

Logistic Regression

Sriram Sankararaman

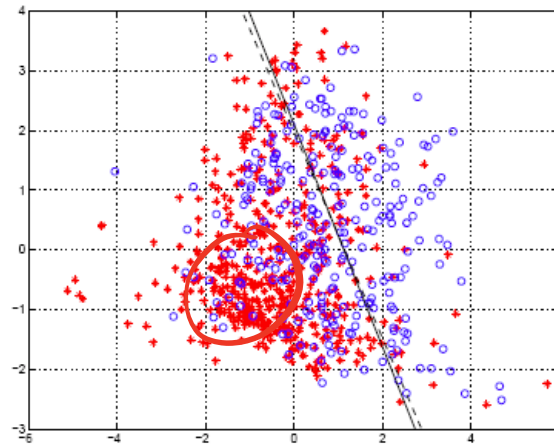
The instructor gratefully acknowledges Fei Sha, Ameet Talwalkar, Eric Eaton, and Jessica Wu whose slides are heavily used, and the many others who made their course material freely available online.

Outline

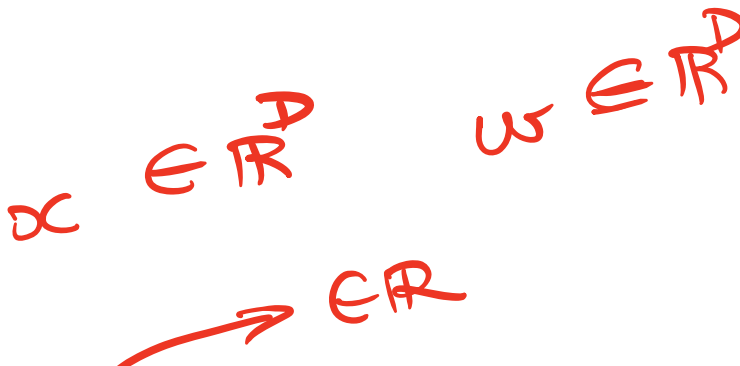
- 1 Logistic regression
 - Prediction in logistic regression
 - Training/Learning a logistic regression model
 - Maximum likelihood estimation
- 2 Optimization
- 3 Stochastic gradient descent

Review

- Instead of predicting the class, predict the probability of instance being in a class
- Perceptron does not produce probability estimates



Prediction in logistic regression


$$P(y = 1|x) = \underbrace{\sigma}_{0 \leq \leq 1}(\underbrace{w^T x + b}_{\in \mathbb{R}}) = \frac{1}{1 + \exp^{-(w^T x + b)}}$$

Compute $\sigma(w^T x + b)$. If this is greater than 0.5, predict 1 else 0.

Decision boundary: Linear or nonlinear?

$$\sigma(a) = \frac{1}{1 + e^{-a}}$$

$a = 0$

$$\sigma(\underline{w^T x + b}) \quad ?$$

$$\begin{aligned} &> 0.5 \Rightarrow 1 \\ &= 0.5 \\ &< 0.5 \Rightarrow 0 \end{aligned}$$

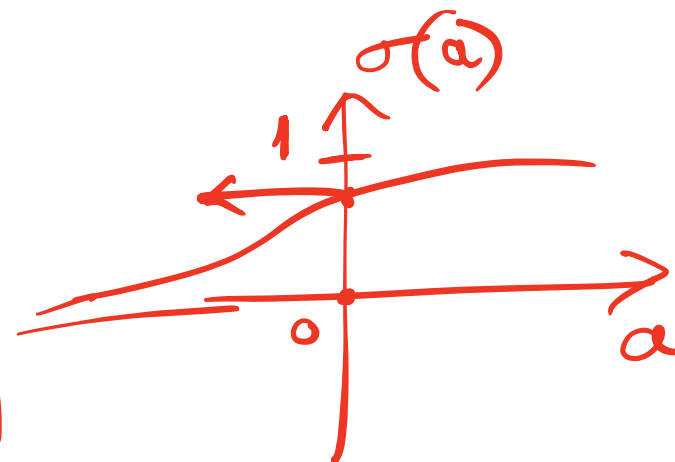
Perceptron

$$\text{sign}(w^T x + b) = +1$$

$$w^T x + b = 0$$

$$-1$$

$$+1$$



$$\underline{\sigma(w^T x + b) = 0.5}$$

$$\underline{w^T x + b = 0}$$

Decision boundary: Linear or nonlinear?

$\sigma(a)$ is **nonlinear**, however, the decision boundary is determined by

$$\sigma(a) = 0.5 \Rightarrow a(\mathbf{x}) = 0 \Rightarrow a(\mathbf{x}) = \underline{\underline{b}} + \underline{\underline{\mathbf{w}^T}} \mathbf{x} = 0$$

which is a **linear** function in \mathbf{x}

As in the case of perceptron, b the **bias or offset or intercept term**.

\mathbf{w} the **weights**.

Logistic regression

Setup for binary classification

- Input: $\underline{x} \in \mathbb{R}^D$
- Output: $\underline{y} \in \{0, 1\}$
- Training data: $\underline{\mathcal{D}} = \{(\underline{x}_n, y_n), n = 1, 2, \dots, N\}$
- Hypotheses/Model:

$$h_{\underline{w}, b}(\underline{x}) = p(y = 1 | \underline{x}; b, \underline{w}) = \sigma(a(\underline{x}))$$

where

Activation
fn.

$$a(\underline{x}) = b + \sum_d w_d x_d = b + \underline{w}^T \underline{x}$$

- Given training data N samples/instances:
 $\mathcal{D}^{\text{TRAIN}} = \{(\underline{x}_1, y_1), (\underline{x}_2, y_2), \dots, (\underline{x}_N, y_N)\}$, train/learn/induce $\underline{h}_{\underline{w}, b}$.
Find values for (\underline{w}, b) .

Example: bag of words

$$P(aDaaa | \text{Bag 1}) \gg P(aDaaa | \text{Bag 2})$$

Which bag of words is more likely to generate : aDaaa ?



1

2



2

?

Example: bag of words

Which bag of words is more likely to generate : aDaaa ?

$$\underline{0.7} \times \underline{0.1} \times 0.7 \times 0.7 \times 0.7 \\ = 2.401 \times 10^{-2}$$



$$0.2 \times 0.1 \times \underline{0.2} \times \underline{0.2} \times \underline{0.2} \\ = 1.6 \times 10^{-4}$$



Example: drawing color cards from an envelope

- An envelope with two colors of cards: yellow and purple.
- Assume if you draw a card at random, probability of drawing a yellow card is θ and probability of drawing a purple card is $1 - \theta$.
- Sample with replacement n times.
- k times we get yellow, $n - k$ times we get purple.
- The joint probability (likelihood) : $\theta^k (1 - \theta)^{n-k}$.
- What is the value of θ that maximizes the joint probability?

$P(k \text{ yellow \& } n-k \text{ purple})$

$$= \underbrace{\theta \theta \dots \theta}_k \underbrace{(1-\theta) \dots (1-\theta)}_{(n-k)}$$

$$= P(\text{yellow in 1st draw}) \times P(\text{yellow in 2nd draw}) \times \dots \times P(\text{yellow in } k\text{th draw}) \times P(\text{purple in } (k+1)\text{th draw}) \times \dots \times P(\text{purple in } n\text{th draw})$$

Example: drawing color cards from an envelope

Solve $\text{argmax}_{\theta} \theta^k (1 - \theta)^{n-k}$

Equivalently, we can solve:

$$\begin{aligned} \text{argmax}_{\theta} \log(\theta^k (1 - \theta)^{n-k}) &= \log(\theta^k) + \log((1 - \theta)^{n-k}) \\ &= \text{argmax}_{\theta} k \log \theta + (n - k) \log(1 - \theta) \\ &\quad \downarrow \\ \frac{d}{d\theta} &= 0 \end{aligned}$$


Example: drawing color cards from an envelope

Solve $\operatorname{argmax}_{\theta} \theta^k (1 - \theta)^{n-k}$

Equivalently, we can solve:

$$\begin{aligned} \operatorname{argmax}_{\theta} \log(\theta^k (1 - \theta)^{n-k}) \\ = \operatorname{argmax}_{\theta} k \log \theta + (n - k) \log(1 - \theta) \end{aligned}$$

At the optimum: $\frac{(k \log \theta + (n - k) \log(1 - \theta))}{d\theta} = 0$



Example: drawing color cards from an envelope

Maximum likelihood estimate (MLE): $\hat{\theta} = \frac{k}{n}$

These are easy examples. We don't always have a closed-form solution for the MLE typically !

