



University of Colorado **Boulder**

Department of Computer Science

CSCI 5622: Machine Learning

Chenhao Tan

Lecture 8: Stochastic Gradient Descent

Slides adapted from Chris Ketelsen, Jordan Boyd-Graber,  
and Noah Smith

# Administrivia

- HW1 grading is done
  - Common issues

# Learning Objectives

- Understand gradient descent
- Understand structural risk minimization
- Understand stochastic gradient descent

# Outline

- Objective function
- Gradient descent
- Structural risk minimization
- Stochastic gradient descent

# Outline

- Objective function
- Gradient descent
- Structural risk minimization
- Stochastic gradient descent

## Reminder: Logistic Regression

---

$$\underline{P(Y = 0 \mid \mathbf{x})} = \frac{1}{1 + \exp \left[ \beta_0 + \sum_j \beta_j \mathbf{x}_j \right]} \quad (1)$$

$$\underline{P(Y = 1 \mid \mathbf{x})} = \frac{\exp \left[ \beta_0 + \sum_j \beta_j \mathbf{x}_j \right]}{1 + \exp \left[ \beta_0 + \sum_j \beta_j \mathbf{x}_j \right]} \quad (2)$$

- Discriminative prediction:  $P(y \mid \mathbf{x})$
- Classification uses: sentiment analysis, spam detection
- What we didn't talk about is how to learn  $\beta$  from data

## Logistic Regression: Objective Function

One idea: find the parameter that maximize the likelihood of observing the training data.

## Logistic Regression: Objective Function

One idea: find the parameter that maximize the likelihood of observing the training data.

Maximize likelihood

$$\text{Obj} \equiv \log P(Y | X, \beta) = \sum_i \log P(y^{(i)} | \mathbf{x}^{(i)}, \beta)$$

$$\begin{aligned} & \underbrace{\sum_i \log P(y^{(i)} | \mathbf{x}^{(i)}, \beta)}_{\text{Likelihood}} \\ &= \sum_i y^{(i)} \left( \beta_0 + \sum_j \beta_j \mathbf{x}_j^{(i)} \right) - \log \left[ 1 + \exp \left( \beta_0 + \sum_j \beta_j \mathbf{x}_j^{(i)} \right) \right] \end{aligned}$$

## Logistic Regression: Objective Function

Minimize negative log likelihood (loss)

$$\begin{aligned}\mathcal{L} &\equiv -\log P(Y \mid X, \beta) = -\sum_i \log P(y^{(i)} \mid \mathbf{x}^{(i)}, \beta)) \\ &= \sum_i -y^{(i)} \left( \beta_0 + \sum_j \beta_j \mathbf{x}_j^{(i)} \right) + \log \left[ 1 + \exp \left( \beta_0 + \sum_j \beta_j \mathbf{x}_j^{(i)} \right) \right]\end{aligned}$$


## Logistic Regression: Objective Function

Minimize negative log likelihood (loss)

$$\begin{aligned}\mathcal{L} &\equiv -\log P(Y \mid X, \beta) = -\sum_i \log P(y^{(i)} \mid \mathbf{x}^{(i)}, \beta)) \\ &= \sum_i -y^{(i)} \left( \beta_0 + \sum_j \beta_j \mathbf{x}_j^{(i)} \right) + \log \left[ 1 + \exp \left( \beta_0 + \sum_j \beta_j \mathbf{x}_j^{(i)} \right) \right]\end{aligned}$$

Training data  $\{(\mathbf{x}, y)\}$  are fixed. Objective function is a function of  $\beta$  ... what values of  $\beta$  give a good value?

## Logistic Regression: Objective Function

Minimize negative log likelihood (loss)

$$\begin{aligned}\mathcal{L} &\equiv -\log P(Y \mid X, \beta) = -\sum_i \log P(y^{(i)} \mid \mathbf{x}^{(i)}, \beta)) \\ &= \sum_i -y^{(i)} \left( \beta_0 + \sum_j \beta_j \mathbf{x}_j^{(i)} \right) + \log \left[ 1 + \exp \left( \beta_0 + \sum_j \beta_j \mathbf{x}_j^{(i)} \right) \right]\end{aligned}$$

Training data  $\{(\mathbf{x}, y)\}$  are fixed. Objective function is a function of  $\beta$  ... what values of  $\beta$  give a good value?

$$\beta^* = \arg \min_{\beta} \mathcal{L}(\beta)$$

## Convexity

$\mathcal{L}(\beta)$  is convex for logistic regression.

Proof.

- Logistic loss  $-yv + \log(1 + \exp(v))$  is convex.
- Composition with linear function maintains convexity.
- Sum of convex functions is convex.

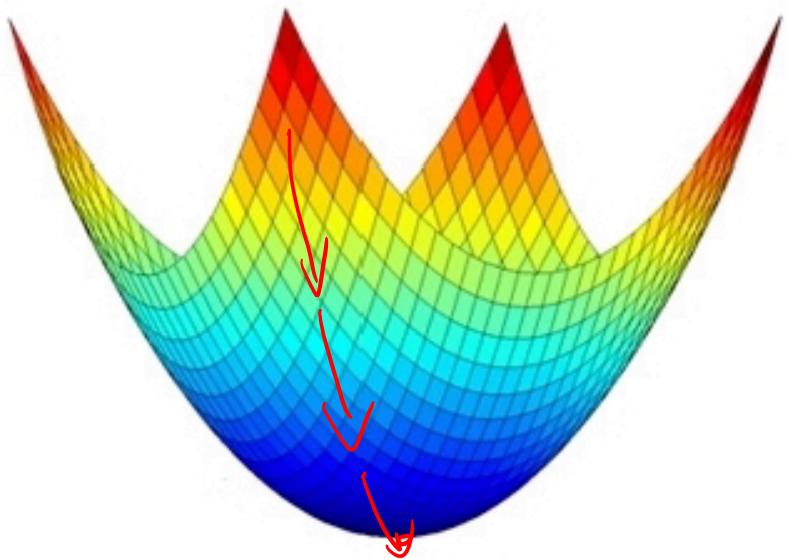
$$\frac{\partial^2 f}{\partial v^2} \geq 0$$



# Outline

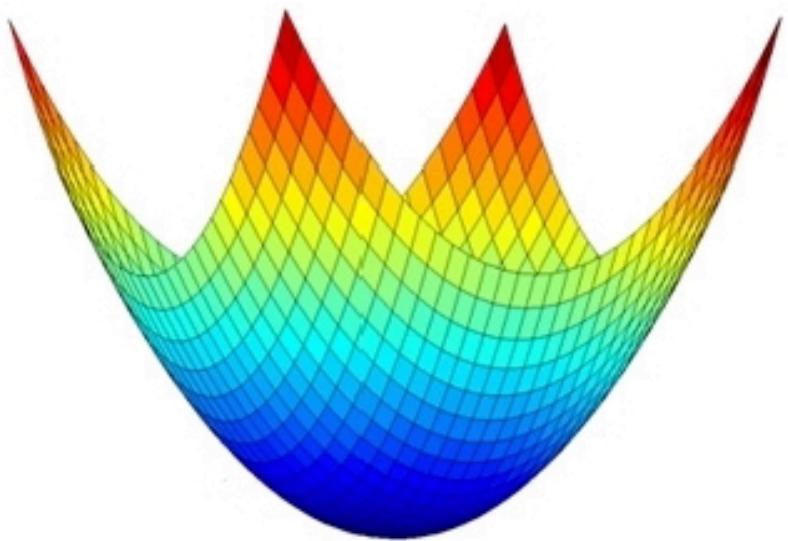
- Objective function
- Gradient descent
- Structural risk minimization
- Stochastic gradient descent

## Convexity



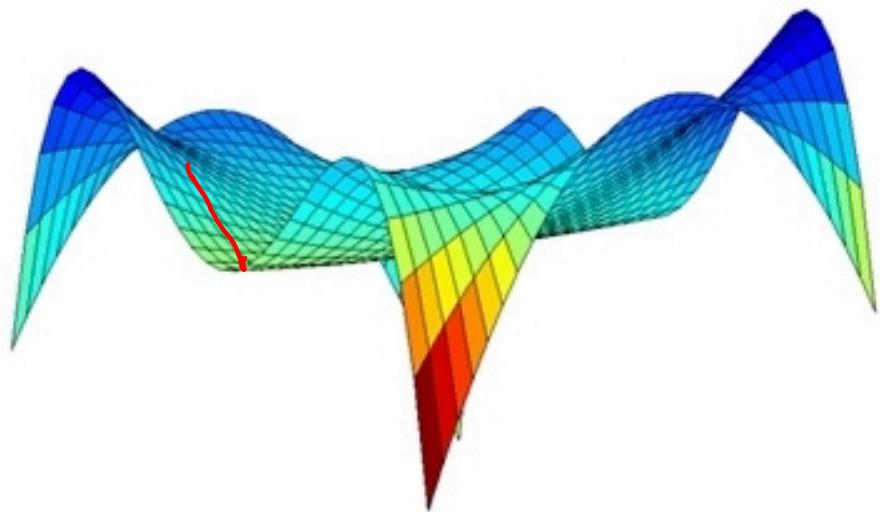
- Convex function
- Doesn't matter where you start, if you go down along the gradient

## Convexity



- Convex function
- Doesn't matter where you start, if you go down along the gradient
- Gradient!

