

10-301/10-601 Introduction to Machine Learning

Machine Learning Department School of Computer Science Carnegie Mellon University

Linear Regression

Matt Gormley & Henry Chai Lecture 7 Feb. 3, 2025

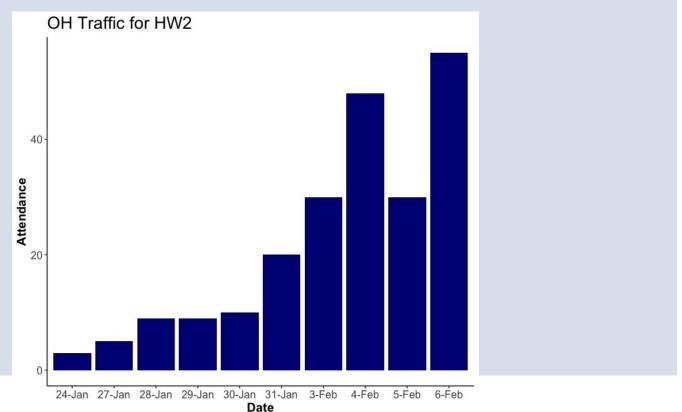
Reminders

- Homework 2: Decision Trees
 - Out: Wed, Jan. 22
 - Due: Mon, Feb. 3 at 11:59pm
- Homework 3: KNN, Perceptron, Lin.Reg.
 - Out: Mon, Feb. 3
 - Due: Mon, Feb. 10 at 11:59pm
 - (only two grace/late days permitted)
- Exam conflicts form

Q&A

Q: How can I get more one-on-one interaction with the course staff?

A: Attend office hours as soon after the homework release as possible!



Q&A

Q: I just asked a question in OH and now my TA is crying quietly -- what did I do wrong?

A: You've just committed the worst of crimes: asking a question that was directly answered in a recitation.

The TA you asked spent hours carefully writing careful recitation notes and solutions, practicing their recitation, responding to criticism / changes from me, etc.

To increase OH efficiency, please review the HW recitation before asking HW questions in OHs.

Q&A

Q: I have a medical emergency or family emergency or disability or other compelling reason and am unable to attend office hours in-person this week. Can an exception be made so I can attend office hours remotely?

A: Yes. Please email the Education Associate(s) and request a period of remote office hours. We will reply with instructions on how to utilize them during the approved time period.

REGRESSION

Regression

Goal:

This is what

differentiates

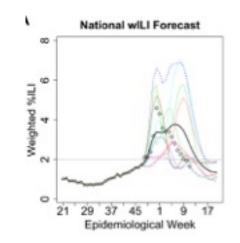
regression from

classification

- Given a training dataset of pairs (x,y) where
 - x is a vector
 - y is a scalar
- Learn a function (aka. curve or line) y' = h(x) that best fits the training data

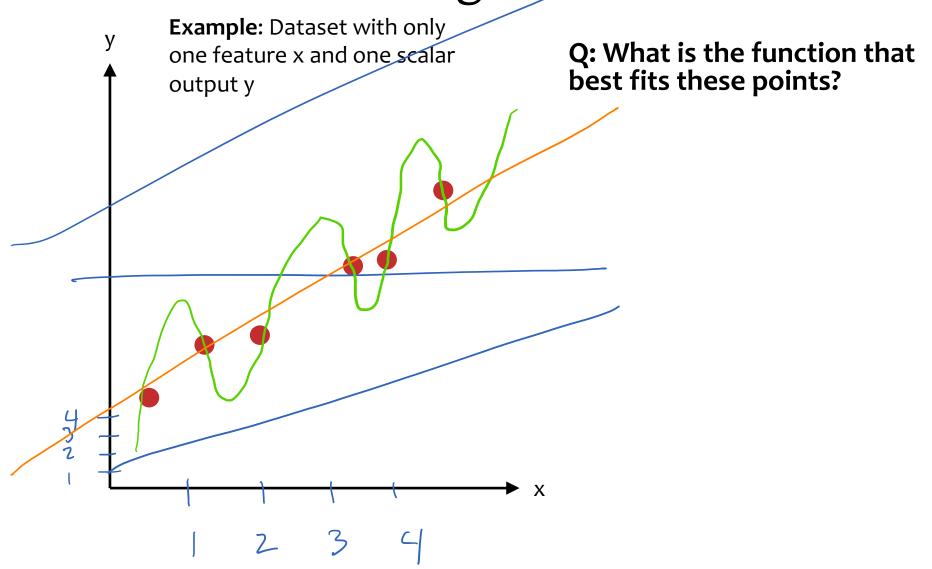
Example Applications:

- Stock price prediction
- Forecasting epidemics
- Speech synthesis
- Generation of images (e.g. Deep Dream)

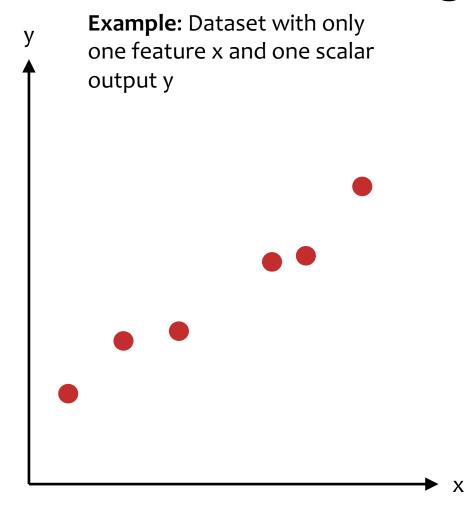




Regression



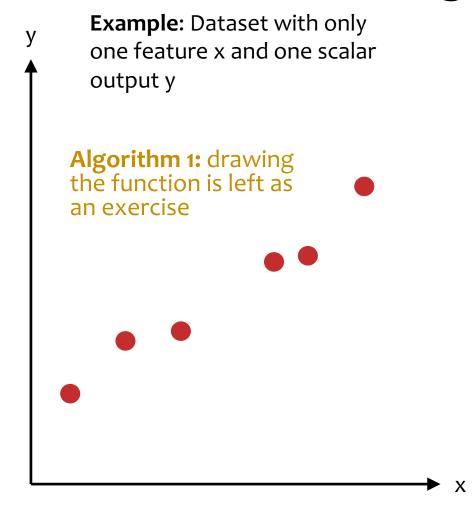
K-NEAREST NEIGHBOR REGRESSION



Algorithm 1: k=1 Nearest Neighbor Regression

- Train: store all (x, y) pairs
- Predict: pick the nearest x in training data and return its y

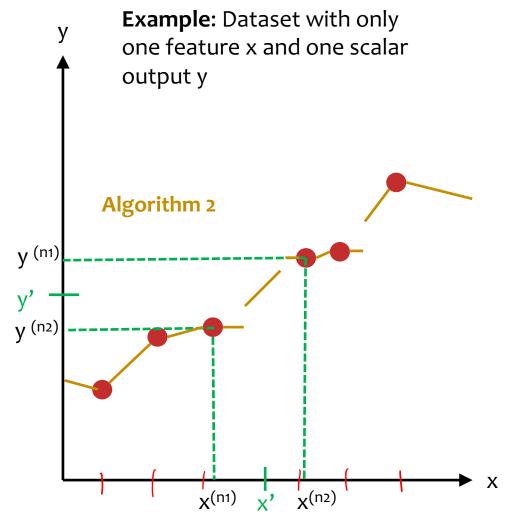
- Train: store all (x, y) pairs
- Predict: pick the nearest two instances $x^{(n_1)}$ and $x^{(n_2)}$ in training data and return the weighted average of their y values



Algorithm 1: k=1 Nearest Neighbor Regression

- Train: store all (x, y) pairs
- Predict: pick the nearest x in training data and return its y

