



10-601 Introduction to Machine Learning

Machine Learning Department School of Computer Science Carnegie Mellon University

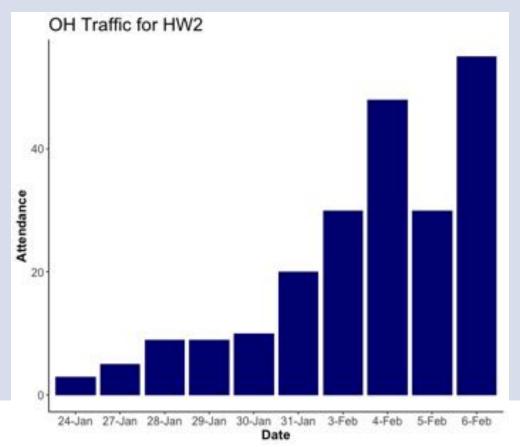
Linear Regression + Optimization for ML

Matt Gormley Lecture 8 Feb. 07, 2020

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!



Reminders

- Homework 3: KNN, Perceptron, Lin.Reg.
 - Out: Wed, Feb. 05 (+ 1 day)
 - Due: Wed, Feb. 12 at 11:59pm
- Today's In-Class Poll
 - http://p8.mlcourse.org

LINEAR REGRESSION

Regression Problems

Chalkboard

- Definition of Regression
- Linear functions
- Residuals
- Notation trick: fold in the intercept

The Big Picture

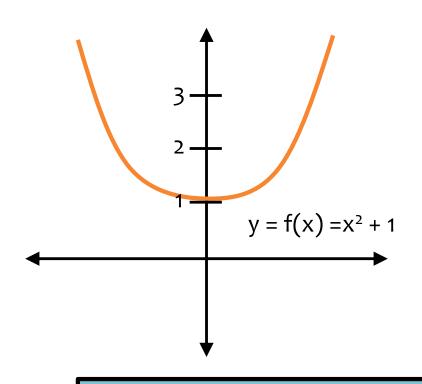
OPTIMIZATION FOR ML

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)

min vs. argmin



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

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

1. Q: What is v*?

v* = 1, the minimum value of the function

2. Q: What is x*?

 $x^* = 0$, the argument that yields the minimum value

Linear Regression as Function $\sum_{\substack{\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^{N} \\ \text{where } \mathbf{x} \in \mathbb{R}^{M} \text{ and } y \in \mathbb{R} } }$ Approximation

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 \mathbf{x} , make prediction \hat{y}

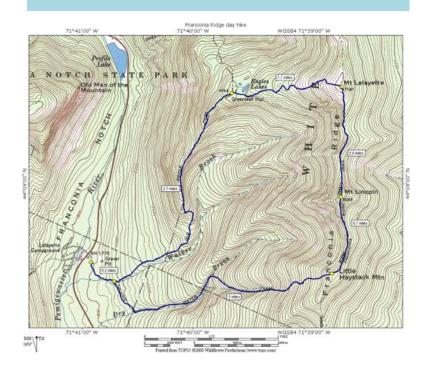
$$\hat{y} = h_{\hat{\boldsymbol{\theta}}}(\mathbf{x}) = \hat{\boldsymbol{\theta}}^T \mathbf{x}$$



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



$$J(\theta) = J(\theta_1, \theta_2) = (10(\theta_1 - 0.5))^2 + (6(\theta_1 - 0.4))^2$$

$$0.8 - \frac{10.000}{0.0}$$

$$0.2 - \frac{10.000}{0.0}$$

$$0.3 - \frac{10.000}{0.0}$$

$$0.4 - \frac{10.000}{0.0}$$

$$0.2 - \frac{10.000}{0.0}$$

$$0.3 - \frac{10.000}{0.0}$$

$$0.4 - \frac{10.000}{0.0}$$

$$0.6 - \frac{10.000}{0.0}$$

$$0.1 - \frac{10.000}{0.0}$$

$$0.2 - \frac{10.000}{0.0}$$

$$0.3 - \frac{10.000}{0.0}$$

$$0.4 - \frac{10.000}{0.0}$$

$$0.6 - \frac{10.000}{0.0}$$

$$0.7 - \frac{10.000}{0.0}$$

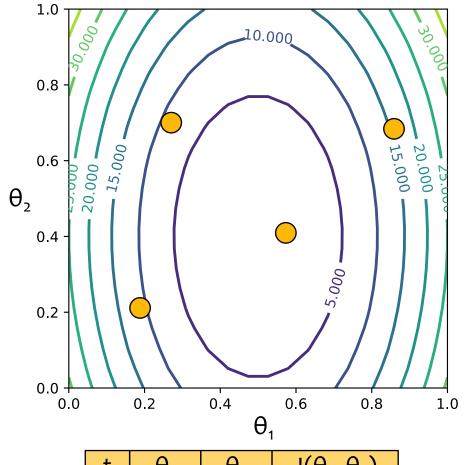
$$0.8 - \frac{10.000}{0.0}$$

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 $J(\theta)$

$$J(\boldsymbol{\theta}) = J(\theta_1, \theta_2) = (10(\theta_1 - 0.5))^2 + (6(\theta_1 - 0.4))^2$$