Likelihood Function

Let X_1, \dots, X_N be IID (independent and identically distributed) random variables with PDF $p(x|\theta)$ (also written as $p(x; \theta)$). The *likelihood function* is defined by $L(\theta)$,

$$L(\theta) = p(X_1, \dots, X_N; \theta).$$

$$= \prod_{i=1}^N p(X_i; \theta).$$

Notes The likelihood function is just the joint density of the data, except that we treat it as a function of the parameter θ .

Maximum Likelihood Estimator

Definition: The maximum likelihood estimator (MLE) $\hat{\theta}$, is the value of θ that maximizes $L(\theta)$.

The log-likelihood function is defined by $l(\theta) = \log L(\theta)$. Its maximum occurs at the same place as that of the likelihood function.

Maximum Likelihood Estimator

$$L(\theta) = \prod_{i=1}^n P(x_i; \theta)$$

$$\begin{aligned} \ell(\theta) &= \log(L(\theta)) = \log\left(\prod_{i=1}^n P(x_i; \theta)\right) \\ &= \sum_{i=1}^n \log P(x_i; \theta) \end{aligned}$$

Definition: The maximum likelihood estimator (MLE) $\hat{\theta}$, is the value of θ that maximizes $L(\theta)$.

The log-likelihood function is defined by $\ell(\theta) = \log L(\theta)$. Its maximum occurs at the same place as that of the likelihood function.

- Using logs simplifies mathematical expressions (converts exponents to products and products to sums)
- Using logs helps with numerical stability

The same is true of the likelihood function times any constant. Thus we shall often drop constants in the likelihood function.

$$\underline{c L(\theta)}$$

$$L(\theta)$$

Likelihood function for logistic regression

$$P(\gamma=1 | x; b, w) = \sigma(w^T x + b)$$
$$P(\gamma=0 | x; b, w) = 1 \Rightarrow P(\gamma=0 | x; b, w) = 1 - \sigma(w^T x + b)$$

Probability of a single training sample (x_n, y_n)

$$p(\underline{y_n} | \underline{x_n}; b, w) = \begin{cases} \underline{h_{w,b}(x_n)} = \sigma(b + w^T x_n) & \text{if } y_n = 1 \\ \underline{= 1 - h_{w,b}(x_n) = 1 - \sigma(b + w^T x_n)} & \text{otherwise} \end{cases}$$

Likelihood function for logistic regression

Probability of a single training sample (\mathbf{x}_n, y_n)

$$p(y_n | \mathbf{x}_n; b, \mathbf{w}) = \begin{cases} h_{\mathbf{w},b}(\mathbf{x}_n) = \sigma(b + \mathbf{w}^T \mathbf{x}_n) & \text{if } y_n = 1 \\ 1 - h_{\mathbf{w},b}(\mathbf{x}_n) = 1 - \sigma(b + \mathbf{w}^T \mathbf{x}_n) & \text{otherwise} \end{cases}$$

Compact expression, exploring that y_n is either 1 or 0

$$p(y_n | \mathbf{x}_n; b; \mathbf{w}) = \underline{h_{\mathbf{w},b}(\mathbf{x}_n)}^{y_n} \underline{[1 - h_{\mathbf{w},b}(\mathbf{x}_n)]}^{1-y_n}$$

$y_n = 1$

$$\begin{array}{cc} \downarrow & \downarrow \\ (h_{\mathbf{w},b}(\mathbf{x}_n))^1 & (1 - h_{\mathbf{w},b}(\mathbf{x}_n))^0 \\ \parallel & \\ h_{\mathbf{w},b}(\mathbf{x}_n) \cdot 1 & \end{array}$$

Log Likelihood

$$\log(L(w, b)) = \log\left(\prod_{i=1}^n P(y_n | x_n, w, b)\right)$$

Log-likelihood of the whole training data \mathcal{D}

$$l(w, b) = \sum_n \{y_n \log h_{w,b}(x_n) + (1 - y_n) \log[1 - h_{w,b}(x_n)]\}$$

$$\underline{l(w, b)} = \sum_{n=1}^N \log \underline{P(y_n | x_n; w, b)}$$

$$= \sum_{n=1}^N \log \left(\overset{y_n}{h_{w,b}(x_n)} \left(1 - h_{w,b}(x_n)\right)^{\overset{(1-y_n)}{}} \right)$$

$$= \sum_{n=1}^N \log \left(\underset{y_n}{h_{w,b}(x_n)} \right) + \log \left(\left(1 - h_{w,b}(x_n)\right)^{\overset{(1-y_n)}{}} \right)$$

$$\log(a^b) = b \log(a)$$

Log Likelihood

Log-likelihood of the whole training data \mathcal{D}

$$\curvearrowright \underline{l(\mathbf{w}, b)} = \sum_n \{y_n \log h_{\mathbf{w}, b}(\mathbf{x}_n) + (1 - y_n) \log[1 - h_{\mathbf{w}, b}(\mathbf{x}_n)]\}$$

It is convenient to work with its negation termed negative log likelihood

$$\underline{J(b, \mathbf{w})} = - \sum_n \{y_n \log h_{\mathbf{w}, b}(\mathbf{x}_n) + (1 - y_n) \log[1 - h_{\mathbf{w}, b}(\mathbf{x}_n)]\}$$

We can ignore the distinction between bias and weights

(b, w)

This is for convenience

- Append 1 to x

$$\underline{x} \leftarrow \underline{[1 \quad x_1 \quad x_2 \quad \cdots \quad x_D]}$$

- Append b to w

$$\underline{\theta} \leftarrow \underline{[b \quad w_1 \quad w_2 \quad \cdots \quad w_D]}$$

-

$$\underline{J(\theta)} = - \sum_n \{y_n \log \underline{h_{\theta}}(\underline{x_n}) + (1 - y_n) \log[1 - \underline{h_{\theta}}(\underline{x_n})]\}$$

- Same trick as in the case of perceptrons
 - ▶ we are rewriting a hyperplane in D dimensions as one in $D + 1$ dimensions that passes through the origin.

How to find the optimal parameters for logistic regression?

$$\begin{aligned} \overset{\text{MLE}}{\theta^*} &= \underset{\theta}{\operatorname{argmax}} \ell(\theta) \\ &= \underset{\theta}{\operatorname{argmin}} -\ell(\theta) \end{aligned}$$

We will minimize the negative log likelihood

$$J(\theta) = - \sum_n \{y_n \log h_{\theta}(\mathbf{x}_n) + (1 - y_n) \log[1 - h_{\theta}(\mathbf{x}_n)]\}$$

How to find the optimal parameters for logistic regression?

We will minimize the negative log likelihood

$$J(\boldsymbol{\theta}) = - \sum_n \{y_n \log h_{\boldsymbol{\theta}}(\mathbf{x}_n) + (1 - y_n) \log[1 - h_{\boldsymbol{\theta}}(\mathbf{x}_n)]\}$$

How do we find its minimum?

$$\boldsymbol{\theta} = (\underbrace{\theta_0, \theta_1, \dots, \theta_D}_{(D+1)})$$
$$\nabla J(\boldsymbol{\theta}) = 0$$
$$\nabla J(\boldsymbol{\theta}) = \begin{pmatrix} \frac{\partial J}{\partial \theta_0} \\ \frac{\partial J}{\partial \theta_1} \\ \vdots \\ \frac{\partial J}{\partial \theta_D} \end{pmatrix} \Rightarrow$$

Outline

- 1 Logistic regression
- 2 Optimization
- 3 Stochastic gradient descent

Maximum Likelihood

$\hat{\theta}$

X_1, \dots, X_n

$X_i \sim P(x|\theta)$

$P(X_i = \text{Yellow}) = 0.2$

$P(X_i = \text{Purple}) = 0.8$

$\theta = 0.1$ (true)

$P(X_1, \dots, X_n | \theta = 0.1)$

$P(X_1, \dots, X_n | \theta = 0.8)$

$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} L(\theta) = \log P(X_1, \dots, X_n | \theta)$

Optimization

Given a function $f(\underline{x})$, find its minimum (or maximum).