## Convexity

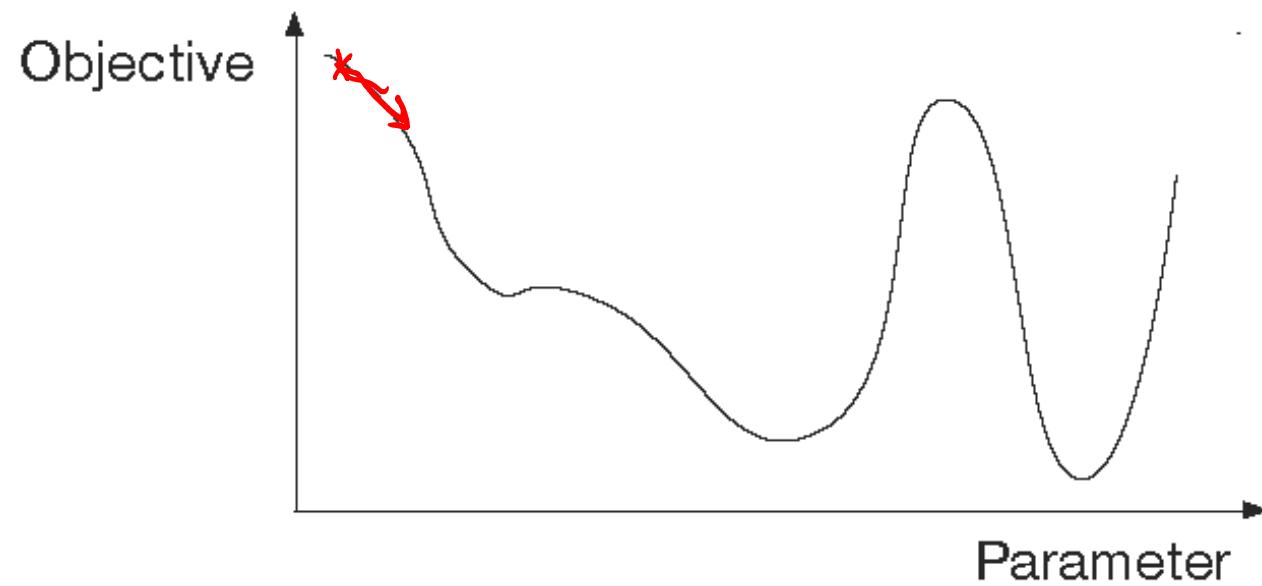


- It would have been much harder if this is not convex.

## Gradient Descent (non-convex)

Goal

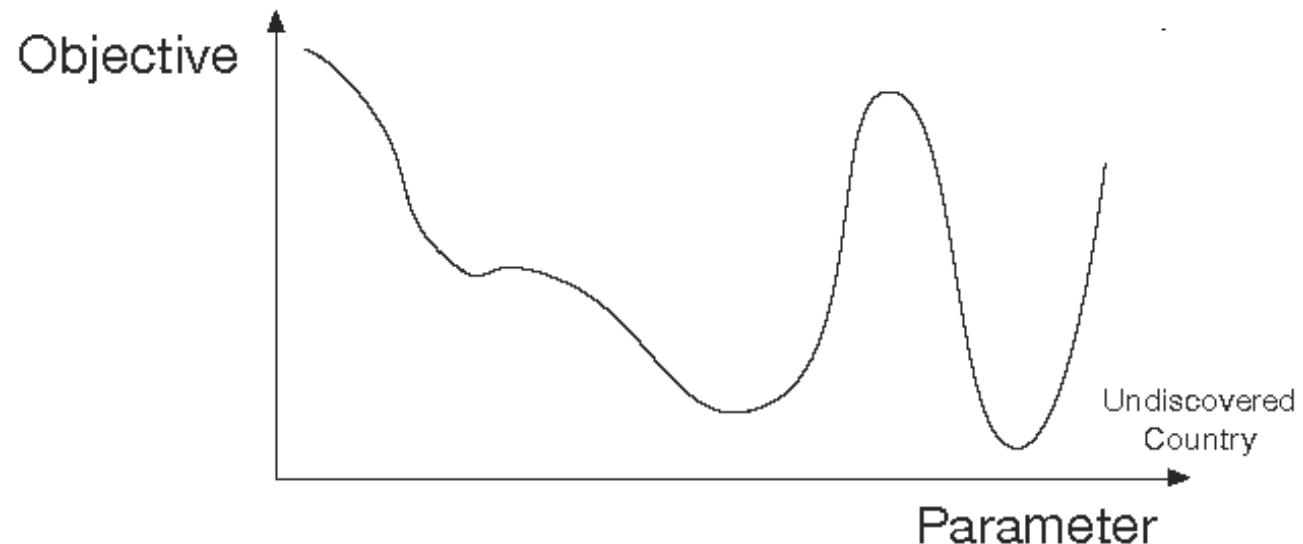
Optimize loss function with respect to variables  $\beta$



## Gradient Descent (non-convex)

Goal

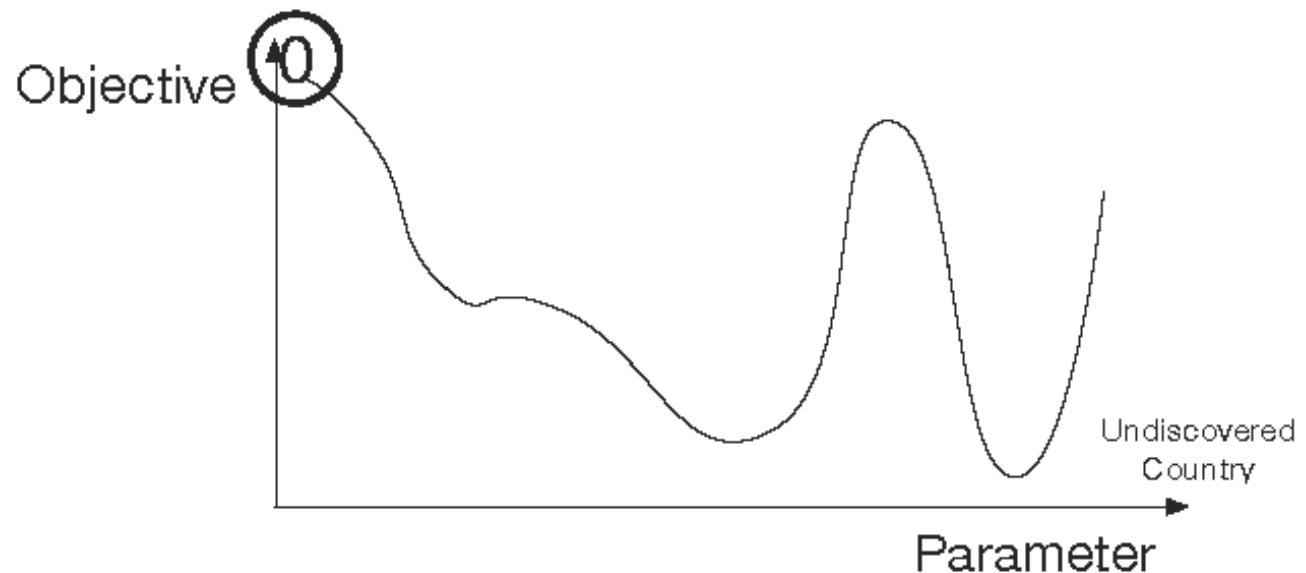
Optimize loss function with respect to variables  $\beta$



