

CSCI567 Machine Learning (Spring 2018)

Michael Shindler

January 29, 2018

Outline

- 1 Administration
- 2 Review of Last Lecture
- 3 Logistic regression
- 4 Summary
- 5 Constrained Optimization (MLE)

Outline

- 1 Administration
- 2 Review of Last Lecture
- 3 Logistic regression
- 4 Summary
- 5 Constrained Optimization (MLE)

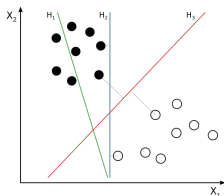
Administrative stuff

- No more D-clearance will be granted.
- Problem Set 1 has been posted
 - No grace days on **problem sets**
 - Will be submitted via USCDen.
- Programming assignment one posted too
 - Grace days allowed (form will be provided)
 - Git + notebook

Outline

- 1 Administration
- 2 Review of Last Lecture**
- 3 Logistic regression
- 4 Summary
- 5 Constrained Optimization (MLE)

Definition of linear classifiers for binary classification



The decision boundary is a linear function of the input \mathbf{x} . The prediction is

$$y = \text{sign}(\mathbf{w}^T \mathbf{x}) = \begin{cases} +1 & \text{if } \mathbf{w}^T \mathbf{x} > 0 \\ -1 & \text{if } \mathbf{w}^T \mathbf{x} \leq 0 \end{cases}$$

Sometimes, we use sgn for sign . Note that a correct classification happens *if and only if*

$$y^* \mathbf{w}^T \mathbf{x} > 0$$

where y^* is the true label.

Perceptron

Iteratively improving one case at a time

- REPEAT
- Pick a data point \mathbf{x}_n randomly
- Make a prediction $y = \text{sign}(\mathbf{w}^T \mathbf{x}_n)$ using the *current* \mathbf{w}
- If $y = y_n$, do nothing. Else,

$$\mathbf{w} \leftarrow \mathbf{w} + y_n \mathbf{x}_n$$

- UNTIL converged.

Note that the new \mathbf{w} is used to make prediction as soon as it is updated.

Properties

- If the training data is linearly separable, the algorithm stops in a finite number of steps.
- The parameter vector is always a linear combination of training instances.

Stochastic gradient descent for linear regression

RSS Loss function

$$RSS(\tilde{\mathbf{w}}) = \sum_n (\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n - y_n)^2$$

Widrow-Hoff rule: update parameters using one example at a time

- Initialize $\tilde{\mathbf{w}}$ to $\tilde{\mathbf{w}}^{(0)}$ (anything reasonable is fine); set $t = 0$; choose $\eta > 0$
- Loop *until convergence*
 - ① random choose a training a sample \mathbf{x}_n
 - ② Compute its contribution to the gradient (ignoring the constant factor)

$$\mathbf{g}_n = (\tilde{\mathbf{x}}_n^T \tilde{\mathbf{w}}^{(t)} - y_n) \tilde{\mathbf{x}}_n$$

- ③ Update the parameters
 $\tilde{\mathbf{w}}^{(t+1)} = \tilde{\mathbf{w}}^{(t)} - \eta \mathbf{g}_n$
- ④ $t \leftarrow t + 1$

Outline

- 1 Administration
- 2 Review of Last Lecture
- 3 Logistic regression
 - Motivation
 - General setup
 - Maximum likelihood estimation
 - Numerical optimization
 - Gradient descent
 - Gradient descent for logistic regression
 - Newton method
 - Stochastic Gradient Descent for Logistic Regression

- 4 Summary

Linear regression method for classification?