- f is called the **objective function**.

Optimization

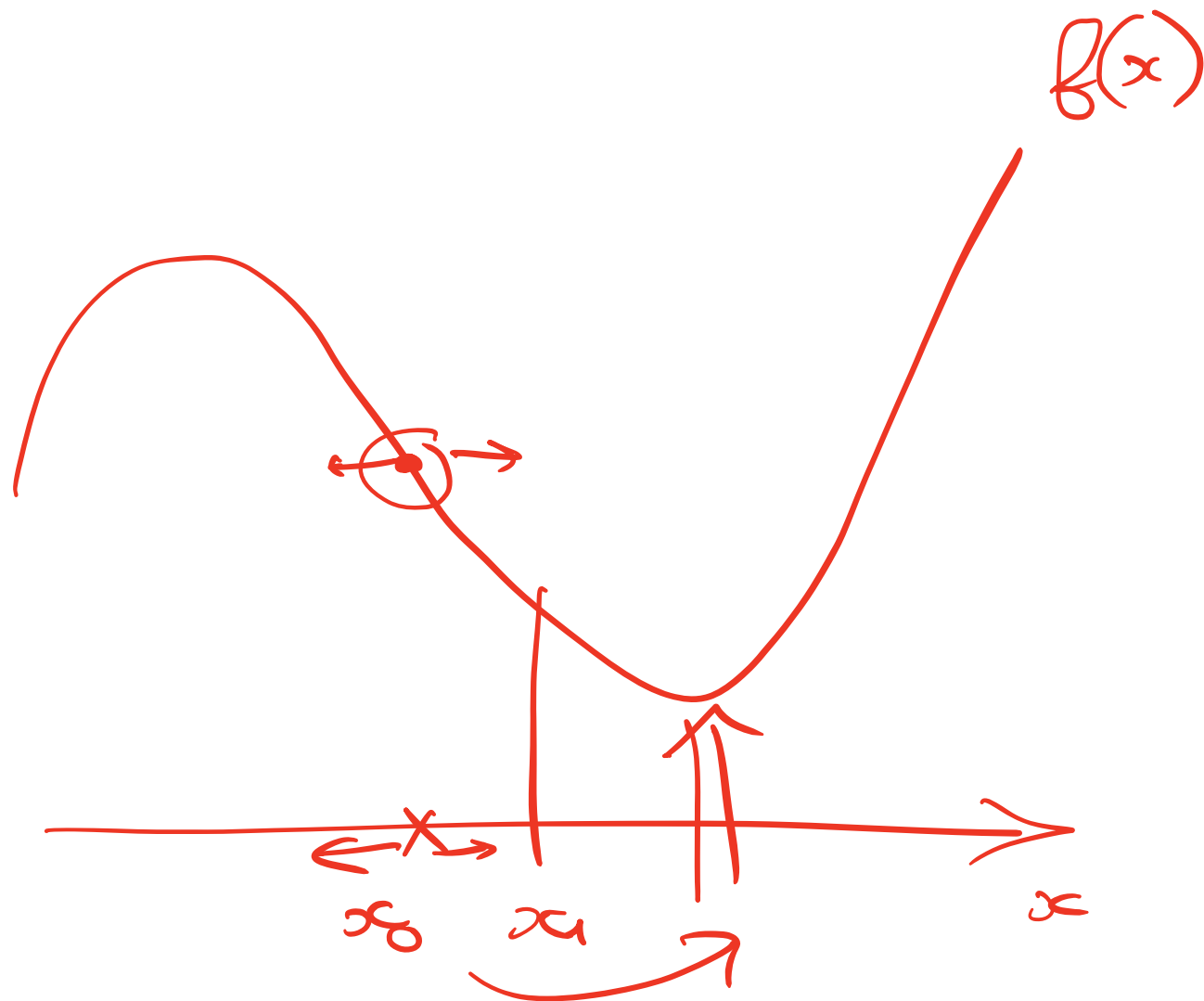
Given a function $f(x)$, find its minimum (or maximum).

- f is called the **objective function**.
- Maximizing f is equivalent to minimizing $-f$.

So we only need to consider minimization problems.

