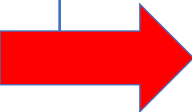# UVA CS 6316:
# Machine Learning

# Lecture 14 Extra: More about Logistic Regression

Dr. Yanjun Qi

University of Virginia

Department of Computer Science

# Today: Extra

✓ Bayes Classifier and MAP Rule?

- Bayes Classifier

- Expected Prediction Error

- 0-1 Loss function for Bayes Classifier

✓ Logistic regression

$$P(y|x) = \frac{e^{\beta x}}{1 + e^{\beta x}}$$

$$\Downarrow$$
$$\beta$$

- Parameter Estimation for LR

# Bayes classifiers

- Treat each feature attribute and the class label as random variables.

- Given a sample **x** with attributes ( $x_1$, $x_2$, … , $x_p$ ):
  - Goal is to predict its class *C.*
  - Specifically, we want to find the value of $C_i$ that maximizes $p( C_i | x_1, x_2, … , x_p )$.

- Can we estimate $p(C_i |\mathbf{x}) = p( C_i | x_1, x_2, … , x_p )$ directly from data?

# Bayes classifiers

- Treat each feature attribute and the class label as random variables.

- Given a sample **x** with attributes ( $x_1$, $x_2$, … , $x_p$ ):
  - Goal is to predict its class $C$.
  - Specifically, we want to find the value of $C_i$ that maximizes $p( C_i | x_1, x_2, … , x_p )$.

- Can we estimate $p(C_i | \mathbf{x}) = p( C_i | x_1, x_2, … , x_p )$ directly from data?

$$\rightarrow \{ C_1, C_2, …, C_L \}$$

# Bayes classifiers
➔ MAP classification rule

- Establishing a probabilistic model for classification
  ➔ **MAP** classification rule
    - **MAP**: **M**aximum **A P**osterior
    - Assign $x$ to $c^*$ if

$$P(C = c^* \mid \mathbf{X} = \mathbf{x}) > P(C = c \mid \mathbf{X} = \mathbf{x})$$
$$\text{for } c \neq c^*, \quad c = c_1, \cdots, c_L$$

Adapt from Prof. Ke Chen NB slides

# Bayes classifiers
➔ MAP classification rule

- Establishing a probabilistic model for classification
  ➔ **MAP** classification rule
    - **MAP**: **M**aximum **A P**osterior
    - Assign $x$ to $c^*$ if

$$P(C = c^* \mid \mathbf{X} = \mathbf{x}) > P(C = c \mid \mathbf{X} = \mathbf{x}) \quad c \neq c^*, \ c = c_1, \cdots, c_L$$

$$\left. \begin{array}{l} P(C = c_1 \mid x) \\ P(C = c_2 \mid x) \\ P(C = c_3 \mid x) \end{array} \right\} \ max \Rightarrow c_i$$

Adapt from Prof. Ke Chen NB slides

# Bayes Classifiers – MAP Rule

*Task*: Classify a new instance *X* based on a tuple of attribute values $X = \langle X_1, X_2, \ldots, X_p \rangle$ into one of the classes

$$c_{MAP} = \underset{c_j \in C}{\operatorname{argmax}} P(c_j \mid x_1, x_2, \ldots, x_p)$$

WHY ?

MAP = Maximum Aposteriori Probability

Adapt From Carols' prob tutorial

# 0-1 LOSS for Classification

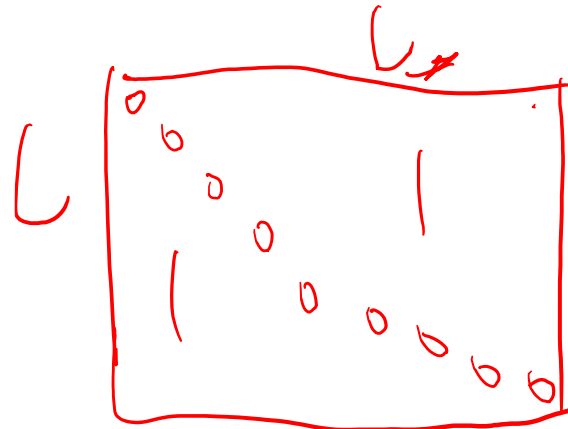if $k = \ell$, $L(k, \ell) = 0$

- Procedure for categorical output variable $C$

if $k \neq \ell$, $L((k, \ell) = 1$

- Frequently,  0-1 loss function used: $L(k, \ell)$

- $L(k, \ell)$ is the price paid for misclassifying an element from class $C_k$ as belonging to class $C_\ell$