Setup for two classes

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

Why not just use RSS?

$$\tilde{\mathbf{w}}_{\text{RSS}} = \arg \min_{\tilde{\mathbf{w}}} \sum_n (\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n - y_n)^2$$

Problems

- $\tilde{\mathbf{w}}_{\text{RSS}}^T \tilde{\mathbf{x}}$ is not guaranteed to be either 0, 1 or even between 0 and 1!
- Unclear why the binary valued y is appropriate as response variable in regression

Strange thing can happen

Toy data

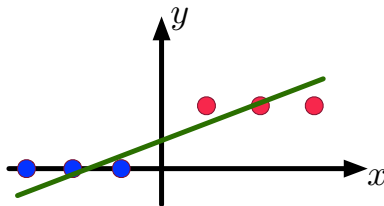
x	label
1	1
2	1
3	1
-1	0
-2	0
-3	0

An intuitive classifier

$$y = \text{sign}(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x < 0 \end{cases}$$

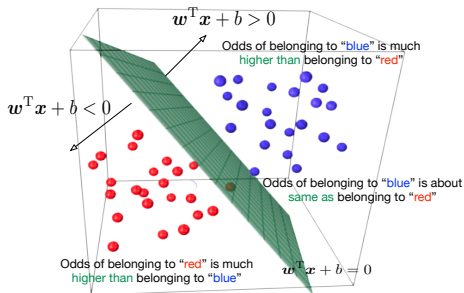
If we use linear regression

We get something not sensible!



Clearly, new perspective is needed!

Reinterpreting the linear decision boundary



The further the point is away from the boundary, the more confident we think it belongs to one of the classes. If a point is on the boundary precisely, then we think it can belong to either class.

We define thus an *odd ratio* to represent how much a point should be interpreted as belonging to class 1 or 0 *over equally possible*.

Link probability to odd ratio

Probability belonging to class 1

$$p(y = 1|\mathbf{x})$$

Probability belonging to class 0

$$p(y = 0|\mathbf{x}) = 1 - p(y = 1|\mathbf{x})$$

Odd ratio of belonging to class 1

$$r = \frac{p(y = 1|\mathbf{x})}{p(y = 0|\mathbf{x})} = \frac{p(y = 1|\mathbf{x})}{1 - p(y = 1|\mathbf{x})}$$

Link probability to odd ratio

Probability belonging to class 1

$$p(y = 1|\mathbf{x})$$

Probability belonging to class 0

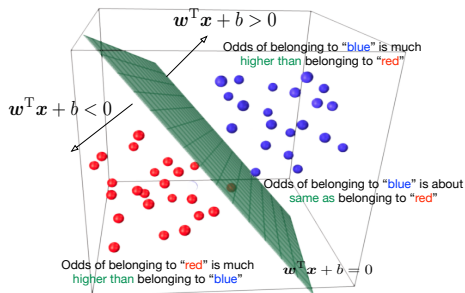
$$p(y = 0|\mathbf{x}) = 1 - p(y = 1|\mathbf{x})$$

Odd ratio of belonging to class 1

$$r = \frac{p(y = 1|\mathbf{x})}{p(y = 0|\mathbf{x})} = \frac{p(y = 1|\mathbf{x})}{1 - p(y = 1|\mathbf{x})}$$

- $r = 1$: equally likely
- $r = +\infty$: absolutely certain class 1
- $r = 0$: absolutely certain class 0

Linking log-odds ratio to decision boundary



$$\log r = \log \frac{p(y = 1|\mathbf{x})}{1 - p(y = 1|\mathbf{x})} = \mathbf{w}^T \mathbf{x} + b$$

Or

$$p(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + b)}}$$

Logistic classification

Setup for two classes

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

$$p(y = 1 | \mathbf{x}; b, \mathbf{w}) = \sigma[g(\mathbf{x})]$$

where

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

and $\sigma[\cdot]$ stands for the *sigmoid* function

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

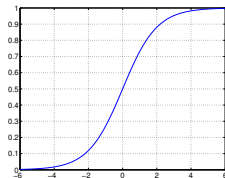
Why the sigmoid function?

What does it look like?

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

where

$$a = b + \mathbf{w}^T \mathbf{x}$$



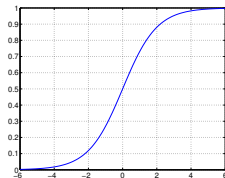
Why the sigmoid function?

What does it look like?

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

where

$$a = b + \mathbf{w}^T \mathbf{x}$$



Properties

- Bounded between 0 and 1 ← thus, interpretable as probability
- Monotonically increasing thus, usable to derive classification rules
 - ① $\sigma(a) > 0.5$, positive (classify as '1')
 - ② $\sigma(a) < 0.5$, negative (classify as '0')
 - ③ $\sigma(a) = 0.5$, undecidable
- Nice computational properties These will unfold in the next few slides

Linear or nonlinear?