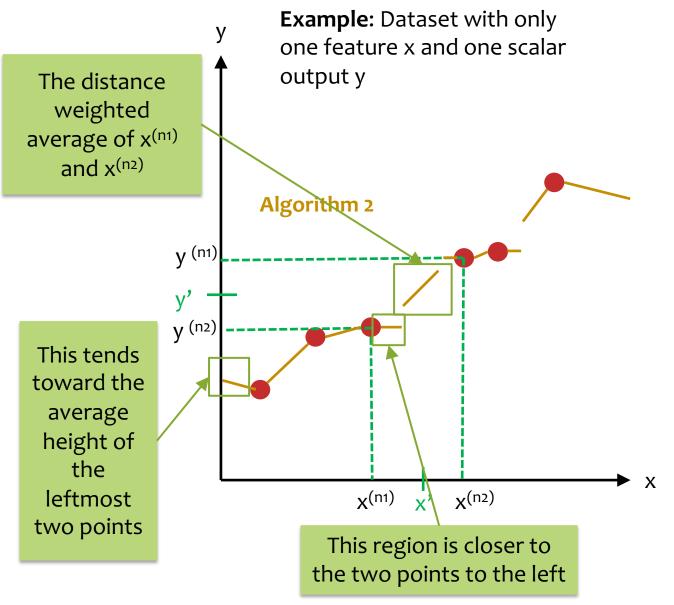
- Train: store all (x, y) pairs
- Predict: pick the nearest two instances $x^{(n_1)}$ and $x^{(n_2)}$ in training data and return the weighted average of their y values



Algorithm 1: k=1 Nearest Neighbor Regression

- Train: store all (x, y) pairs
- Predict: pick the nearest x in training data and return its y

- Train: store all (x, y) pairs
- Predict: pick the nearest two instances $x^{(n_1)}$ and $x^{(n_2)}$ in training data and return the weighted average of their y values



Algorithm 1: k=1 Nearest Neighbor Regression

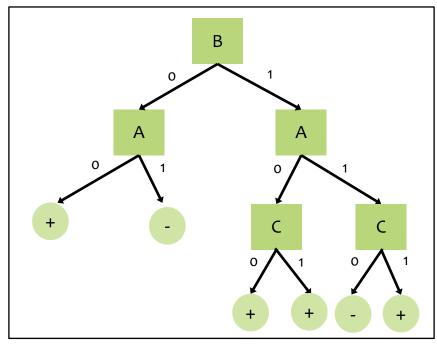
- Train: store all (x, y) pairs
- Predict: pick the nearest x in training data and return its y

- Train: store all (x, y) pairs
- Predict: pick the nearest two instances $x^{(n_1)}$ and $x^{(n_2)}$ in training data and return the weighted average of their y values

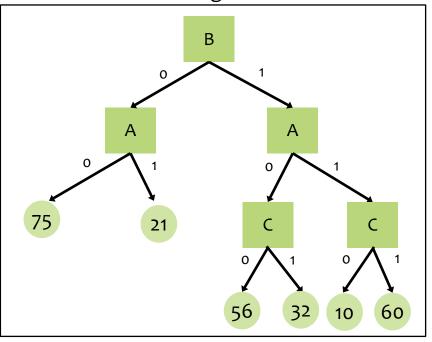
DECISION TREE REGRESSION

Decision Tree Regression

Decision Tree for Classification



Decision Tree for Regression

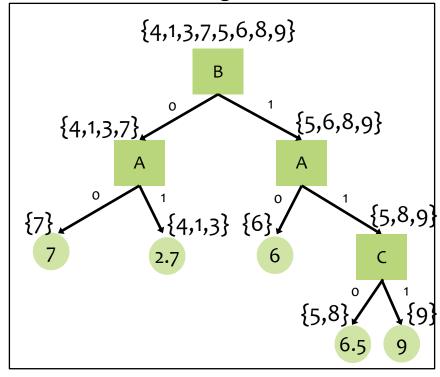


Decision Tree Regression

Dataset for Regression

Y	А	В	C
4	1	0	0
1	1	0	1
3	1	0	0
7	0	0	1
5	1	1	0
5 6	0	1	1
8	1	1	0
9	1	1	1

Decision Tree for Regression



During learning, choose the attribute that minimizes an appropriate splitting criterion (e.g. mean squared error, mean absolute error)

LINEAR FUNCTIONS, RESIDUALS, AND MEAN SQUARED ERROR

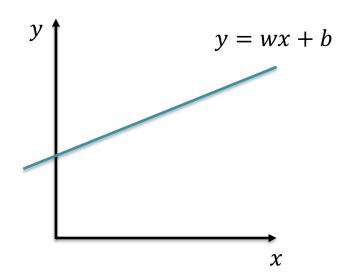
Linear Functions

<u>Def</u>: Regression is predicting real-valued outputs

$$\mathcal{D} = \left\{ \left(\mathbf{x}^{(i)}, y^{(i)} \right) \right\}_{i=1}^{n} \text{ with } \mathbf{x}^{(i)} \in \mathbb{R}^{M}, y^{(i)} \in \mathbb{R}$$

Common Misunderstanding:

Linear functions ≠ Linear decision boundaries



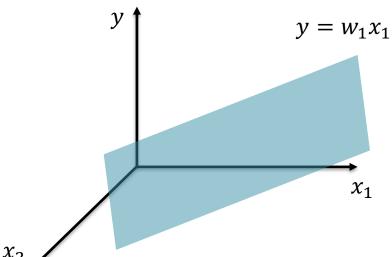
Linear Functions

<u>Def</u>: Regression is predicting real-valued outputs

$$\mathcal{D} = \left\{ \left(\mathbf{x}^{(i)}, y^{(i)} \right) \right\}_{i=1}^{n} \text{ with } \mathbf{x}^{(i)} \in \mathbb{R}^{M}, y^{(i)} \in \mathbb{R}$$

Common Misunderstanding:

Linear functions ≠ Linear decision boundaries



$$y = w_1 x_1 + w_2 x_2 + b$$
 $\overrightarrow{w} = [\omega_1, \ldots, \omega_M]$

- A general linear function is $y = \mathbf{w}^T \mathbf{x} + b$
- A general linear decision boundary is $y = sign(\mathbf{w}^T \mathbf{x} + b)$

Key Idea of Linear Regression

Residuals

Def: a residual is the vertical distance

Find the linear function observed value
$$y''(i)$$

to predicted value $y''(i)$

For some training set

 $y''(i) = y''(i) - h(x''(i))$

Mean squared error

 $y''(i) = y''(i) - h(x''(i))$

Mean squared error

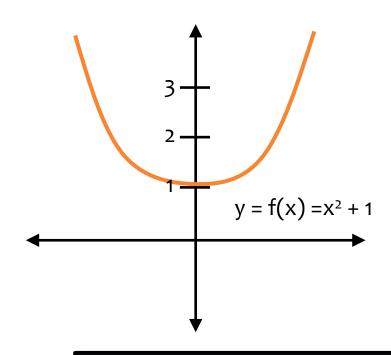
Key Idea of Linear Regression

$$\frac{\langle v \rangle \langle o \rangle}{\langle v \rangle \langle v \rangle} \times (v) \times (v)$$

The Big Picture

OPTIMIZATION FOR ML

min vs. argmin



$$v* = min_x f(x)$$

$$x^* = \operatorname{argmin}_x f(x)$$

1. Poll Question 1: What is v*?

2. Poll Question 2: What is x*?

Unconstrained Optimization

• *Def*: In **unconstrained optimization**, we try minimize (or maximize) a function with *no constraints* on the inputs to the function

Given a function

$$J(\boldsymbol{ heta}), J: \mathbb{R}^M
ightarrow \mathbb{R}$$

Our goal is to find

$$\hat{m{ heta}} = \operatornamewithlimits{argmin}_{m{ heta} \in \mathbb{R}^M} J(m{ heta})$$
 For ML, this is the objective function

Optimization for ML

Not quite the same setting as other fields...