➔ $L*L$ matrix

$C_1, C_2, \cdots, C_L$

# Expected prediction error (EPE)

- Expected prediction error (EPE), with expectation taken w.r.t. the joint distribution Pr(C,X)

  - $\Pr(C,X) = \Pr(C \mid X)\Pr(X)$

$\nearrow$ e.g. 0-1 loss

$E_x(x)$

$E_x(g(x))$

$$\text{EPE}(f) = E_{X,C}(L(C, f(X)))$$

$$= E_X \sum_{k=1}^{L} L[C_k, f(X)]\Pr(C_k \mid X)$$

Consider sample population distribution

$$EPE(f) = E_{X,C}\left(L(C, f(X))\right)$$

$$= E_X E_{C|X}\left[L(C, f(X)) \mid X\right]$$

$$E_C(C) = \sum_{i=1}^{L} c_i \hat{P}(c_i)$$

Discrete RV's Expectation

$$= E_X \sum_{k=1}^{L} L[C_k, f(X)] \, Pr(C_k \mid X)$$

$$\underset{f}{argmin} \, EPE(f(X))$$

$\Rightarrow$ Pointwise minimization   when   $X = x$

$$\Rightarrow \hat{f}(X=x) = \underset{f(x) \in C}{argmin} \sum_{k=1}^{L} L(C_k, f(x)) \, Pr(C_k \mid X=x)$$

$$\begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}$$

$$\Rightarrow \hat{f}(x) = \underset{C_K \in C_i}{argmax} \, Pr(C_k \mid X=x)$$

$$\begin{cases} C_2 \\ C_3 \\ \vdots \\ C_L \end{cases}$$

$$\begin{cases} P(C_1 \mid x) \\ P(C_2 \mid x) \\ \vdots \\ P(C_L \mid x) \end{cases}$$

# Expected prediction error (EPE)

$$\text{EPE}(f) = E_{X,C}(L(C, f(X))) = E_X \sum_{k=1}^{K} L[C_k, f(X)]\Pr(C_k \mid X)$$

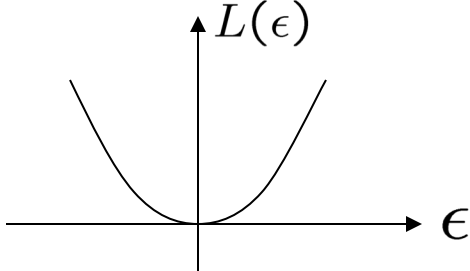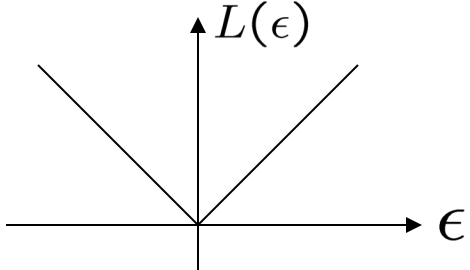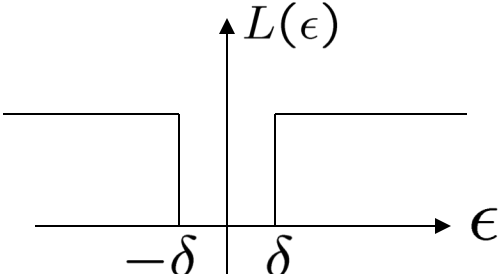Consider sample population distribution

- Pointwise minimization suffices

- ➜ simply $$\hat{f}(X) = \operatorname{argmin}_{g \in C} \sum_{k=1}^{K} L(C_k, g)\Pr(C_k \mid X = x)$$

$$\hat{f}(X) = C_k \text{ if}$$

$$\Pr(C_k \mid X = x) = \max_{g \in C} \Pr(g \mid X = x)$$

Bayes Classifier

# SUMMARY: WHEN Expected prediction error (EPE) USES DIFFERENT LOSS

| Loss Function | Estimator $\hat{f}(x)$ |
|---|---|
| $L_2$ <br> $L(\epsilon)$ vs $\epsilon$ (parabola) | $EPE = E_{X,Y}\left(Y - f(X)\right)^2$ <br> $\hat{f}(x) = E[Y \mid X = x]$ |
| $L_1$ <br> $L(\epsilon)$ vs $\epsilon$ (V-shape) | $\hat{f}(x) = \text{median}(Y \mid X = x)$ |
| 0-1 <br> $L(\epsilon)$ vs $\epsilon$ (step at $-\delta$, $\delta$) | $\hat{f}(x) = \arg\max_{Y} P(Y \mid X = x)$ <br><br> (Bayes classifier / MAP) |

# **Today: Extra**

✓ Why Bayes Classification – MAP Rule?

