

CS189/CS289A
Introduction to Machine Learning
Lecture 11: Logistic Regression

Peter Bartlett

February 24, 2015

First: A Bayesian view of linear regression

Logistic Regression:

- Gaussian generative to logistic discriminative models.
- Parameter estimates for logistic models.
 - Maximum likelihood: coupled non-linear equations.
 - Gradient ascent.
 - Stochastic gradient.
 - (Detour: stochastic gradient for linear regression.)
 - Newton's method: iteratively reweighted least squares.

Bayesian linear regression

Linear model

Bayesian linear regression

Linear model

$$P(Y|X = x) = \mathcal{N}(x'\beta, \sigma^2).$$

Bayesian linear regression

Linear model

$$P(Y|X = x) = \mathcal{N}(x'\beta, \sigma^2).$$

Equivalently: $Y = x'\beta + \epsilon$, where $\epsilon \sim \mathcal{N}(0, \sigma^2)$.

Bayesian linear regression

Linear model

$$P(Y|X = x) = \mathcal{N}(x'\beta, \sigma^2).$$

Equivalently: $Y = x'\beta + \epsilon$, where $\epsilon \sim \mathcal{N}(0, \sigma^2)$.

Bayesian Analysis

Bayesian linear regression

Linear model

$$P(Y|X = x) = \mathcal{N}(x'\beta, \sigma^2).$$

Equivalently: $Y = x'\beta + \epsilon$, where $\epsilon \sim \mathcal{N}(0, \sigma^2)$.

Bayesian Analysis

- Model β as a random variable.

Bayesian linear regression

Linear model

$$P(Y|X = x) = \mathcal{N}(x'\beta, \sigma^2).$$

Equivalently: $Y = x'\beta + \epsilon$, where $\epsilon \sim \mathcal{N}(0, \sigma^2)$.

Bayesian Analysis

- Model β as a random variable.
- e.g., $\beta \sim \mathcal{N}(0, \tau^2 I)$.

Bayesian linear regression

Linear model

$$P(Y|X = x) = \mathcal{N}(x'\beta, \sigma^2).$$

Equivalently: $Y = x'\beta + \epsilon$, where $\epsilon \sim \mathcal{N}(0, \sigma^2)$.

Bayesian Analysis

- Model β as a random variable.
- e.g., $\beta \sim \mathcal{N}(0, \tau^2 I)$. \leftarrow prior distribution of β .

Bayesian linear regression

Linear model

$$P(Y|X = x) = \mathcal{N}(x'\beta, \sigma^2).$$

Equivalently: $Y = x'\beta + \epsilon$, where $\epsilon \sim \mathcal{N}(0, \sigma^2)$.

Bayesian Analysis

- Model β as a random variable.
- e.g., $\beta \sim \mathcal{N}(0, \tau^2 I)$. \leftarrow prior distribution of β .
- Then compute posterior distribution $P(\beta|X, Y)$:

Bayesian linear regression

Linear model

$$P(Y|X = x) = \mathcal{N}(x'\beta, \sigma^2).$$

Equivalently: $Y = x'\beta + \epsilon$, where $\epsilon \sim \mathcal{N}(0, \sigma^2)$.

Bayesian Analysis

- Model β as a random variable.
- e.g., $\beta \sim \mathcal{N}(0, \tau^2 I)$. \leftarrow prior distribution of β .
- Then compute posterior distribution $P(\beta|X, Y)$:

$$P(\beta|X, Y) = \frac{P(Y|\beta, X)P(\beta)}{P(Y|X)}$$

Bayesian linear regression

Linear model

$$P(Y|X = x) = \mathcal{N}(x'\beta, \sigma^2).$$

Equivalently: $Y = x'\beta + \epsilon$, where $\epsilon \sim \mathcal{N}(0, \sigma^2)$.

Bayesian Analysis

- Model β as a random variable.
- e.g., $\beta \sim \mathcal{N}(0, \tau^2 I)$. \leftarrow prior distribution of β .
- Then compute posterior distribution $P(\beta|X, Y)$:

$$P(\beta|X, Y) = \frac{P(Y|\beta, X)P(\beta)}{P(Y|X)} \propto \underbrace{P(Y|\beta, X)}_{\text{likelihood}} \underbrace{P(\beta)}_{\text{prior}}.$$

Bayesian linear regression

Linear model

$$P(Y|X = x) = \mathcal{N}(x'\beta, \sigma^2).$$

Equivalently: $Y = x'\beta + \epsilon$, where $\epsilon \sim \mathcal{N}(0, \sigma^2)$.

Bayesian Analysis

- Model β as a random variable.
- e.g., $\beta \sim \mathcal{N}(0, \tau^2 I)$. \leftarrow prior distribution of β .
- Then compute posterior distribution $P(\beta|X, Y)$:

$$P(\beta|X, Y) = \frac{P(Y|\beta, X)P(\beta)}{P(Y|X)} \propto \underbrace{P(Y|\beta, X)}_{\text{likelihood}} \underbrace{P(\beta)}_{\text{prior}}.$$

$$P(\beta|X_1, Y_1, X_2, Y_2) \propto \underbrace{P(Y_2|\beta, X_2)P(Y_1|\beta, X_1)}_{\text{likelihood}} P(\beta)$$

Bayesian linear regression

Linear model

$$P(Y|X = x) = \mathcal{N}(x'\beta, \sigma^2).$$

Equivalently: $Y = x'\beta + \epsilon$, where $\epsilon \sim \mathcal{N}(0, \sigma^2)$.

Bayesian Analysis

- Model β as a random variable.
- e.g., $\beta \sim \mathcal{N}(0, \tau^2 I)$. \leftarrow prior distribution of β .
- Then compute posterior distribution $P(\beta|X, Y)$:

$$P(\beta|X, Y) = \frac{P(Y|\beta, X)P(\beta)}{P(Y|X)} \propto \underbrace{P(Y|\beta, X)}_{\text{likelihood}} \underbrace{P(\beta)}_{\text{prior}}.$$

$$\begin{aligned} P(\beta|X_1, Y_1, X_2, Y_2) &\propto \underbrace{P(Y_2|\beta, X_2)P(Y_1|\beta, X_1)}_{\text{likelihood}} P(\beta) \\ &= P(Y_2|\beta, X_2) \underbrace{P(Y_1|\beta, X_1)P(\beta)}_{\text{prior}}. \end{aligned}$$

Bayesian Analysis

- Prior $\beta \sim \mathcal{N}(0, \tau^2 I)$.
- Posterior distribution $P(\beta|X, Y)$:

$$P(\beta|X, Y) \propto P(Y|\beta, X)P(\beta)$$

Bayesian Analysis

- Prior $\beta \sim \mathcal{N}(0, \tau^2 I)$.
- Posterior distribution $P(\beta|X, Y)$:

$$\begin{aligned} P(\beta|X, Y) &\propto P(Y|\beta, X)P(\beta) \\ &\propto \exp\left(-\frac{(Y - X'\beta)^2}{2\sigma^2} - \frac{\|\beta\|^2}{2\tau^2}\right) \end{aligned}$$