## Gradient Descent (non-convex)

Goal

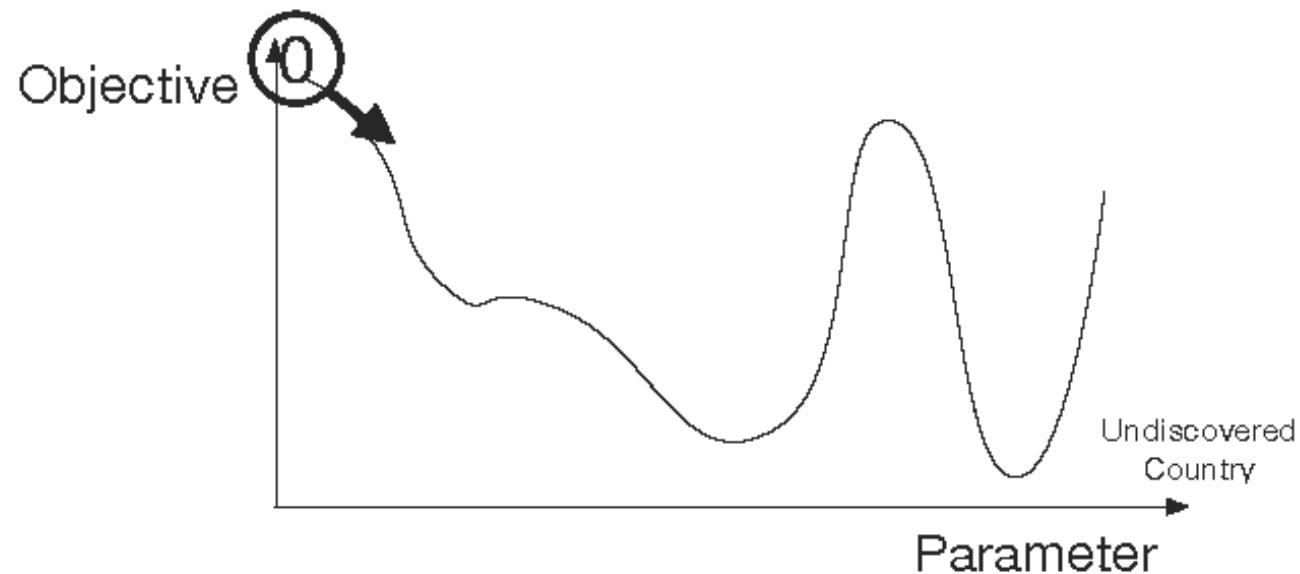
Optimize loss function with respect to variables  $\beta$



## Gradient Descent (non-convex)

Goal

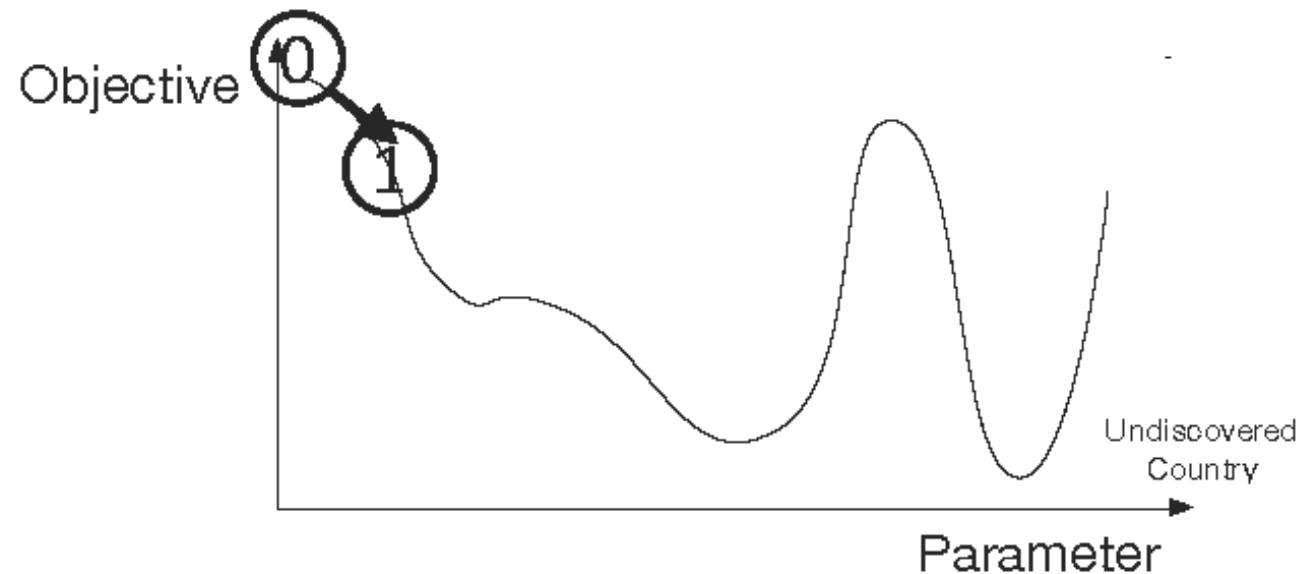
Optimize loss function with respect to variables  $\beta$



## Gradient Descent (non-convex)

Goal

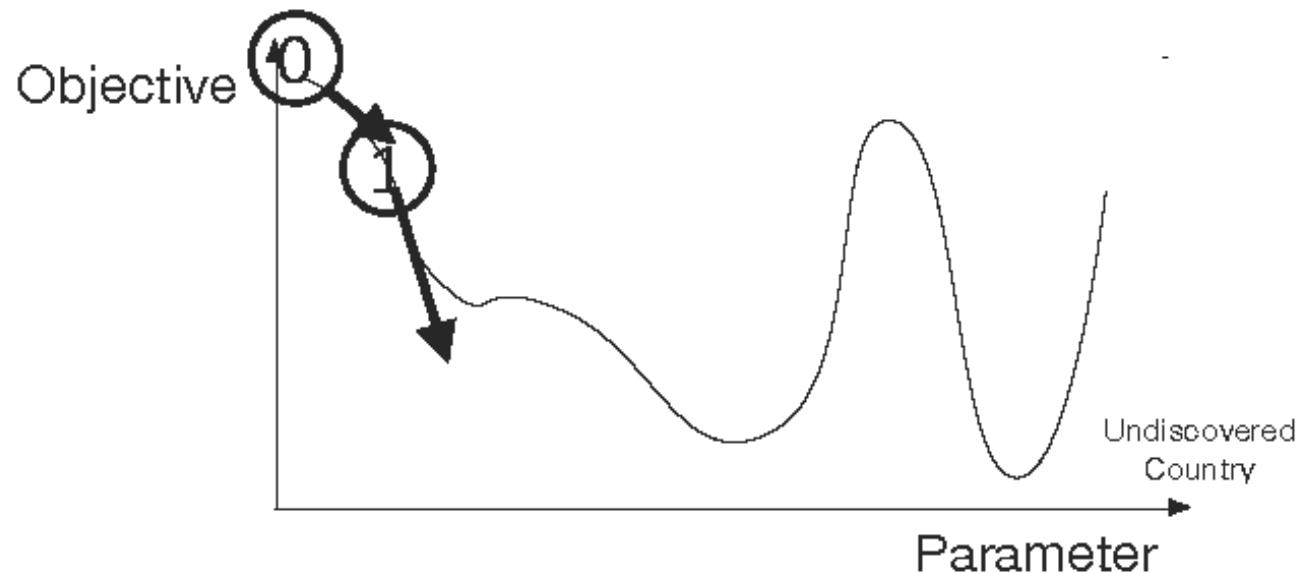
Optimize loss function with respect to variables  $\beta$