- Function we are optimizing might not be the true goal (e.g. likelihood vs generalization error)
- Precision might not matter
 (e.g. data is noisy, so optimal up to 1e-16 might not help)
- Stopping early can help generalization error
 (i.e. "early stopping" is a technique for regularization discussed more next time)

OPTIMIZATION METHOD #0: RANDOM GUESSING

Notation Trick: Folding in the Intercept Term

$$\mathbf{x}' = [1, x_1, x_2, \dots, x_M]^T$$

 $\boldsymbol{\theta} = [b, w_1, \dots, w_M]^T$

Notation Trick: fold the bias b and the weights w into a single vector $\boldsymbol{\theta}$ by prepending a constant to x and increasing dimensionality by one!

$$h_{\mathbf{w},b}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$$

 $h_{\boldsymbol{\theta}}(\mathbf{x}') = \boldsymbol{\theta}^T \mathbf{x}'$

This convenience trick allows us to more compactly talk about linear functions as a simple dot product (without explicitly writing out the intercept term every time).

Linear Regression as Function Approximation

$$\mathcal{D}=\{\mathbf{x}^{(i)},y^{(i)}\}_{i=1}^N$$
 where $\mathbf{x}\in\mathbb{R}^M$ and $y\in\mathbb{R}$

$$\mathbf{x}^{(i)} \sim p^*(\cdot)$$
$$y^{(i)} = h^*(\mathbf{x}^{(i)})$$

2. Choose hypothesis space, \mathcal{H} : all linear functions in M-dimensional space

$$\mathcal{H} = \{h_{\boldsymbol{\theta}} : h_{\boldsymbol{\theta}}(\mathbf{x}) = \boldsymbol{\theta}^T \mathbf{x}, \boldsymbol{\theta} \in \mathbb{R}^M \}$$

3. Choose an objective function: mean squared error (MSE)

$$J(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^{N} e_i^2$$

$$= \frac{1}{N} \sum_{i=1}^{N} \left(y^{(i)} - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) \right)^2$$

$$= \frac{1}{N} \sum_{i=1}^{N} \left(y^{(i)} - \boldsymbol{\theta}^T \mathbf{x}^{(i)} \right)^2$$

- 4. Solve the unconstrained optimization problem via favorite method:
 - gradient descent
 - closed form
 - stochastic gradient descent

$$\hat{m{ heta}} = \operatorname*{argmin}_{m{ heta}} J(m{ heta})$$

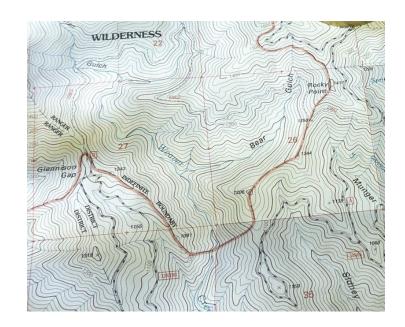
5. Test time: given a new x, make prediction \hat{y}

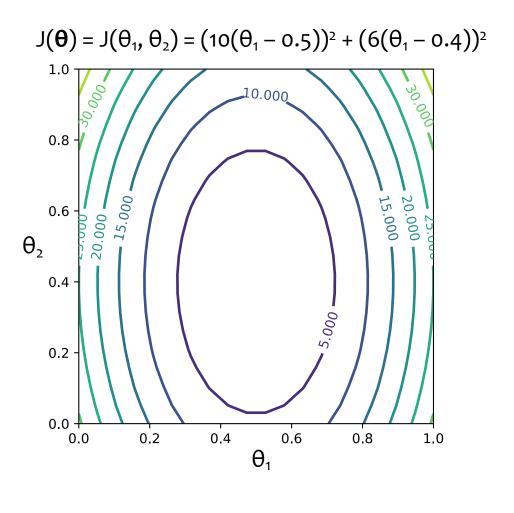
$$\hat{y} = h_{\hat{oldsymbol{ heta}}}(\mathbf{x}) = \hat{oldsymbol{ heta}}^T \mathbf{x}$$

Contour Plots

Contour Plots

- Each level curve labeled with value
- value label indicates the value of the function for all points lying on that level curve
- 3. Just like a topographical map, but for a function



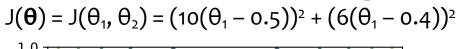


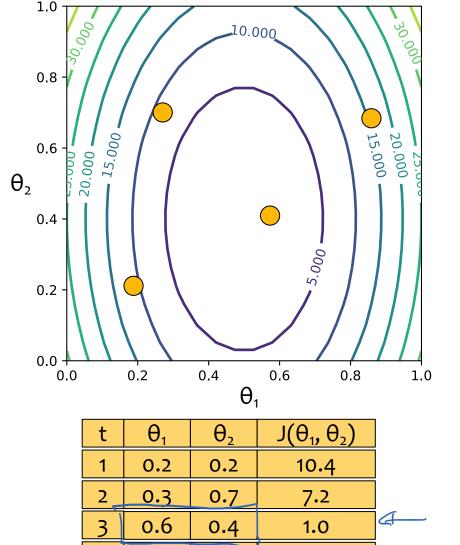
Optimization by Random Guessing



Optimization Method #0: Random Guessing

- 1. Pick a random $\boldsymbol{\theta}$
- 2. Evaluate $J(\theta)$
- 3. Repeat steps 1 and 2 many times
- 4. Return θ that gives smallest $J(\theta)$





0.7

0.9

16.2

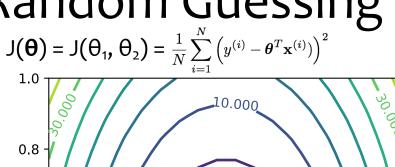
Optimization by Random Guessing

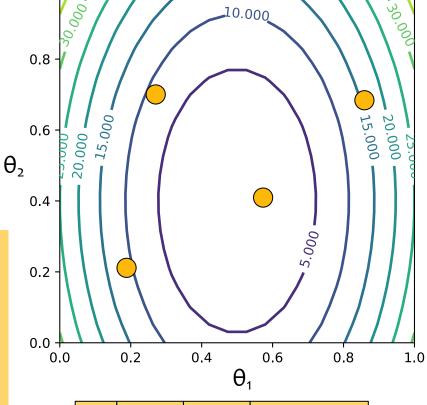
Optimization Method #0: Random Guessing

- 1. Pick a random θ
- 2. Evaluate $J(\theta)$
- 3. Repeat steps 1 and 2 many times
- 4. Return θ that gives smallest $J(\theta)$

For Linear Regression:

- objective function is Mean Squared Error (MSE)
- MSE = J(w, b) = J(θ_1 , θ_2) = $\frac{1}{N} \sum_{i=1}^{N} (y^{(i)} - \boldsymbol{\theta}^T \mathbf{x}^{(i)})^2$
- contour plot: each line labeled with MSE – lower means a better fit
- **minimum** corresponds to parameters (w,b) = (θ_1, θ_2) that **best fit** some training dataset