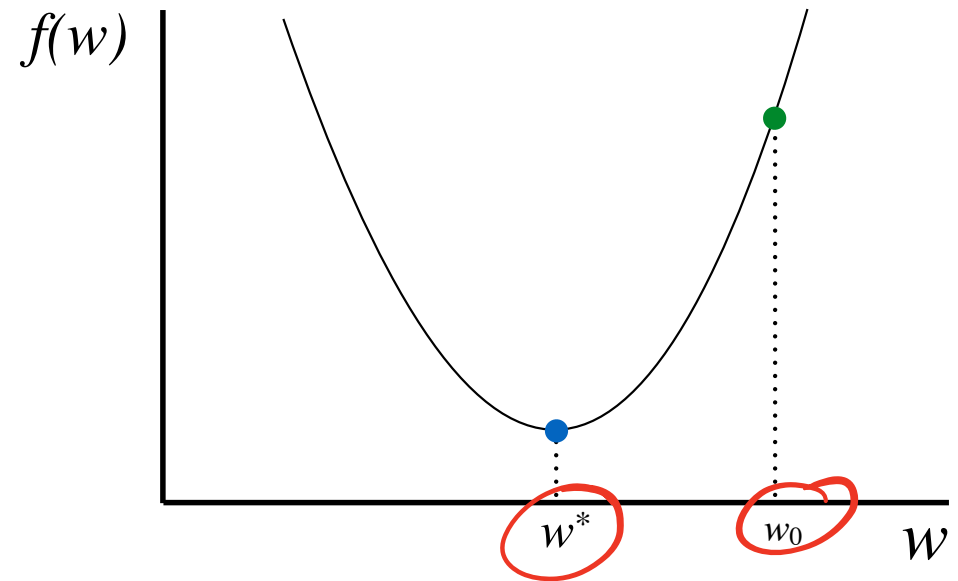
One way to minimize a function f

Gradient descent



Gradient Descent

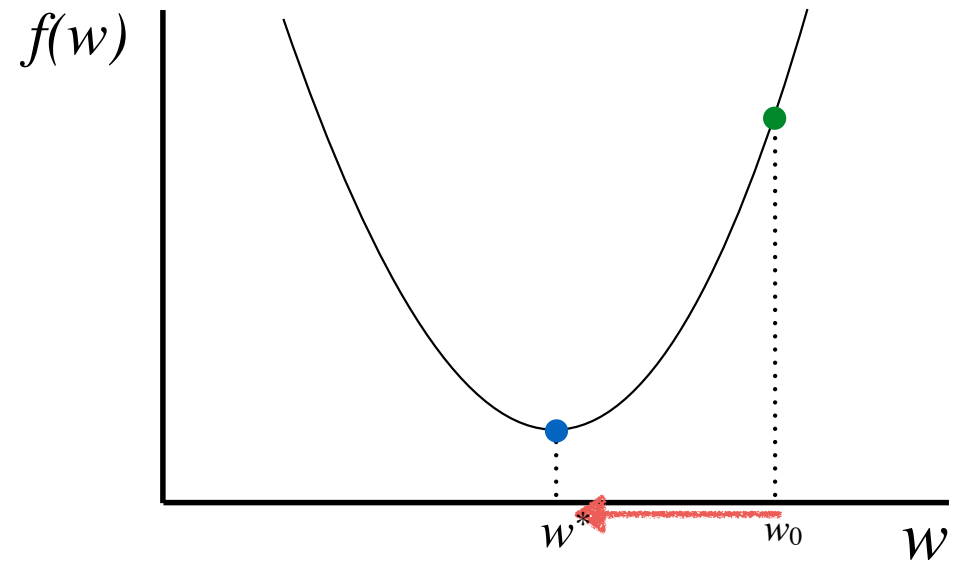
Start at a random point



Gradient Descent

Start at a random point

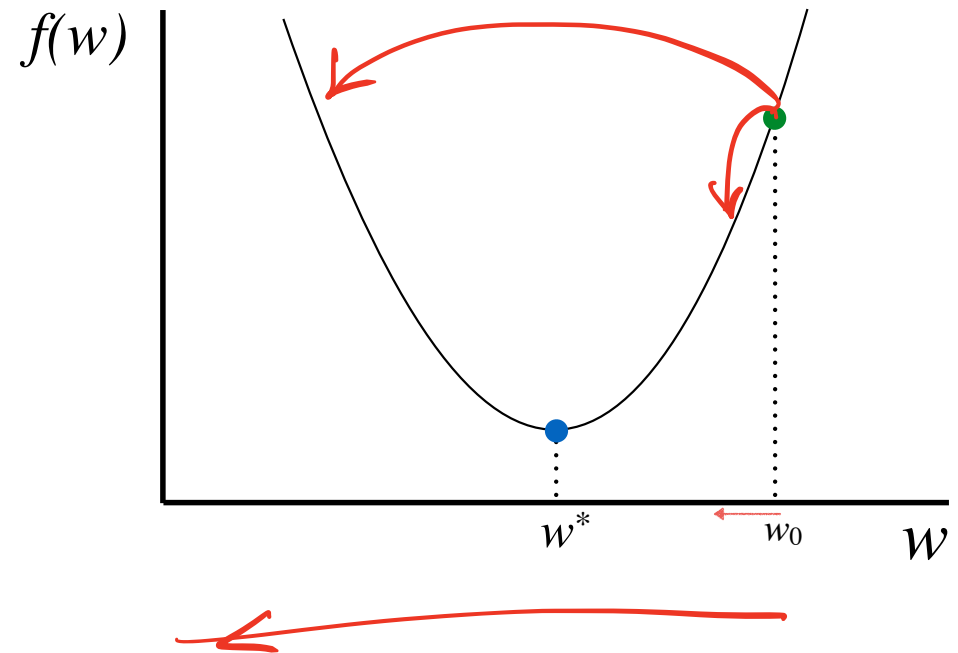
Determine a descent direction



Gradient Descent

Start at a random point

Determine a descent direction
Choose a step size



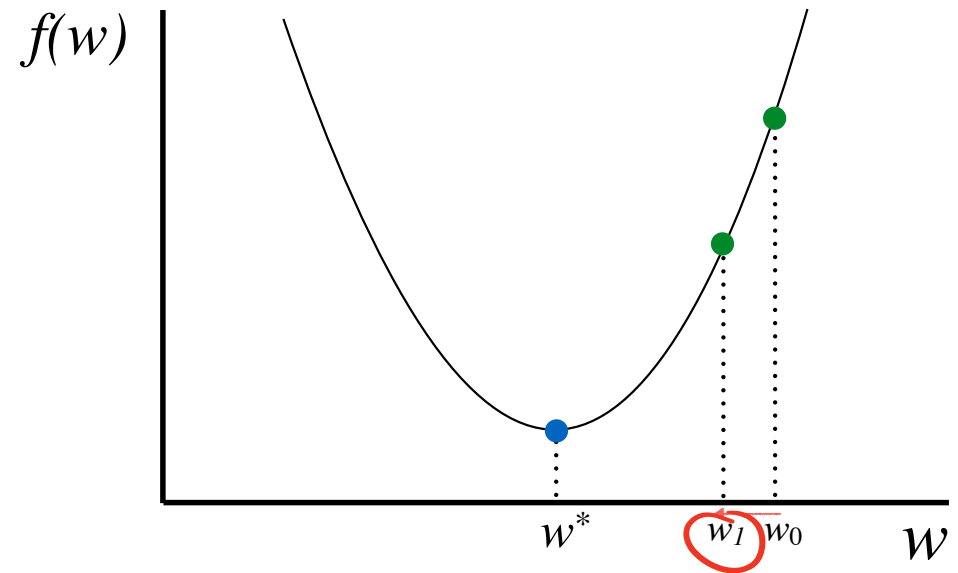
Gradient Descent

Start at a random point

Determine a descent direction

Choose a step size

Update



Gradient Descent

Start at a random point

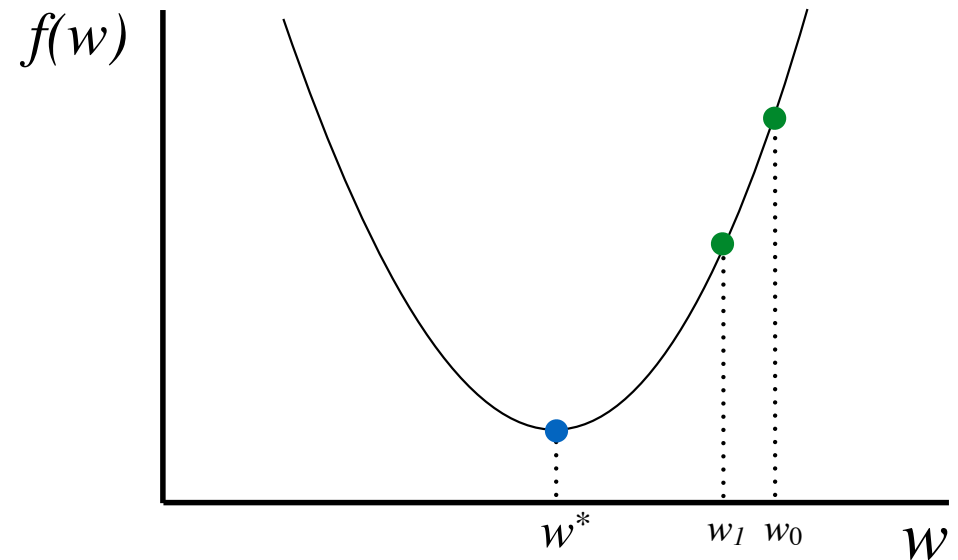
Repeat

Determine a descent direction

Choose a step size

Update

Until stopping criterion is satisfied



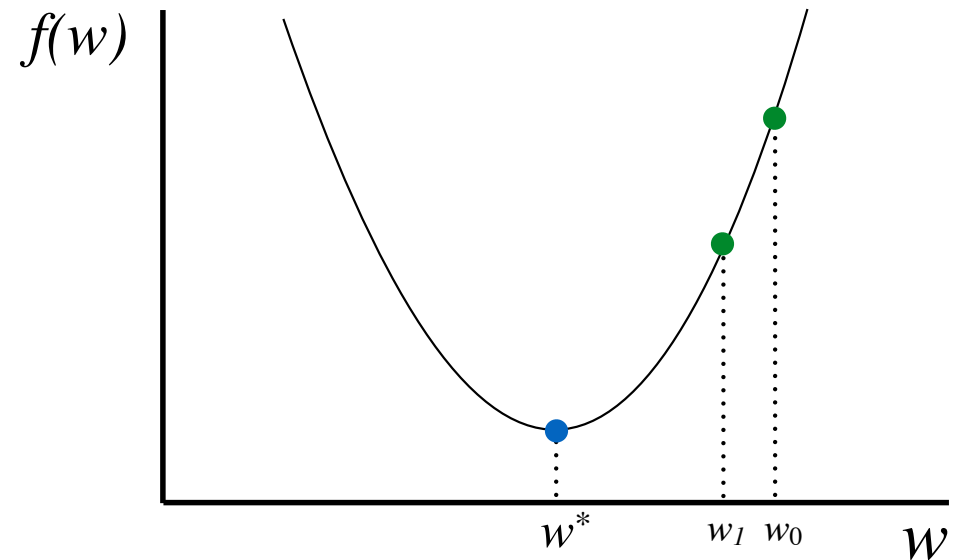
Gradient Descent

Start at a random point

Repeat