## Gradient Descent (non-convex)

Goal

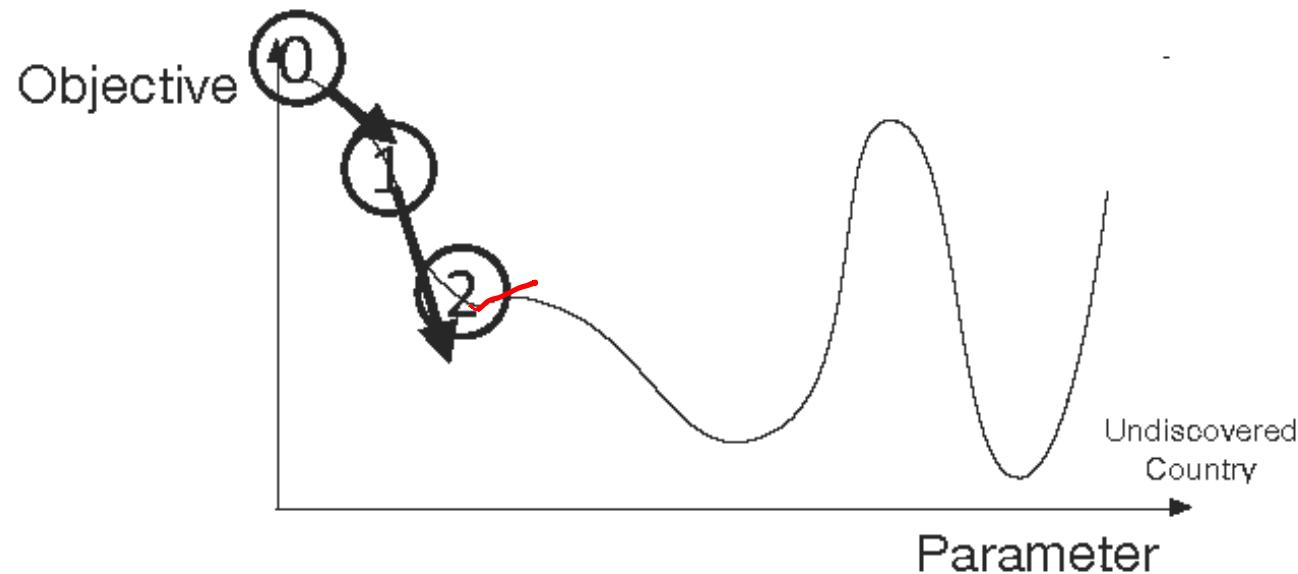
Optimize loss function with respect to variables  $\beta$



## Gradient Descent (non-convex)

Goal

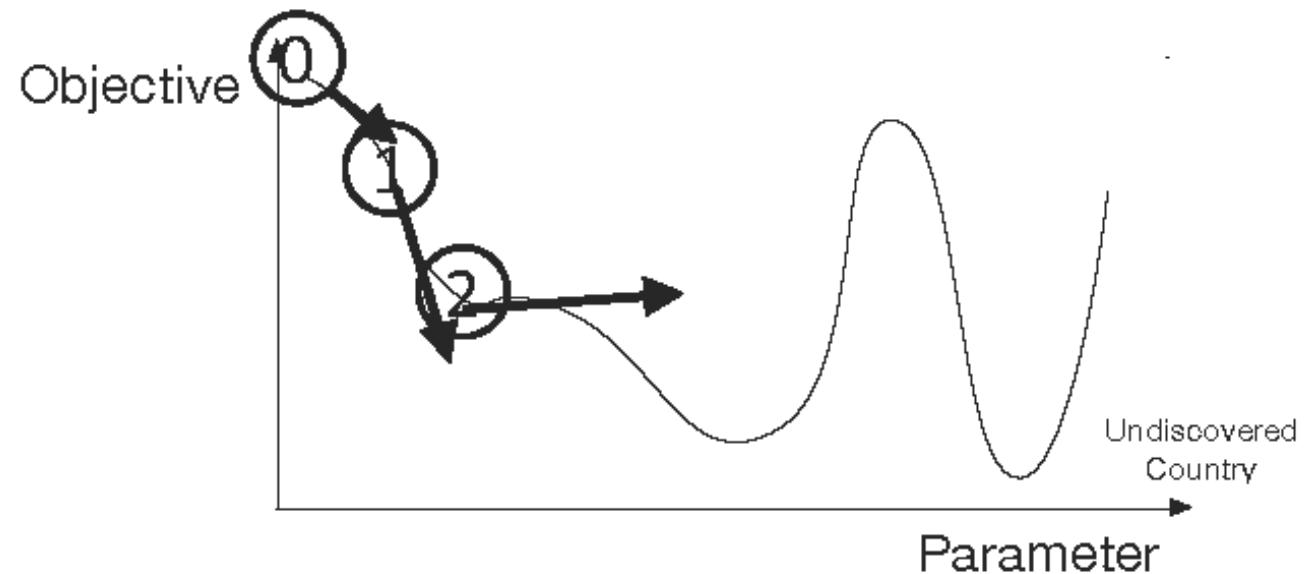
Optimize loss function with respect to variables  $\beta$



## Gradient Descent (non-convex)

Goal

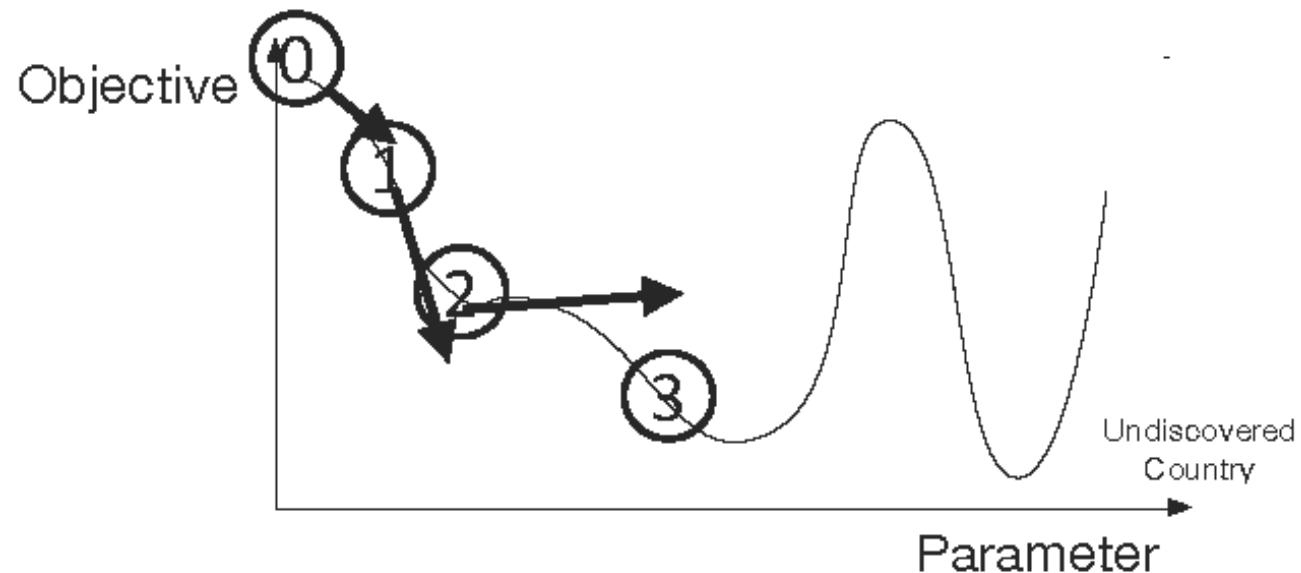
Optimize loss function with respect to variables  $\beta$



## Gradient Descent (non-convex)

Goal

Optimize loss function with respect to variables  $\beta$



## Gradient Descent (non-convex)

---

Goal

Optimize loss function with respect to variables  $\beta$

$$\underline{\beta_j^{l+1}} = \beta_j^l - \eta \frac{\partial \mathcal{L}}{\partial \underline{\beta_j}}$$

## Gradient Descent (non-convex)

---

Goal

Optimize loss function with respect to variables  $\beta$

$$\beta_j^{l+1} = \beta_j^l - \eta \frac{\partial \mathcal{L}}{\partial \beta_j}$$

Luckily, (vanilla) logistic regression is convex

## Gradient for Logistic Regression

To ease notation, let's define

$$\begin{aligned}
 \frac{\partial \pi_i}{\partial \beta} &= \frac{\partial \left( 1 - \frac{1}{1 + \exp(\beta^T x_i)} \right)}{\partial \beta} = - \cdot \left( -\frac{1}{(1 + \exp(\beta^T x))^2} \right) \cdot \frac{\partial (\exp(\beta^T x))}{\partial \beta} \\
 &= \frac{1}{(1 + \exp(\beta^T x))^2} \cdot \exp(\beta^T x) \cdot x \\
 &= \pi_i (1 - \pi_i) \cdot x
 \end{aligned} \tag{3}$$

Our objective function is

$$\mathcal{L} = - \sum_i \log p(y_i | x_i) = \sum_i \mathcal{L}_i = \sum_i \begin{cases} -\log \pi_i & \text{if } \underline{y_i = 1} \\ -\log(1 - \pi_i) & \text{if } \underline{y_i = 0} \end{cases} \tag{4}$$

$$\frac{\partial \mathcal{L}}{\partial \pi_i} = \begin{cases} -\frac{1}{\pi_i} & \text{if } y_i = 1 \\ \frac{1}{1 - \pi_i} & \text{if } y_i = 0 \end{cases}$$

$$\begin{aligned}
 \frac{\partial \mathcal{L}}{\partial \beta} &= \frac{\partial \mathcal{L}}{\partial \pi_i} \frac{\partial \pi_i}{\partial \beta} = \begin{cases} -(1 - \pi_i)x & \text{if } y_i = 1 \\ \pi_i \cdot x & \text{if } y_i = 0 \end{cases} \\
 &= -(y_i - \pi_i)x
 \end{aligned}$$