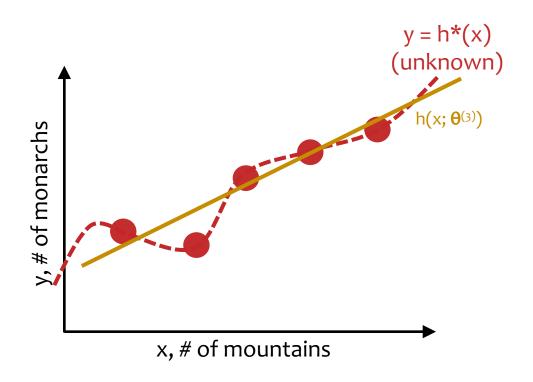


t	θ_1	θ_2	$J(\theta_1, \theta_2)$
1	0.2	0.2	10.4
2	0.3	0.7	7.2
3	0.6	0.4	1.0
4	0.9	0.7	16.2

Linear Regression: Running Example



Counting Butterflies



breed only in places where milkweed grows. When the last ice age ended, and

MIGRATION ROUTES OF MONARCH BUTTERFLIES CANADA UNITED STATES ATLANTIC OCEAN Summer breeding area Spring breeding area Wintering area PACIFIC Corn Belt region OCEAN → Spring migration route - Fall migration route

This map shows migration routes of fall and spring migrations, both east and west of the Rocky Mountains.

the cold and glaciers retreated, milkweed may have gradually spread northward, and monarchs may have followed. But the monarch butterfly remained a tropical creature, unable to survive the severe northern winters. So every year as winter approached, monarchs left their summer fields of milkweed and flew south again. To this day, every spring and summer, monarchs travel north to their breeding grounds across the eastern United States and Canada. Every winter, they return to Mexico.

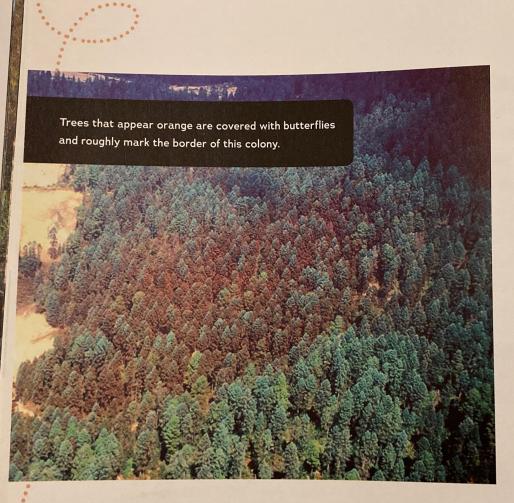
Researchers began taking measurements in 1993. The highest year on record came in 1997, when the colonies covered about 45 acres (18 ha), an area equal to about thirty-four football fields. Scientists aren't sure exactly how many butterflies

LOCATION OF MONARCH BUTTERFLY COLONIES WINTERING IN MEXICO UNITED STATES Contapec bulf of Mexico Mexico City Toluca Angangueo Cumbres MEXICO Chivati-Huacal omas de Aparicio PACIFIC Zitácuaro OCEAN erro Pelon Monarch butterfly colony Capital City (Mexican location detail) Town

The eastern monarchs migrate to just twelve mountaintops, all located in central Mexico.

that represented, but one estimate is that there were one billion monarchs in the colonies that winter.

But as researchers measured the colonies year after year, they noticed that the colonies were shrinking. By 2014 the colonies measured just 1.7 acres (0.7 ha), or less than one and a half football fields. That year there may have been only about thirty-five million monarchs in the colonies.



Many scientists were worried. The population of eastern monarchs had dropped more than 90 percent in just seventeen years.

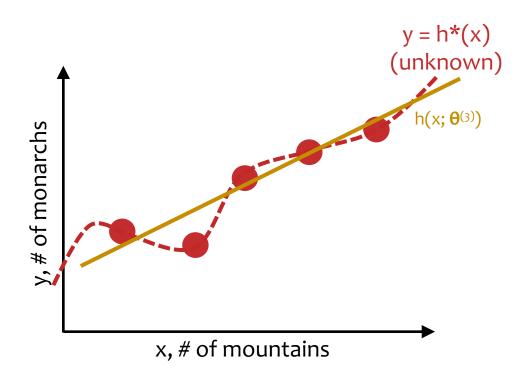
At the same time, scientists in California reported that the number of western monarchs was dropping as well. From 1997 to 2014, the number of monarchs overwintering along the California coast had fallen by 74 percent.

Populations of overwintering monarchs were falling fast. By 2014 their numbers had fallen so far that people wondered

whether the monarch butterfly should be listed as an endangered species—a species in danger of becoming extinct, or disappearing forever.

Losing monarchs could be bad for our world because monarchs play an important part in the food web. Despite the milkweed toxins in their bodies, they are food for songbirds, spiders, and insects. Monarchs visit many flowers and act as pollinators.

Counting Butterflies



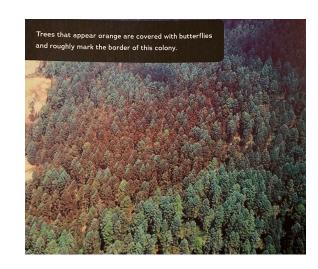
Linear Regression in High Dimensions

- In our discussions of linear regression, we will always assume there is just one output,
- But our inputs will usually have many features:

$$\mathbf{x} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_M]^\mathsf{T}$$

- For example:
 - suppose we had a drone take pictures of each section of forest
 - each feature could correspond to a pixel in this image such that $x_m = 1$ if the pixel is orange and $x_m = 0$ otherwise
 - the output y would be the number of butterflies in each picture

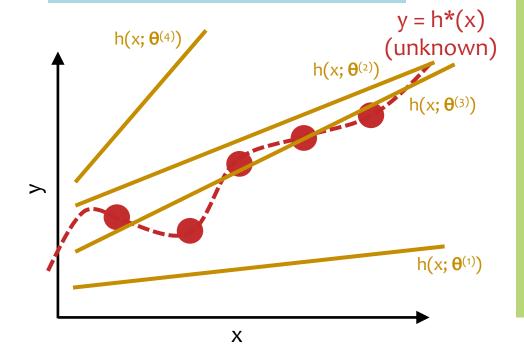
Q: How would you obtain ground truth data?



Linear Regression by Rand. Guessing

Optimization Method #0: Random Guessing

- 1. Pick a random $\boldsymbol{\theta}$
- 2. Evaluate $J(\theta)$
- 3. Repeat steps 1 and 2 many times
- 4. Return θ that gives smallest $J(\theta)$



For Linear Regression:

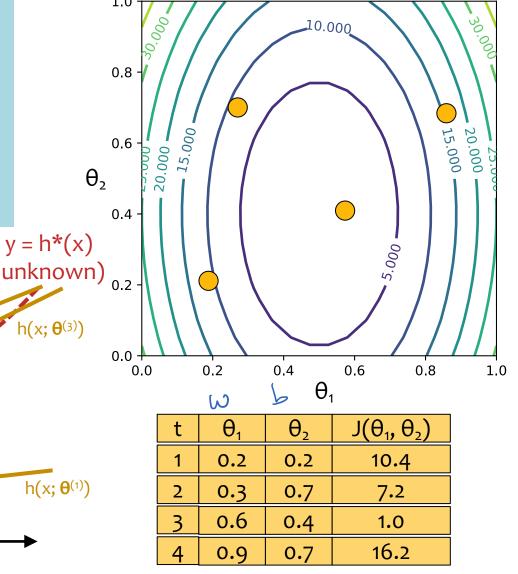
- target function h*(x) is unknown
- only have access to h*(x) through training examples (x⁽ⁱ⁾,y⁽ⁱ⁾)
- want $h(x; \theta^{(t)})$ that **best** approximates $h^*(x)$
- enable generalization w/inductive bias that restricts hypothesis class to linear functions