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

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

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

We often call b the bias term.

Likelihood function

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

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

Likelihood function

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

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

Compact expression, exploring that y_n is either 1 or 0

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

Or its logarithm

$$\log p(y_n|\mathbf{x}_n; b, \mathbf{w}) = \{y_n \log \sigma(b + \mathbf{w}^T \mathbf{x}_n) + (1 - y_n) \log [1 - \sigma(b + \mathbf{w}^T \mathbf{x}_n)]\}$$

Cross-entropy error

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

$$\begin{aligned}\log P(\mathcal{D}) &= \log \prod_n p(y_n | \mathbf{x}_n; b, \mathbf{w}) \\ &= \sum_n \{y_n \log \sigma(b + \mathbf{w}^T \mathbf{x}_n) + (1 - y_n) \log [1 - \sigma(b + \mathbf{w}^T \mathbf{x}_n)]\}\end{aligned}\tag{1}$$

Cross-entropy error

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

$$\begin{aligned}\log P(\mathcal{D}) &= \log \prod_n p(y_n | \mathbf{x}_n; b, \mathbf{w}) \\ &= \sum_n \{y_n \log \sigma(b + \mathbf{w}^T \mathbf{x}_n) + (1 - y_n) \log [1 - \sigma(b + \mathbf{w}^T \mathbf{x}_n)]\}\end{aligned}\tag{1}$$

It is convenient to work with its negation, which is called *cross-entropy error function*

$$\mathcal{E}(b, \mathbf{w}) = - \sum_n \{y_n \log \sigma(b + \mathbf{w}^T \mathbf{x}_n) + (1 - y_n) \log [1 - \sigma(b + \mathbf{w}^T \mathbf{x}_n)]\}$$

Shorthand notation

This is for convenience

- Append 1 to \mathbf{x}

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

- Append b to \mathbf{w}

$$\mathbf{w} \leftarrow [b \quad w_1 \quad w_2 \quad \cdots \quad w_D]$$

- Cross-entropy is then

$$\mathcal{E}(\mathbf{w}) = - \sum_n \{y_n \log \sigma(\mathbf{w}^T \mathbf{x}_n) + (1 - y_n) \log [1 - \sigma(\mathbf{w}^T \mathbf{x}_n)]\}$$

NB. We are not using the $\tilde{\mathbf{x}}$ and $\tilde{\mathbf{w}}$ (as in several textbooks) for cosmetic reasons.

How to find the optimal parameters for logistic regression?

We will minimize the error function

$$\mathcal{E}(\mathbf{w}) = - \sum_n \{y_n \log \sigma(\mathbf{w}^T \mathbf{x}_n) + (1 - y_n) \log[1 - \sigma(\mathbf{w}^T \mathbf{x}_n)]\}$$

However, this function is *highly nonlinear* and we will not be able to find the simple closed-form solution. So we need to use *numerical* methods.

- Numerical methods are messier, in contrast to cleaner analytic solutions.
- In practice, we often have to tune a few optimization parameters — patience is necessary.

An overview of numerical methods

We describe two

- Gradient descent (our focus in lecture): simple, especially effective for large-scale problems
- Newton method: classical and powerful method

Gradient descent is often referred to as a *first-order* method as it requires only to compute the gradients (i.e., the first-order derivative) of the function.

In contrast, Newton method is often referred as to a *second-order* method.

Gradient descent

General form for minimizing $f(\theta)$

$$\theta^{t+1} \leftarrow \theta^t - \eta \frac{\partial f}{\partial \theta}$$

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

$$\frac{\partial f}{\partial \theta} = 0$$

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

How do we do this for logistic regression?

Simple fact: derivatives of $\sigma(a)$

$$\frac{d\sigma(a)}{da} = \frac{d}{da} \left(\frac{1}{1 + e^{-a}} \right) = \frac{-(1 + e^{-a})'}{(1 + e^{-a})^2}$$

How do we do this for logistic regression?