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	19.2

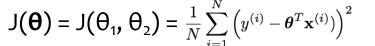
Optimization by Random Guessing

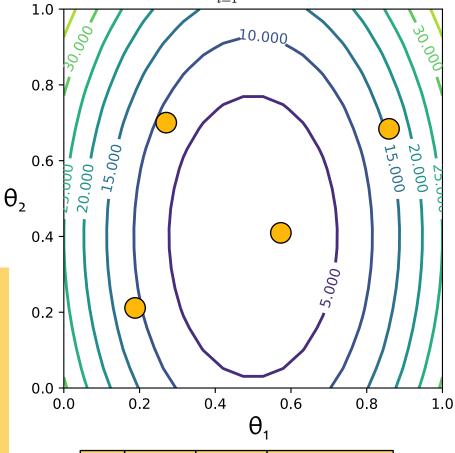
Optimization Method #0: Random Guessing

- 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



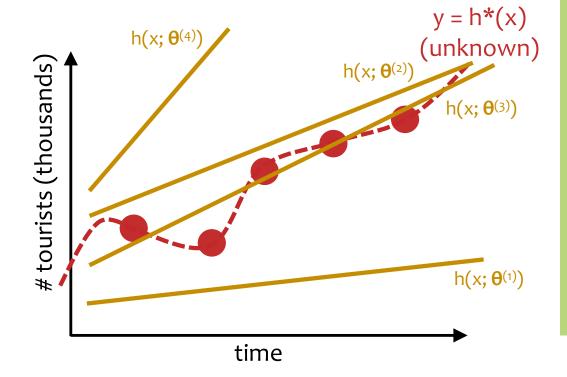


t	θ_1	θ_2	$J(\theta_1, \theta_2)$
1	0.2	0.2	10.4
2	0.3	0.7	7.2
3	0.6	0.4	1.0
4	0.9	0.7	19.2

Linear Regression by Rand. 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 $J(\theta)$



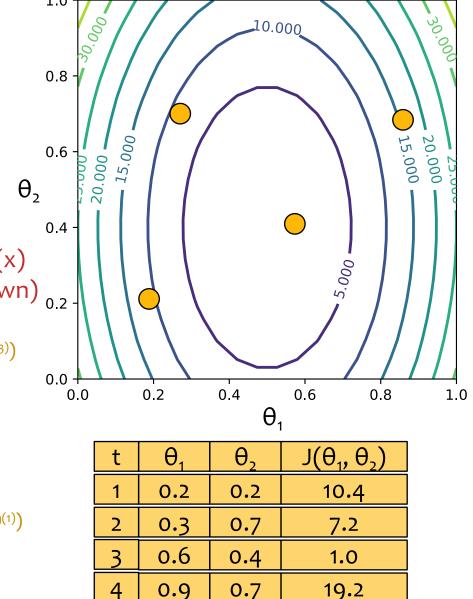
For Linear Regression:

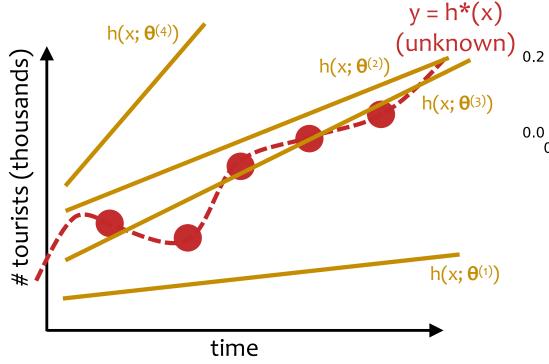
- target function h*(x) is unknown
- only have access to h*(x) through training examples (x⁽ⁱ⁾,y⁽ⁱ⁾)
- want $h(x; \theta^{(t)})$ that **best** approximates $h^*(x)$
- enable generalization w/inductive bias that restricts hypothesis class to linear functions