- ▪ Expected Prediction Error
- ▪ 0-1 Loss function for Bayes Classifier

✓ Logistic regression

➡ ▪ Parameter Estimation for LR

$$P(y|x) = \frac{e^{\beta x}}{1 + e^{\beta x}}$$

$$\Downarrow$$
$$\beta$$

# Newton's method for optimization

- The most basic second-order optimization algorithm
- Updating parameter with

$$GD: \theta_{k+1} = \theta_k - \alpha g_k$$

$$\text{Newton:} \quad \boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \mathbf{H}_K^{-1}\mathbf{g}_k$$

$$\underbrace{P \times P \quad P \times 1}_{P \times 1}$$

# Review: Hessian Matrix / n==2 case

Singlevariate        → multivariate            $f(x, y)$

- 1$^{st}$ derivative to gradient,

$$g = \nabla f = \begin{pmatrix} \dfrac{\partial f}{\partial x} \\[2ex] \dfrac{\partial f}{\partial y} \end{pmatrix}$$

- 2$^{nd}$ derivative to Hessian

$$H = \begin{pmatrix} \dfrac{\partial^2 f}{\partial x^2} & \dfrac{\partial^2 f}{\partial x \partial y} \\[3ex] \dfrac{\partial^2 f}{\partial x \partial y} & \dfrac{\partial^2 f}{\partial y^2} \end{pmatrix}$$

# Review: Hessian Matrix

Suppose that $f : \mathbb{R}^n \to \mathbb{R}$ is a function that takes a vector in $\mathbb{R}^n$ and returns a real number. Then the **Hessian** matrix with respect to $x$, written $\nabla_x^2 f(x)$ or simply as $H$ is the $n \times n$ matrix of partial derivatives,

$$\nabla_x^2 f(x) \in \mathbb{R}^{n \times n} = \begin{bmatrix} \frac{\partial^2 f(x)}{\partial x_1^2} & \frac{\partial^2 f(x)}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f(x)}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f(x)}{\partial x_2 \partial x_1} & \frac{\partial^2 f(x)}{\partial x_2^2} & \cdots & \frac{\partial^2 f(x)}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f(x)}{\partial x_n \partial x_1} & \frac{\partial^2 f(x)}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f(x)}{\partial x_n^2} \end{bmatrix}.$$

# Newton's method for optimization

- Making a quadratic/second-order Taylor series approximation

$$\widehat{f_{quad}}(\boldsymbol{\theta}) = f(\boldsymbol{\theta}_k) + \mathbf{g}_k^T(\boldsymbol{\theta} - \boldsymbol{\theta}_k) + \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_k)^T \mathbf{H}_k (\boldsymbol{\theta} - \boldsymbol{\theta}_k)$$

Finding the minimum solution of the above right quadratic approximation (quadratic function minimization is easy !)