Simple fact: derivatives of $\sigma(a)$

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

How do we do this for logistic regression?

Simple fact: derivatives of $\sigma(a)$

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

Gradients of the cross-entropy error function

Gradients

$$\frac{\partial \mathcal{E}(\mathbf{w})}{\partial \mathbf{w}} = - \sum_n \{ y_n [1 - \sigma(\mathbf{w}^T \mathbf{x}_n)] \mathbf{x}_n - (1 - y_n) \sigma(\mathbf{w}^T \mathbf{x}_n) \mathbf{x}_n \} \quad (2)$$

$$= \sum_n \{ \sigma(\mathbf{w}^T \mathbf{x}_n) - y_n \} \mathbf{x}_n \quad (3)$$

Remarks

- $e_n = \{ \sigma(\mathbf{w}^T \mathbf{x}_n) - y_n \}$ is called **error** for the n th training sample.
- Stationary point (in this case, the optimum):

$$\sum_n \{ \sigma(\mathbf{w}^T \mathbf{x}_n) - y_n \} \mathbf{x}_n = 0$$

Intuition: *on average, the error is zero.*

Numerical optimization

Gradient descent

- Choose a proper step size $\eta > 0$

Numerical optimization

Gradient descent

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

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

Numerical optimization

Gradient descent

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

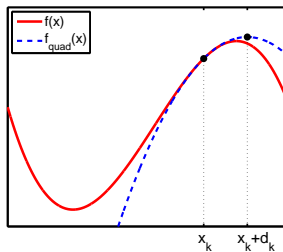
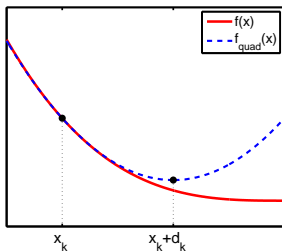
$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \sum_n \{ \sigma(\mathbf{w}^T \mathbf{x}_n) - y_n \} \mathbf{x}_n$$

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 – we will discuss in a bit.

Intuition for Newton method

Approximate the true function with an easy-to-solve optimization problem



Approximation

Taylor expansion of the cross-entropy function

$$\mathcal{E}(\mathbf{w}) \approx \mathcal{E}(\mathbf{w}^{(t)}) + (\mathbf{w} - \mathbf{w}^{(t)})^T \nabla \mathcal{E}(\mathbf{w}^{(t)}) + \frac{1}{2}(\mathbf{w} - \mathbf{w}^{(t)})^T \mathbf{H}^{(t)}(\mathbf{w} - \mathbf{w}^{(t)})$$

where

- $\nabla \mathcal{E}(\mathbf{w}^{(t)})$ is the gradient
- $\mathbf{H}^{(t)}$ is the Hessian matrix evaluated at $\mathbf{w}^{(t)}$

Example: a scalar function

$$\sin(\theta) \approx \sin(0) + \theta \cos(\theta = 0) + \frac{1}{2}\theta^2[-\sin(\theta = 0)] = \theta$$

where $\nabla \sin(\theta) = \cos(\theta)$ and $\mathbf{H} = \nabla \cos(\theta) = -\sin(\theta)$

So what is the Hessian matrix?

The matrix of second-order derivatives

$$\mathbf{H} = \frac{\partial^2 \mathcal{E}(\mathbf{w})}{\partial \mathbf{w} \mathbf{w}^T}$$

In other words,

$$H_{ij} = \frac{\partial}{\partial w_j} \left(\frac{\partial \mathcal{E}(\mathbf{w})}{\partial w_i} \right)$$

So the Hessian matrix is $\mathbb{R}^{D \times D}$, where $\mathbf{w} \in \mathbb{R}^D$.