Bayesian Analysis

- Prior $\beta \sim \mathcal{N}(0, \tau^2 I)$.
- Posterior distribution $P(\beta|X, Y)$:

$$\begin{aligned} P(\beta|X, Y) &\propto P(Y|\beta, X)P(\beta) \\ &\propto \exp\left(-\frac{(Y - X'\beta)^2}{2\sigma^2} - \frac{\|\beta\|^2}{2\tau^2}\right) \end{aligned}$$

$$P(\beta|X_1, Y_1, \dots, X_n, Y_n) \propto P(Y_1, \dots, Y_n|\beta, X_1, \dots, X_n)P(\beta)$$

Bayesian Analysis

- Prior $\beta \sim \mathcal{N}(0, \tau^2 I)$.
- Posterior distribution $P(\beta|X, Y)$:

$$\begin{aligned} P(\beta|X, Y) &\propto P(Y|\beta, X)P(\beta) \\ &\propto \exp\left(-\frac{(Y - X'\beta)^2}{2\sigma^2} - \frac{\|\beta\|^2}{2\tau^2}\right) \end{aligned}$$

$$\begin{aligned} P(\beta|X_1, Y_1, \dots, X_n, Y_n) &\propto P(Y_1, \dots, Y_n|\beta, X_1, \dots, X_n)P(\beta) \\ &\propto \exp\left(-\frac{1}{2}\left(\sum_{i=1}^n \frac{(Y_i - X_i'\beta)^2}{\sigma^2} + \frac{1}{\tau^2}\|\beta\|^2\right)\right). \end{aligned}$$

Bayesian Analysis

Bayesian Analysis

- Maximum a posteriori probability estimate: estimate β as the mode of the posterior distribution.

Bayesian Analysis

- Maximum a posteriori probability estimate: estimate β as the mode of the posterior distribution.
- Ridge regression (least squares with a squared Euclidean norm penalty) is equivalent to a MAP estimate with a Gaussian prior.

Bayesian Analysis

- Maximum a posteriori probability estimate: estimate β as the mode of the posterior distribution.
- Ridge regression (least squares with a squared Euclidean norm penalty) is equivalent to a MAP estimate with a Gaussian prior.
- Consider a Laplace prior:

$$P(\beta) \propto \exp(-\lambda \|\beta\|_1).$$

Bayesian Analysis

- Maximum a posteriori probability estimate: estimate β as the mode of the posterior distribution.
- Ridge regression (least squares with a squared Euclidean norm penalty) is equivalent to a MAP estimate with a Gaussian prior.
- Consider a Laplace prior:

$$P(\beta) \propto \exp(-\lambda \|\beta\|_1).$$

- Lasso (least squares with a squared one-norm penalty) is equivalent to a MAP estimate with a Laplace prior.

Logistic Regression:

Logistic Regression:

- Gaussian generative to logistic discriminative models.

Logistic Regression:

- **Gaussian generative to logistic discriminative models.**
- Parameter estimates for logistic models.

Logistic Regression:

- **Gaussian generative to logistic discriminative models.**
- Parameter estimates for logistic models.
 - Maximum likelihood: coupled non-linear equations.

Logistic Regression:

- **Gaussian generative to logistic discriminative models.**
- Parameter estimates for logistic models.
 - Maximum likelihood: coupled non-linear equations.
 - Gradient ascent.

Logistic Regression:

- **Gaussian generative to logistic discriminative models.**
- Parameter estimates for logistic models.
 - Maximum likelihood: coupled non-linear equations.
 - Gradient ascent.
 - Stochastic gradient.

Logistic Regression:

- **Gaussian generative to logistic discriminative models.**
- Parameter estimates for logistic models.
 - Maximum likelihood: coupled non-linear equations.
 - Gradient ascent.
 - Stochastic gradient.
 - (Detour: stochastic gradient for linear regression.)

Logistic Regression:

- **Gaussian generative to logistic discriminative models.**
- Parameter estimates for logistic models.
 - Maximum likelihood: coupled non-linear equations.
 - Gradient ascent.
 - Stochastic gradient.
 - (Detour: stochastic gradient for linear regression.)
 - Newton's method: iteratively reweighted least squares.

Logistic regression

Class conditionals to posterior

For Gaussian class conditional densities $P(X|Y = 1)$, $P(X|Y = 0)$ (with the same variance), the posterior probability is logistic

Logistic regression

Class conditionals to posterior

For Gaussian class conditional densities $P(X|Y=1)$, $P(X|Y=0)$ (with the same variance), the posterior probability is logistic:

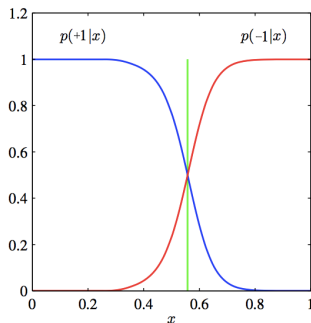
$$P(Y=1|x) = \frac{1}{1 + \exp(-x \cdot \beta - \beta_0)}.$$

Logistic regression

Class conditionals to posterior

For Gaussian class conditional densities $P(X|Y=1)$, $P(X|Y=0)$ (with the same variance), the posterior probability is logistic:

$$P(Y=1|x) = \frac{1}{1 + \exp(-x \cdot \beta - \beta_0)}.$$



Gaussian generative to logistic discriminative models

- Suppose the class conditional distributions are Gaussian:

$$p(x|y = 1) = \mathcal{N}(\mu_1, \Sigma), \quad p(x|y = 0) = \mathcal{N}(\mu_0, \Sigma)$$

Gaussian generative to logistic discriminative models

- Suppose the class conditional distributions are Gaussian:

$$p(x|y=1) = \mathcal{N}(\mu_1, \Sigma), \quad p(x|y=0) = \mathcal{N}(\mu_0, \Sigma)$$

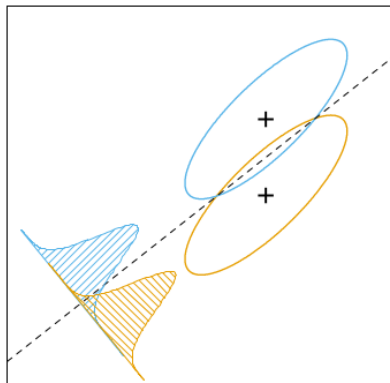
$$p(x|y=1) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left(-\frac{1}{2} (x - \mu_1)' \Sigma^{-1} (x - \mu_1) \right).$$

Gaussian generative to logistic discriminative models

- Suppose the class conditional distributions are Gaussian:

$$p(x|y=1) = \mathcal{N}(\mu_1, \Sigma), \quad p(x|y=0) = \mathcal{N}(\mu_0, \Sigma)$$

$$p(x|y=1) = \frac{1}{(2\pi)^{d/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_1)' \Sigma^{-1}(x - \mu_1)\right).$$



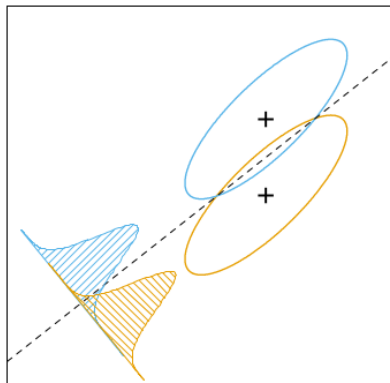
Gaussian generative to logistic discriminative models

- Suppose the class conditional distributions are Gaussian:

$$p(x|y=1) = \mathcal{N}(\mu_1, \Sigma), \quad p(x|y=0) = \mathcal{N}(\mu_0, \Sigma)$$

$$p(x|y=1) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left(-\frac{1}{2} (x - \mu_1)' \Sigma^{-1} (x - \mu_1) \right).$$

- Class 1 has mean μ_1 .
Class 0 has mean μ_0 .



