# Classification: Logistic Regression from Data

Machine Learning: Jordan Boyd-Graber
University of Colorado Boulder
LECTURE 3

Slides adapted from Emily Fox

$$P(Y = 0|X) = \frac{1}{1 + \exp\left[\beta_0 + \sum_i \beta_i X_i\right]} \tag{1}$$

$$P(Y = 1|X) = \frac{\exp\left[\beta_0 + \sum_i \beta_i X_i\right]}{1 + \exp\left[\beta_0 + \sum_i \beta_i X_i\right]} \tag{2}$$
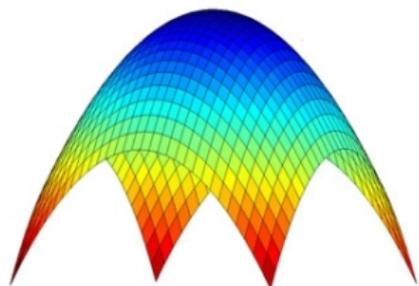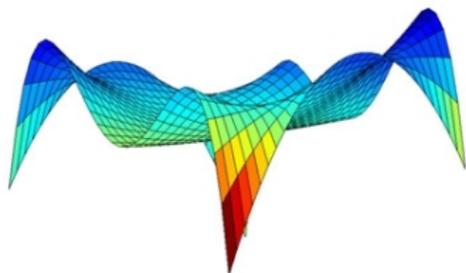
- Discriminative prediction: $p(y|x)$
- Classification uses: ad placement, spam detection
- What we didn't talk about is how to learn $\beta$ from data

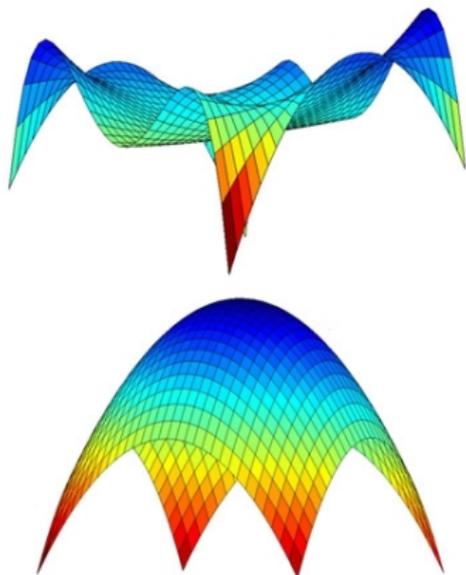$$\mathcal{L} \equiv \ln p(Y \mid X, \beta) = \sum_j \ln p(y^{(j)} \mid x^{(j)}, \beta) \tag{3}$$

$$= \sum_j y^{(j)} \left( \beta_0 + \sum_i \beta_i x_i^{(j)} \right) - \ln \left[ 1 + \exp \left( \beta_0 + \sum_i \beta_i x_i^{(j)} \right) \right]$$

$$\tag{4}$$

- Convex function
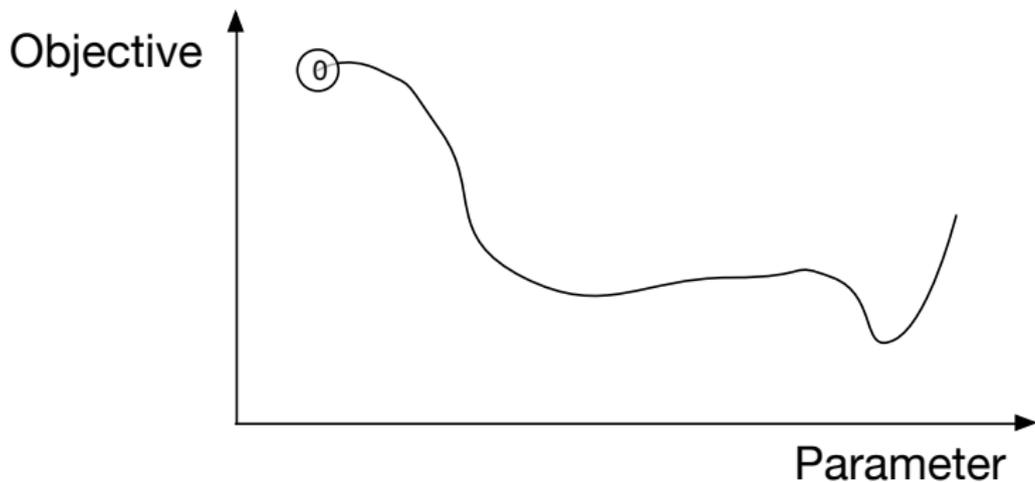- Doesn't matter where you start, if you "walk up" objective

- Convex function
- Doesn't matter where you start, if you "walk up" objective
- Gradient!