$$\hat{f}(\theta) = f(\theta_k) + g_k^T(\theta - \theta_k) +$$

$$\frac{1}{2}(\theta - \theta_k)^T H_k (\theta - \theta_k)$$

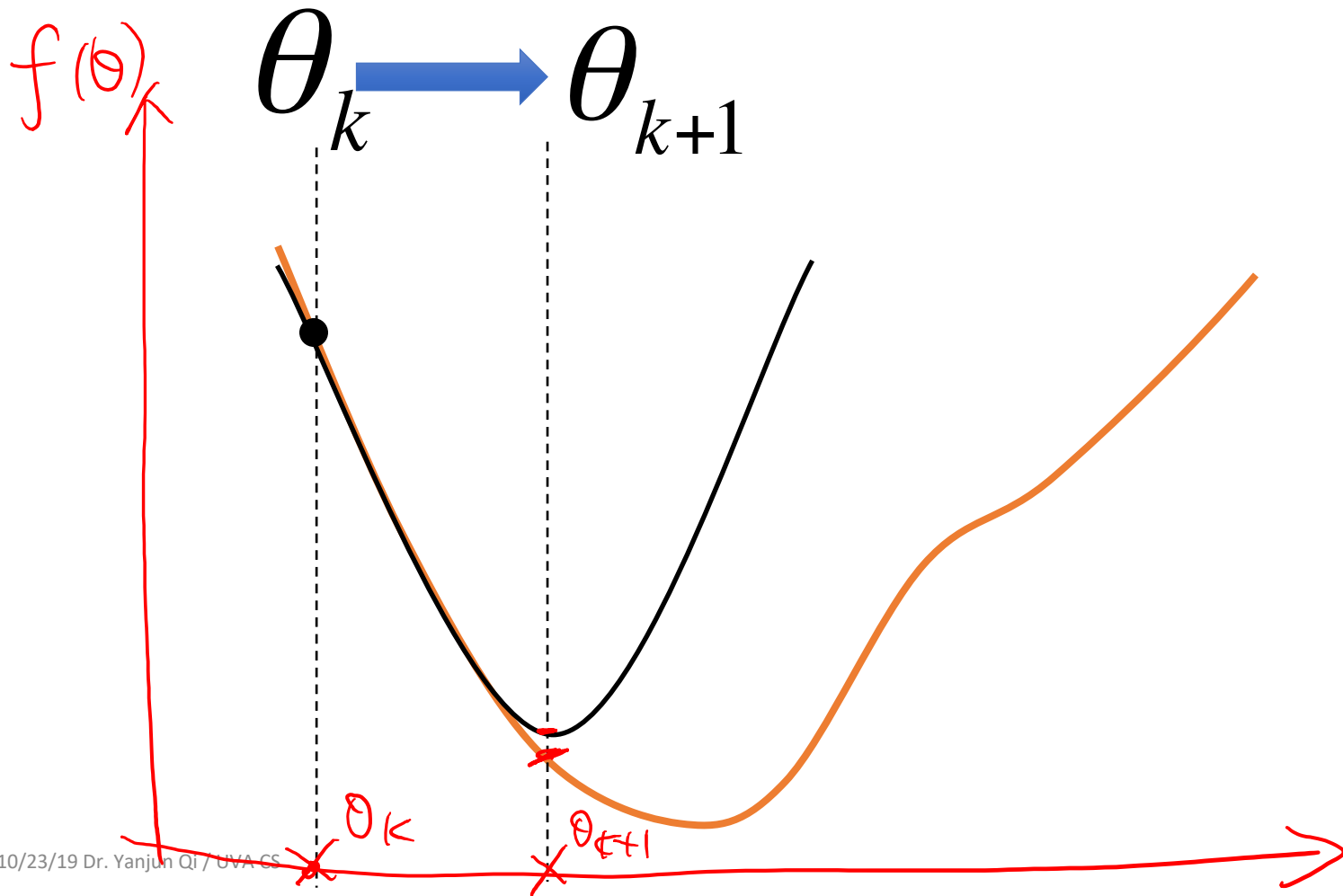$$\frac{1}{2}\left(\theta^T H_k \theta - 2\theta^T H_k \theta_k + \theta_k^T H_k \theta_k\right)$$

$$\frac{\partial \hat{f}(\theta)}{\partial \theta} = 0 + g_k + \frac{2}{2} H_k \theta - \frac{2}{2} H_k \theta_k \quad := = 0$$

See p24 handout

$$g_k + H_k(\theta - \theta_k) = 0$$

$$\Rightarrow \quad \theta = \theta_k - H_k^{-1} g_k$$

where $\begin{cases} H_k \in \mathbb{R}^{P \times P} \\ g_k \in \mathbb{R}^{P} \end{cases}$