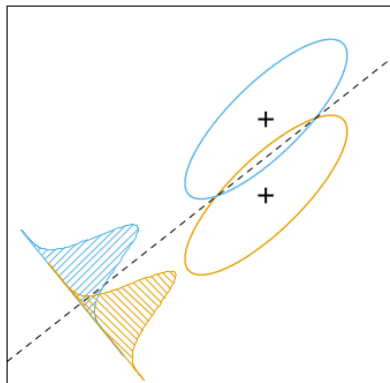
Gaussian generative to logistic discriminative models

- Suppose the class conditional distributions are Gaussian:

$$p(x|y=1) = \mathcal{N}(\mu_1, \Sigma), \quad p(x|y=0) = \mathcal{N}(\mu_0, \Sigma)$$

$$p(x|y=1) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left(-\frac{1}{2} (x - \mu_1)' \Sigma^{-1} (x - \mu_1) \right).$$

- Class 1 has mean μ_1 .
Class 0 has mean μ_0 .
- Both have covariance matrix Σ .



Gaussian generative to logistic discriminative models

$$p(x|Y=1) = \frac{1}{(2\pi)^{d/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_1)'\Sigma^{-1}(x - \mu_1)\right)$$

Gaussian generative to logistic discriminative models

$$p(x|Y=1) = \frac{1}{(2\pi)^{d/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_1)'\Sigma^{-1}(x - \mu_1)\right)$$

$$\log \frac{P(Y=1|x)}{P(Y=0|x)} = \log \frac{P(x|Y=1)P(Y=1)}{P(x|Y=0)P(Y=0)}$$

Gaussian generative to logistic discriminative models

$$p(x|Y=1) = \frac{1}{(2\pi)^{d/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_1)'\Sigma^{-1}(x - \mu_1)\right)$$

$$\begin{aligned} \log \frac{P(Y=1|x)}{P(Y=0|x)} &= \log \frac{P(x|Y=1)P(Y=1)}{P(x|Y=0)P(Y=0)} \\ &= \frac{1}{2} [(x - \mu_0)'\Sigma^{-1}(x - \mu_0) \\ &\quad - (x - \mu_1)'\Sigma^{-1}(x - \mu_1)] + \log \frac{P(Y=1)}{P(Y=0)} \end{aligned}$$

Gaussian generative to logistic discriminative models

$$p(x|Y=1) = \frac{1}{(2\pi)^{d/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_1)'\Sigma^{-1}(x - \mu_1)\right)$$

$$\begin{aligned}\log \frac{P(Y=1|x)}{P(Y=0|x)} &= \log \frac{P(x|Y=1)P(Y=1)}{P(x|Y=0)P(Y=0)} \\&= \frac{1}{2} \left[(x - \mu_0)'\Sigma^{-1}(x - \mu_0) \right. \\&\quad \left. - (x - \mu_1)'\Sigma^{-1}(x - \mu_1) \right] + \log \frac{P(Y=1)}{P(Y=0)} \\&= \frac{1}{2} (\mu_0'\Sigma^{-1}\mu_0 - \mu_1'\Sigma^{-1}\mu_1) + \log \frac{P(Y=1)}{P(Y=0)} \\&\quad + (\mu_1 - \mu_0)'\Sigma^{-1}x\end{aligned}$$

Gaussian generative to logistic discriminative models

$$p(x|Y=1) = \frac{1}{(2\pi)^{d/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_1)'\Sigma^{-1}(x - \mu_1)\right)$$

$$\begin{aligned}\log \frac{P(Y=1|x)}{P(Y=0|x)} &= \log \frac{P(x|Y=1)P(Y=1)}{P(x|Y=0)P(Y=0)} \\&= \frac{1}{2} \left[(x - \mu_0)'\Sigma^{-1}(x - \mu_0) \right. \\&\quad \left. - (x - \mu_1)'\Sigma^{-1}(x - \mu_1) \right] + \log \frac{P(Y=1)}{P(Y=0)} \\&= \frac{1}{2} (\mu_0'\Sigma^{-1}\mu_0 - \mu_1'\Sigma^{-1}\mu_1) + \log \frac{P(Y=1)}{P(Y=0)} \\&\quad + (\mu_1 - \mu_0)'\Sigma^{-1}x \\&= \beta_0 + \beta'x.\end{aligned}$$

Gaussian generative to logistic discriminative models

$$\log \frac{P(Y = 1|x)}{P(Y = 0|x)} = \beta_0 + \beta'x.$$

Gaussian generative to logistic discriminative models

$$\log \frac{P(Y = 1|x)}{P(Y = 0|x)} = \beta_0 + \beta'x.$$

For $p = P(Y = 1|x)$, we have

$$\log \frac{p}{1-p} = \beta_0 + \beta'x,$$

Gaussian generative to logistic discriminative models

$$\log \frac{P(Y = 1|x)}{P(Y = 0|x)} = \beta_0 + \beta'x.$$

For $p = P(Y = 1|x)$, we have

$$\begin{aligned}\log \frac{p}{1-p} &= \beta_0 + \beta'x, \\ \frac{p}{1-p} &= \exp(\beta_0 + \beta'x),\end{aligned}$$

Gaussian generative to logistic discriminative models

$$\log \frac{P(Y = 1|x)}{P(Y = 0|x)} = \beta_0 + \beta'x.$$

For $p = P(Y = 1|x)$, we have

$$\log \frac{p}{1-p} = \beta_0 + \beta'x,$$

$$\frac{p}{1-p} = \exp(\beta_0 + \beta'x),$$

$$P(Y = 1|X) = p = \frac{1}{1 + \exp(-(\beta_0 + \beta'x))}.$$

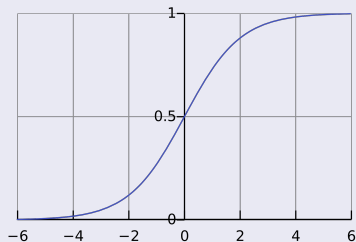
Gaussian generative to logistic discriminative models

Logistic model:
$$P(Y = 1|x) = \frac{1}{1 + \exp(-\beta \cdot x - \beta_0)}.$$

Gaussian generative to logistic discriminative models

Logistic model:

$$P(Y = 1|x) = \frac{1}{1 + \exp(-\beta \cdot x - \beta_0)}.$$

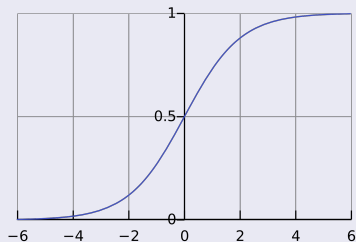


The logistic function $\frac{1}{1 + e^{-\alpha}}$.