- Determine a descent direction
- Choose a step size
- Update

Until stopping criterion is satisfied



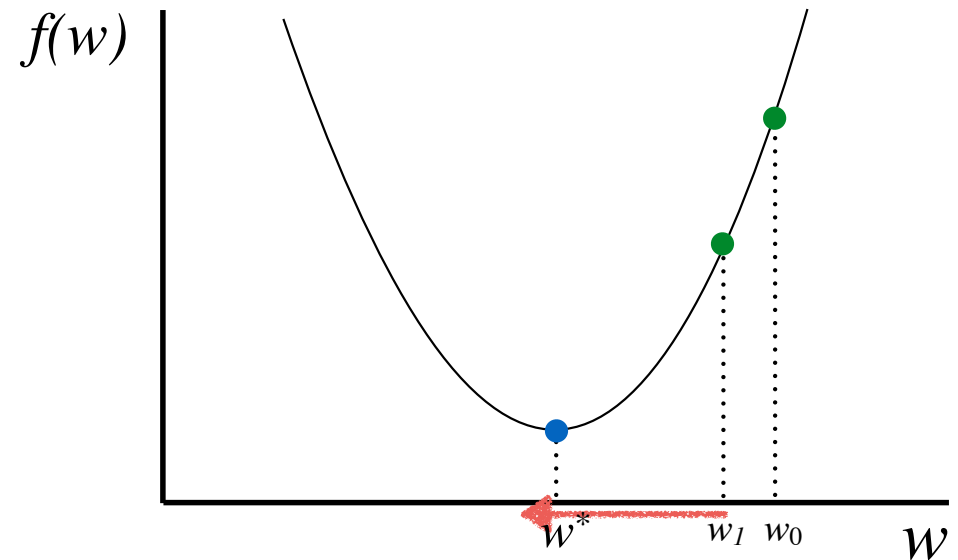
Gradient Descent

Start at a random point

Repeat

- I Determine a descent direction
- Choose a step size
- Update

Until stopping criterion is satisfied



Gradient Descent

Start at a random point

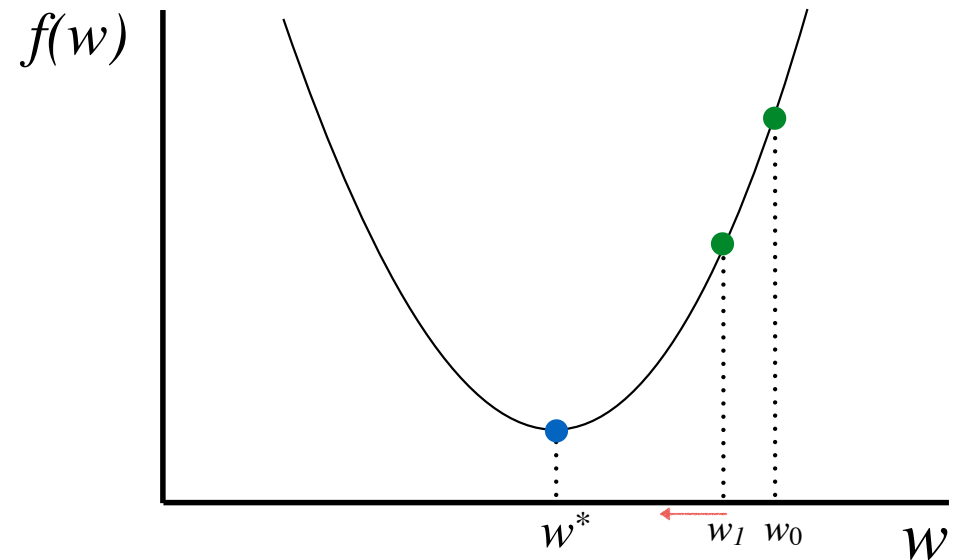
Repeat

Determine a descent direction

■ Choose a step size

Update

Until stopping criterion is satisfied



Gradient Descent

Start at a random point

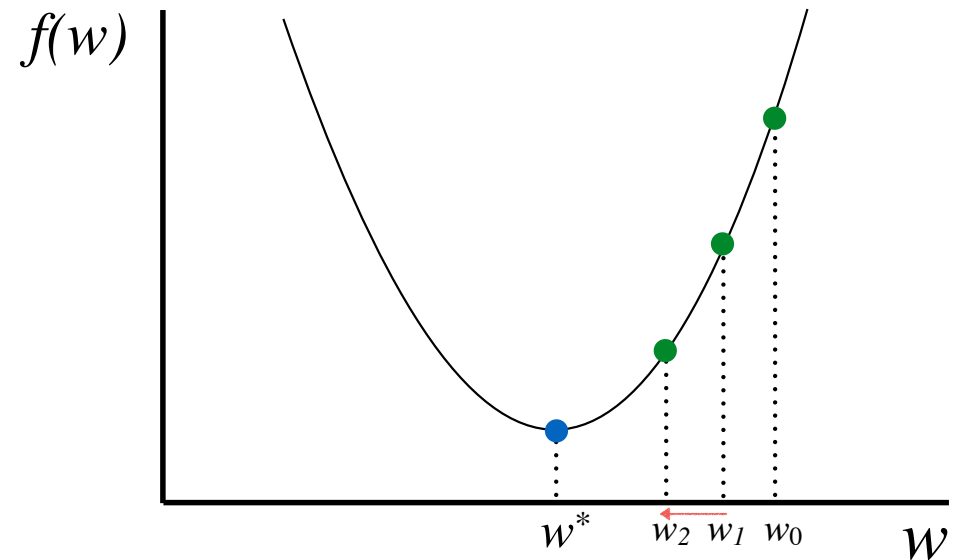
Repeat

Determine a descent direction

Choose a step size

Update

Until stopping criterion is satisfied



Gradient Descent

Start at a random point

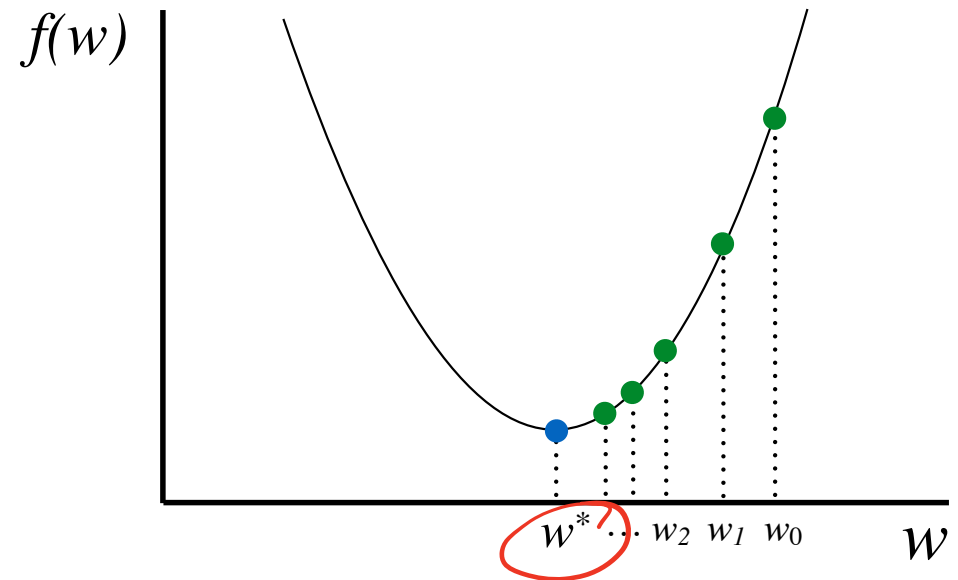
Repeat

Determine a descent direction