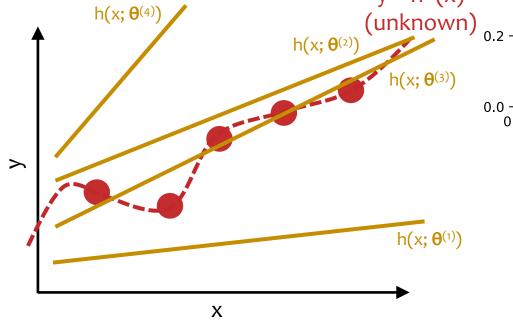
Linear Regression by Rand. Guessing

J(θ) = J(θ_1 , θ_2) = $\frac{1}{N} \sum_{i=1}^{N} (y^{(i)} - \theta^T \mathbf{x}^{(i)})^2$

Optimization Method #0: Random Guessing

- Pick a random θ
- Evaluate $J(\theta)$
- Repeat steps 1 and 2 many times
- Return θ that gives smallest $J(\theta)$





OPTIMIZATION METHOD #1: GRADIENT DESCENT

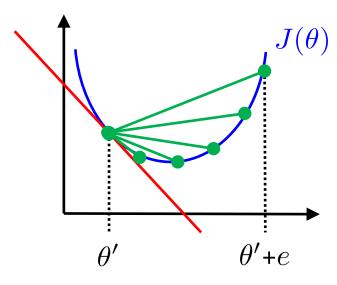
Derivatives

1. Derivative as a Slope

$J(\theta)$

slope =
$$\frac{\partial J(\theta)}{\partial \theta}$$
of the tangent line

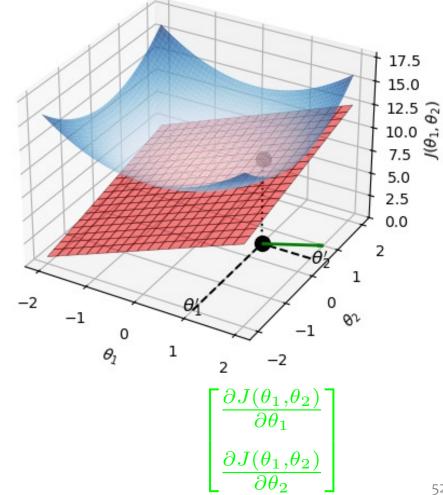
2. Derivative as a Limit



$$\frac{\partial J(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \lim_{e \to 0} \frac{J(\boldsymbol{\theta} + e) - J(\boldsymbol{\theta})}{e}$$

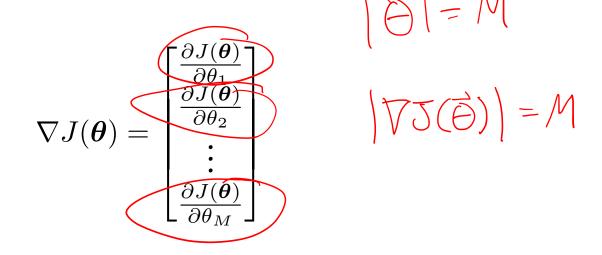
The limit of the secants is the tangent

3. Derivative as a Tangent Plane



Gradient

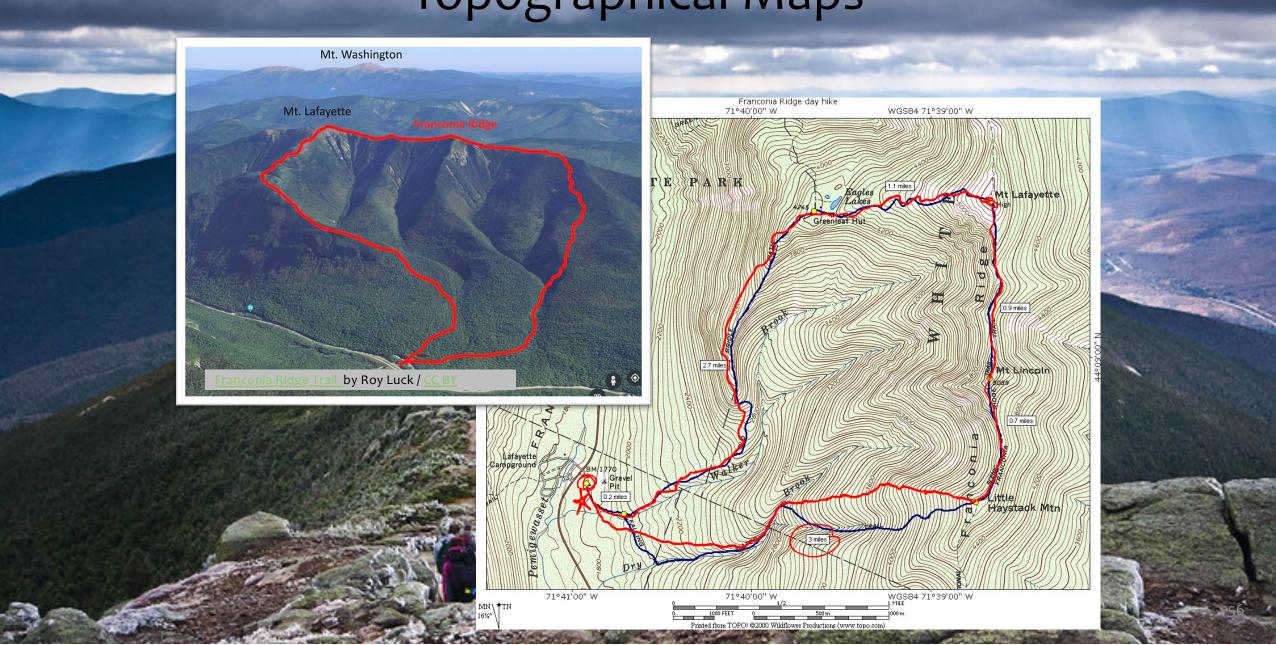
Def: The gradient of $J:\mathbb{R}^M \to \mathbb{R}$ is



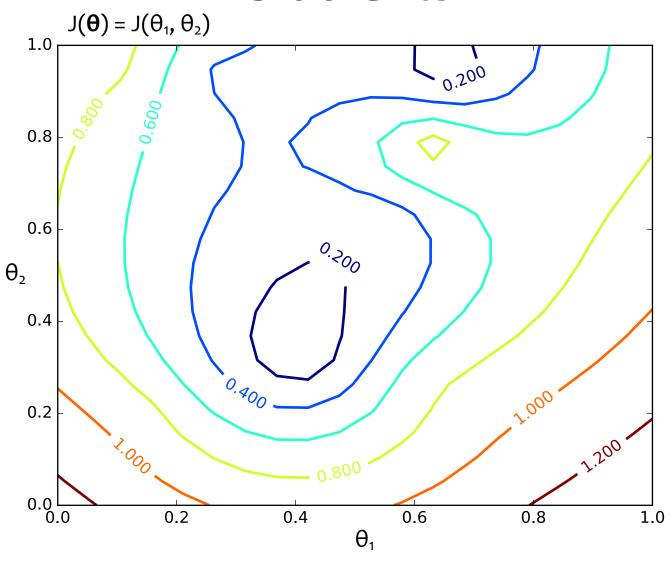
Each entry is a first-order partial derivative