Gaussian generative to logistic discriminative models

Logistic model:
$$P(Y = 1|x) = \frac{1}{1 + \exp(-\beta \cdot x - \beta_0)}.$$

$$\beta = \Sigma^{-1}(\mu_1 - \mu_0), \quad \beta_0 = \frac{\mu_0' \Sigma^{-1} \mu_0 - \mu_1' \Sigma^{-1} \mu_1}{2} + \log \frac{P(Y = 1)}{P(Y = 0)}.$$



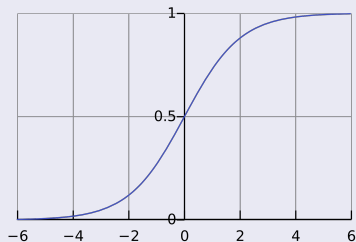
The logistic function $\frac{1}{1 + e^{-\alpha}}.$

Gaussian generative to logistic discriminative models

Logistic model:
$$P(Y = 1|x) = \frac{1}{1 + \exp(-\beta \cdot x - \beta_0)}.$$

$$\beta = \Sigma^{-1}(\mu_1 - \mu_0), \quad \beta_0 = \frac{\mu_0' \Sigma^{-1} \mu_0 - \mu_1' \Sigma^{-1} \mu_1}{2} + \log \frac{P(Y = 1)}{P(Y = 0)}.$$

- $P(Y = 1|x)$ increases as



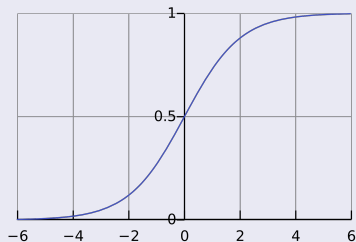
The logistic function $\frac{1}{1 + e^{-\alpha}}.$

Gaussian generative to logistic discriminative models

Logistic model:
$$P(Y = 1|x) = \frac{1}{1 + \exp(-\beta \cdot x - \beta_0)}.$$

$$\beta = \Sigma^{-1}(\mu_1 - \mu_0), \quad \beta_0 = \frac{\mu_0' \Sigma^{-1} \mu_0 - \mu_1' \Sigma^{-1} \mu_1}{2} + \log \frac{P(Y = 1)}{P(Y = 0)}.$$

- $P(Y = 1|x)$ increases as
 - x moves from μ_0 to μ_1 ,



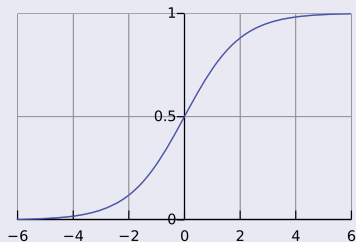
The logistic function $\frac{1}{1 + e^{-\alpha}}.$

Gaussian generative to logistic discriminative models

Logistic model:
$$P(Y = 1|x) = \frac{1}{1 + \exp(-\beta \cdot x - \beta_0)}.$$

$$\beta = \Sigma^{-1}(\mu_1 - \mu_0), \quad \beta_0 = \frac{\mu_0' \Sigma^{-1} \mu_0 - \mu_1' \Sigma^{-1} \mu_1}{2} + \log \frac{P(Y = 1)}{P(Y = 0)}.$$

- $P(Y = 1|x)$ increases as
 - x moves from μ_0 to μ_1 ,
 - $P(Y = 1)$ increases.



The logistic function
$$\frac{1}{1 + e^{-\alpha}}.$$

Gaussian generative to logistic discriminative models

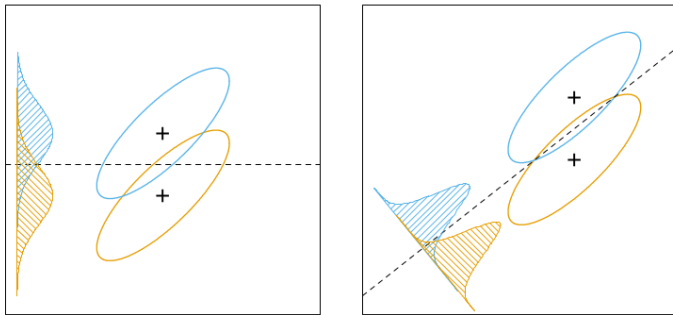


FIGURE 4.9. *Although the line joining the centroids defines the direction of greatest centroid spread, the projected data overlap because of the covariance (left panel). The discriminant direction minimizes this overlap for Gaussian data (right panel).*

Gaussian generative to logistic discriminative models

- This is one motivation for logistic regression.

Gaussian generative to logistic discriminative models

- This is one motivation for logistic regression.
But logistic regression does not *require* an assumption of Gaussian class conditionals!

Gaussian generative to logistic discriminative models

- This is one motivation for logistic regression.
But logistic regression does not *require* an assumption of Gaussian class conditionals!
- Logistic regression: Model log odds ($\log p/(1 - p)$) as an affine function of x .

Logistic Regression:

- Gaussian generative to logistic discriminative models.
- **Parameter estimates for logistic models.**
 - Maximum likelihood: coupled non-linear equations.
 - Gradient ascent.
 - Stochastic gradient.
 - (Detour: stochastic gradient for linear regression.)
 - Newton's method: iteratively reweighted least squares.

Logistic Model

$$P(Y = 1|x) = \frac{1}{1 + \exp(-\beta'x)}.$$

(incorporate β_0 into β)

Logistic Model

$$P(Y = 1|x) = \frac{1}{1 + \exp(-\beta'x)}. \quad (\text{incorporate } \beta_0 \text{ into } \beta)$$

Logistic Regression

Given data $(X_1, Y_1), \dots, (X_n, Y_n) \in \mathbb{R}^p \times \{0, 1\}$,
estimate β

Logistic Model

$$P(Y = 1|x) = \frac{1}{1 + \exp(-\beta'x)}. \quad (\text{incorporate } \beta_0 \text{ into } \beta)$$

Logistic Regression

Given data $(X_1, Y_1), \dots, (X_n, Y_n) \in \mathbb{R}^p \times \{0, 1\}$,
estimate β : maximum likelihood.

Maximum likelihood estimation

Maximum likelihood estimation

Log likelihood: $\ell(\beta) = \log P(y_1, \dots, y_n | x_1, \dots, x_n; \beta)$

Maximum likelihood estimation

Log likelihood:

$$\begin{aligned}\ell(\beta) &= \log P(y_1, \dots, y_n | x_1, \dots, x_n; \beta) \\ &= \sum_{i=1}^n y_i \log \mu_i(\beta) + (1 - y_i) \log(1 - \mu_i(\beta)),\end{aligned}$$

Maximum likelihood estimation

Log likelihood: $\ell(\beta) = \log P(y_1, \dots, y_n | x_1, \dots, x_n; \beta)$

$$= \sum_{i=1}^n y_i \log \mu_i(\beta) + (1 - y_i) \log(1 - \mu_i(\beta)),$$

where $\mu_i(\beta) = P(Y = 1 | X = x_i, \beta) = \frac{1}{1 + \exp(-\beta' x_i)}.$

Maximum likelihood estimation

Log likelihood: $\ell(\beta) = \log P(y_1, \dots, y_n | x_1, \dots, x_n; \beta)$

$$= \sum_{i=1}^n y_i \log \mu_i(\beta) + (1 - y_i) \log(1 - \mu_i(\beta)),$$

where $\mu_i(\beta) = P(Y = 1 | X = x_i, \beta) = \frac{1}{1 + \exp(-\beta' x_i)}.$