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

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 $J(\theta)$



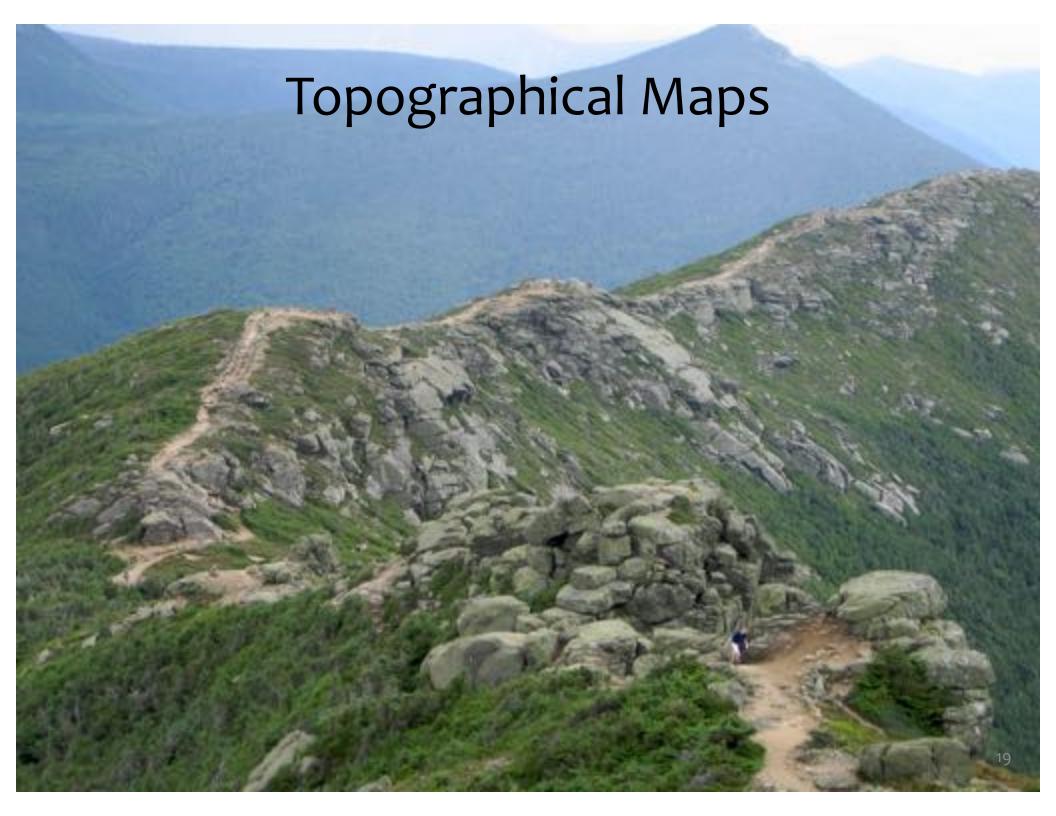


OPTIMIZATION METHOD #1: GRADIENT DESCENT

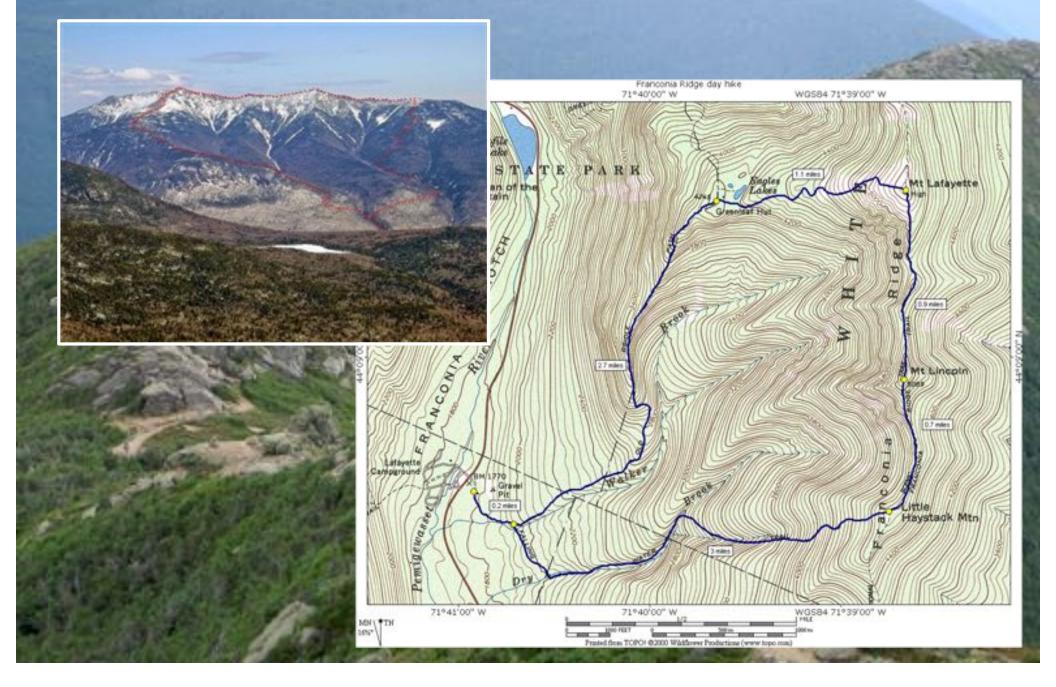
Optimization for ML

Chalkboard

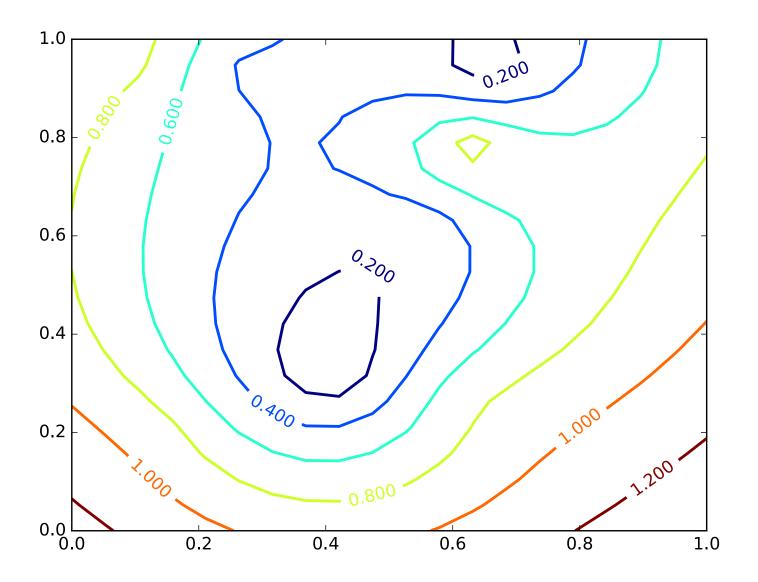
- Unconstrained optimization
- Derivatives
- Gradient