# Newton's Method / second-order Taylor series approximation

$$\theta_k \longrightarrow \theta_{k+1}$$
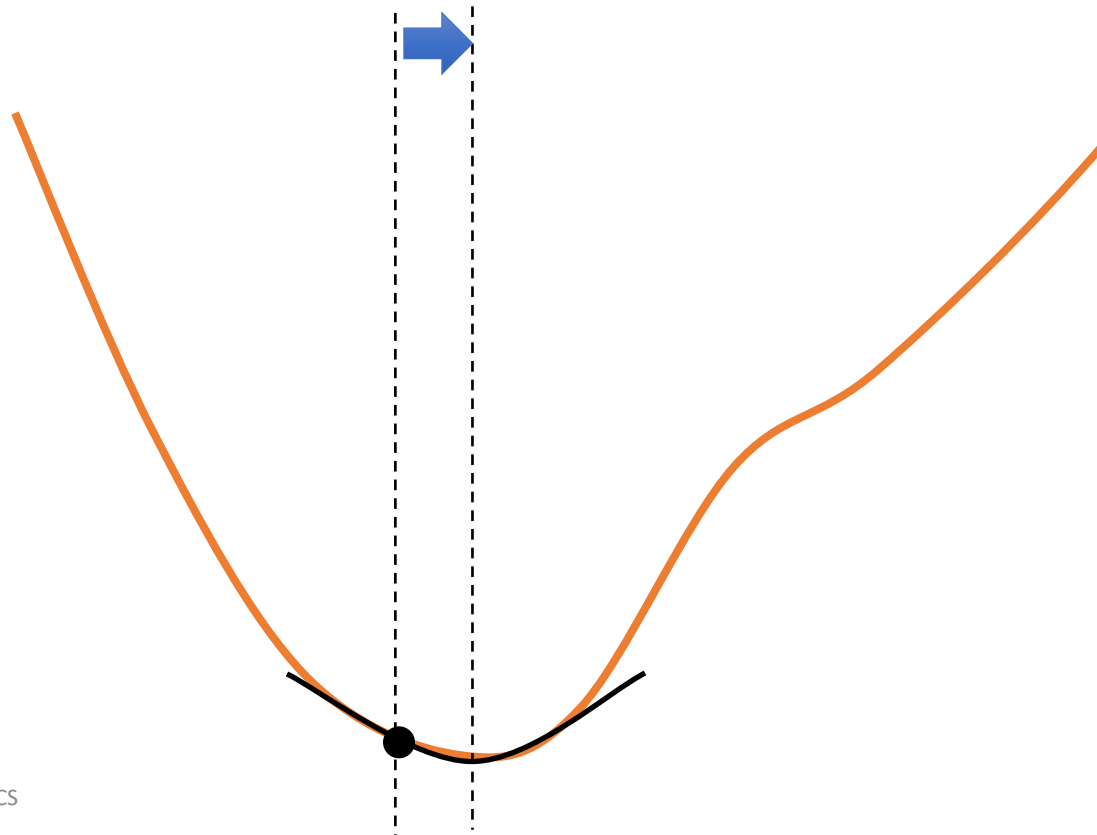
$f(\theta)$
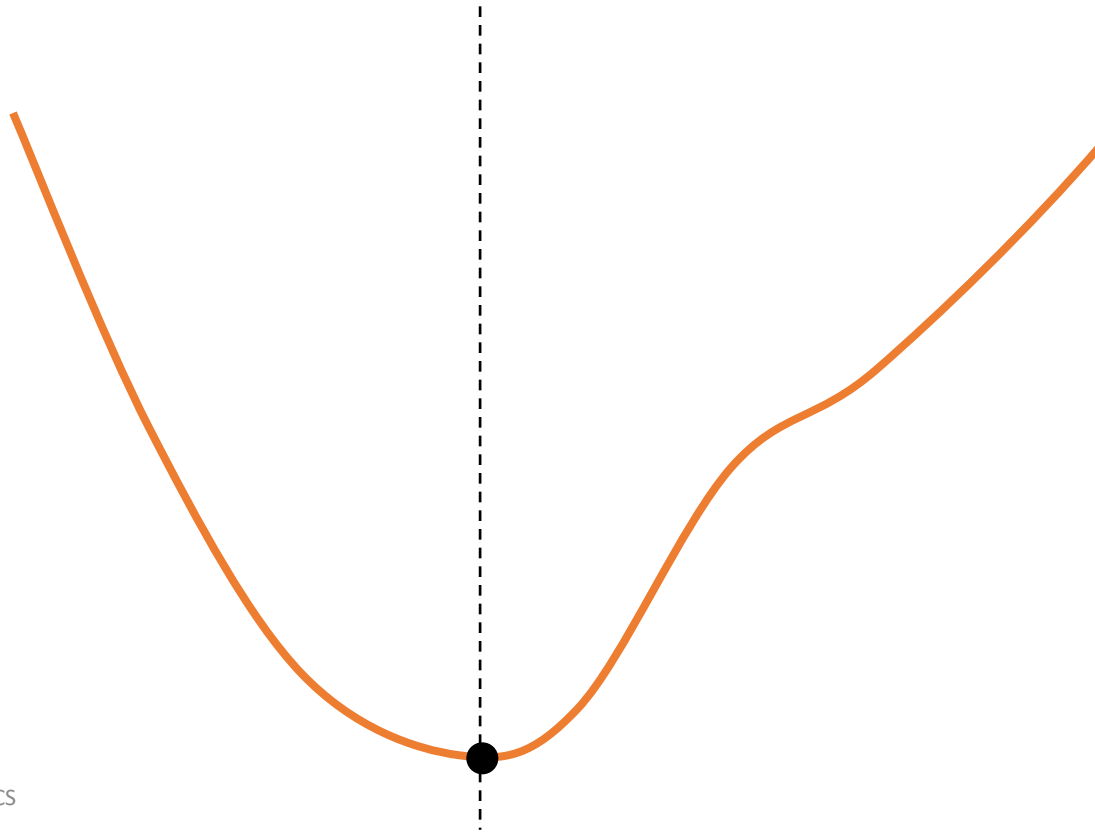
$\theta_k$