## Taking the Derivative

---

Apply chain rule:

$$\frac{\partial \mathcal{L}}{\partial \beta_j} = \sum_i \frac{\partial \mathcal{L}_i(\vec{\beta})}{\partial \beta_j} = \sum_i \begin{cases} -\frac{1}{\pi_i} \frac{\partial \pi_i}{\partial \beta_j} & \text{if } y_i = 1 \\ -\frac{1}{1-\pi_i} \left( -\frac{\partial \pi_i}{\partial \beta_j} \right) & \text{if } y_i = 0 \end{cases} \quad (5)$$

If we plug in the derivative,

$$\frac{\partial \pi_i}{\partial \beta_j} = \pi_i(1 - \pi_i)x_j, \quad (6)$$

we can merge these two cases

$$\frac{\partial \mathcal{L}_i}{\partial \underline{\beta_j}} = -(y_i - \pi_i)\underline{x_j}. \quad (7)$$

## Gradient for Logistic Regression

### Gradient

$$\nabla_{\beta} \mathcal{L}(\vec{\beta}) = \left[ \frac{\partial \mathcal{L}(\vec{\beta})}{\partial \beta_0}, \dots, \frac{\partial \mathcal{L}(\vec{\beta})}{\partial \beta_n} \right] \quad (8)$$

### Update

$$\Delta \beta \equiv \eta \nabla_{\beta} \mathcal{L}(\vec{\beta}) \quad (9)$$

$$\beta'_i \leftarrow \beta_i - \eta \frac{\partial \mathcal{L}(\vec{\beta})}{\partial \beta_i} \quad (10)$$

## Gradient for Logistic Regression

### Gradient

$$\nabla_{\beta} \mathcal{L}(\vec{\beta}) = \left[ \frac{\partial \mathcal{L}(\vec{\beta})}{\partial \beta_0}, \dots, \frac{\partial \mathcal{L}(\vec{\beta})}{\partial \beta_n} \right] \quad (8)$$

### Update

$$\Delta \beta \equiv \eta \nabla_{\beta} \mathcal{L}(\vec{\beta}) \quad (9)$$

$$\beta'_i \leftarrow \beta_i - \eta \frac{\partial \mathcal{L}(\vec{\beta})}{\partial \beta_i} \quad (10)$$

$\eta$ : step size, must be greater than zero

## Gradient for Logistic Regression

### Gradient

$$\nabla_{\beta} \mathcal{L}(\vec{\beta}) = \left[ \frac{\partial \mathcal{L}(\vec{\beta})}{\partial \beta_0}, \dots, \frac{\partial \mathcal{L}(\vec{\beta})}{\partial \beta_n} \right] \quad (8)$$

### Update

$$\Delta \beta \equiv \eta \nabla_{\beta} \mathcal{L}(\vec{\beta}) \quad (9)$$

$$\beta'_i \leftarrow \beta_i - \eta \frac{\partial \mathcal{L}(\vec{\beta})}{\partial \beta_i} \quad (10)$$

## Overfitting

- It is not ideal to maximize the likelihood of training data

## Overfitting

- It is not ideal to maximize the likelihood of training data
  - When to stop?
  - Simple models (avoid  $\beta$  to get too big)

## Overfitting

- It is not ideal to maximize the likelihood of training data
  - When to stop?
  - Simple models (avoid  $\beta$  to get too big)

**Regularization**

$$\sum \beta^2$$

# Outline

- Objective function
- Gradient descent
- Structural risk minimization
- Stochastic gradient descent

## Regularized Conditional Log Likelihood

Unregularized

$$\beta^* = \underbrace{\arg \min_{\beta} -\ln [p(y^{(j)} | x^{(j)}, \beta)]}_{(11)}$$

Regularized

$$\beta^* = \underbrace{\arg \min_{\beta} -\ln [p(y^{(j)} | x^{(j)}, \beta)] + \frac{1}{2}\lambda \sum_i \beta_i^2}_{(12)}$$

## Regularized Conditional Log Likelihood

$$\max \ln \frac{P(y^c | x^{cd}, \beta)}{\tau}$$

Unregularized

$$\beta^* = \arg \min_{\beta} -\ln [p(y^j | x^{(j)}, \beta)] \quad (11)$$

Regularized

$$\beta^* = \arg \min_{\beta} -\ln [p(y^j | x^{(j)}, \beta)] + \frac{1}{2}\lambda \sum_i \beta_i^2 \quad (12)$$

$\lambda$  is the “regularization” parameter (a hyperparameter) that trades off between likelihood and having small parameters

## Alternative view of regularization

Can also get to regularization by putting prior beliefs on parameters

$$\underline{p(\beta \mid \mathcal{D})} \propto \underline{p(\mathcal{D} \mid \beta)} \underline{p(\beta)}$$

Then MAP estimate for  $\beta$  is  $\hat{\beta}$  which maximizes posterior

## Alternative view of regularization

---

Can also get to regularization by putting prior beliefs on parameters

$$p(\beta \mid \mathcal{D}) \propto p(\mathcal{D} \mid \beta)p(\beta)$$

Then MAP estimate for  $\beta$  is  $\hat{\beta}$  which maximizes posterior

Ridge: Assume Gaussian prior  $p(\beta_j) = \mathcal{N}(\beta_j \mid 0, \tau^2)$ , we will obtain the same regularized objective function

You can learn more about this view in “Bayesian statistics”

## Risk minimization

---

$$\min_{\beta} \sum_i \ell(y^{(i)}, h_{\beta}(\mathbf{x}^{(i)})) + \lambda R(\beta) \geq 0$$

## Risk minimization

---

$$\min_{\beta} \sum_i \ell(y^{(i)}, h_{\beta}(\mathbf{x}^{(i)})) + \lambda R(\beta)$$

### Loss functions ( $\ell$ )

Describe how well the model fits the training data

- $-y\hat{y} + \log(1 + \exp(\hat{y}))$

### Regularization ( $R$ )

Control the complexity of the model

- $||\beta||^2 = \sum_j \beta_j^2$

## Risk minimization

$$\min_{\beta} \sum_i \ell(y^{(i)}, h_{\beta}(\mathbf{x}^{(i)})) + \lambda \underline{R(\beta)}$$

### Loss functions ( $\ell$ )

Describe how well the model fits the training data

- $-y\hat{y} + \log(1 + \exp(\hat{y}))$
- $(y - \hat{y})^2$
- $\max\{0, 1 - y\hat{y}\}$

### Regularization ( $R$ )

Control the complexity of the model

- $\|\beta\|^2 = \sum_j \beta_j^2$
- $\|\beta\|_p = \left( \sum_j |\beta_j|^p \right)^{\frac{1}{p}}$ 
  - $\ell_1$ -regularization:  $\sum_j |\beta_j|$

# Outline

- Objective function
- Gradient descent
- Structural risk minimization
- Stochastic gradient descent

## Approximating the Gradient

---

- Our datasets are big (to fit into memory)
- ... or data are changing / streaming

## Approximating the Gradient

---

- Our datasets are big (to fit into memory)
- ... or data are changing / streaming
- Hard to compute true gradient

$$\mathcal{L}(\beta) \equiv \mathbb{E}_{\mathbf{x}} [\nabla \mathcal{L}(\beta, \mathbf{x})] \quad (1)$$

- Average over all observations

## Approximating the Gradient

---

- Our datasets are big (to fit into memory)
- ... or data are changing / streaming
- Hard to compute true gradient

$$\mathcal{L}(\beta) \equiv \mathbb{E}_{\mathbf{x}} [\nabla \mathcal{L}(\beta, \mathbf{x})] \quad (1)$$

- Average over all observations
- What if we compute an update just from one observation?