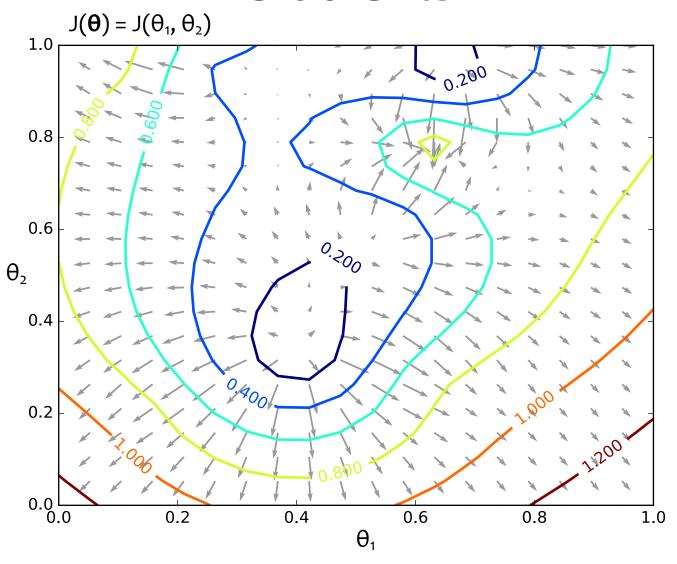
Topographical Maps



Gradients

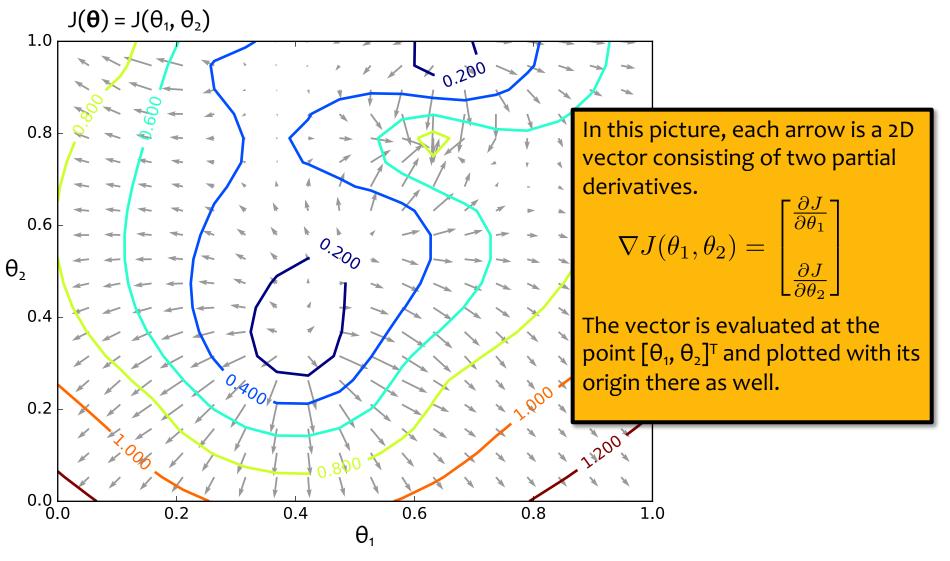


Gradients



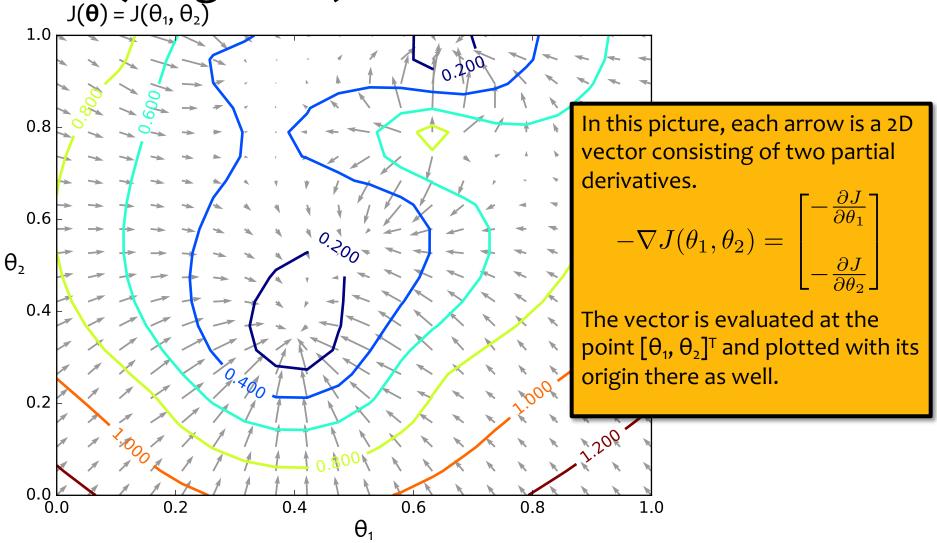
These are the **gradients** that Gradient **Ascent** would follow.

Gradients



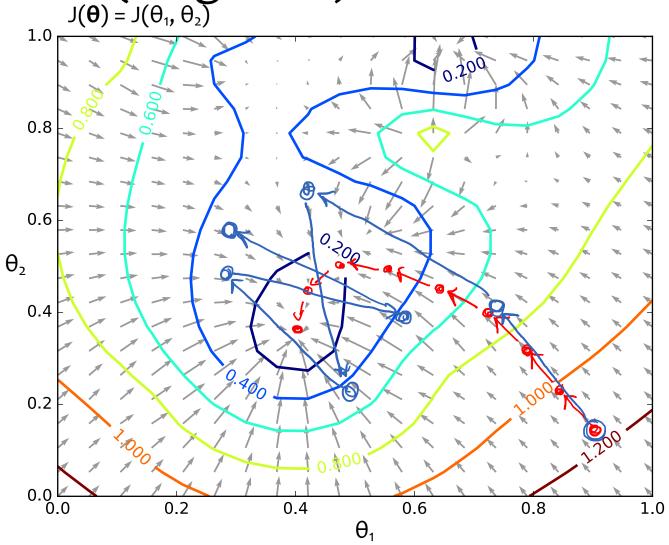
These are the **gradients** that Gradient **Ascent** would follow.

(Negative) Gradients $J(\theta) = J(\theta_1, \theta_2)$



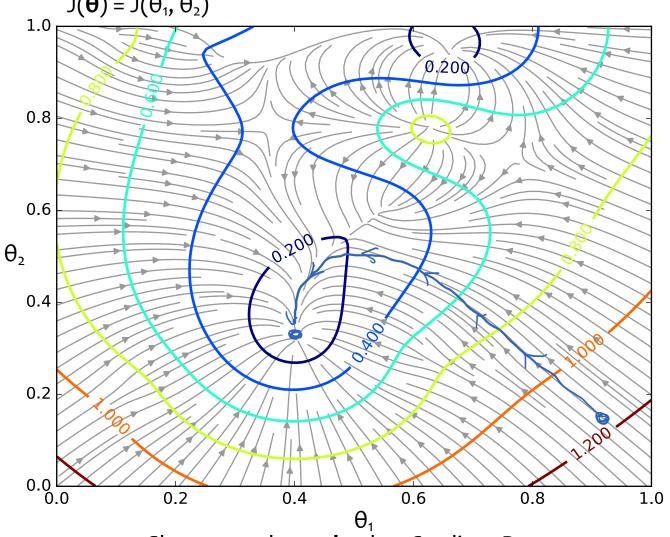
These are the **negative** gradients that Gradient **Descent** would follow.

(Negative) Gradients $J(\theta) = J(\theta_1, \theta_2)$



These are the **negative** gradients that Gradient **Descent** would follow.

(Negative) Gradient Paths $J(\theta) = J(\theta_1, \theta_2)$



Shown are the **paths** that Gradient Descent would follow if it were making **infinitesimally small steps**.