$\theta_{k+1}$

# Newton's Method / second-order Taylor series approximation

# Newton's Method / second-order Taylor series approximation

# Newton's Method / second-order Taylor series approximation

# Newton's Method

- At each step:

$$\theta_{k+1} = \theta_k - \frac{f'(\theta_k)}{f''(\theta_k)}$$

$$\theta_{k+1} = \theta_k - H^{-1}(\theta_k)\nabla f(\theta_k)$$

- Requires 1st and 2nd derivatives
- Quadratic convergence
- ➔ However, finding the inverse of the Hessian matrix is often expensive

# Newton vs. GD for optimization

- Newton: a quadratic/second-order Taylor series approximation

$$\theta_{k+1} = \theta_k - \frac{1}{H(\theta_k)} g(\theta_k)$$

$$\widehat{f}_{quad}(\boldsymbol{\theta}) = f(\boldsymbol{\theta}_k) + \mathbf{g}_k^T(\boldsymbol{\theta} - \boldsymbol{\theta}_k) + \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_k)^T \mathbf{H}_k(\boldsymbol{\theta} - \boldsymbol{\theta}_k)$$

Finding the minimum solution of the above right quadratic approximation (quadratic function minimization is easy !)

$$\widehat{f}_{quad}(\boldsymbol{\theta}) = f(\boldsymbol{\theta}_k) + \mathbf{g}_k^T(\boldsymbol{\theta} - \boldsymbol{\theta}_k) + \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_k)^T \frac{1}{\alpha}(\boldsymbol{\theta} - \boldsymbol{\theta}_k)$$

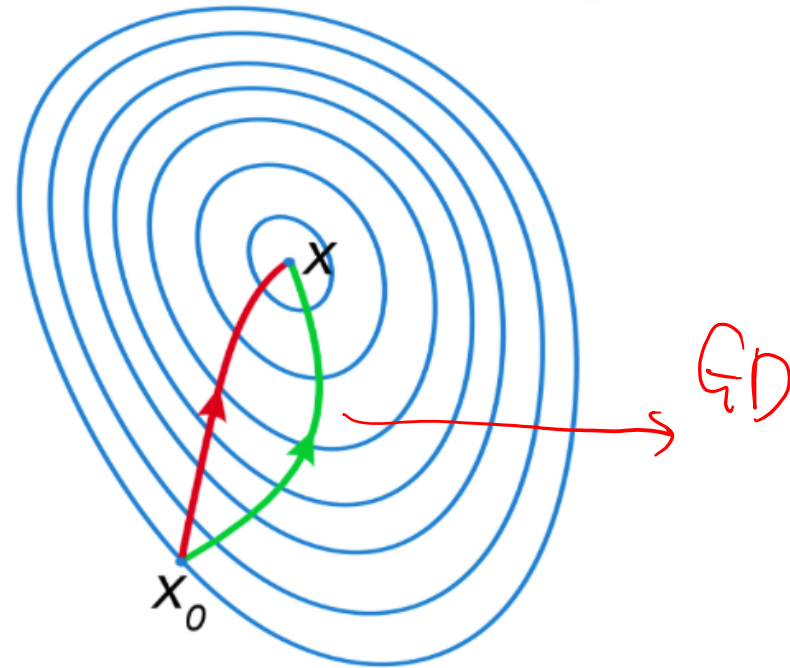$$\theta_{k+1} = \theta_k - \alpha\, g(\theta_k)$$

# Comparison

- Newton's method vs. Gradient descent

A comparison of gradient descent (green) and Newton's method (red) for minimizing a function (with small step sizes).