**Goal**

Optimize log likelihood with respect to variables $W$ and $b$

**Goal**

Optimize log likelihood with respect to variables $W$ and $b$

**Goal**

Optimize log likelihood with respect to variables *W* and *b*

**Goal**
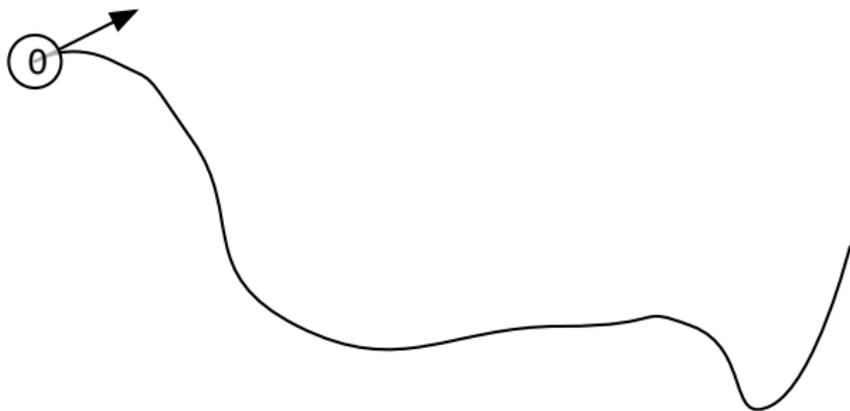
Optimize log likelihood with respect to variables *W* and *b*

**Goal**

Optimize log likelihood with respect to variables *W* and *b*