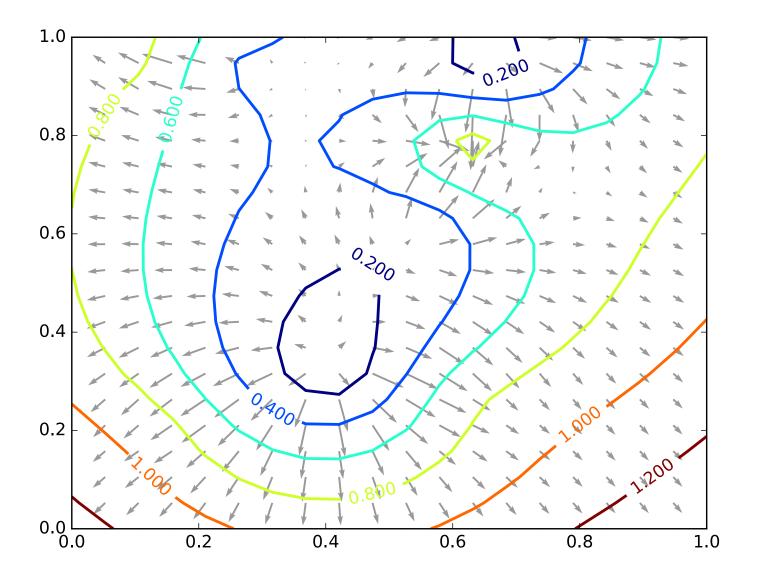
Topographical Maps



Gradients

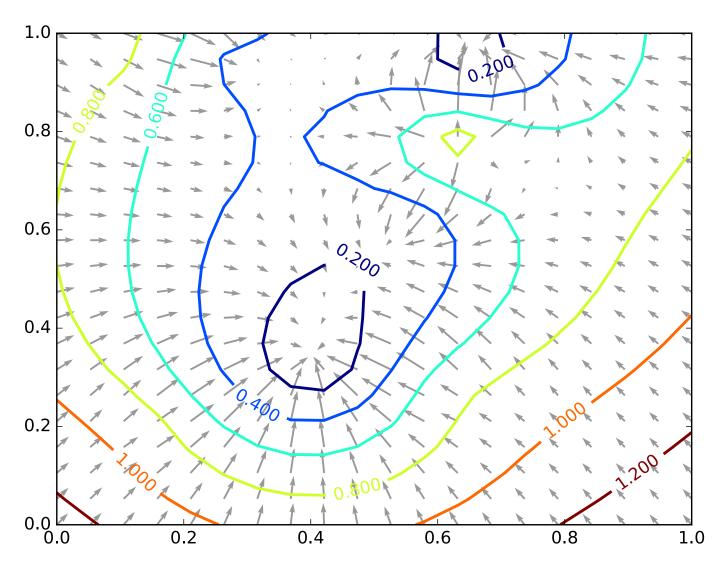


Gradients



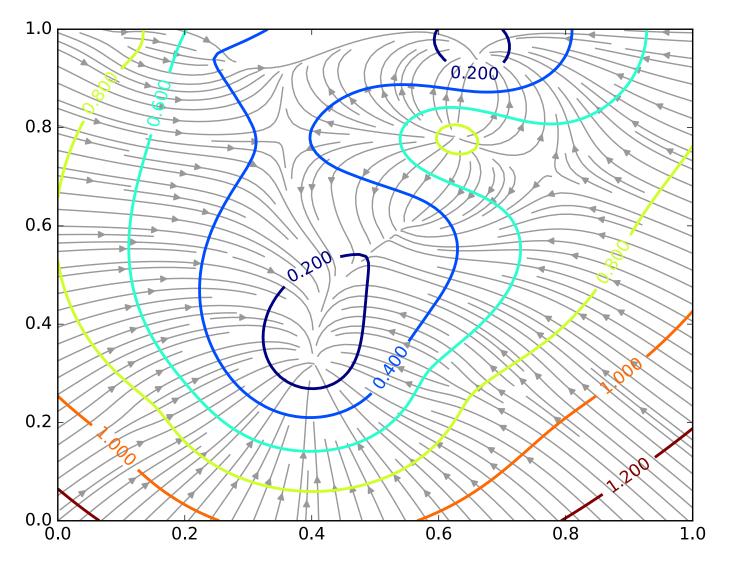
These are the **gradients** that Gradient **Ascent** would follow.

(Negative) Gradients



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

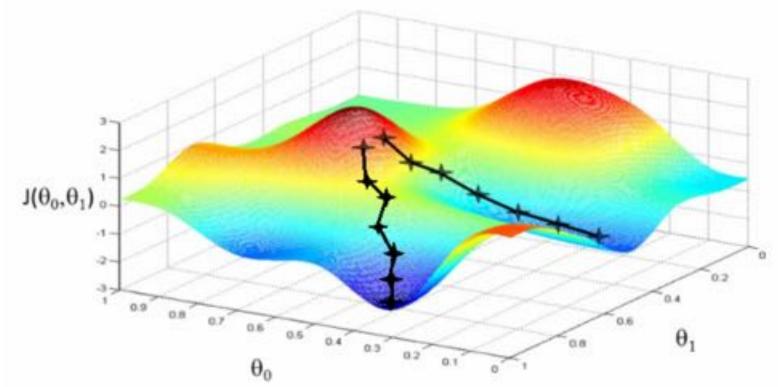
(Negative) Gradient Paths



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

Pros and cons of gradient descent

- Simple and often quite effective on ML tasks
- Often very scalable
- Only applies to smooth functions (differentiable)
- Might find a local minimum, rather than a global one



Gradient Descent

Chalkboard

- Gradient Descent Algorithm
- Details: starting point, stopping criterion, line search

Gradient Descent

Algorithm 1 Gradient Descent

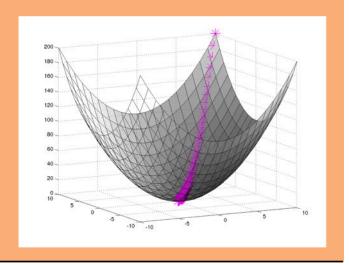
1: **procedure** $GD(\mathcal{D}, \boldsymbol{\theta}^{(0)})$

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).

$$abla_{m{ heta}} J(m{ heta}) = egin{bmatrix} rac{d heta_1}{d heta_2} J(m{ heta}) \ dots \ rac{d}{d heta_M} J(m{ heta}) \end{bmatrix}$$

Gradient Descent

Algorithm 1 Gradient Descent

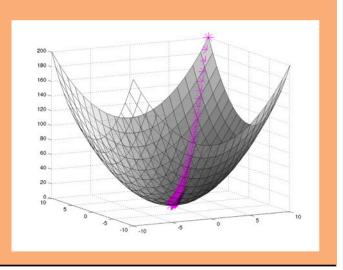
1: **procedure** $GD(\mathcal{D}, \boldsymbol{\theta}^{(0)})$