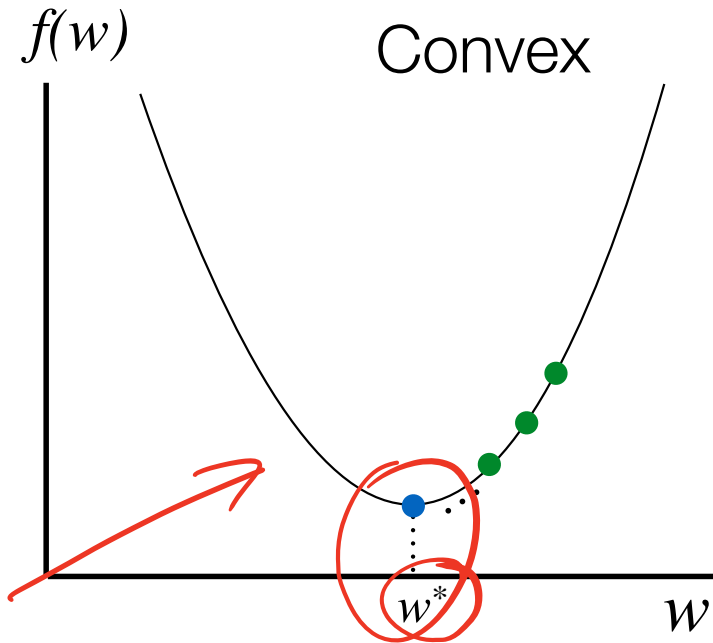
Choose a step size

Update

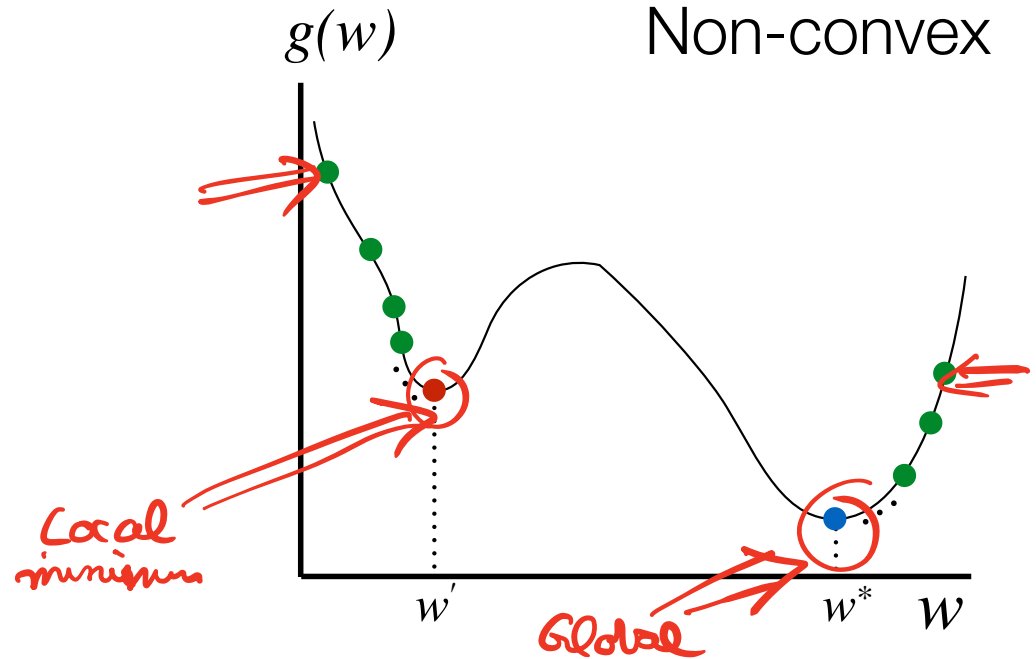
Until stopping criterion is satisfied



Where Will We Converge?



Any local minimum is a global minimum



Multiple local minima may exist

**Least Squares, Ridge Regression and
Logistic Regression are all convex!**

Convex functions

A function $f(x)$ is convex if

$$f(\lambda a + (1 - \lambda)b) \leq \lambda f(a) + (1 - \lambda)f(b)$$

for

$$x_\lambda = \lambda a + (1 - \lambda)b$$

$$0 \leq \lambda \leq 1$$

$$f(\lambda a + (1 - \lambda)b) \leq$$

$$0 \leq \lambda \leq 1$$

$$\lambda = 1$$

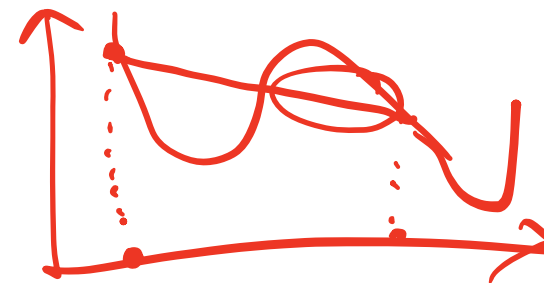
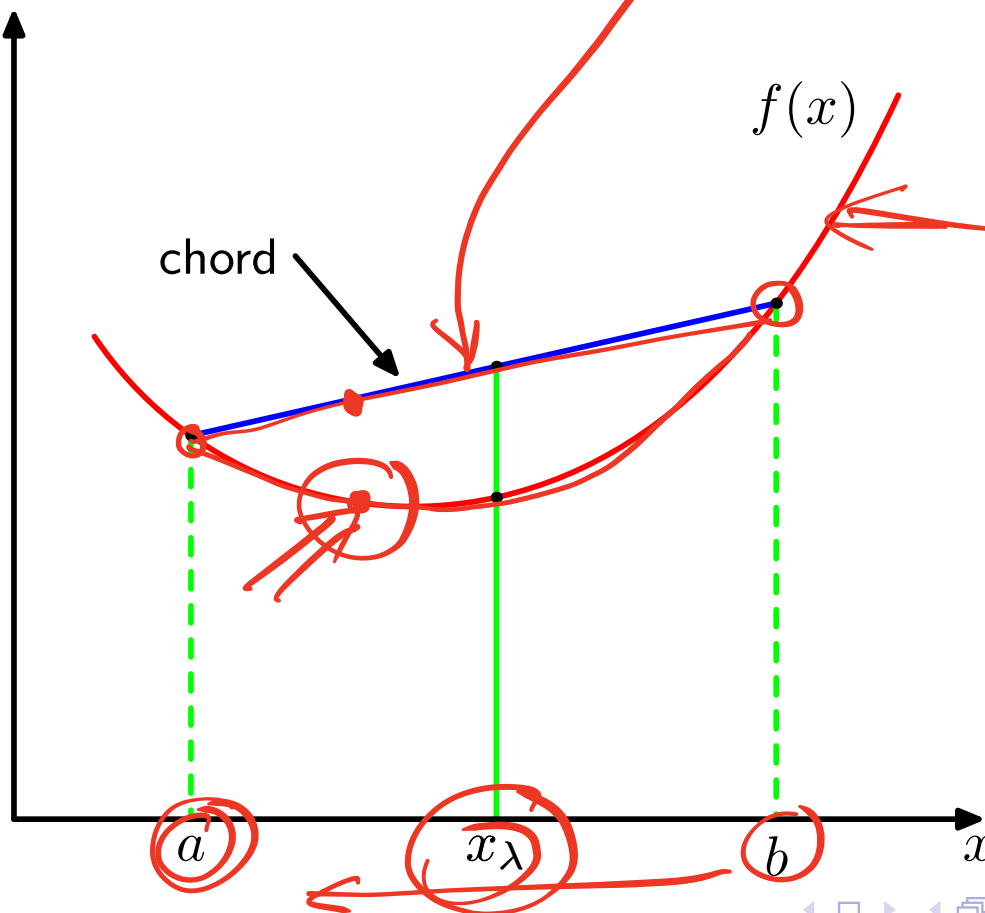
$$a = a$$

$$\lambda = 0$$

$$x = b$$

$$\lambda = \frac{1}{2}$$

$$x = \frac{1}{2}a + \frac{1}{2}b$$



How to determine convexity?