**Goal**

Optimize log likelihood with respect to variables *W* and *b*

**Goal**

Optimize log likelihood with respect to variables *W* and *b*

**Goal**

Optimize log likelihood with respect to variables *W* and *b*
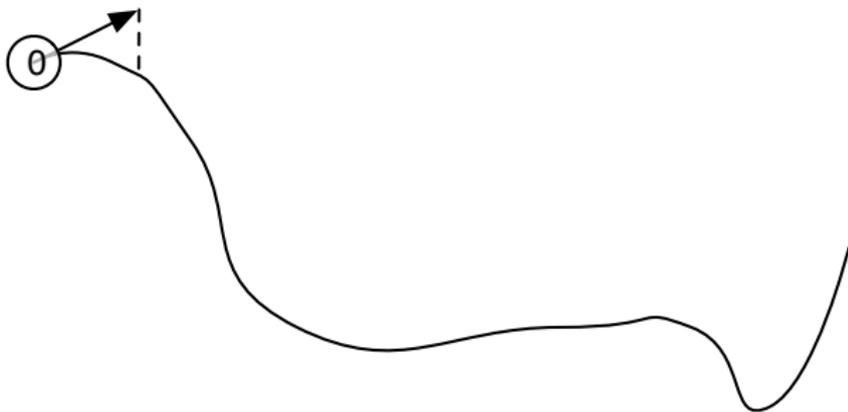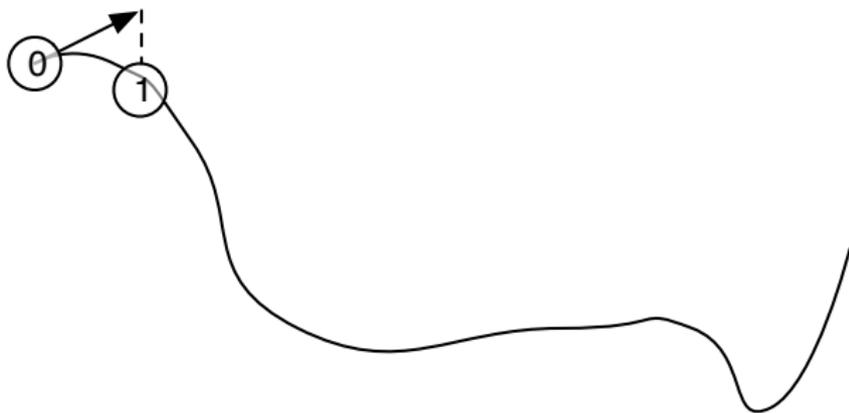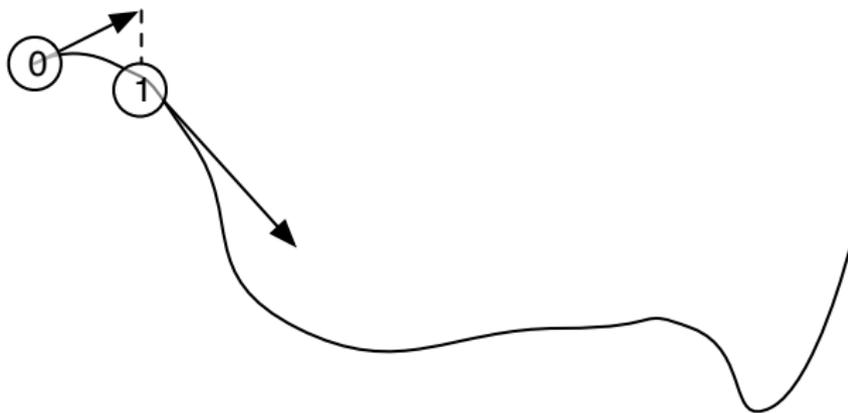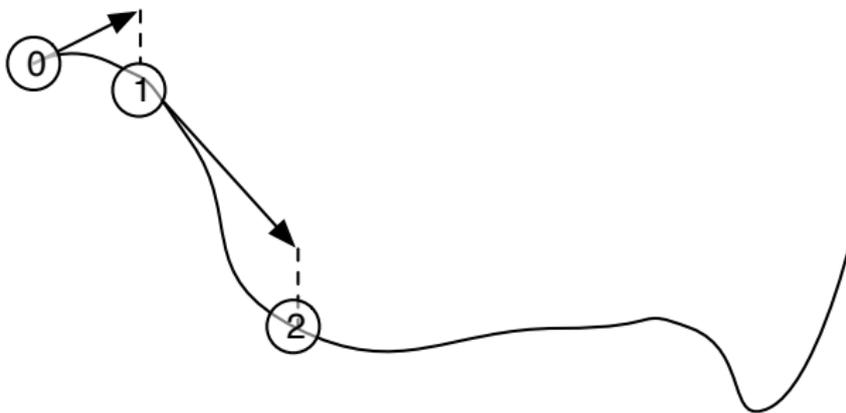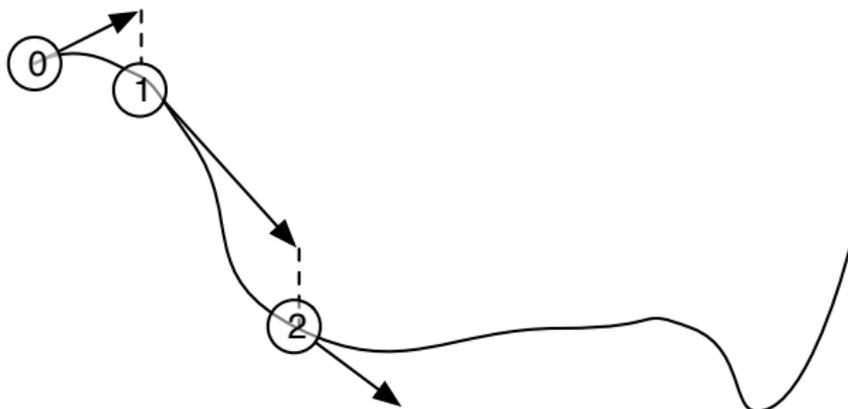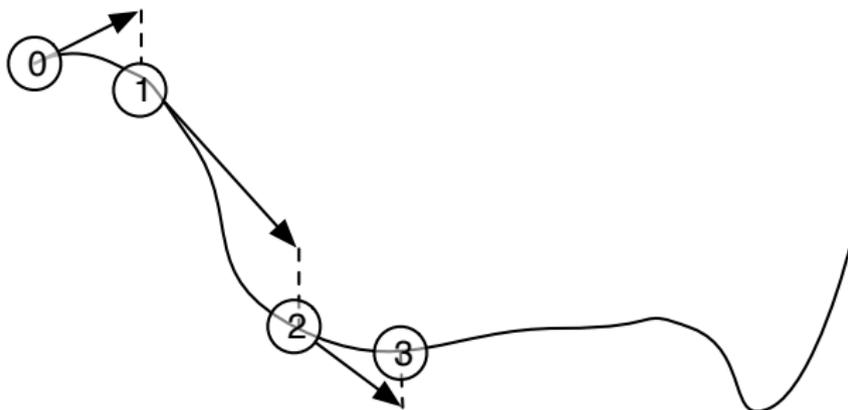
**Goal**