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

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

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

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
 - ...

$$\hat{\boldsymbol{\theta}} = \operatorname*{argmin}_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$$

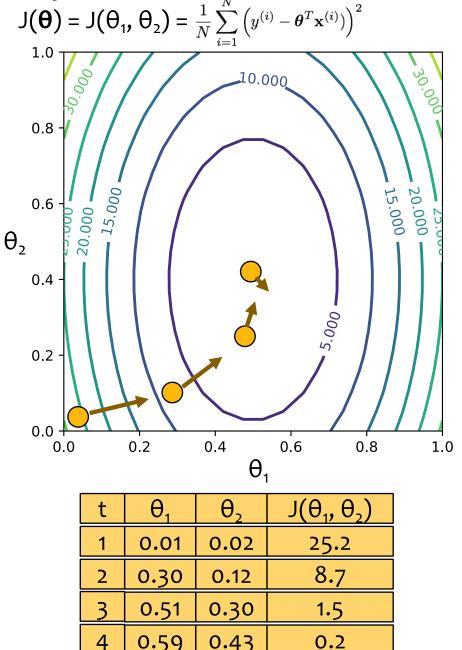
5. Test time: given a new x, make prediction \hat{y}

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

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

Optimization Method #1: Gradient Descent

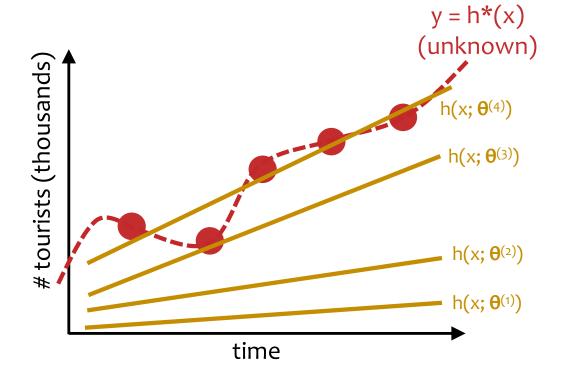
- 1. Pick a random $\boldsymbol{\theta}$
- 2. Repeat:
 - a. Evaluate gradient $\nabla J(\boldsymbol{\theta})$
 - b. Step opposite gradient
- 3. Return θ that gives smallest $J(\theta)$



Linear Regression by Gradient Desc.

Optimization Method #1: Gradient Descent

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

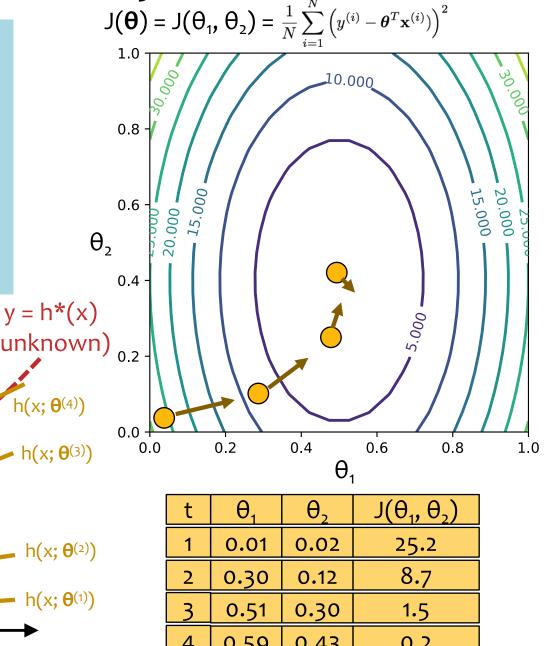


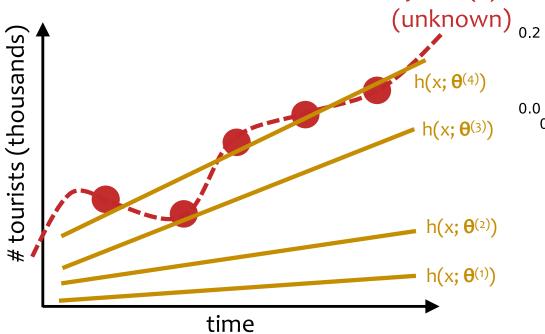
t	θ_1	θ_2	$J(\theta_1, \theta_2)$
1	0.01	0.02	25.2
2	0.30	0.12	8.7
3	0.51	0.30	1.5
4	0.59	0.43	0.2

Linear Regression by Gradient Desc.

Optimization Method #1: Gradient Descent

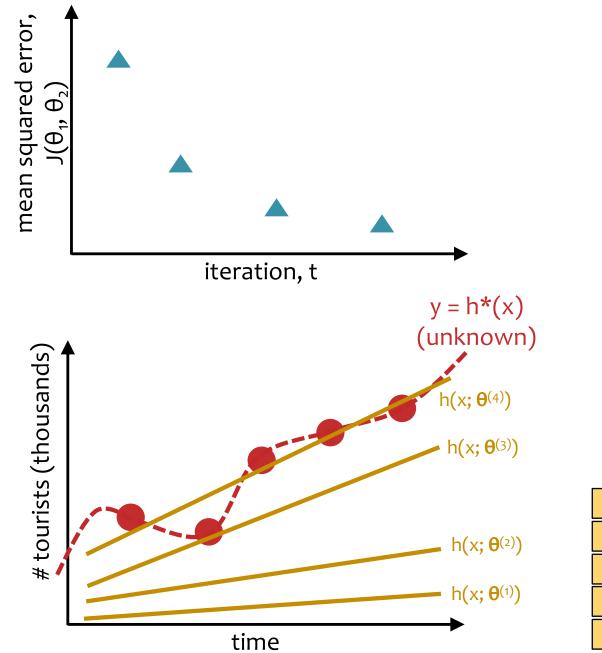
- Pick a random θ
- Repeat: 2.
 - a. Evaluate gradient $\nabla J(\boldsymbol{\theta})$
 - b. Step opposite gradient
- Return θ that gives 3. smallest $J(\theta)$