## Getting to Union Station

Pretend it's a pre-smartphone world and you want to get to Union Station



## Stochastic Gradient for Regularized Regression

---

$$\mathcal{L} = -\log p(y \mid \mathbf{x}; \boldsymbol{\beta}) + \frac{1}{2}\lambda \sum_j \beta_j^2 \quad (2)$$

## Stochastic Gradient for Regularized Regression

---

$$\mathcal{L} = -\log p(y | \mathbf{x}; \beta) + \frac{1}{2} \lambda \sum_j \beta_j^2 \quad (2)$$

Taking the derivative (with respect to example  $\underline{x_i}$ )

$$\frac{\partial \mathcal{L}}{\partial \beta_j} = - \underline{(y_i - \pi_i)x_{ij}} + \lambda \beta_j \quad (3)$$

## Stochastic Gradient for Logistic Regression

---

Given a **single observation**  $x_i$  chosen at random from the dataset,

$$\beta_j \leftarrow \beta'_j - \eta \underbrace{(\lambda\beta'_j - x_{ij} [y_i - \pi_i])}_{(4)}$$

## Example Documents

$$\beta_j = \beta_j + \eta(y_i - \pi_i)x_{ij}$$
$$\vec{\beta} = \langle \beta_{bias} = 0, \beta_A = 0, \beta_B = 0, \beta_C = 0, \beta_D = 0 \rangle$$
$$\frac{\exp(\vec{\beta}^T \vec{x})}{1 + \exp(\vec{\beta}^T \vec{x})}$$

$y_1 = 1$

A A A A B B B C

(Assume step size  $\eta = 1.0$ )

$y_2 = 0$

B C C C D D D D

You first see the positive example. First, compute  $\pi_1 = 0.5$

$\beta_{bias} = 0.5$

## Example Documents

---

$$\beta_j = \beta_j + \eta(y_i - \pi_i)x_{ij}$$

$$\vec{\beta} = \langle 0, 0, 0, 0, 0 \rangle$$

$$y_1 = 1$$

A A A A B B B C

(Assume step size  $\eta = 1.0$ .)

$$y_2 = 0$$

B C C C D D D D

You first see the positive example. First, compute  $\pi_1$

$$\pi_1 = \Pr(y_1 = 1 | \mathbf{x}_1) = \frac{\exp \beta^T \mathbf{x}_i}{1 + \exp \beta^T \mathbf{x}_i} =$$

## Example Documents

---

$$\beta_j = \beta_j + \eta(y_i - \pi_i)x_{ij}$$

$$\vec{\beta} = \langle 0, 0, 0, 0, 0 \rangle$$

$y_1 = 1$

A A A A B B B C

(Assume step size  $\eta = 1.0.$ )

$y_2 = 0$

B C C C D D D D

You first see the positive example. First, compute  $\pi_1$

$$\pi_1 = \Pr(y_1 = 1 | \mathbf{x}_1) = \frac{\exp \beta^T \mathbf{x}_i}{1 + \exp \beta^T \mathbf{x}_i} = \frac{\exp 0}{\exp 0 + 1} = 0.5$$

## Example Documents

---

$$\begin{aligned}\beta_j &= \beta_j + \eta(y_i - \pi_i)x_{ij} \\ \vec{\beta} &= \langle 0, 0, 0, 0, 0 \rangle\end{aligned}$$

$y_1 = 1$

A A A A B B B C

(Assume step size  $\eta = 1.0$ .)

$y_2 = 0$

B C C C D D D D

$\pi_1 = 0.5$  What's the update for  $\beta_{bias}$ ?

## Example Documents

---

$$\beta_j = \beta_j + \eta(y_i - \pi_i)x_{ij}$$

$$\vec{\beta} = \langle 0, 0, 0, 0, 0 \rangle$$

$y_1 = 1$

A A A A B B B C

(Assume step size  $\eta = 1.0$ .)

$y_2 = 0$

B C C C D D D D

What's the update for  $\beta_{bias}$ ?

$$\beta_{bias} = \beta'_{bias} + \eta \cdot (y_1 - \pi_1) \cdot x_{1,bias} = 0.0 + 1.0 \cdot (1.0 - 0.5) \cdot 1.0$$

## Example Documents

---

$$\beta_j = \beta_j + \eta(y_i - \pi_i)x_{ij}$$

$$\vec{\beta} = \langle 0, 0, 0, 0, 0 \rangle$$

$$y_1 = 1$$

AAAAA BBBB C

(Assume step size  $\eta = 1.0$ .)

$$\begin{matrix} A:4 \\ B:3 \\ C:1 \end{matrix}$$

$$y_2 = 0$$

B C C C D D D D

What's the update for  $\beta_{bias}$ ?

$$\beta_{bias} = \beta'_{bias} + \eta \cdot (y_1 - \pi_1) \cdot x_{1,bias} = 0.0 + 1.0 \cdot (1.0 - 0.5) \cdot 1.0 = 0.5$$

## Example Documents

---

$$\begin{aligned}\beta_j &= \beta_j + \eta(y_i - \pi_i)x_{ij} \\ \vec{\beta} &= \langle 0, 0, 0, 0, 0 \rangle\end{aligned}$$

$y_1 = 1$

A A A A B B B C

(Assume step size  $\eta = 1.0$ .)

$y_2 = 0$

B C C C D D D D

What's the update for  $\beta_A$ ?

## Example Documents

---

$$\begin{aligned}\beta_j &= \beta_j + \eta(y_i - \pi_i)x_{ij} \\ \vec{\beta} &= \langle 0, 0, 0, 0, 0 \rangle\end{aligned}$$

$y_1 = 1$

A A A A B B B C

(Assume step size  $\eta = 1.0$ .)

$y_2 = 0$

B C C C D D D D

What's the update for  $\beta_A$ ?     $\beta_A = \beta'_A + \eta \cdot (y_1 - \pi_1) \cdot x_{1,A} = 0.0 + 1.0 \cdot (1.0 - 0.5) \cdot 4.0$

## Example Documents

---

$$\begin{aligned}\beta_j &= \beta_j + \eta(y_i - \pi_i)x_{ij} \\ \vec{\beta} &= \langle 0, 0, 0, 0, 0 \rangle\end{aligned}$$

$y_1 = 1$

A A A A B B B C

(Assume step size  $\eta = 1.0$ .)

$y_2 = 0$

B C C C D D D D

What's the update for  $\beta_A$ ?     $\beta_A = \beta'_A + \eta \cdot (y_1 - \pi_1) \cdot x_{1,A} = 0.0 + 1.0 \cdot (1.0 - 0.5) \cdot 4.0 = 2.0$

## Example Documents

---

$$\begin{aligned}\beta_j &= \beta_j + \eta(y_i - \pi_i)x_{ij} \\ \vec{\beta} &= \langle 0, 0, 0, 0, 0 \rangle\end{aligned}$$

$y_1 = 1$

A A A A B B B C

(Assume step size  $\eta = 1.0$ .)

$y_2 = 0$

B C C C D D D D

What's the update for  $\beta_B$ ?

## Example Documents

---

$$\begin{aligned}\beta_j &= \beta_j + \eta(y_i - \pi_i)x_{ij} \\ \vec{\beta} &= \langle 0, 0, 0, 0, 0 \rangle\end{aligned}$$

$y_1 = 1$

A A A A B B B C

(Assume step size  $\eta = 1.0$ .)

$y_2 = 0$

B C C C D D D D

What's the update for  $\beta_B$ ?  $\beta_B = \beta'_B + \eta \cdot (y_1 - \pi_1) \cdot x_{1,B} = 0.0 + 1.0 \cdot (1.0 - 0.5) \cdot 3.0$

## Example Documents

---

$$\begin{aligned}\beta_j &= \beta_j + \eta(y_i - \pi_i)x_{ij} \\ \vec{\beta} &= \langle 0, 0, 0, 0, 0 \rangle\end{aligned}$$

$y_1 = 1$

A A A A B B B C

(Assume step size  $\eta = 1.0$ .)

$y_2 = 0$

B C C C D D D D

What's the update for  $\beta_B$ ?     $\beta_B = \beta'_B + \eta \cdot (y_1 - \pi_1) \cdot x_{1,B} = 0.0 + 1.0 \cdot (1.0 - 0.5) \cdot 3.0 = 1.5$

## Example Documents

---

$$\begin{aligned}\beta_j &= \beta_j + \eta(y_i - \pi_i)x_{ij} \\ \vec{\beta} &= \langle 0, 0, 0, 0, 0 \rangle\end{aligned}$$

$y_1 = 1$

A A A A B B B C

(Assume step size  $\eta = 1.0$ .)

$y_2 = 0$

B C C C D D D D

What's the update for  $\beta_C$ ?

## Example Documents

---

$$\beta_j = \beta_j + \eta(y_i - \pi_i)x_{ij}$$

$$\vec{\beta} = \langle 0, 0, 0, 0, 0 \rangle$$

$y_1 = 1$

A A A A B B B C

(Assume step size  $\eta = 1.0$ .)