NB: $\nabla_{\beta} \mu_i(\beta) = \mu_i(\beta)(1 - \mu_i(\beta))x_i.$

Maximum likelihood estimation

Log likelihood:
$$\ell(\beta) = \sum_{i=1}^n y_i \log \mu_i(\beta) + (1 - y_i) \log(1 - \mu_i(\beta)).$$

Maximum likelihood estimation

Log likelihood:
$$\ell(\beta) = \sum_{i=1}^n y_i \log \mu_i(\beta) + (1 - y_i) \log(1 - \mu_i(\beta)).$$

First derivative:
$$\nabla_{\beta} \ell(\beta) = \sum_{i=1}^n \left(\left(\frac{y_i}{\mu_i(\beta)} - \frac{1 - y_i}{1 - \mu_i(\beta)} \right) \nabla_{\beta} \mu_i(\beta) \right)$$

Maximum likelihood estimation

Log likelihood:
$$\ell(\beta) = \sum_{i=1}^n y_i \log \mu_i(\beta) + (1 - y_i) \log(1 - \mu_i(\beta)).$$

First derivative:
$$\begin{aligned} \nabla_{\beta} \ell(\beta) &= \sum_{i=1}^n \left(\left(\frac{y_i}{\mu_i(\beta)} - \frac{1 - y_i}{1 - \mu_i(\beta)} \right) \nabla_{\beta} \mu_i(\beta) \right) \\ &= \sum_{i=1}^n (y_i - \mu_i(\beta)) x_i. \end{aligned}$$

Maximum likelihood estimation

Log likelihood:
$$\ell(\beta) = \sum_{i=1}^n y_i \log \mu_i(\beta) + (1 - y_i) \log(1 - \mu_i(\beta)).$$

First derivative:
$$\begin{aligned}\nabla_{\beta} \ell(\beta) &= \sum_{i=1}^n \left(\left(\frac{y_i}{\mu_i(\beta)} - \frac{1 - y_i}{1 - \mu_i(\beta)} \right) \nabla_{\beta} \mu_i(\beta) \right) \\ &= \sum_{i=1}^n (y_i - \mu_i(\beta)) x_i.\end{aligned}$$

Second derivative:
$$\nabla_{\beta}^2 \ell(\beta) = \sum_{i=1}^n -\mu_i(\beta)(1 - \mu_i(\beta)) x_i x_i'.$$

Maximum likelihood estimate

$$\hat{\beta}^{ml} \text{ solves: } 0 = \sum_{i=1}^n (y_i - \mu_i(\beta)) x_i$$

Maximum likelihood estimate

$$\hat{\beta}^{ml} \text{ solves: } 0 = \sum_{i=1}^n (y_i - \mu_i(\beta)) x_i = \sum_{i=1}^n \left(y_i - \frac{1}{1 + \exp(-\beta' x_i)} \right) x_i$$

Maximum likelihood estimate

$$\hat{\beta}^{ml} \text{ solves: } 0 = \sum_{i=1}^n (y_i - \mu_i(\beta)) x_i = \sum_{i=1}^n \left(y_i - \frac{1}{1 + \exp(-\beta' x_i)} \right) x_i$$
$$\Leftrightarrow \sum_{i=1}^n y_i x_i = \sum_{i=1}^n \mu_i(\beta) x_i.$$

Maximum likelihood estimate

$$\hat{\beta}^{ml} \text{ solves: } 0 = \sum_{i=1}^n (y_i - \mu_i(\beta)) x_i = \sum_{i=1}^n \left(y_i - \frac{1}{1 + \exp(-\beta' x_i)} \right) x_i$$
$$\Leftrightarrow \sum_{i=1}^n y_i x_i = \sum_{i=1}^n \mu_i(\beta) x_i. \quad (\text{moment matching})$$

Maximum likelihood estimate

$$\hat{\beta}^{ml} \text{ solves: } 0 = \sum_{i=1}^n (y_i - \mu_i(\beta)) x_i = \sum_{i=1}^n \left(y_i - \frac{1}{1 + \exp(-\beta' x_i)} \right) x_i$$
$$\Leftrightarrow \sum_{i=1}^n y_i x_i = \sum_{i=1}^n \mu_i(\beta) x_i. \quad (\text{moment matching})$$

- System of p coupled non-linear equations.

Maximum likelihood estimate

$$\hat{\beta}^{ml} \text{ solves: } 0 = \sum_{i=1}^n (y_i - \mu_i(\beta)) x_i = \sum_{i=1}^n \left(y_i - \frac{1}{1 + \exp(-\beta' x_i)} \right) x_i$$
$$\Leftrightarrow \sum_{i=1}^n y_i x_i = \sum_{i=1}^n \mu_i(\beta) x_i. \quad (\text{moment matching})$$

- System of p coupled non-linear equations.
- No closed-form solution.

Maximum likelihood estimate

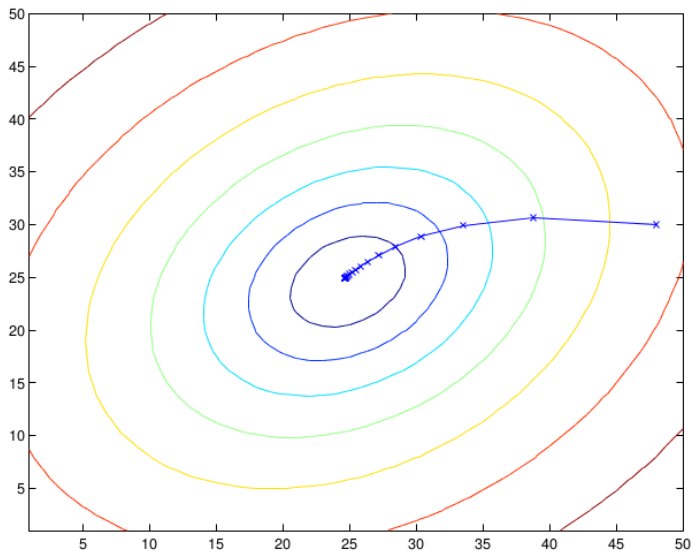
$$\hat{\beta}^{ml} \text{ solves: } 0 = \sum_{i=1}^n (y_i - \mu_i(\beta)) x_i = \sum_{i=1}^n \left(y_i - \frac{1}{1 + \exp(-\beta' x_i)} \right) x_i$$
$$\Leftrightarrow \sum_{i=1}^n y_i x_i = \sum_{i=1}^n \mu_i(\beta) x_i. \quad (\text{moment matching})$$

- System of p coupled non-linear equations.
- No closed-form solution.
- But ℓ is *concave*, so we can find the solution.

Logistic Regression:

- Gaussian generative to logistic discriminative models.
- Parameter estimates for logistic models.
 - Maximum likelihood: coupled non-linear equations.
 - **Gradient ascent.**
 - Stochastic gradient.
 - (Detour: stochastic gradient for linear regression.)
 - Newton's method: iteratively reweighted least squares.