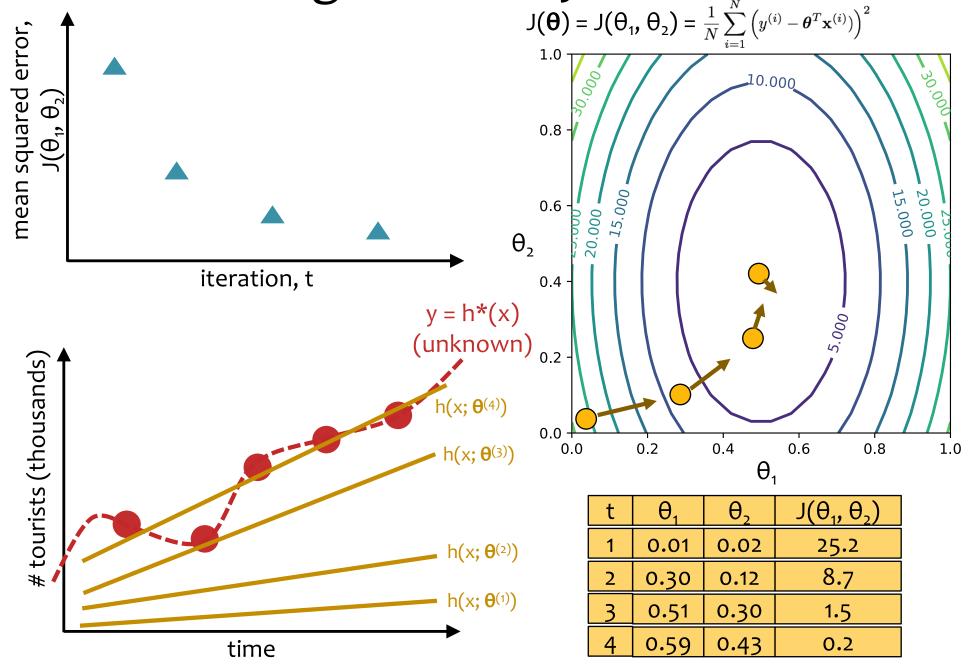
t	θ_1	θ_2	$J(\theta_1, \theta_2)$
1	0.01	0.02	25.2
2	0.30	0.12	8.7
3	0.51	0.30	1.5
4	0.59	0.43	0.2

Linear Regression by Gradient Desc.



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 for Linear Regression

Chalkboard

- Computing the gradient for Linear Regression
- Gradient Descent for Linear Regression

Gradient Calculation for Linear Regression

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

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

Derivative of $J(\theta)$:

$$egin{aligned} rac{d}{d heta_k}J(oldsymbol{ heta}) &= \sum_{i=1}^N rac{d}{d heta_k}J^{(i)}(oldsymbol{ heta}) \ &= \sum_{i=1}^N (oldsymbol{ heta}^T\mathbf{x}^{(i)} - y^{(i)})x_k^{(i)} \end{aligned}$$

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

GD for Linear Regression

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

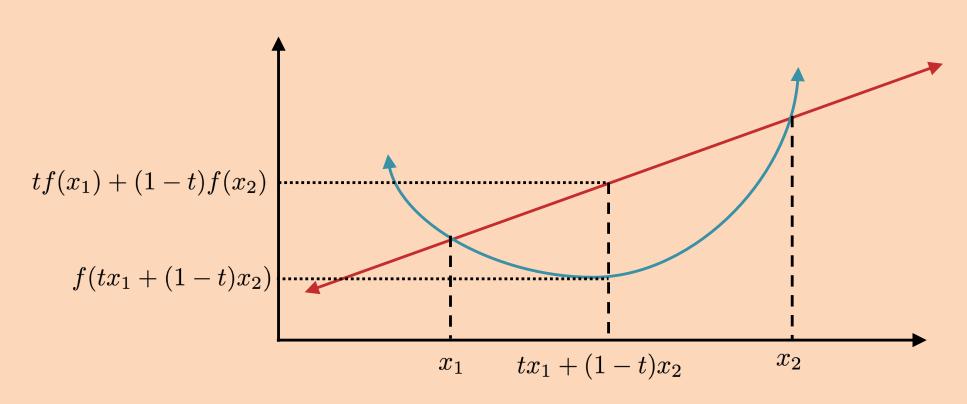
Algorithm 1 GD for Linear Regression 1: procedure GDLR(\mathcal{D} , $\theta^{(0)}$) 2: $\theta \leftarrow \theta^{(0)}$ > Initialize parameters 3: while not converged do 4: $\mathbf{g} \leftarrow \sum_{i=1}^{N} (\theta^T \mathbf{x}^{(i)} - y^{(i)}) \mathbf{x}^{(i)}$ > Compute gradient 5: $\theta \leftarrow \theta - \gamma \mathbf{g}$ > Update parameters 6: return θ

CONVEXITY

Convexity

Function $f: \mathbb{R}^M \to \mathbb{R}$ is **convex** if $\forall \ \mathbf{x}_1 \in \mathbb{R}^M, \mathbf{x}_2 \in \mathbb{R}^M, 0 \leq t \leq 1$:

$$f(t\mathbf{x}_1 + (1-t)\mathbf{x}_2) \le tf(\mathbf{x}_1) + (1-t)f(\mathbf{x}_2)$$

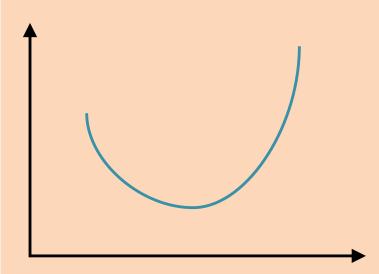


Convexity

Suppose we have a function $f(x): \mathcal{X} \to \mathcal{Y}$.

