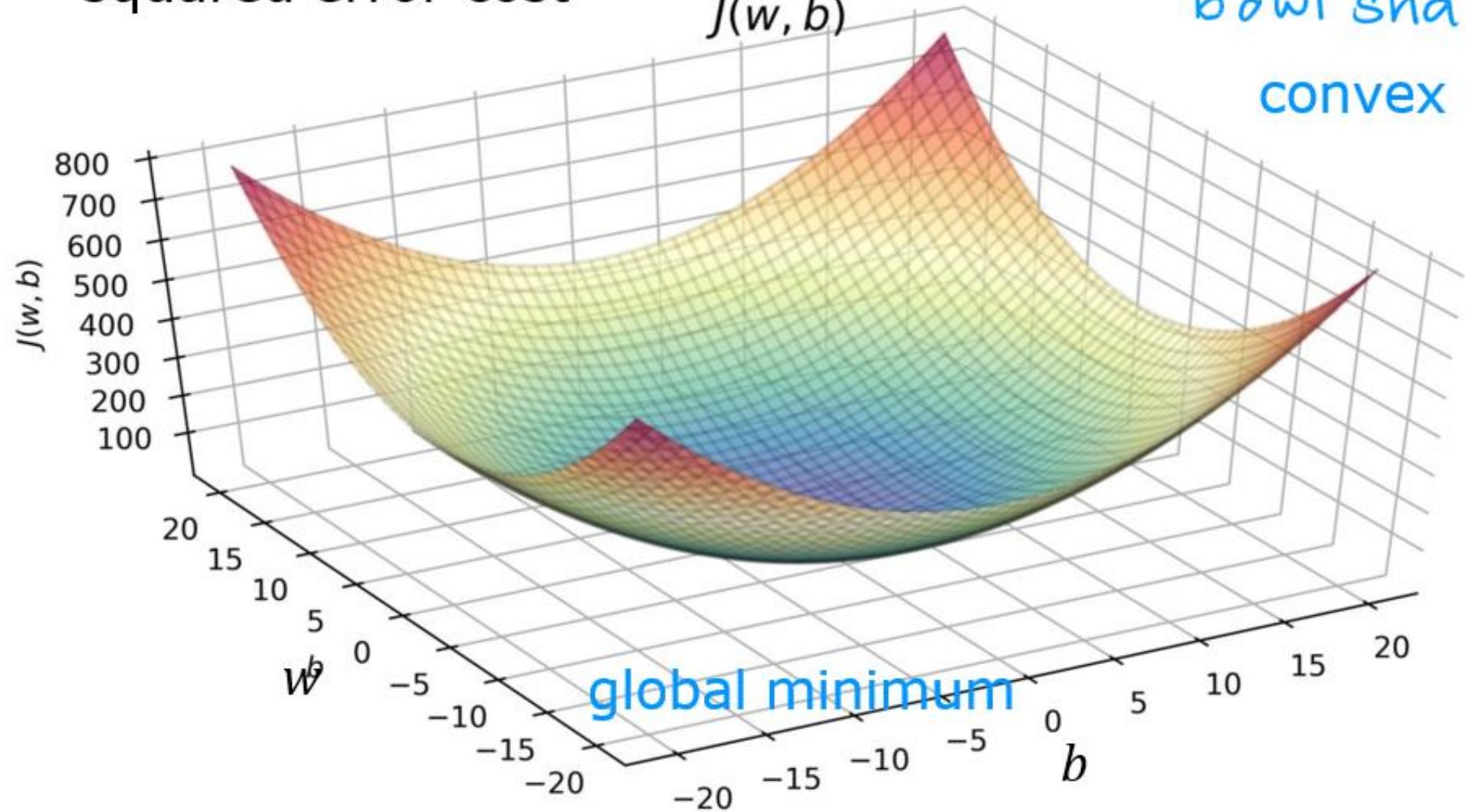
squared error cost

$J(w, b)$

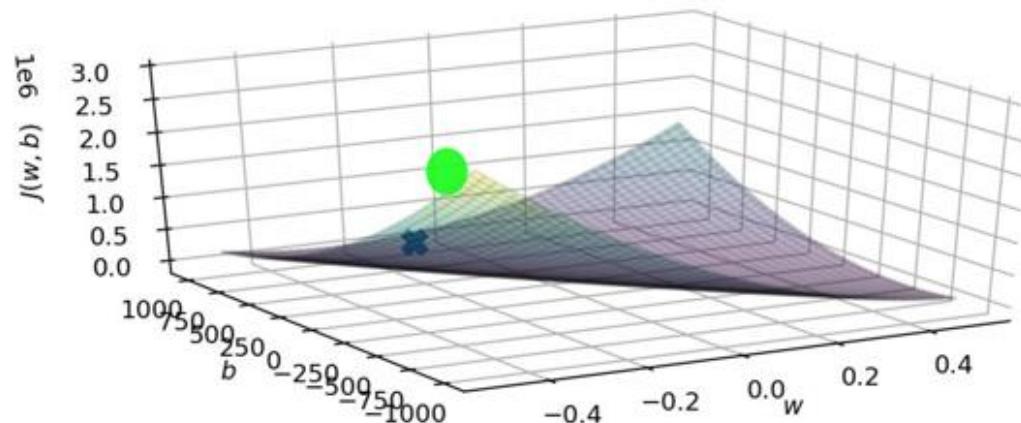