Optimizing the approximation

Minimize the approximation

$$\mathcal{E}(\mathbf{w}) \approx \mathcal{E}(\mathbf{w}^{(t)}) + (\mathbf{w} - \mathbf{w}^{(t)})^T \nabla \mathcal{E}(\mathbf{w}^{(t)}) + \frac{1}{2}(\mathbf{w} - \mathbf{w}^{(t)})^T \mathbf{H}^{(t)}(\mathbf{w} - \mathbf{w}^{(t)})$$

and use the solution as the new estimate of the parameters

$$\mathbf{w}^{(t+1)} \leftarrow \min_{\mathbf{w}} (\mathbf{w} - \mathbf{w}^{(t)})^T \nabla \mathcal{E}(\mathbf{w}^{(t)}) + \frac{1}{2}(\mathbf{w} - \mathbf{w}^{(t)})^T \mathbf{H}^{(t)}(\mathbf{w} - \mathbf{w}^{(t)})$$

Optimizing the approximation

Minimize the approximation

$$\mathcal{E}(\mathbf{w}) \approx \mathcal{E}(\mathbf{w}^{(t)}) + (\mathbf{w} - \mathbf{w}^{(t)})^T \nabla \mathcal{E}(\mathbf{w}^{(t)}) + \frac{1}{2} (\mathbf{w} - \mathbf{w}^{(t)})^T \mathbf{H}^{(t)} (\mathbf{w} - \mathbf{w}^{(t)})$$

and use the solution as the new estimate of the parameters

$$\mathbf{w}^{(t+1)} \leftarrow \min_{\mathbf{w}} (\mathbf{w} - \mathbf{w}^{(t)})^T \nabla \mathcal{E}(\mathbf{w}^{(t)}) + \frac{1}{2} (\mathbf{w} - \mathbf{w}^{(t)})^T \mathbf{H}^{(t)} (\mathbf{w} - \mathbf{w}^{(t)})$$

The quadratic function minimization has a *closed* form, thus, we have

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \left(\mathbf{H}^{(t)} \right)^{-1} \nabla \mathcal{E}(\mathbf{w}^{(t)})$$

i.e., the Newton method.

Contrast gradient descent and Newton method

Similar

Both are iterative procedures.

Difference

- Newton method requires second-order derivatives.
- Newton method does not have the magic η to be set.

Other important things about Hessian

Our cross-entropy error function is convex

$$\frac{\partial \mathcal{E}(\mathbf{w})}{\partial \mathbf{w}} = \sum_n \{ \sigma(\mathbf{w}^T \mathbf{x}_n) - y_n \} \mathbf{x}_n \quad (4)$$

$$\Rightarrow \mathbf{H} = \frac{\partial^2 \mathcal{E}(\mathbf{w})}{\partial \mathbf{w} \mathbf{w}^T} \quad (5)$$

Other important things about Hessian

Our cross-entropy error function is convex

$$\frac{\partial \mathcal{E}(\mathbf{w})}{\partial \mathbf{w}} = \sum_n \{ \sigma(\mathbf{w}^T \mathbf{x}_n) - y_n \} \mathbf{x}_n \quad (4)$$

$$\Rightarrow \mathbf{H} = \frac{\partial^2 \mathcal{E}(\mathbf{w})}{\partial \mathbf{w} \mathbf{w}^T} \quad (5)$$

For any vector \mathbf{v} ,

$$\mathbf{v}^T \mathbf{H} \mathbf{v} \geq 0$$

Thus, positive semidefinite. Thus, the cross-entropy error function is convex, with only one global optimum.

Good about Newton method

Fast!

Suppose we want to minimize $f(x) = x^2 + 2x$ and we have its current estimate at $x^{(t)} \neq -1$. So what is the next estimate?

$$x^{(t+1)} \leftarrow x^{(t)} - [f''(x)]^{-1} f'(x) = x^{(t)} - \frac{1}{2}(2x^{(t)} + 2) = -1$$

Namely, the next step (of iteration) immediately tells us the global optimum! (In optimization, this is called *superlinear convergence rate*).

In general, the better our approximation, the fast the Newton method is in solving our optimization problem.

Bad about Newton method

Not scalable!

- Computing and inverting Hessian matrix can be very expensive for large-scale problems where the dimensionally D is very large.
- Newton method does not guarantee convergence if your starting point is far away from the optimum

NB. There are fixes and alternatives, such as Quasi-Newton/Quasi-second order method.