$f(x)$ is convex if

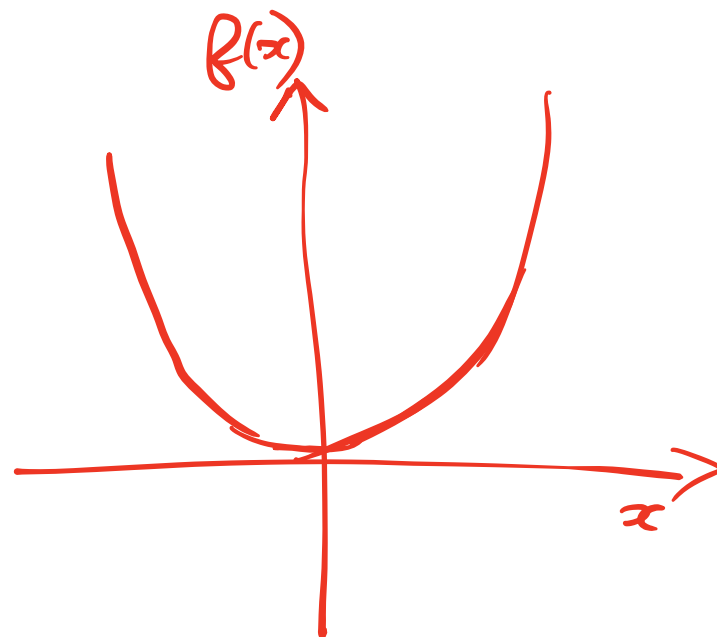
$$\underline{\underline{f''(x) \geq 0}}$$

Examples:

$$\underline{\underline{f(x) = x^2, f''(x) = 2 > 0}}$$

$$f'(x) = 2x$$

$$f''(x) = 2$$



Examples

Convex functions

$$f(x) = ax + b$$

$$f(x) = x^2$$

$$f(x) = e^x$$

$$f(x) = \frac{1}{x}, x \geq 0$$

$$f''(x) = 0 \geq 0$$

$$f''(x) = e^x \geq 0$$

Examples

Nonconvex functions

$$f(x) = \cos(x)$$

$$f(x) = e^x - x^2$$

$$f(x) = \log(x)$$

Multi-variate functions

Definition

$f(\mathbf{x})$ is convex

$$f(\lambda \mathbf{a} + (1 - \lambda) \mathbf{b}) \leq \lambda f(\mathbf{a}) + (1 - \lambda) f(\mathbf{b})$$

for all \mathbf{a}, \mathbf{b} , $0 \leq \lambda \leq 1$

Multi-variate functions

$$\frac{\partial^2 f}{\partial x_i \partial x_j}$$

How to determine convexity in this case?

Matrix of second-order derivatives (**Hessian**)

$$f''(x) = \frac{d^2 f(x)}{dx^2} \geq 0$$

$$H = \begin{pmatrix} \frac{\partial^2 f(x)}{\partial x_1^2} & \frac{\partial^2 f(x)}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f(x)}{\partial x_1 \partial x_D} \\ \frac{\partial^2 f(x)}{\partial x_1 \partial x_2} & \frac{\partial^2 f(x)}{\partial x_2^2} & \cdots & \frac{\partial^2 f(x)}{\partial x_2 \partial x_D} \\ \cdots & \cdots & \cdots & \cdots \\ \frac{\partial^2 f(x)}{\partial x_1 \partial x_D} & \frac{\partial^2 f(x)}{\partial x_2 \partial x_D} & \cdots & \frac{\partial^2 f(x)}{\partial x_D^2} \end{pmatrix}$$

(DxD)

Multi-variate functions

How to determine convexity in this case?

If the Hessian is positive semi-definite $\underline{H \succeq 0}$ then f is convex.

A matrix \mathbf{H} is positive semi-definite if and only if

$$\underbrace{z^T \mathbf{H} z}_{\text{red bracket}} = \sum_{j,k} H_{j,k} z_j z_k \geq 0$$

for all z .

$$\begin{array}{cc} \mathbf{H} & \mathbf{z} \\ D \times D & D \end{array}$$

$$\underbrace{\begin{array}{ccc} & \mathbf{z}^T & \\ \mathbf{z} & \mathbf{H} & \mathbf{z} \\ 1 \times D & D \times D & D \times 1 \end{array}}_{\text{red oval}} \geq 0$$

Multi-variate functions

Example

$$x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$

$$\underline{f(x)} = x_1^2 + 2x_2^2$$

symmetric \swarrow

$$\underline{H} = \begin{pmatrix} 2 & 0 \\ 0 & 4 \end{pmatrix}$$

$$z^T H z = \underline{2z_1^2} + \underline{4z_2^2}$$

$$\neq 0$$

Multi-variate functions

Example

$$f(\mathbf{x}) = x_1^2 + 2x_2^2$$

$$\mathbf{H} = \begin{pmatrix} 2 & 0 \\ 0 & 4 \end{pmatrix}$$

$$\mathbf{z}^T \mathbf{H} \mathbf{z} = 2z_1^2 + 4z_2^2 \geq 0$$

Example: $\min f(\boldsymbol{\theta}) = 0.5(\theta_1^2 - \theta_2)^2 + 0.5(\theta_1 - 1)^2$

- We compute the gradients

$$\left\{ \begin{array}{l} \frac{\partial f}{\partial \theta_1} = 2(\theta_1^2 - \theta_2)\theta_1 + \theta_1 - 1 \\ \frac{\partial f}{\partial \theta_2} = -(\theta_1^2 - \theta_2) \end{array} \right. \quad (1)$$

$$(2)$$

- Use the following *iterative* procedure for *gradient descent*

- 1 Initialize $\theta_1^{(0)}$ and $\theta_2^{(0)}$, and $t = 0$
- 2 do

$$\theta_1^{(t+1)} \leftarrow \theta_1^{(t)} - \eta [2(\theta_1^{(t)^2} - \theta_2^{(t)})\theta_1^{(t)} + \theta_1^{(t)} - 1] \quad (3)$$

$$\theta_2^{(t+1)} \leftarrow \theta_2^{(t)} - \eta [-(\theta_1^{(t)^2} - \theta_2^{(t)})] \quad (4)$$

$$t \leftarrow t + 1 \quad (5)$$

- 3 until $f(\boldsymbol{\theta}^{(t)})$ does not change much

Gradient descent

General form for minimizing $f(\theta)$

$$\underline{\underline{\theta^{t+1}}} \leftarrow \underline{\underline{\theta^t}} - \eta \nabla f(\underline{\underline{\theta^t}})$$

Remarks

- η is often called *step size* – literally, how far our update will go along the the direction of the negative gradient
- Note that this is for *minimizing* a function, hence the subtraction $(-\eta)$
- With a *suitable* choice of η , the iterative procedure converges to a stationary point where

$$\underline{\underline{\nabla f(\theta)}} = 0$$

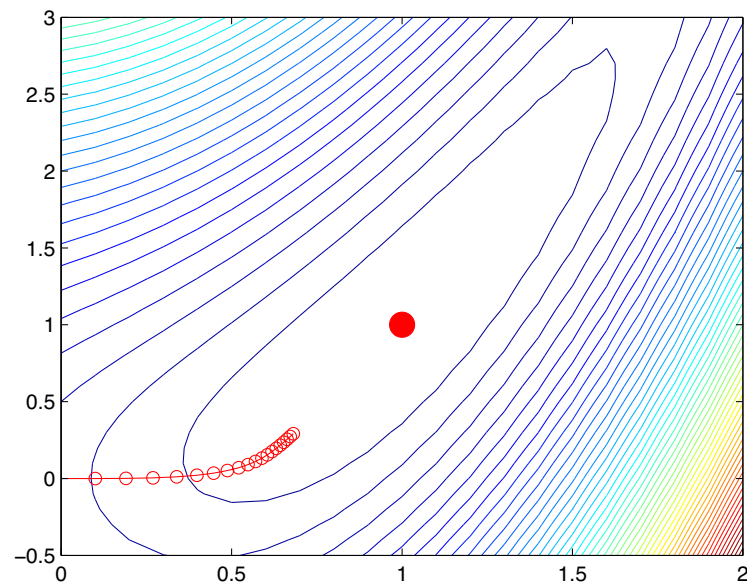
- A stationary point is only necessary for being the minimum.