Gradient Ascent



Gradient ascent

$$\beta^{(t+1)} = \beta^{(t)} + \eta \nabla_{\beta} \ell(\beta^{(t)})$$

Gradient ascent

$$\begin{aligned}\beta^{(t+1)} &= \beta^{(t)} + \eta \nabla_{\beta} \ell \left(\beta^{(t)} \right) \\ &= \beta^{(t)} + \eta \sum_{i=1}^n \left(y_i - \mu_i(\beta^{(t)}) \right) x_i.\end{aligned}$$

Gradient ascent

$$\begin{aligned}\beta^{(t+1)} &= \beta^{(t)} + \eta \nabla_{\beta} \ell \left(\beta^{(t)} \right) \\ &= \beta^{(t)} + \eta \sum_{i=1}^n \left(y_i - \mu_i(\beta^{(t)}) \right) x_i.\end{aligned}$$

- η is a step-size parameter.

Gradient ascent

$$\begin{aligned}\beta^{(t+1)} &= \beta^{(t)} + \eta \nabla_{\beta} \ell \left(\beta^{(t)} \right) \\ &= \beta^{(t)} + \eta \sum_{i=1}^n \left(y_i - \mu_i(\beta^{(t)}) \right) x_i.\end{aligned}$$

- η is a step-size parameter.
- Since the μ_i depend on β , the gradient changes as β changes.

Gradient ascent

$$\begin{aligned}\beta^{(t+1)} &= \beta^{(t)} + \eta \nabla_{\beta} \ell(\beta^{(t)}) \\ &= \beta^{(t)} + \eta \sum_{i=1}^n (y_i - \mu_i(\beta^{(t)})) x_i.\end{aligned}$$

- η is a step-size parameter.
- Since the μ_i depend on β , the gradient changes as β changes.
- Each gradient calculation takes $O(np)$ time.

Logistic regression

- Instead of computing $\nabla_{\beta} \ell(\beta^{(t)})$, we compute a random approximation that is on average in the direction of $\nabla_{\beta} \ell(\beta^{(t)})$.

Logistic regression

- Instead of computing $\nabla_{\beta} \ell(\beta^{(t)})$, we compute a random approximation that is on average in the direction of $\nabla_{\beta} \ell(\beta^{(t)})$.
- Here, gradient is

$$\nabla_{\beta} \ell(\beta) = \sum_{i=1}^n (y_i - \mu_i(\beta)) x_i.$$

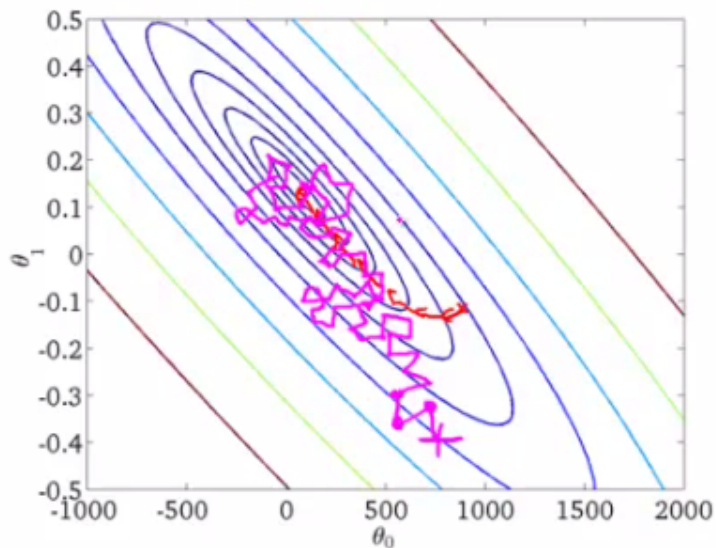
Logistic regression

- Instead of computing $\nabla_{\beta} \ell(\beta^{(t)})$, we compute a random approximation that is on average in the direction of $\nabla_{\beta} \ell(\beta^{(t)})$.
- Here, gradient is

$$\nabla_{\beta} \ell(\beta) = \sum_{i=1}^n (y_i - \mu_i(\beta)) x_i.$$

- We need an (easy-to-compute) quantity that is in this direction on average.

Stochastic Gradient



Logistic regression

- Instead of computing $\nabla_{\beta} \ell(\beta^{(t)})$, we compute a random approximation that is on average in the direction of $\nabla_{\beta} \ell(\beta^{(t)})$.

Logistic regression

- Instead of computing $\nabla_{\beta} \ell(\beta^{(t)})$, we compute a random approximation that is on average in the direction of $\nabla_{\beta} \ell(\beta^{(t)})$.
- At step t , we'll choose a random index $i_t \in \{1, \dots, n\}$, and use

$$g_{i_t} = \left(y_{i_t} - \mu_{i_t}(\beta^{(t)}) \right) x_{i_t}.$$

Logistic regression

- Instead of computing $\nabla_{\beta} \ell(\beta^{(t)})$, we compute a random approximation that is on average in the direction of $\nabla_{\beta} \ell(\beta^{(t)})$.
- At step t , we'll choose a random index $i_t \in \{1, \dots, n\}$, and use

$$g_{i_t} = \left(y_{i_t} - \mu_{i_t}(\beta^{(t)}) \right) x_{i_t}.$$

For instance, if we assume that the training examples are in a random order, we can cycle through the training data and interpret the choice of (x_i, y_i) as uniform.

Logistic regression

- Instead of computing $\nabla_{\beta} \ell(\beta^{(t)})$, we compute a random approximation that is on average in the direction of $\nabla_{\beta} \ell(\beta^{(t)})$.
- At step t , we'll choose a random index $i_t \in \{1, \dots, n\}$, and use

$$g_{i_t} = \left(y_{i_t} - \mu_{i_t}(\beta^{(t)}) \right) x_{i_t}.$$

For instance, if we assume that the training examples are in a random order, we can cycle through the training data and interpret the choice of (x_i, y_i) as uniform.

- Each stochastic gradient calculation takes $O(p)$ time.

Logistic regression

- Instead of computing $\nabla_{\beta} \ell(\beta^{(t)})$, we compute a random approximation that is on average in the direction of $\nabla_{\beta} \ell(\beta^{(t)})$.
- At step t , we'll choose a random index $i_t \in \{1, \dots, n\}$, and use

$$g_{i_t} = \left(y_{i_t} - \mu_{i_t}(\beta^{(t)}) \right) x_{i_t}.$$

For instance, if we assume that the training examples are in a random order, we can cycle through the training data and interpret the choice of (x_i, y_i) as uniform.

- Each stochastic gradient calculation takes $O(p)$ time.
- If n is large, we might obtain a good estimate long before we have seen n training examples.

Logistic regression

- Instead of computing $\nabla_{\beta} \ell(\beta^{(t)})$, we compute a random approximation that is on average in the direction of $\nabla_{\beta} \ell(\beta^{(t)})$.
- At step t , we'll choose a random index $i_t \in \{1, \dots, n\}$, and use

$$g_{i_t} = \left(y_{i_t} - \mu_{i_t}(\beta^{(t)}) \right) x_{i_t}.$$

For instance, if we assume that the training examples are in a random order, we can cycle through the training data and interpret the choice of (x_i, y_i) as uniform.

- Each stochastic gradient calculation takes $O(p)$ time.
- If n is large, we might obtain a good estimate long before we have seen n training examples.