Gradient Descent

Gradient Descent Algorithm

(2) Repeat
$$t = 1, 2, ..., T$$

a) Compute gradient $g = \nabla J(G)$

b) Select step size S_t

c) Update params $G = G - V_2 G$

(3) Return $G = G - V_2 G$

Remarks

Initial point

- Ô randomly

- Ô - Ô all zeros

Step Size

- fixed value
$$Y = 0.1$$

- set a schedule

 $X_t = \frac{X_0}{(t-1)X_0 + 1}$

- line search

Stopping Criterian

- $\hat{g} \approx \hat{O}$

- $\|\nabla J(\hat{O})\|_2 \leq \epsilon$ where $\epsilon = |O^{-8}|$

Gradient Descent: Step Size

Poll Question 3:

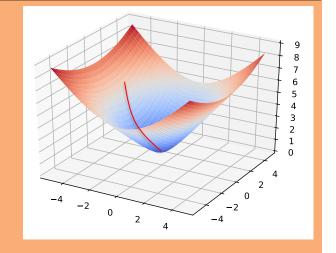
In gradient descent, what could go wrong if we *always* use the same step size (or step size schedule) for every problem we encounter?

Answer:			

Gradient Descent

Algorithm 1 Gradient Descent

- 1: **procedure** $GD(\mathcal{D}, \boldsymbol{\theta}^{(0)})$
- 2: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta}^{(0)}$
- 3: **while** not converged **do**
- 4: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} \gamma \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$
- 5: return θ



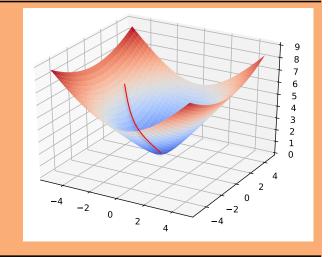
In order to apply GD to Linear Regression all we need is the gradient of the objective function (i.e. vector of partial derivatives).

$$abla_{m{ heta}} J(m{ heta}) = egin{bmatrix} rac{d}{d heta_2} J(m{ heta}) \ dots \ rac{d}{d heta_M} J(m{ heta}) \end{bmatrix}$$

Gradient Descent

Algorithm 1 Gradient Descent

- 1: **procedure** $\mathrm{GD}(\mathcal{D}, \boldsymbol{\theta}^{(0)})$
- 2: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta}^{(0)}$
- 3: while not converged do
- 4: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} \gamma \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$
- 5: return θ



There are many possible ways to detect **convergence**. For example, we could check whether the L2 norm of the gradient is below some small tolerance.

$$||\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})||_2 \leq \epsilon$$

Alternatively we could check that the reduction in the objective function from one iteration to the next is small.

GRADIENT DESCENT FOR LINEAR REGRESSION

Linear Regression as Function Approximation

$$\mathcal{D}=\{\mathbf{x}^{(i)},y^{(i)}\}_{i=1}^N$$
 where $\mathbf{x}\in\mathbb{R}^M$ and $y\in\mathbb{R}$ 1. Assume \mathcal{D} generated as:

$$\mathbf{x}^{(i)} \sim p^*(\cdot)$$
$$y^{(i)} = h^*(\mathbf{x}^{(i)})$$

2. Choose hypothesis space, \mathcal{H} : all linear functions in M-dimensional space

$$\mathcal{H} = \{h_{\boldsymbol{\theta}} : h_{\boldsymbol{\theta}}(\mathbf{x}) = \boldsymbol{\theta}^T \mathbf{x}, \boldsymbol{\theta} \in \mathbb{R}^M \}$$

3. Choose an objective function: mean squared error (MSE)

$$J(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^{N} e_i^2$$

$$= \frac{1}{N} \sum_{i=1}^{N} \left(y^{(i)} - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) \right)^2$$

$$= \frac{1}{N} \sum_{i=1}^{N} \left(y^{(i)} - \boldsymbol{\theta}^T \mathbf{x}^{(i)} \right)^2$$

- 4. Solve the unconstrained optimization problem via favorite method:
 - gradient descent
 - closed form
 - stochastic gradient descent

$$\hat{m{ heta}} = \operatorname*{argmin}_{m{ heta}} J(m{ heta})$$

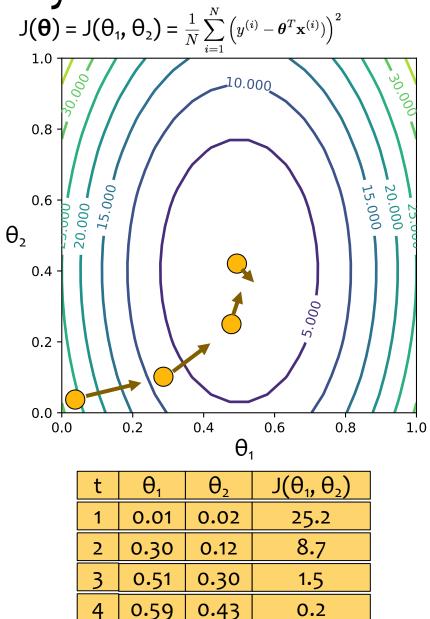
5. Test time: given a new x, make prediction \hat{y}

$$\hat{y} = h_{\hat{oldsymbol{ heta}}}(\mathbf{x}) = \hat{oldsymbol{ heta}}^T \mathbf{x}$$

Linear Regression by Gradient Desc. Optimization Method #1: $J(\theta) = J(\theta_1, \theta_2) = \frac{1}{N} \sum_{i=1}^{N} (y^{(i)} - \theta^T \mathbf{x}^{(i)})^2$

Optimization Method #1: Gradient Descent

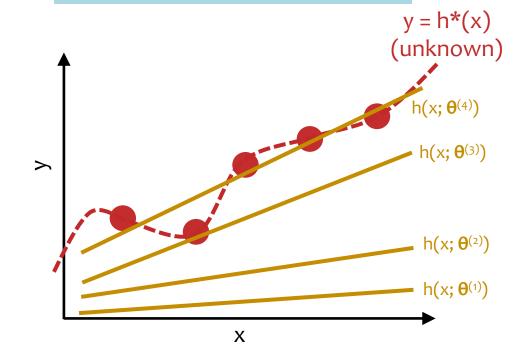
- Pick a random θ
- Repeat:
 - a. Evaluate gradient $\nabla J(\boldsymbol{\theta})$
 - b. Step opposite gradient
- Return θ that gives smallest $J(\theta)$



Linear Regression by Gradient Desc.

Optimization Method #1: Gradient Descent

- 1. Pick a random θ
- 2. Repeat:
 - a. Evaluate gradient $\nabla J(\boldsymbol{\theta})$
 - b. Step opposite gradient
- 3. Return θ that gives smallest $J(\theta)$

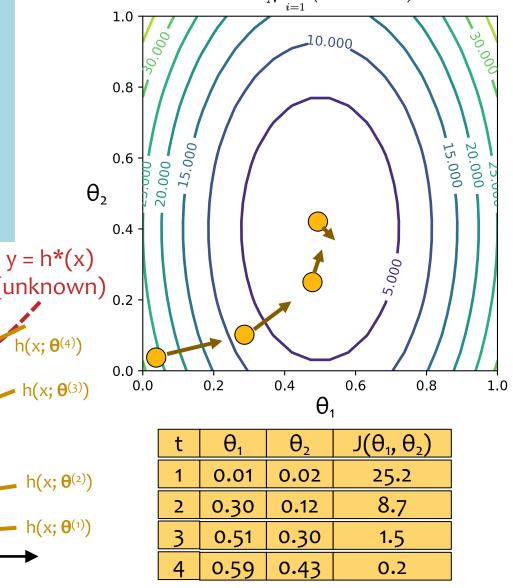


t	θ_1	θ_2	$J(\theta_1, \theta_2)$
1	0.01	0.02	25.2
2	0.30	0.12	8.7
3	0.51	0.30	1.5
4	0.59	0.43	0.2