Newton's method uses curvature information to get a more direct route …

# MLE for Logistic Regression Training

Let's fit the logistic regression model for $K$=2, i.e., number of classes is 2

Training set: $(x_i, y_i)$, i=1,…,$N$

For Bernoulli distribution

$$p(y \mid x)^y (1-p)^{1-y}$$

(conditional )
Log-likelihood:

$$l(\beta) = \sum_{i=1}^{N} \{\log \Pr(Y = y_i \mid X = x_i)\}$$

**How?**

$$= \sum_{i=1}^{N} y_i \log(\Pr(Y = 1 \mid X = x_i)) + (1 - y_i)\log(\Pr(Y = 0 \mid X = x_i))$$

$$= \sum_{i=1}^{N} (y_i \log \frac{\exp(\beta^T x_i)}{1 + \exp(\beta^T x_i)}) + (1 - y_i)\log \frac{1}{1 + \exp(\beta^T x_i)})$$

$$= \sum_{i=1}^{N} (y_i \beta^T x_i - \log(1 + \exp(\beta^T x_i)))$$

$x_i$ are ($p$+1)-dimensional input vector with leading entry 1
\beta is a ($p$+1)-dimensional vector

We want to maximize the log-likelihood in order to estimate \beta

$$l(\beta) = \sum_{i=1}^{N} \{\log \Pr(Y = y_i \mid X = x_i)\}$$

$$\mathcal{Y}_i$$

$$P(\mathcal{Y}_i = 1 \mid x)$$

$$\log \{ \Pr(Y = y_i \mid X = x_i) = P(y_i \mid x_i) \} \Rightarrow \begin{cases} \mathcal{Y}_i = 1 \\ \mathcal{Y}_i = 0 \end{cases}$$

$$1 - P(\mathcal{Y}_i = 1 \mid x)$$

$$= \log \{ P(\mathcal{Y}_i = 1 \mid x)^{y_i} (1 - P(\mathcal{Y}_i = 1 \mid x_i))^{1 - y_i} \}$$

$$= y_i \log P(\mathcal{Y}_i = 1 \mid x) + (1 - y_i) \log (1 - P(\mathcal{Y}_i = 1 \mid x))$$

# Newton-Raphson for LR (optional)

$$\frac{\partial l(\beta)}{\partial \beta} = \sum_{i=1}^{N} (y_i - \frac{\exp(\beta^T x)}{1 + \exp(\beta^T x)}) x_i = 0$$

*Vector*

(*p*+1) Non-linear equations to solve for $\boxed{(p+1)}$ unknowns $\beta$

Solve by Newton-Raphson method:

$$\beta^{new} \leftarrow \beta^{old} - [(\frac{\partial^2 l(\beta)}{\partial \beta \partial \beta^T})]^{-1} \frac{\partial l(\beta)}{\partial \beta},$$

$$\text{where, } (\frac{\partial^2 l(\beta)}{\partial \beta \partial \beta^T}) = -\sum_{i=1}^{N} x_i x_i^T (\frac{\exp(\beta^T x_i)}{1 + \exp(\beta^T x_i)})(\frac{1}{1 + \exp(\beta^T x_i)})$$

**minimizes a quadratic approximation to the function we are really interested in.**

$$\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \mathbf{H}_K^{-1} \mathbf{g}_k$$

p(x_i ; β)

1 - p(x_i ; β)

# Newton-Raphson for LR…

$$\frac{\partial l(\beta)}{\partial \beta} = \sum_{i=1}^{N}(y_i - \frac{\exp(\beta^T x)}{1+\exp(\beta^T x)})x_i = X^T(y-p)$$


$$p(y=1|x) = \frac{e^{\beta x}}{1+e^{\beta x}}$$

$$(\frac{\partial^2 l(\beta)}{\partial \beta \partial \beta^T}) = -X^T W X$$

So, NR rule becomes:
$$\beta^{new} \leftarrow \beta^{old} + (X^T W X)^{-1} X^T(y-p),$$

$$X = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{bmatrix}_{N-by-(p+1)}, \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}_{N-by-1}, \quad p = \begin{bmatrix} \exp(\beta^T x_1)/(1+\exp(\beta^T x_1)) \\ \exp(\beta^T x_2)/(1+\exp(\beta^T x_2)) \\ \vdots \\ \exp(\beta^T x_N)/(1+\exp(\beta^T x_N)) \end{bmatrix}_{N-by-1},$$