Stochastic gradient descent

Cross entropy error function

$$\mathcal{E}(\mathbf{w}) = - \sum_n \{y_n \log \sigma(\mathbf{w}^T \mathbf{x}_n) + (1 - y_n) \log[1 - \sigma(\mathbf{w}^T \mathbf{x}_n)]\}$$

Update rules

- Initialize \mathbf{w} to $\mathbf{w}^{(0)}$ (anything reasonable is fine); set $t = 0$; choose $\eta > 0$

Stochastic gradient descent

Cross entropy error function

$$\mathcal{E}(\mathbf{w}) = - \sum_n \{y_n \log \sigma(\mathbf{w}^T \mathbf{x}_n) + (1 - y_n) \log[1 - \sigma(\mathbf{w}^T \mathbf{x}_n)]\}$$

Update rules

- Initialize \mathbf{w} to $\mathbf{w}^{(0)}$ (anything reasonable is fine); set $t = 0$; choose $\eta > 0$
- Loop *until convergence*
 - 1 random choose a training a sample \mathbf{x}_n
 - 2 Compute its contribution to the gradient

$$\mathbf{g}_n = (\sigma(\mathbf{w}^T \mathbf{x}_n) - y_n) \mathbf{x}_n$$

Stochastic gradient descent

Cross entropy error function

$$\mathcal{E}(\mathbf{w}) = - \sum_n \{y_n \log \sigma(\mathbf{w}^T \mathbf{x}_n) + (1 - y_n) \log[1 - \sigma(\mathbf{w}^T \mathbf{x}_n)]\}$$

Update rules

- Initialize \mathbf{w} to $\mathbf{w}^{(0)}$ (anything reasonable is fine); set $t = 0$; choose $\eta > 0$
- Loop *until convergence*
 - 1 random choose a training a sample \mathbf{x}_n
 - 2 Compute its contribution to the gradient

$$\mathbf{g}_n = (\sigma(\mathbf{w}^T \mathbf{x}_n) - y_n) \mathbf{x}_n$$

- 3 Update the parameters
$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \mathbf{g}_n$$

Stochastic gradient descent

Cross entropy error function

$$\mathcal{E}(\mathbf{w}) = - \sum_n \{y_n \log \sigma(\mathbf{w}^T \mathbf{x}_n) + (1 - y_n) \log[1 - \sigma(\mathbf{w}^T \mathbf{x}_n)]\}$$

Update rules

- Initialize \mathbf{w} to $\mathbf{w}^{(0)}$ (anything reasonable is fine); set $t = 0$; choose $\eta > 0$
- Loop *until convergence*
 - ① random choose a training a sample \mathbf{x}_n
 - ② Compute its contribution to the gradient

$$\mathbf{g}_n = (\sigma(\mathbf{w}^T \mathbf{x}_n) - y_n) \mathbf{x}_n$$

- ③ Update the parameters
$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \mathbf{g}_n$$
- ④ $t \leftarrow t + 1$

So why $(\sigma(\mathbf{w}^T \mathbf{x}_n) - y_n)$ is called error?

Recall $\sigma(\mathbf{w}^T \mathbf{x}_n)$ is the probability

$$p(y = 1 | \mathbf{x}_n)$$

Given the update rule,

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta g_n = \mathbf{w}^{(t)} - \eta (\sigma(\mathbf{w}^T \mathbf{x}_n) - y_n) \mathbf{x}_n$$

The effect is:

- If $y_n = 1$ and $\sigma(\mathbf{w}^T \mathbf{x}_n) \approx 1$, we do nothing
- If $y_n = 1$ and $\sigma(\mathbf{w}^T \mathbf{x}_n) \approx 0$, we **add** \mathbf{x}_n to $\mathbf{w}^{(t)}$
- If $y_n = 0$ and $\sigma(\mathbf{w}^T \mathbf{x}_n) \approx 1$, we **subtract** \mathbf{x}_n to $\mathbf{w}^{(t)}$
- If $y_n = 0$ and $\sigma(\mathbf{w}^T \mathbf{x}_n) \approx 0$, we do nothing

This rule is similar to perceptron update rule except “softer” as $\sigma(\cdot)$ is often between 0 and 1

Outline

- 1 Administration
- 2 Review of Last Lecture
- 3 Logistic regression
- 4 Summary**
- 5 Constrained Optimization (MLE)