Optimize log likelihood with respect to variables *W* and *b*



Undiscovered
Country

**Goal**

Optimize log likelihood with respect to variables *W* and *b*



$$\alpha \frac{\partial}{\partial W} J$$

$$W_{ij}^{(l)} \quad W_{ij}^{(l)}$$

**Gradient for Logistic Regression**

**Gradient**

$$\nabla_\beta \mathscr{L}(\vec{\beta}) = \left[ \frac{\partial \mathscr{L}(\vec{\beta})}{\partial \beta_0}, \ldots, \frac{\partial \mathscr{L}(\vec{\beta})}{\partial \beta_n} \right] \tag{5}$$

**Update**

$$\Delta\beta \equiv \eta \nabla_\beta \mathscr{L}(\vec{\beta}) \tag{6}$$

$$\beta_i \leftarrow \beta_i' + \eta \frac{\partial \mathscr{L}(\vec{\beta})}{\partial \beta_i} \tag{7}$$

**Gradient**

$$\nabla_\beta \mathscr{L}(\vec{\beta}) = \left[ \frac{\partial \mathscr{L}(\vec{\beta})}{\partial \beta_0}, \ldots, \frac{\partial \mathscr{L}(\vec{\beta})}{\partial \beta_n} \right] \tag{5}$$

**Update**

$$\Delta \beta \equiv \eta \nabla_\beta \mathscr{L}(\vec{\beta}) \tag{6}$$

$$\beta_i \leftarrow \beta_i' + \eta \frac{\partial \mathscr{L}(\vec{\beta})}{\partial \beta_i} \tag{7}$$

Why are we adding? What would well do if we wanted to do **descent**?

**Gradient**

$$\nabla_\beta \mathscr{L}(\vec{\beta}) = \left[ \frac{\partial \mathscr{L}(\vec{\beta})}{\partial \beta_0}, \ldots, \frac{\partial \mathscr{L}(\vec{\beta})}{\partial \beta_n} \right] \tag{5}$$

**Update**

$$\Delta\beta \equiv \eta \nabla_\beta \mathscr{L}(\vec{\beta}) \tag{6}$$

$$\beta_i \leftarrow \beta_i' + \eta \frac{\partial \mathscr{L}(\vec{\beta})}{\partial \beta_i} \tag{7}$$

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

**Gradient for Logistic Regression**

**Gradient**

$$\nabla_\beta \mathscr{L}(\vec{\beta}) = \left[ \frac{\partial \mathscr{L}(\vec{\beta})}{\partial \beta_0}, \ldots, \frac{\partial \mathscr{L}(\vec{\beta})}{\partial \beta_n} \right] \tag{5}$$
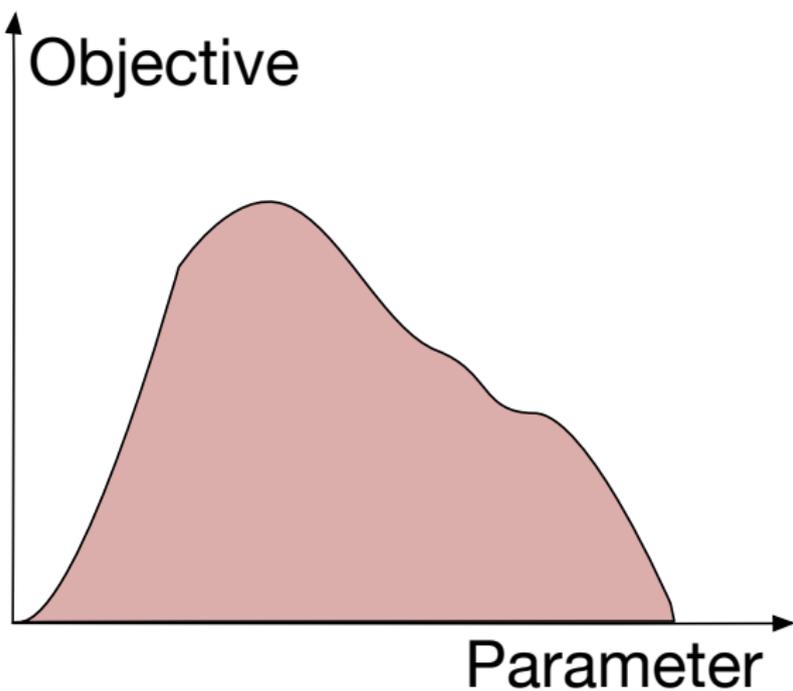
**Update**