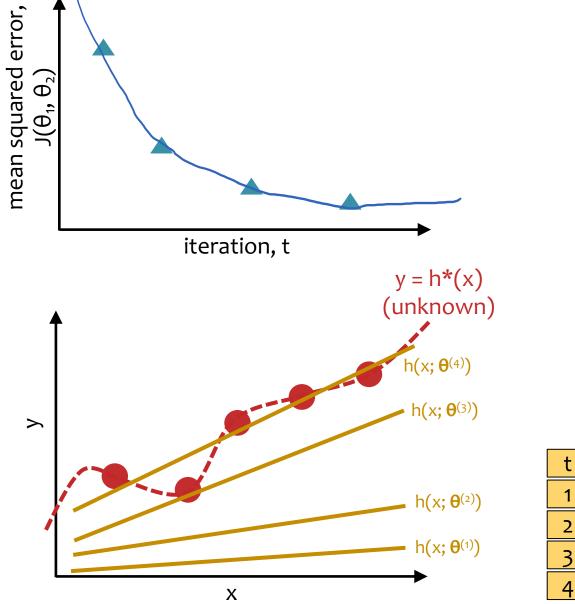
Linear Regression by Gradient Desc. $J(\theta) = J(\theta_1, \theta_2) = \frac{1}{N} \sum_{i=1}^{N} (y^{(i)} - \theta^T \mathbf{x}^{(i)})^2$

Optimization Method #1: Gradient Descent

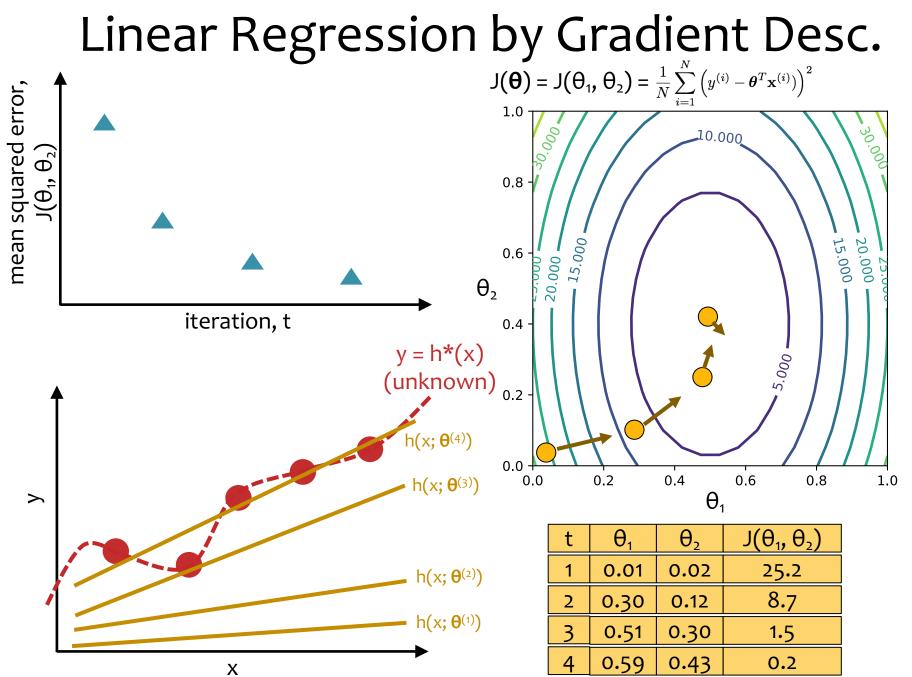
- Pick a random θ
- Repeat:
 - a. Evaluate gradient $\nabla J(\boldsymbol{\theta})$
 - b. Step opposite gradient
- Return **θ** that gives smallest $J(\theta)$



Linear Regression by Gradient Desc.



t	θ_1	θ_2	$J(\theta_1, \theta_2)$
1	0.01	0.02	25.2
2	0.30	0.12	8.7
3	0.51	0.30	1.5
4	0.59	0.43	0.2



Gradient Calculation for Linear Regression

$$\frac{MSE:}{J(\Theta)} = \frac{1}{N} \sum_{i=1}^{N} \overline{J^{(i)}}(\Theta) \quad \text{where} \quad \overline{J^{(i)}}(\Theta) = \frac{1}{2} (y^{(i)} - \overrightarrow{D}^{T} \overrightarrow{X}^{(i)})^{2}$$
Partial Derivy
$$\frac{3J^{(i)}(\overrightarrow{\Theta})}{\delta\Theta_{j}} = \frac{\delta}{\delta\Theta_{j}} \left(\frac{1}{2} (y^{(i)} - \overrightarrow{D}^{T} \overrightarrow{X}^{(i)})^{2} \right)$$

$$= \frac{1}{2} 2 (y^{(i)} - \overrightarrow{D}^{T} \overrightarrow{X}^{(i)}) \frac{\delta}{\delta\Theta_{j}} (y^{(i)} - \overrightarrow{D}^{T} \overrightarrow{X}^{(i)})$$

$$= - (y^{(i)} - \overrightarrow{D}^{T} \overrightarrow{X}^{(i)}) \times_{j} G_{j}$$

$$\frac{\delta J(\Theta)}{\delta\Theta_{j}} = \frac{1}{N} \sum_{i=1}^{N} - (y^{(i)} - \overrightarrow{D}^{T} \overrightarrow{X}^{(i)}) \times_{j} G_{j}$$

$$\frac{\delta J(\Theta)}{\delta\Theta_{j}} = \frac{1}{N} \sum_{i=1}^{N} - (y^{(i)} - \overrightarrow{D}^{T} \overrightarrow{X}^{(i)}) \times_{j} G_{j}$$

Gradient Calculation for Linear Regression

Derivative of $J^{(i)}(\boldsymbol{\theta})$:

$$\begin{split} \frac{d}{d\theta_k} J^{(i)}(\boldsymbol{\theta}) &= \frac{d}{d\theta_k} \frac{1}{2} (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)})^2 \\ &= \frac{1}{2} \frac{d}{d\theta_k} (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)})^2 \\ &= (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) \frac{d}{d\theta_k} (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) \\ &= (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) \frac{d}{d\theta_k} \left(\sum_{j=1}^K \theta_j x_j^{(i)} - y^{(i)} \right) \\ &= (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) x_k^{(i)} \end{split}$$

Derivative of $J(\theta)$:

$$\frac{d}{d\theta_k} J(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N \frac{d}{d\theta_k} J^{(i)}(\boldsymbol{\theta})$$
$$= \frac{1}{N} \sum_{i=1}^N (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) x_k^{(i)}$$

Gradient of
$$J(\boldsymbol{\theta})$$
 [used by Gradient Descent]
$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \begin{bmatrix} \frac{d}{d\theta_1} J(\boldsymbol{\theta}) \\ \frac{d}{d\theta_2} J(\boldsymbol{\theta}) \\ \vdots \\ \frac{d}{d\theta_M} J(\boldsymbol{\theta}) \end{bmatrix} = \begin{bmatrix} \frac{1}{N} \sum_{i=1}^N (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) x_1^{(i)} \\ \frac{1}{N} \sum_{i=1}^N (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) x_2^{(i)} \end{bmatrix} \\ = \frac{1}{N} \sum_{i=1}^N (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) \mathbf{x}_M^{(i)} \end{bmatrix}$$

GD for Linear Regression

Gradient Descent for Linear Regression repeatedly takes steps opposite the gradient of the objective function

 $\theta \leftarrow \theta - \gamma \mathbf{g}$

return θ

6:

Algorithm 1 GD for Linear Regression 1: procedure GDLR(\mathcal{D} , $\theta^{(0)}$) 2: $\theta \leftarrow \theta^{(0)}$ \triangleright Initialize parameters 3: while not converged do 4: $\mathbf{g} \leftarrow \frac{1}{N} \sum_{i=1}^{N} (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) \mathbf{x}^{(i)}$ \triangleright Compute gradient

□ Update parameters