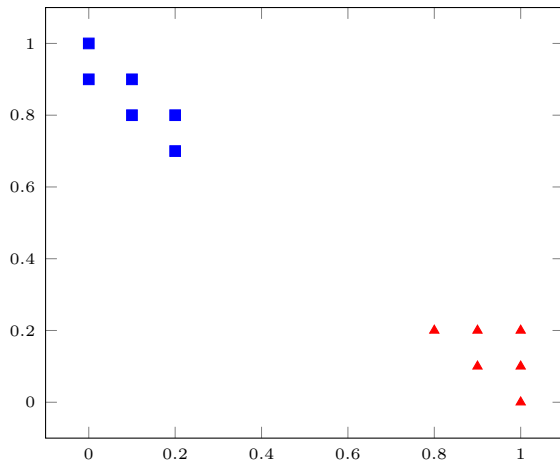
Important concepts and tools

- Logistic regression:
probabilistic modeling of the (log) odds ratio
cross-entropy error function
- Numerical optimization:
First-order method: gradient descent, batch update, stochastic gradient descent
Second-order method: Newton method

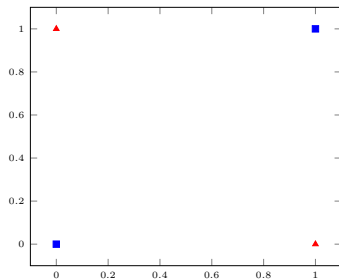
When to use perceptron vs logistic regression?

- Perceptron minimizes a loss function.
- Logistic is same, but probabilistic.

Example



How do we know when data is linearly separable?



Outline

- 1 Administration
- 2 Review of Last Lecture
- 3 Logistic regression
- 4 Summary
- 5 Constrained Optimization (MLE)

Constrained Optimization

General Case

- Minimize $f(x)$
- such that $g(x) = 0$

Method of Lagrange Multipliers

- $L(x, \lambda) = f(x) + \lambda g(x)$

Lagrange multipliers

- 1 Set derivative to zero

$$\frac{\partial L(x, \lambda)}{\partial x} = f'(x) + \lambda g'(x) = 0$$

- 2 Solve x in terms of λ

$$x = h(\lambda)$$

- 3 Substitute into constraint, solve λ , then x

$$g(h(\lambda)) = 0$$

Example: Dice rolls

Model

Probability of seeing a number k between 1 and 6 is $P(X = k) = \Theta_k$

Observations

$$\mathcal{D} = \{x_1, x_2, \dots, x_n\} \quad x_n \in \{1, 2, \dots, 6\}$$

Likelihood

$$L(\theta) = \prod_{n=1}^N P(X = x_n) = \prod_{k=1}^6 \Theta_k^{n_k}$$

Optimization

Objective function (log-likelihood)

$$\max \sum_k n_k \log \theta_k$$

Constraints

$$\sum_k \theta_k = 1 \quad \theta_k \geq 0$$

Lagrangian (ignoring non-negative constraint)

$$L(\theta, \lambda) = \sum_k n_k \log \theta_k + \lambda \left(\sum_k \theta_k - 1 \right)$$

Finding both multiplier and the parameters

Derivatives

$$\frac{\partial L(\theta, \lambda)}{\partial \theta_k} = \frac{n_k}{\theta_k} + \lambda$$

Setting them to zero

$$\theta_k = -\frac{1}{\lambda} n_k$$

Solving the multiplier by using the constraint

$$\sum_k \theta_k = -\frac{1}{\lambda} \sum_k n_k = 1 \rightarrow \lambda = -\sum_k n_k$$

Finding both multiplier and the parameters

Derivatives

$$\frac{\partial L(\theta, \lambda)}{\partial \theta_k} = \frac{n_k}{\theta_k} + \lambda$$

Setting them to zero

$$\theta_k = -\frac{1}{\lambda} n_k$$

Solving the multiplier by using the constraint

$$\sum_k \theta_k = -\frac{1}{\lambda} \sum_k n_k = 1 \rightarrow \lambda = -\sum_k n_k$$

Finally

$$\theta_k = \frac{n_k}{\sum_k n_k}$$