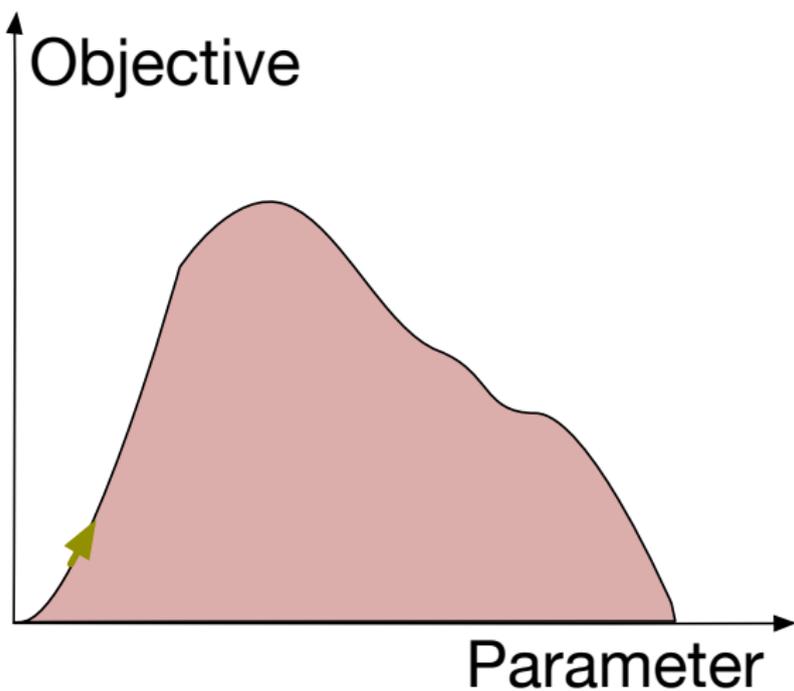$$\Delta \beta \equiv \eta \nabla_\beta \mathscr{L}(\vec{\beta}) \tag{6}$$

$$\beta_i \leftarrow \beta_i' + \eta \frac{\partial \mathscr{L}(\vec{\beta})}{\partial \beta_i} \tag{7}$$

NB: Conjugate gradient is usually better, but harder to implement

Objective

Parameter

Objective

Parameter

- When to stop?
- What if $\beta$ keeps getting bigger?

**Unregularized**

$$\beta^* = \arg\max_\beta \ln\left[p(y^{(j)}\,|\,x^{(j)}, \beta)\right] \tag{8}$$

**Regularized**

$$\beta^* = \arg\max_\beta \ln\left[p(y^{(j)}\,|\,x^{(j)}, \beta)\right] - \mu \sum_i \beta_i^2 \tag{9}$$

**Unregularized**

$$\beta^* = \arg\max_{\beta} \ln\left[p(y^{(j)}|x^{(j)},\beta)\right] \qquad (8)$$

**Regularized**

$$\beta^* = \arg\max_{\beta} \ln\left[p(y^{(j)}|x^{(j)},\beta)\right] - \mu\sum_i \beta_i^2 \qquad (9)$$

$\mu$ is "regularization" parameter that trades off between likelihood and having small parameters

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

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

$$\mathcal{L}(\beta) \equiv \mathbb{E}_x [\nabla \mathcal{L}(\beta, x)] \tag{10}$$

- Average over all observations

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

$$\mathscr{L}(\beta) \equiv \mathbb{E}_x \left[ \nabla \mathscr{L}(\beta, x) \right] \tag{10}$$

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

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

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

$$\beta_i \leftarrow \beta_i' + \eta \left( -\mu \beta_i' + x_i \left[ y - p(y = 1 \,|\, \vec{x}, \vec{\beta}') \right] \right) \tag{11}$$

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

$$\beta_i \leftarrow \beta_i' + \eta \left( -\mu \beta_i' + x_i \left[ y - p(y = 1 \,|\, \vec{x}, \vec{\beta}') \right] \right) \tag{11}$$

Examples in class.

$$\mathcal{L} = \log p(y \mid x; \beta) - \mu \sum_j \beta_j^2 \qquad (12)$$

$$\mathcal{L} = \log p(y \,|\, x; \beta) - \mu \sum_j \beta_j^2 \tag{12}$$

Taking the derivative (with respect to example $x_i$)

$$\frac{\partial \mathcal{L}}{\partial \beta_j} = (y_i - p(y_i = 1 \,|\, \vec{x}_i; \beta)) x_j - 2\mu \beta_j \tag{13}$$

**Proofs about Stochastic Gradient**

- Depends on convexity of objective and how close $\epsilon$ you want to get to actual answer
- Best bounds depend on changing $\eta$ over time and **per dimension** (not all features created equal)

- Your questions!
- Working through simple example
- Prepared for logistic regression homework