$y_2 = 0$

B C C C D D D D

What's the update for  $\beta_C$ ?

$$\beta_C = \beta'_C + \eta \cdot (y_1 - \pi_1) \cdot x_{1,C} = 0.0 + 1.0 \cdot (1.0 - 0.5) \cdot 1.0$$

## Example Documents

---

$$\beta_j = \beta_j + \eta(y_i - \pi_i)x_{ij}$$

$$\vec{\beta} = \langle 0, 0, 0, 0, 0 \rangle$$

$y_1 = 1$

A A A A B B B C

(Assume step size  $\eta = 1.0$ .)

$y_2 = 0$

B C C C D D D D

What's the update for  $\beta_C$ ?

$$\beta_C = \beta'_C + \eta \cdot (y_1 - \pi_1) \cdot x_{1,C} = 0.0 + 1.0 \cdot (1.0 - 0.5) \cdot 1.0 = 0.5$$

## Example Documents

---

$$\begin{aligned}\beta_j &= \beta_j + \eta(y_i - \pi_i)x_{ij} \\ \vec{\beta} &= \langle 0, 0, 0, 0, 0 \rangle\end{aligned}$$

$y_1 = 1$

A A A A B B B C

(Assume step size  $\eta = 1.0$ .)

$y_2 = 0$

B C C C D D D D

What's the update for  $\beta_D$ ?

## Example Documents

---

$$\beta_j = \beta_j + \eta(y_i - \pi_i)x_{ij}$$

$$\vec{\beta} = \langle 0, 0, 0, 0, 0 \rangle$$

$y_1 = 1$

A A A A B B B C

(Assume step size  $\eta = 1.0$ .)

$y_2 = 0$

B C C C D D D D

What's the update for  $\beta_D$ ?

$$\beta_D = \beta'_D + \eta \cdot (y_1 - \pi_1) \cdot x_{1,D} = 0.0 + 1.0 \cdot (1.0 - 0.5) \cdot 0.0$$

## Example Documents

---

$$\beta_j = \beta_j + \eta(y_i - \pi_i)x_{ij}$$

$$\vec{\beta} = \langle 0, 0, 0, 0, 0 \rangle$$

$y_1 = 1$

A A A A B B B C

(Assume step size  $\eta = 1.0$ .)

$y_2 = 0$

B C C C D D D D

What's the update for  $\beta_D$ ?

$$\beta_D = \beta'_D + \eta \cdot (y_1 - \pi_1) \cdot x_{1,D} = 0.0 + 1.0 \cdot (1.0 - 0.5) \cdot 0.0 = 0.0$$

## Example Documents

---

$$\beta_j = \beta_j + \eta(y_i - \pi_i)x_{ij}$$

$$\vec{\beta} = \langle .5, 2, 1.5, 0.5, 0 \rangle$$

$$y_1 = 1$$

A A A A B B B C

(Assume step size  $\eta = 1.0$ .)

$$y_2 = 0$$

B C C C D D D D

Now you see the negative example. What's  $\pi_2$ ?

## Example Documents

---

$$\beta_j = \beta_j + \eta(y_i - \pi_i)x_{ij}$$

$$\vec{\beta} = \langle .5, 2, 1.5, 0.5, 0 \rangle$$

$$y_1 = 1$$

A A A A B B B C

(Assume step size  $\eta = 1.0$ .)

$$y_2 = 0$$

B C C C D D D D

Now you see the negative example. What's  $\pi_2$ ?

$$\pi_2 = \Pr(y_2 = 1 | \vec{x}_2) = \frac{\exp \beta^T x_i}{1 + \exp \beta^T x_i} = \frac{\exp\{.5+1.5+1.5+0\}}{\exp\{.5+1.5+1.5+0\}+1} =$$

## Example Documents

---

$$\beta_j = \beta_j + \eta(y_i - \pi_i)x_{ij}$$

$$\vec{\beta} = \langle .5, 2, 1.5, 0.5, 0 \rangle$$

$$y_1 = \underline{1}$$

A A A A B B B C

(Assume step size  $\eta = 1.0$ .)

$$y_2 = \underline{0}$$

B C C C D D D D

Now you see the negative example. What's  $\pi_2$ ?

$$\pi_2 = \Pr(y_2 = 1 \mid \vec{x}_2) = \frac{\exp \beta^T x_i}{1 + \exp \beta^T x_i} = \frac{\exp\{.5+1.5+1.5+0\}}{\exp\{.5+1.5+1.5+0\}+1} = 0.97$$

## Example Documents

---

$$\beta_j = \beta_j + \eta(y_i - \pi_i)x_{ij}$$

$$\vec{\beta} = \langle .5, 2, 1.5, 0.5, 0 \rangle$$

$$y_1 = 1$$

A A A A B B B C

(Assume step size  $\eta = 1.0.$ )

$$y_2 = 0$$

B C C C D D D D

Now you see the negative example. What's  $\pi_2$ ?

$$\pi_2 = 0.97$$

What's the update for  $\beta_{bias}$ ?

## Example Documents

---

$$\beta_j = \beta_j + \eta(y_i - \pi_i)x_{ij}$$

$$\vec{\beta} = \langle .5, 2, 1.5, 0.5, 0 \rangle$$

$$y_1 = 1$$

A A A A B B B C

(Assume step size  $\eta = 1.0$ .)

$$y_2 = 0$$

B C C C D D D D

What's the update for  $\beta_{bias}$ ?

$$\beta_{bias} = \beta'_{bias} + \eta \cdot (y_2 - \pi_2) \cdot x_{2,bias} = 0.5 + 1.0 \cdot (0.0 - 0.97) \cdot 1.0$$

## Example Documents

---

$$\beta_j = \beta_j + \eta(y_i - \pi_i)x_{ij}$$

$$\vec{\beta} = \langle .5, 2, 1.5, 0.5, 0 \rangle$$

$$y_1 = 1$$

A A A A B B B C

(Assume step size  $\eta = 1.0$ .)

$$y_2 = 0$$

B C C C D D D D

What's the update for  $\beta_{bias}$ ?

-0.97

$$\beta_{bias} = \beta'_{bias} + \eta \cdot (y_2 - \pi_2) \cdot x_{2,bias} = 0.5 + 1.0 \cdot (0.0 - 0.97) \cdot 1.0 = -0.47$$

## Example Documents

---

$$\begin{aligned}\beta_j &= \beta_j + \eta(y_i - \pi_i)x_{ij} \\ \vec{\beta} &= \langle .5, 2, 1.5, 0.5, 0 \rangle\end{aligned}$$

$y_1 = 1$

A A A A B B B C

(Assume step size  $\eta = 1.0$ .)

$y_2 = 0$

B C C C D D D D

What's the update for  $\beta_A$ ?

## Example Documents

---

$$\beta_j = \beta_j + \eta(y_i - \pi_i)x_{ij}$$

$$\vec{\beta} = \langle .5, 2, 1.5, 0.5, 0 \rangle$$

$$y_1 = 1$$

A A A A B B B C

(Assume step size  $\eta = 1.0$ .)

$$y_2 = 0$$

B C C C D D D D

What's the update for  $\beta_A$ ?

$$\beta_A = \beta'_A + \eta \cdot (y_2 - \pi_2) \cdot x_{2,A} = 2.0 + 1.0 \cdot (0.0 - 0.97) \cdot 0.0$$

## Example Documents

---

$$\beta_j = \beta_j + \eta(y_i - \pi_i)x_{ij}$$

$$\vec{\beta} = \langle .5, 2, 1.5, 0.5, 0 \rangle$$

$$y_1 = 1$$

A A A A B B B C

(Assume step size  $\eta = 1.0$ .)

$$y_2 = 0$$

B C C C D D D D

What's the update for  $\beta_A$ ?

$$\beta_A = \beta'_A + \eta \cdot (y_2 - \pi_2) \cdot x_{2,A} = 2.0 + 1.0 \cdot (0.0 - 0.97) \cdot 0.0 = 2.0$$

## Example Documents

---

$$\begin{aligned}\beta_j &= \beta_j + \eta(y_i - \pi_i)x_{ij} \\ \vec{\beta} &= \langle .5, 2, 1.5, 0.5, 0 \rangle\end{aligned}$$

$y_1 = 1$

A A A A B B B C

(Assume step size  $\eta = 1.0$ .)

$y_2 = 0$

B C C C D D D D

What's the update for  $\beta_B$ ?