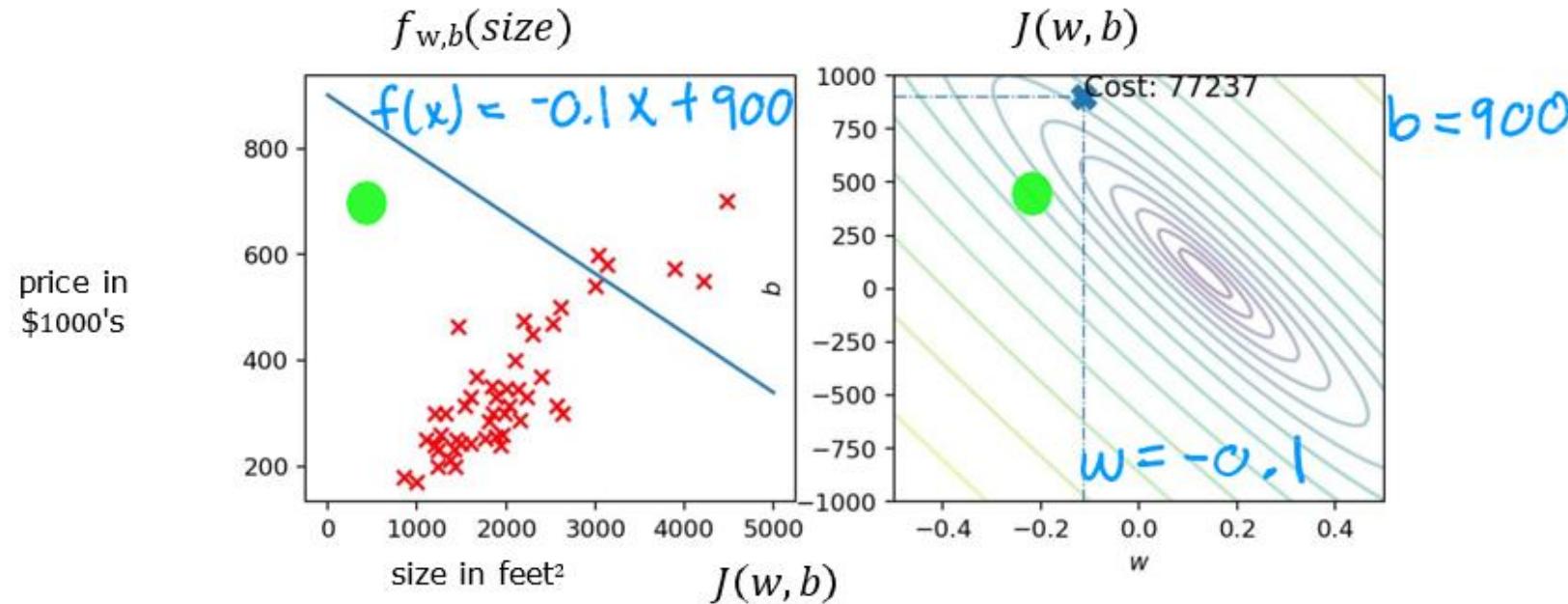
bowl shape

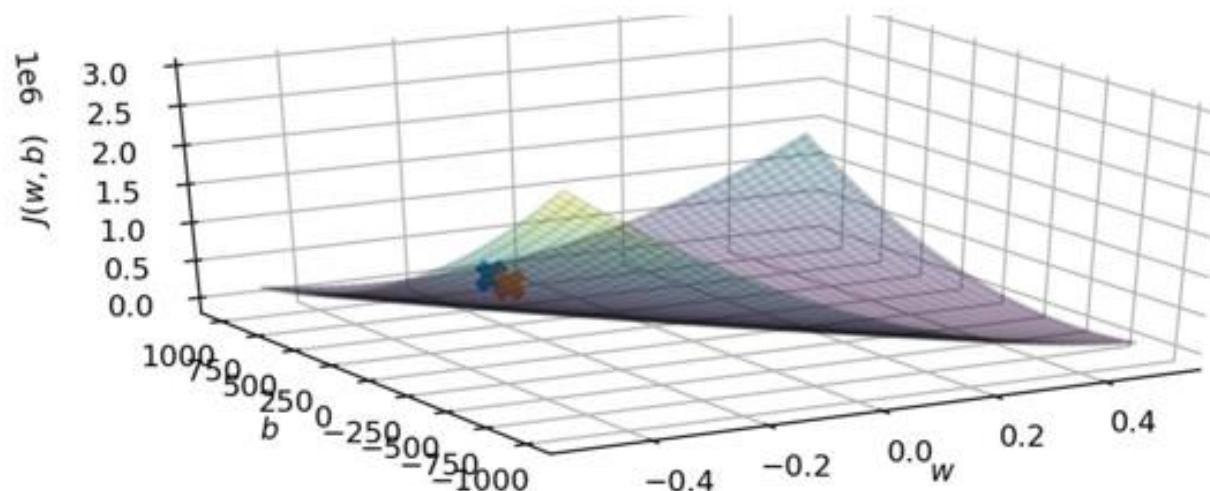