Stochastic gradient

$$\beta^{(t+1)} = \beta^{(t)} + \eta \left(y_{i_t} - \mu_{i_t}(\beta^{(t)}) \right) x_{i_t}.$$

Logistic Regression:

- Gaussian generative to logistic discriminative models.
- Parameter estimates for logistic models.
 - Maximum likelihood: coupled non-linear equations.
 - Gradient ascent.
 - Stochastic gradient.
 - **(Detour: stochastic gradient for linear regression.)**
 - Newton's method: iteratively reweighted least squares.

Detour: Back to linear regression

Detour: Back to linear regression

Residual sum of squares

Detour: Back to linear regression

Residual sum of squares

$$\nabla_{\beta} RSS(\beta) = X'X\beta - X'y$$

Detour: Back to linear regression

Residual sum of squares

$$\nabla_{\beta} RSS(\beta) = X'X\beta - X'y = \sum_{i=1}^n x_i(x_i'\beta - y_i).$$

Detour: Back to linear regression

Residual sum of squares

$$\nabla_{\beta} RSS(\beta) = X'X\beta - X'y = \sum_{i=1}^n x_i(x_i'\beta - y_i).$$

- This gave us the normal equations.

Detour: Back to linear regression

Residual sum of squares

$$\nabla_{\beta} RSS(\beta) = X'X\beta - X'y = \sum_{i=1}^n x_i(x_i'\beta - y_i).$$

- This gave us the normal equations.
- We could also consider a gradient descent approach.

Detour: Back to linear regression

Residual sum of squares

$$\nabla_{\beta} RSS(\beta) = X'X\beta - X'y = \sum_{i=1}^n x_i(x'_i\beta - y_i).$$

- This gave us the normal equations.
- We could also consider a gradient descent approach.
- Or a stochastic gradient approach:

$$\beta^{(t+1)} = \beta^{(t)} + \eta x_{it} \left(y_{it} - x'_{it}\beta^{(t)} \right).$$

Detour: Back to linear regression

Residual sum of squares

$$\nabla_{\beta} RSS(\beta) = X'X\beta - X'y = \sum_{i=1}^n x_i(x'_i\beta - y_i).$$

- This gave us the normal equations.
- We could also consider a gradient descent approach.
- Or a stochastic gradient approach:

$$\beta^{(t+1)} = \beta^{(t)} + \eta x_{i_t} \left(y_{i_t} - x'_{i_t} \beta^{(t)} \right).$$

- This might be simpler for large scale problems.

Detour: Back to linear regression

Residual sum of squares

$$\nabla_{\beta} RSS(\beta) = X'X\beta - X'y = \sum_{i=1}^n x_i(x_i'\beta - y_i).$$

- This gave us the normal equations.
- We could also consider a gradient descent approach.
- Or a stochastic gradient approach:

$$\beta^{(t+1)} = \beta^{(t)} + \eta x_{it} (y_{it} - x_{it}'\beta^{(t)}).$$

- This might be simpler for large scale problems.
(e.g., suppose each x encodes a bag-of-words representation of a document.)

Detour: Back to linear regression

Residual sum of squares

$$\nabla_{\beta} RSS(\beta) = X'X\beta - X'y = \sum_{i=1}^n x_i(x'_i\beta - y_i).$$

- This gave us the normal equations.
- We could also consider a gradient descent approach.
- Or a stochastic gradient approach:

$$\beta^{(t+1)} = \beta^{(t)} + \eta x_{i_t} (y_{i_t} - x'_{i_t}\beta^{(t)}).$$

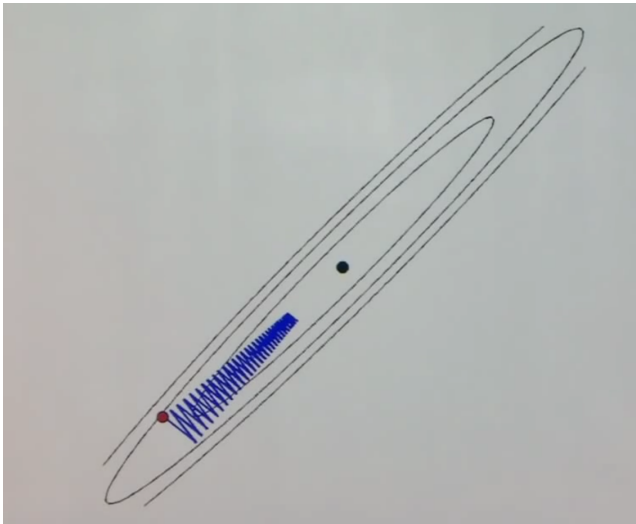
- This might be simpler for large scale problems.
(e.g., suppose each x encodes a bag-of-words representation of a document.)
- Convex regularization terms (like $\|\beta\|_2^2$ and $\|\beta\|_1$) can be easily incorporated.

Logistic Regression:

- Gaussian generative to logistic discriminative models.
- Parameter estimates for logistic models.
 - Maximum likelihood: coupled non-linear equations.
 - Gradient ascent.
 - Stochastic gradient.
 - (Detour: stochastic gradient for linear regression.)
 - **Newton's method: iteratively reweighted least squares.**

Newton's method

- Incorporating second derivative information can significantly speed up gradient methods.



Finding roots

Newton's method

Finding roots

- To solve $f(x) = 0$:

Finding roots

- To solve $f(x) = 0$:
- Start at x_0 .

Finding roots

- To solve $f(x) = 0$:
- Start at x_0 .
- Calculate the linear (first order Taylor series) approximation of f at x_0 :

$$f(x_0) + f'(x_0)(x - x_0).$$

Finding roots

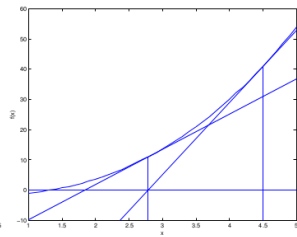
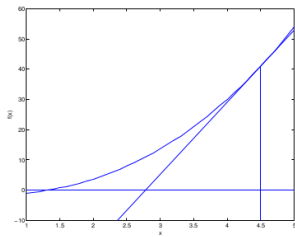
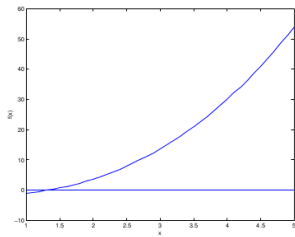
- To solve $f(x) = 0$:
- Start at x_0 .
- Calculate the linear (first order Taylor series) approximation of f at x_0 :

$$f(x_0) + f'(x_0)(x - x_0).$$

- Solve for $f(x) = 0$:

$$x_1 = x_0 - \frac{f(x_0)}{f'(x_0)}.$$

Newton's method



Finding maxima

Finding maxima

- To solve $\nabla f(x) = 0$:

Finding maxima

- To solve $\nabla f(x) = 0$:
- Start at x_0 .