## Example Documents

---

$$\beta_j = \beta_j + \eta(y_i - \pi_i)x_{ij}$$

$$\vec{\beta} = \langle .5, 2, 1.5, 0.5, 0 \rangle$$

$$y_1 = 1$$

A A A A B B B C

(Assume step size  $\eta = 1.0$ .)

$$y_2 = 0$$

B C C C D D D D

What's the update for  $\beta_B$ ?

$$\beta_B = \beta'_B + \eta \cdot (y_2 - \pi_2) \cdot x_{2,B} = 1.5 + 1.0 \cdot (0.0 - 0.97) \cdot 1.0$$

## Example Documents

---

$$\beta_j = \beta_j + \eta(y_i - \pi_i)x_{ij}$$

$$\vec{\beta} = \langle .5, 2, 1.5, 0.5, 0 \rangle$$

$$y_1 = 1$$

A A A A B B B C

(Assume step size  $\eta = 1.0$ .)

$$y_2 = 0$$

B C C C D D D D

What's the update for  $\beta_B$ ?

$$\beta_B = \beta'_B + \eta \cdot (y_2 - \pi_2) \cdot x_{2,B} = 1.5 + 1.0 \cdot (0.0 - 0.97) \cdot 1.0 = 0.53$$

## Example Documents

---

$$\begin{aligned}\beta_j &= \beta_j + \eta(y_i - \pi_i)x_{ij} \\ \vec{\beta} &= \langle .5, 2, 1.5, 0.5, 0 \rangle\end{aligned}$$

$y_1 = 1$

A A A A B B B C

(Assume step size  $\eta = 1.0$ .)

$y_2 = 0$

B C C C D D D D

What's the update for  $\beta_C$ ?

## Example Documents

---

$$\beta_j = \beta_j + \eta(y_i - \pi_i)x_{ij}$$

$$\vec{\beta} = \langle .5, 2, 1.5, 0.5, 0 \rangle$$

$$y_1 = 1$$

A A A A B B B C

(Assume step size  $\eta = 1.0$ .)

$$y_2 = 0$$

B C C C D D D D

What's the update for  $\beta_C$ ?

$$\beta_C = \beta'_C + \eta \cdot (y_2 - \pi_2) \cdot x_{2,C} = 0.5 + 1.0 \cdot (0.0 - 0.97) \cdot 3.0$$

## Example Documents

---

$$\beta_j = \beta_j + \eta(y_i - \pi_i)x_{ij}$$

$$\vec{\beta} = \langle .5, 2, 1.5, 0.5, 0 \rangle$$

$$y_1 = 1$$

A A A A B B B C

(Assume step size  $\eta = 1.0$ .)

$$y_2 = 0$$

B C C C D D D D

What's the update for  $\beta_C$ ?

$$\beta_C = \beta'_C + \eta \cdot (y_2 - \pi_2) \cdot x_{2,C} = 0.5 + 1.0 \cdot (0.0 - 0.97) \cdot 3.0 = -2.41$$

## Example Documents

---

$$\begin{aligned}\beta_j &= \beta_j + \eta(y_i - \pi_i)x_{ij} \\ \vec{\beta} &= \langle .5, 2, 1.5, 0.5, 0 \rangle\end{aligned}$$

$y_1 = 1$

A A A A B B B C

(Assume step size  $\eta = 1.0$ .)

$y_2 = 0$

B C C C D D D D

What's the update for  $\beta_D$ ?

## Example Documents

---

$$\beta_j = \beta_j + \eta(y_i - \pi_i)x_{ij}$$

$$\vec{\beta} = \langle .5, 2, 1.5, 0.5, 0 \rangle$$

$$y_1 = 1$$

A A A A B B B C

(Assume step size  $\eta = 1.0$ .)

$$y_2 = 0$$

B C C C D D D D

What's the update for  $\beta_D$ ?

$$\beta_D = \beta'_D + \eta \cdot (y_2 - \pi_2) \cdot x_{2,D} = 0.0 + 1.0 \cdot (0.0 - 0.97) \cdot 4.0$$

## Example Documents

---

$$\beta_j = \beta_j + \eta(y_i - \pi_i)x_{ij}$$
$$\vec{\beta} = \langle .5, 2, 1.5, 0.5, 0 \rangle$$

$$y_1 = 1$$

AAAAA BBB C

(Assume step size  $\eta = 1.0$ .)

$$y_2 = 0$$

B C C C D D D D

What's the update for  $\beta_D$ ?

$$\beta_D = \beta'_D + \eta \cdot (y_2 - \pi_2) \cdot x_{2,D} = 0.0 + 1.0 \cdot (0.0 - 0.97) \cdot 4.0 = -3.88$$

## Algorithm

---

1. Initialize a vector  $\beta$  to be all zeros
2. For  $t = 1, \dots, T$ 
  - For each example  $x_i, y_i$  and feature  $j$ :
    - Compute  $\pi_i \equiv \Pr(y_i = 1 | x_i)$
    - Set  $\beta_j = \beta'_j - \eta(\lambda\beta'_j + (y_i - \pi_i)x_i)$
3. Output the parameters  $\beta_1, \dots, \beta_d$ .

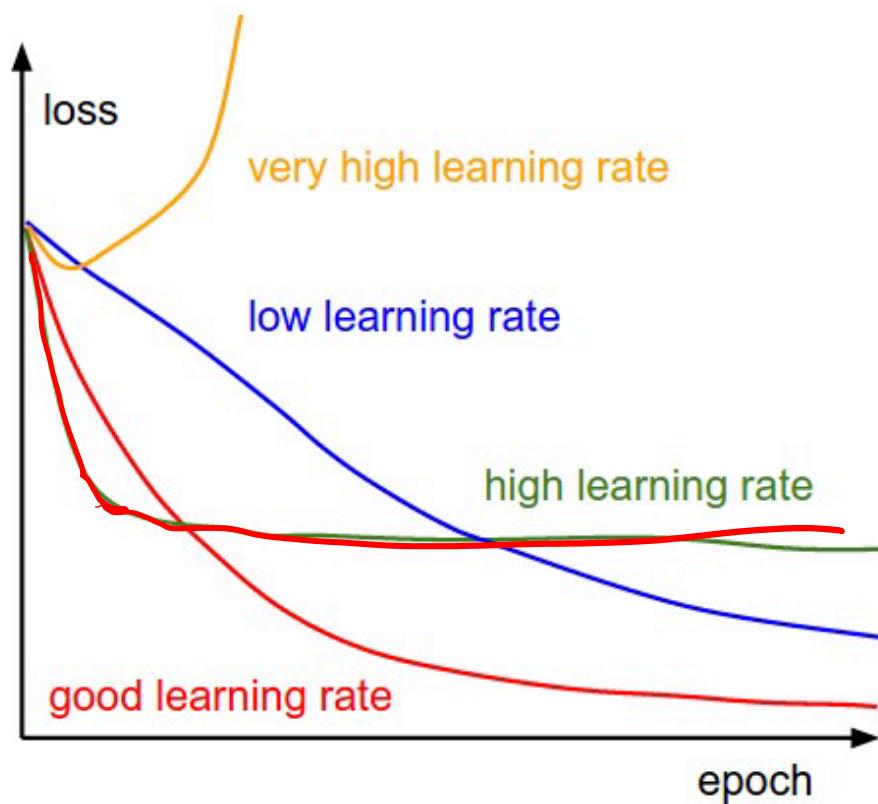
## Algorithm

---

1. Initialize a vector  $\beta$  to be all zeros
2. For  $t = 1, \dots, T$ 
  - For each example  $x_i, y_i$  and feature  $j$ :
    - Compute  $\pi_i \equiv \Pr(y_i = 1 | x_i)$
    - Set  $\beta_j = \beta'_j - \eta(\lambda\beta'_j - (y_i - \pi_i)x_i)$
3. Output the parameters  $\beta_1, \dots, \beta_d$ .

How to decide  $\eta$ ?

## Choosing learning rate



## Learning rate decay

- Decay after each epoch (e.g.,  $\frac{\eta_0}{t^2}$ ,  $\eta_0 e^{-kt}$ )
- Decay after each example (e.g.,  $\frac{\eta_0}{1+kn}$ )

Decay schedule can be seen as a hyperparameter too.



## Learning rate decay

- Decay after each epoch (e.g.,  $\frac{\eta_0}{t^2}$ ,  $\eta_0 e^{-kt}$ )
- Decay after each example (e.g.,  $\frac{\eta_0}{1+kn}$ )

Decay schedule can be seen as a hyperparameter too.

Advanced stochastic gradient descent:

<http://ruder.io/optimizing-gradient-descent/>

*Adam*

# Recap

- Follow the gradient to fit the logistic regression model
- Most machine learning methods fall into the framework of (loss + regularization)
- Stochastic gradient descent allows for approximating the gradient from one instance