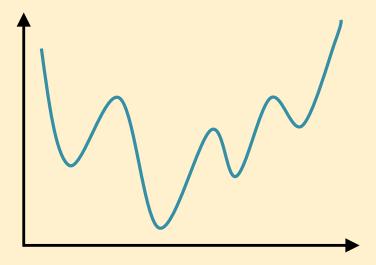
- The value x^* is a **global minimum** of f iff $f(x^*) \leq f(x), \forall x \in \mathcal{X}$.
- The value x^* is a **local minimum** of f iff $\exists \epsilon$ s.t. $f(x^*) \leq f(x), \forall x \in [x^* \epsilon, x^* + \epsilon]$.

Convex Function



Each local minimum is a global minimum

Nonconvex Function

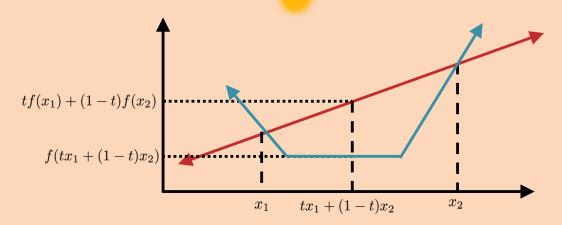


- A nonconvex function is not convex
- Each local minimum is not necessarily a global minimum

Convexity

Function $f: \mathbb{R}^M \to \mathbb{R}$ is **convex** if $\forall \mathbf{x}_1 \in \mathbb{R}^M, \mathbf{x}_2 \in \mathbb{R}^M, 0 \leq t \leq 1$:

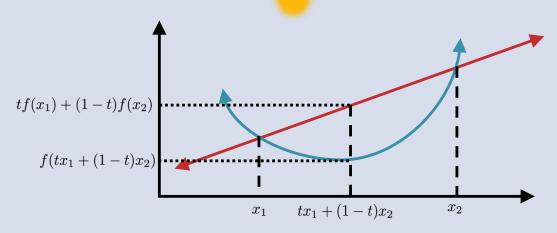
$$f(t\mathbf{x}_1 + (1-t)\mathbf{x}_2) \le tf(\mathbf{x}_1) + (1-t)f(\mathbf{x}_2)$$



Each local
minimum of a
convex function is
also a global
minimum.

Function $f: \mathbb{R}^M \to \mathbb{R}$ is **strictly convex** if $\forall \mathbf{x}_1 \in \mathbb{R}^M, \mathbf{x}_2 \in \mathbb{R}^M, 0 \le t \le 1$:

$$f(t\mathbf{x}_1 + (1-t)\mathbf{x}_2) < tf(\mathbf{x}_1) + (1-t)f(\mathbf{x}_2)$$



A strictly convex function has a unique global minimum.

CONVEXITY AND LINEAR REGRESSION

Convexity and Linear Regression

The Mean Squared Error function, which we minimize for learning the parameters of Linear Regression, is convex!

... but in the general case it is **not** strictly convex.

Regression Loss Functions

In-Class Exercise:

Which of the following could be used as loss functions for training a linear regression model?

Select all that apply.

A.
$$\ell(\hat{y}, y) = ||\hat{y} - y||_2$$

B.
$$\ell(\hat{y}, y) = |\hat{y} - y|$$

C.
$$\ell(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2$$

D.
$$\ell(\hat{y}, y) = \frac{1}{4}(\hat{y} - y)^4$$

$$\text{E. } \ell(\hat{y},y) = \begin{cases} \frac{1}{2}(\hat{y}-y)^2 & \text{if } |\hat{y}-y| \leq \delta \\ \delta|\hat{y}-y| - \frac{1}{2}\delta^2 & \text{otherwise} \end{cases}$$

F.
$$\ell(\hat{y}, y) = \log(\cosh(\hat{y} - y))$$

Solving Linear Regression

Question:

True or False: If Mean Squared Error (i.e. $\frac{1}{N} \sum_{i=1}^{N} (y^{(i)} - h(\mathbf{x}^{(i)}))^2$) has a unique minimizer (i.e. argmin), then Mean Absolute Error (i.e. $\frac{1}{N} \sum_{i=1}^{N} |y^{(i)} - h(\mathbf{x}^{(i)})|$) must also have a unique minimizer.

Answer:

OPTIMIZATION METHOD #2: CLOSED FORM SOLUTION

Calculus and Optimization

In-Class Exercise

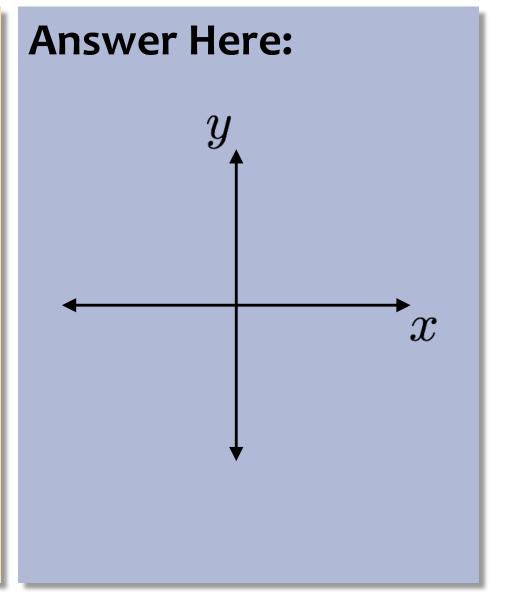
Plot three functions:

1.
$$f(x) = x^3 - x$$

2.
$$f'(x) = \frac{\partial y}{\partial x}$$

2.
$$f'(x)=rac{\partial y}{\partial x}$$

3. $f''(x)=rac{\partial^2 y}{\partial x^2}$



Optimization: Closed form solutions

Chalkboard

- Zero Derivatives
- Example: 1-D function
- Example: higher dimensions

CLOSED FORM SOLUTION FOR LINEAR REGRESSION

Linear Regression as Function $\sum_{\substack{\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^{N} \\ \text{where } \mathbf{x} \in \mathbb{R}^{M} \text{ and } y \in \mathbb{R} } }$ Approximation

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

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

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
 - ...

$$\hat{\boldsymbol{\theta}} = \operatorname*{argmin}_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$$

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

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

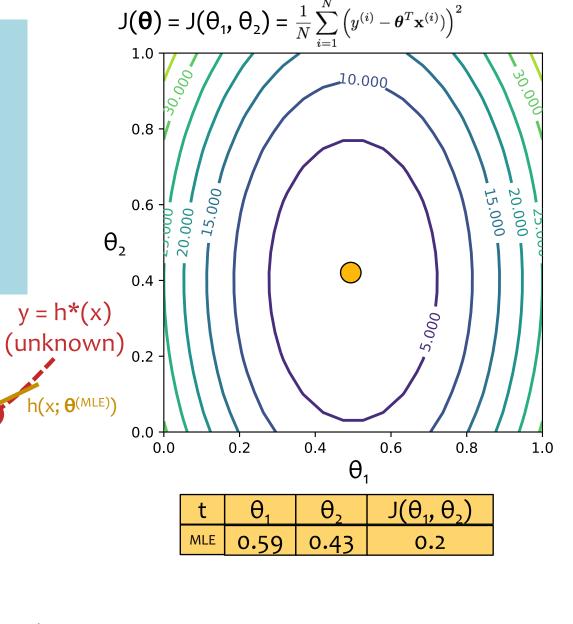
Linear Regression: Closed Form

Optimization Method #2: Closed Form

1. Evaluate

$$\boldsymbol{\theta}^{\mathsf{MLE}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

2. Return θ^{MLE}



Optimization for Linear Regression

Chalkboard

– Closed-form (Normal Equations)

Computational Complexity of OLS

To solve the Ordinary Least Squares problem we compute:

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} = \frac{1}{N} \sum_{i=1}^{N} \frac{1}{2} (y^{(i)} - (\boldsymbol{\theta}^{T} \mathbf{x}^{(i)}))^{2}$$
$$= (\mathbf{X}^{T} \mathbf{X})^{-1} (\mathbf{X}^{T} \mathbf{Y})$$

The resulting shape of the matrices:

$$(\mathbf{X}^{T} \mathbf{X})^{-1} (\mathbf{X}^{T} \mathbf{Y})$$

$$M \times N N \times M$$

$$M \times N N \times 1$$

$$M \times M$$

$$M \times 1$$

Background: Matrix Multiplication Given matrices **A** and **B**

- If **A** is $q \times r$ and **B** is $r \times s$, computing **AB** takes O(qrs)
- If **A** and **B** are $q \times q$, computing **AB** takes $O(q^{2.373})$
- If **A** is $q \times q$, computing A^{-1} takes $O(q^{2.373})$.

Computational Complexity of OLS:

$$\begin{array}{cccc} \mathbf{X}^T\mathbf{X} & O(M^2N) \\ (&)^{-1} & O(M^{2.373}) \\ & \mathbf{X}^T\mathbf{Y} & O(MN) \\ (&)^{-1}(&) & O(M^2) \\ \hline \text{total} & O(M^2N+M^{2.373}) \end{array}$$

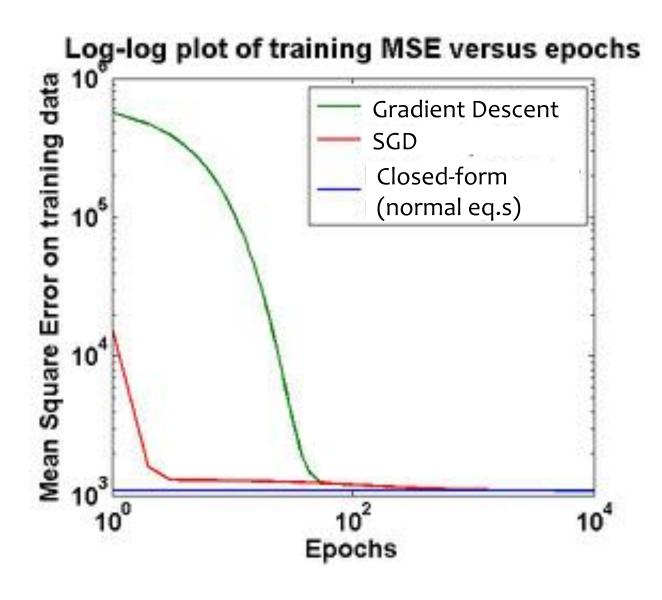
Linear in # of examples, N
Polynomial in # of features, M

Gradient Descent

Cases to consider gradient descent:

- 1. What if we can not find a closed-form solution?
- 2. What if we **can**, but it's inefficient to compute?
- 3. What if we can, but it's numerically unstable to compute?

Convergence Curves



- Def: an epoch is a single pass through the training data
- For GD, only one update per epoch
- For SGD, N updatesper epochN = (# train examples)
- SGD reduces MSE much more rapidly than GD
- For GD / SGD, training MSE is initially large due to uninformed initialization