Finding maxima

- To solve $\nabla f(x) = 0$:
- Start at x_0 .
- Calculate the linear (first order Taylor series) approximation of ∇f at x_0 :

$$\nabla f(x_0) + \nabla^2 f(x_0)(x - x_0).$$

Finding maxima

- To solve $\nabla f(x) = 0$:
- Start at x_0 .
- Calculate the linear (first order Taylor series) approximation of ∇f at x_0 :

$$\nabla f(x_0) + \nabla^2 f(x_0)(x - x_0).$$

- Solve for $\nabla f(x) = 0$:

$$x_1 = x_0 - [\nabla^2 f(x_0)]^{-1} \nabla f(x_0).$$

Newton's method

Finding maxima

- To solve $\nabla f(x) = 0$:
- Start at x_0 .
- Calculate the linear (first order Taylor series) approximation of ∇f at x_0 :

$$\nabla f(x_0) + \nabla^2 f(x_0)(x - x_0).$$

- Solve for $\nabla f(x) = 0$:

$$x_1 = x_0 - [\nabla^2 f(x_0)]^{-1} \nabla f(x_0).$$

- **Newton-Raphson:** $x_{t+1} = x_t - [\nabla^2 f(x_t)]^{-1} \nabla f(x_t).$

Newton-Raphson method for logistic regression

Logistic regression

Log likelihood:
$$\ell(\beta) = \sum_{i=1}^n y_i \log \mu_i(\beta) + (1 - y_i) \log(1 - \mu_i(\beta)),$$

Newton-Raphson method for logistic regression

Logistic regression

Log likelihood:
$$\ell(\beta) = \sum_{i=1}^n y_i \log \mu_i(\beta) + (1 - y_i) \log(1 - \mu_i(\beta)),$$

First derivative:
$$\nabla_{\beta} \ell(\beta) = \sum_{i=1}^n (y_i - \mu_i(\beta)) x_i$$

Newton-Raphson method for logistic regression

Logistic regression

Log likelihood:
$$\ell(\beta) = \sum_{i=1}^n y_i \log \mu_i(\beta) + (1 - y_i) \log(1 - \mu_i(\beta)),$$

First derivative:
$$\nabla_{\beta} \ell(\beta) = \sum_{i=1}^n (y_i - \mu_i(\beta)) x_i$$

Second derivative:
$$\nabla_{\beta}^2 \ell(\beta) = \sum_{i=1}^n -\mu_i(\beta)(1 - \mu_i(\beta)) x_i x_i'$$

Newton-Raphson method for logistic regression

Logistic regression

Log likelihood:
$$\ell(\beta) = \sum_{i=1}^n y_i \log \mu_i(\beta) + (1 - y_i) \log(1 - \mu_i(\beta)),$$

First derivative:
$$\begin{aligned}\nabla_{\beta} \ell(\beta) &= \sum_{i=1}^n (y_i - \mu_i(\beta)) x_i \\ &= X'(y - \mu).\end{aligned}$$

Second derivative:
$$\nabla_{\beta}^2 \ell(\beta) = \sum_{i=1}^n -\mu_i(\beta)(1 - \mu_i(\beta)) x_i x_i'$$

Newton-Raphson method for logistic regression

Logistic regression

Log likelihood:
$$\ell(\beta) = \sum_{i=1}^n y_i \log \mu_i(\beta) + (1 - y_i) \log(1 - \mu_i(\beta)),$$

First derivative:
$$\begin{aligned}\nabla_{\beta} \ell(\beta) &= \sum_{i=1}^n (y_i - \mu_i(\beta)) x_i \\ &= X'(y - \mu).\end{aligned}$$

Second derivative:
$$\begin{aligned}\nabla_{\beta}^2 \ell(\beta) &= \sum_{i=1}^n -\mu_i(\beta)(1 - \mu_i(\beta)) x_i x_i' \\ &= -X' \text{diag}(\mu(1 - \mu)) X.\end{aligned}$$

Newton-Raphson method for logistic regression

Logistic regression

$$\beta^{(t+1)} = \beta^{(t)} - \left[\nabla^2 \ell(\beta^{(t)}) \right]^{-1} \nabla \ell(\beta^{(t)})$$

Newton-Raphson method for logistic regression

Logistic regression

$$\begin{aligned}\beta^{(t+1)} &= \beta^{(t)} - \left[\nabla^2 \ell(\beta^{(t)}) \right]^{-1} \nabla \ell(\beta^{(t)}) \\ &= \beta^{(t)} \quad \left[\quad \quad \quad \right]^{-1} \quad .\end{aligned}$$

Newton-Raphson method for logistic regression

Logistic regression

$$\begin{aligned}\beta^{(t+1)} &= \beta^{(t)} - \left[\nabla^2 \ell(\beta^{(t)}) \right]^{-1} \nabla \ell(\beta^{(t)}) \\ &= \beta^{(t)} + \left[X' \text{diag}(\mu(1 - \mu)) X \right]^{-1} \cdot\end{aligned}$$

Newton-Raphson method for logistic regression

Logistic regression

$$\begin{aligned}\beta^{(t+1)} &= \beta^{(t)} - \left[\nabla^2 \ell(\beta^{(t)}) \right]^{-1} \nabla \ell(\beta^{(t)}) \\ &= \beta^{(t)} + \left[X' \text{diag}(\mu(1 - \mu)) X \right]^{-1} X'(y - \mu).\end{aligned}$$

Newton-Raphson method for logistic regression

Logistic regression

$$\begin{aligned}\beta^{(t+1)} &= \beta^{(t)} - \left[\nabla^2 \ell(\beta^{(t)}) \right]^{-1} \nabla \ell(\beta^{(t)}) \\ &= \beta^{(t)} + \left[X' \text{diag}(\mu(1 - \mu)) X \right]^{-1} X' (y - \mu).\end{aligned}$$

- Each iteration is like a least squares problem, except that the $x_i x_i'$ entries are weighted by $\mu_i(1 - \mu_i)$.

Newton-Raphson method for logistic regression

Logistic regression

$$\begin{aligned}\beta^{(t+1)} &= \beta^{(t)} - \left[\nabla^2 \ell(\beta^{(t)}) \right]^{-1} \nabla \ell(\beta^{(t)}) \\ &= \beta^{(t)} + \left[X' \text{diag}(\mu(1 - \mu)) X \right]^{-1} X' (y - \mu).\end{aligned}$$

- Each iteration is like a least squares problem, except that the $x_i x_i'$ entries are weighted by $\mu_i(1 - \mu_i)$.
- *Iteratively reweighted least squares.*

Newton-Raphson method for logistic regression

Logistic regression

$$\begin{aligned}\beta^{(t+1)} &= \beta^{(t)} - \left[\nabla^2 \ell(\beta^{(t)}) \right]^{-1} \nabla \ell(\beta^{(t)}) \\ &= \beta^{(t)} + \left[X' \text{diag}(\mu(1 - \mu)) X \right]^{-1} X' (y - \mu).\end{aligned}$$

- Each iteration is like a least squares problem, except that the $x_i x_i'$ entries are weighted by $\mu_i(1 - \mu_i)$.
- *Iteratively reweighted least squares.*
- Reweighting is according to how far μ_i is from 0 or 1.

Logistic Regression:

- Gaussian generative to logistic discriminative models.
- Parameter estimates for logistic models.
 - Maximum likelihood: coupled non-linear equations.
 - Gradient ascent.
 - Stochastic gradient.
 - (Detour: stochastic gradient for linear regression.)
 - Newton's method: iteratively reweighted least squares.