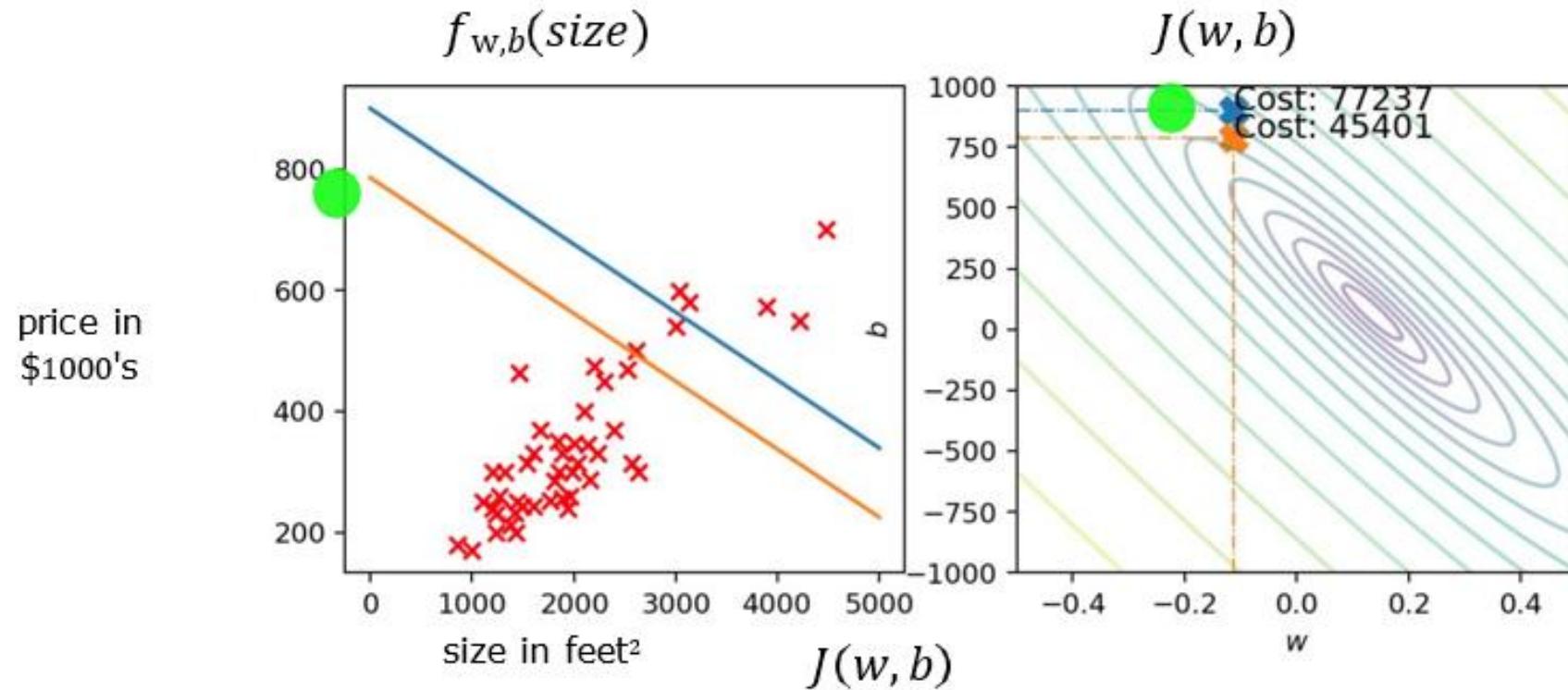


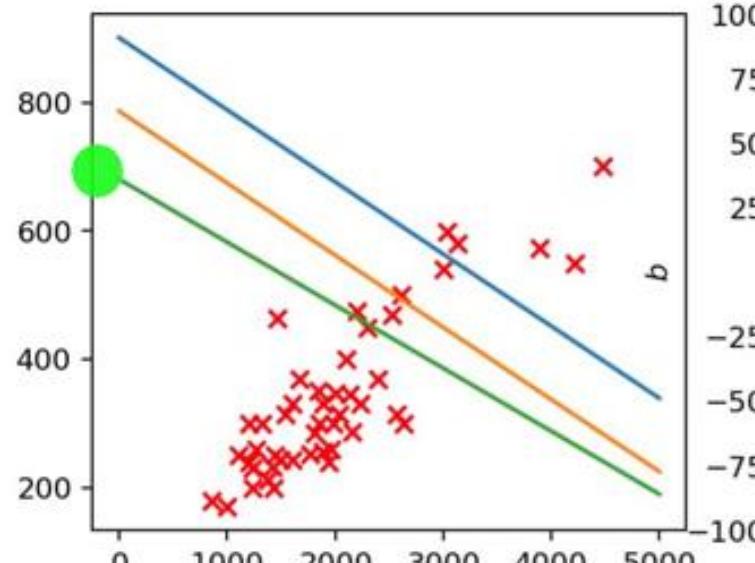
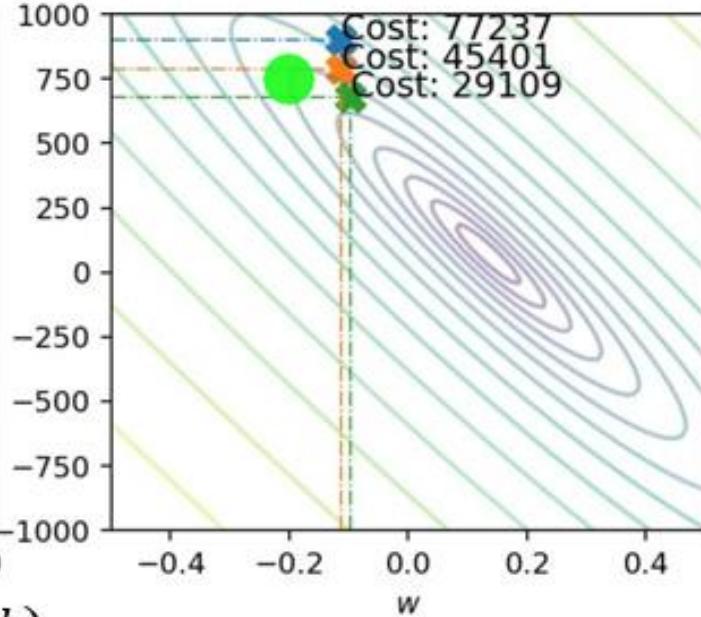
convex function



Running Gradient Descent Algorithm





$f_{w,b}(\text{size})$  $J(w, b)$  $J(w, b)$ 