$X : N\times(p+1)$ matrix of $x_i$

$y : N\times 1$ matrix of $y_i$

$p : N\times 1$ matrix of $p(x_i; \beta^{old})$

$W : N\times N$ diagonal matrix of $p(x_i; \beta^{old})(1-p(x_i; \beta^{old}))$

$$(\frac{\exp(\beta^T x_i)}{(1+\exp(\beta^T x_i))})(1-\frac{1}{(1+\exp(\beta^T x_i))})$$

# Newton-Raphson for LR…

- Newton-Raphson

  - $\beta^{new} = \beta^{old} + (X^T W X)^{-1} X^T (y - p)$

    $= (X^T W X)^{-1} X^T W (X \beta^{old} + W^{-1}(y - p))$

    $= (X^T W X)^{-1} X^T W z$

  Re expressing Newton step as weighted least square step

  $(X^T X)^{-1} X^T Y$

  $W$     $W z$

  - Adjusted response
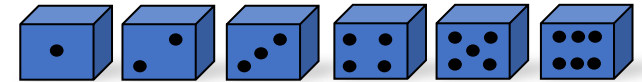
    $z = X \beta^{old} + W^{-1}(y - p)$

  - Iteratively reweighted least squares (IRLS)

    $\beta^{new} \leftarrow \arg\min_{\beta} (z - X\beta^T)^T W (z - X\beta^T)$

    $\leftarrow \arg\min_{\beta} (y - p)^T W^{-1} (y - p)$

# Binary ➜ Multinoulli Logistic Regression Model

$p(G|x)$

Directly models the posterior probabilities as the output of regression

$$\Pr(G=k|X=x)=\frac{\exp(\beta_{k0}+\beta_k^T x)}{1+\sum_{l=1}^{K-1}\exp(\beta_{l0}+\beta_l^T x)}, \; k=1,\dots,K-1$$

$$\Pr(G=K|X=x)=\frac{1}{1+\sum_{l=1}^{K-1}\exp(\beta_{l0}+\beta_l^T x)}$$

$x$ is $p$-dimensional input vector
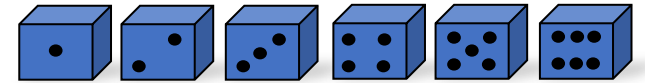
$\beta_k^T$ is a $p$-dimensional vector for each class $k$

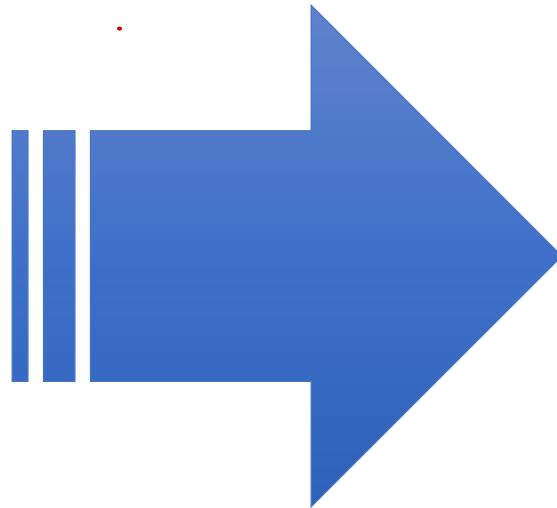Total number of parameters is $(K$-$1)(p+1)$ $\quad \beta_{k0}, \vec{\beta_k}, \; k=1,2,\dots,k-1$

Note that the class boundaries are linear

# Binary ➜ Multinoulli Logistic Regression Model

$$p(y=1|x)$$
$$\Rightarrow \frac{e^{\beta x}}{1+e^{\beta x}}$$

$$p(y=0|x)$$
$$\Rightarrow \frac{1}{1+e^{\beta x}}$$

$$\frac{e^{\beta_k^T x}}{1+e^{\beta_1^T x}+e^{\beta_2^T x}\ldots}$$

e.g.

$$\ln \frac{P(G=k|x)}{P(G=K|x)} = 0 \Rightarrow linear$$

$$\parallel$$

$$\beta_{k0}+\beta_k^T x$$

Note that the class boundaries are linear

# References

❑ Prof. Tan, Steinbach, Kumar's "Introduction to Data Mining" slide

❑ Prof. Andrew Moore's slides

❑ Prof. Eric Xing's slides

❑ Prof. Ke Chen NB slides

❑ Hastie, Trevor, et al. *The elements of statistical learning*. Vol. 2. No. 1. New York: Springer, 2009.