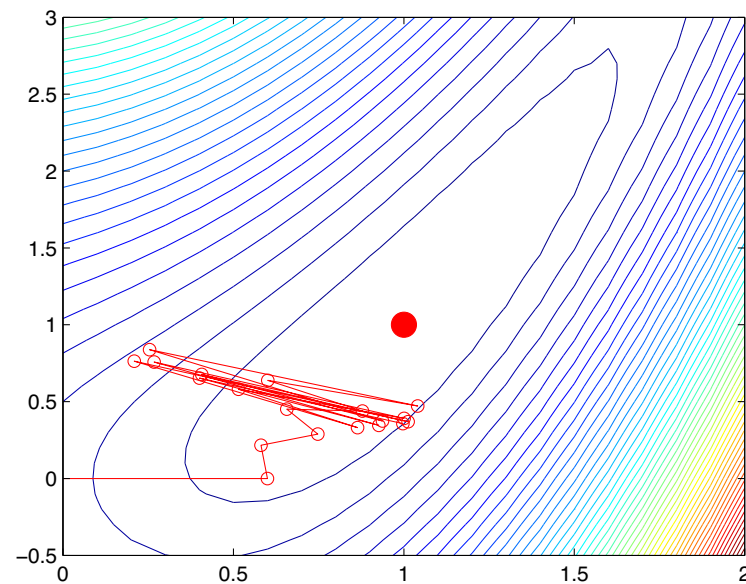
Seeing in action

Choosing the right η is important

small η is too slow?



large η is too unstable?



Gradient descent

$$\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$$

Algorithm 1 Gradient descent

```
1:  $\theta \leftarrow \mathbf{0}$ .  
2: for  $epoch = 1 \dots T$  do  
3:    $\theta \leftarrow \theta - \eta \nabla J(\theta)$   
4: end for  
5: return  $\theta$ 
```

Gradient Descent Update for Logistic Regression

Derivatives of $\sigma(a)$

$$\begin{aligned}\frac{d\sigma(a)}{da} &= \frac{d}{da} (1 + e^{-a})^{-1} = \frac{-(1 + e^{-a})'}{(1 + e^{-a})^2} \\ &= \frac{e^{-a}}{(1 + e^{-a})^2} = \frac{1}{1 + e^{-a}} \frac{e^{-a}}{1 + e^{-a}} \\ &= \sigma(a)[1 - \sigma(a)]\end{aligned}$$

Gradients of the negative log likelihood

Negative log likelihood

$$J(\boldsymbol{\theta}) = - \sum_n \{y_n \log h_{\boldsymbol{\theta}}(\mathbf{x}_n) + (1 - y_n) \log[1 - h_{\boldsymbol{\theta}}(\mathbf{x}_n)]\}$$

Gradients

$$\frac{\partial J(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = - \sum_n \{y_n [1 - \sigma(\boldsymbol{\theta}^T \mathbf{x}_n)] \mathbf{x}_n - (1 - y_n) \sigma(\boldsymbol{\theta}^T \mathbf{x}_n) \mathbf{x}_n\} \quad (6)$$

$$= \sum_n \{\sigma(\boldsymbol{\theta}^T \mathbf{x}_n) - y_n\} \mathbf{x}_n \quad (7)$$

$$= \sum_n \{h_{\boldsymbol{\theta}}(\mathbf{x}_n) - y_n\} \mathbf{x}_n \quad (8)$$

Remark

Gradients of the negative log likelihood

Negative log likelihood

$$J(\boldsymbol{\theta}) = - \sum_n \{y_n \log h_{\boldsymbol{\theta}}(\mathbf{x}_n) + (1 - y_n) \log[1 - h_{\boldsymbol{\theta}}(\mathbf{x}_n)]\}$$

Gradients

$$\frac{\partial J(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = - \sum_n \{y_n [1 - \sigma(\boldsymbol{\theta}^T \mathbf{x}_n)] \mathbf{x}_n - (1 - y_n) \sigma(\boldsymbol{\theta}^T \mathbf{x}_n) \mathbf{x}_n\} \quad (6)$$

$$= \sum_n \{\sigma(\boldsymbol{\theta}^T \mathbf{x}_n) - y_n\} \mathbf{x}_n \quad (7)$$

$$= \sum_n \{h_{\boldsymbol{\theta}}(\mathbf{x}_n) - y_n\} \mathbf{x}_n \quad (8)$$

Remark

- $e_n = \{h_{\boldsymbol{\theta}}(\mathbf{x}_n) - y_n\}$ is called **error** for the n th training sample.

Numerical optimization

Gradient descent

- Choose a proper step size $\eta > 0$
- Iteratively update the parameters following the negative gradient to minimize the error function

$$\boldsymbol{\theta}^{(t+1)} \leftarrow \boldsymbol{\theta}^{(t)} - \eta \sum_n \{ \sigma(\boldsymbol{\theta}^T \mathbf{x}_n) - y_n \} \mathbf{x}_n$$

Remarks

- The step size needs to be chosen carefully to ensure convergence.
- The step size can be adaptive (i.e. varying from iteration to iteration). For example, we can use techniques such as *line search*
- There is a variant called *stochastic* gradient descent, also popularly used.

Outline

- 1 Logistic regression
- 2 Optimization
- 3 Stochastic gradient descent

Stochastic gradient descent

$$\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$$

$$J(\boldsymbol{\theta}) = -\sum_n \{y_n \log h_{\boldsymbol{\theta}}(\mathbf{x}_n) + (1 - y_n) \log[1 - h_{\boldsymbol{\theta}}(\mathbf{x}_n)]\}$$

$$J(\boldsymbol{\theta}) = \frac{1}{N} \sum_n J_n(\boldsymbol{\theta})$$

$$\nabla J(\boldsymbol{\theta}) = \frac{1}{N} \sum_n \nabla J_n(\boldsymbol{\theta}) = \mathbb{E}_{n \sim \mathcal{D}} \nabla J_n(\boldsymbol{\theta})$$

Approximate the gradient by the gradient computed on a single example at a time.

- Repeat until convergence
 - ▶ Randomly pick one example (\mathbf{x}_n, y_n) .
 - ▶ Update $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \nabla J_n(\boldsymbol{\theta})$.

Stochastic Gradient descent (SGD)

$$\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$$

Algorithm 2 Stochastic Gradient descent

```
1:  $\theta \leftarrow \mathbf{0}$ .  
2: for  $epoch = 1 \dots T$  do  
3:   for  $(\mathbf{x}, y) \in \mathcal{D}$  do  
4:      $\theta \leftarrow \theta - \eta \nabla J_{(\mathbf{x}, y)}(\theta)$   
5:   end for  
6: end for  
7: return  $\theta$ 
```

Compare to perceptron learning

$$\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$$

Algorithm 3 PerceptronTrain

```
1:  $\theta \leftarrow \mathbf{0}$ 
2: for  $iter = 1 \dots MaxIter$  do
3:   for  $(\mathbf{x}, y) \in \mathcal{D}$  do
4:      $a \leftarrow \theta^T \mathbf{x}$ 
5:     if  $ay \leq 0$  then
6:        $\theta \leftarrow \theta + y\mathbf{x}$ 
7:     end if
8:   end for
9: end for
10: return  $\theta$ 
```

What is the objective function?

Compare to perceptron learning

$$\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$$

Algorithm 4 PerceptronTrain

```
1:  $\theta \leftarrow \mathbf{0}$ 
2: for  $iter = 1 \dots MaxIter$  do
3:   for  $(\mathbf{x}, y) \in \mathcal{D}$  do
4:      $a \leftarrow \theta^T \mathbf{x}$ 
5:     if  $ay \leq 0$  then
6:        $\theta \leftarrow \theta + y\mathbf{x}$ 
7:     end if
8:   end for
9: end for
10: return  $\theta$ 
```

Perceptron effectively minimizes:

$$J(\theta) = \sum_n \max(0, 1 - y_n \theta^T \mathbf{x}_n)$$

Summary

Setup for binary classification

- Logistic Regression models conditional distribution as:
 $p(y = 1|\mathbf{x}; \boldsymbol{\theta}) = \sigma[a(\mathbf{x})]$ where $a(\mathbf{x}) = \boldsymbol{\theta}^T \mathbf{x}$
- Linear decision boundary: $a(\mathbf{x}) = \boldsymbol{\theta}^T \mathbf{x} = 0$

Minimizing the negative log-likelihood

- $J(\boldsymbol{\theta}) = -\sum_n \{y_n \log \sigma(\boldsymbol{\theta}^T \mathbf{x}_n) + (1 - y_n) \log[1 - \sigma(\boldsymbol{\theta}^T \mathbf{x}_n)]\}$
- No closed form solution; must rely on iterative solvers

Numerical optimization

- Gradient descent: simple, scalable to large-scale problems
 - ▶ move in direction opposite of gradient!
 - ▶ gradient of logistic function takes nice form