# Warm-up as You Walk In

Bernouli distribution:

$$Y \sim Bern(z)$$

$$p(y) = \begin{cases} z, & y = 1 \\ 1 - z, & y = 0 \end{cases}$$

What is the log likelihood for three i.i.d. samples, given parameter $z$:

$$\mathcal{D} = \{y^{(1)} = 1, y^{(2)} = 1, y^{(3)} = 0\}$$

$$L(z) =$$

$$\ell(z) =$$

# Introduction to Machine Learning

## Logistic Regression

Instructor: Pat Virtue

# Announcements

## Assignments:

- HW2 (written & programming)
  - Due Tue 2/4, 11:59 pm

## Early Feedback

- More mathematical rigor
- Consolidated course notes
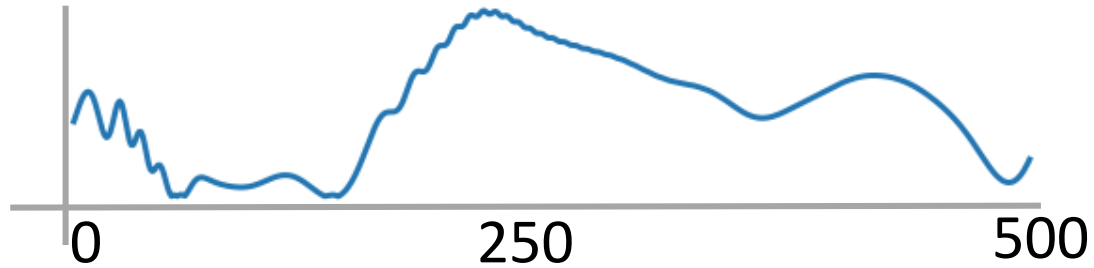- Lots of concepts, how does it all fit together?
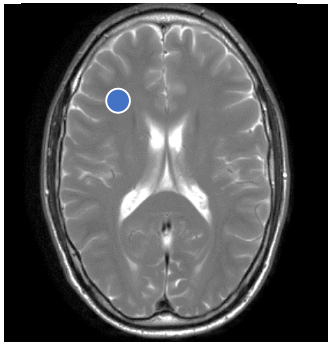
# Plan

## Last time

- Likelihood
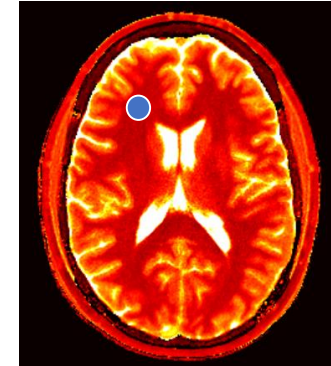- Density Estimation
- MLE for Density Estimation

## Today

- Wrap up MLE for linear regression
- Classification models
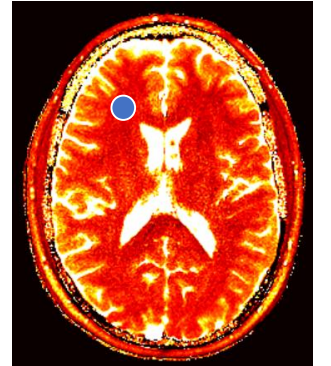- MLE for logistic regression

# MR Fingerprinting Assumptions

Forgot a really important assumption!!
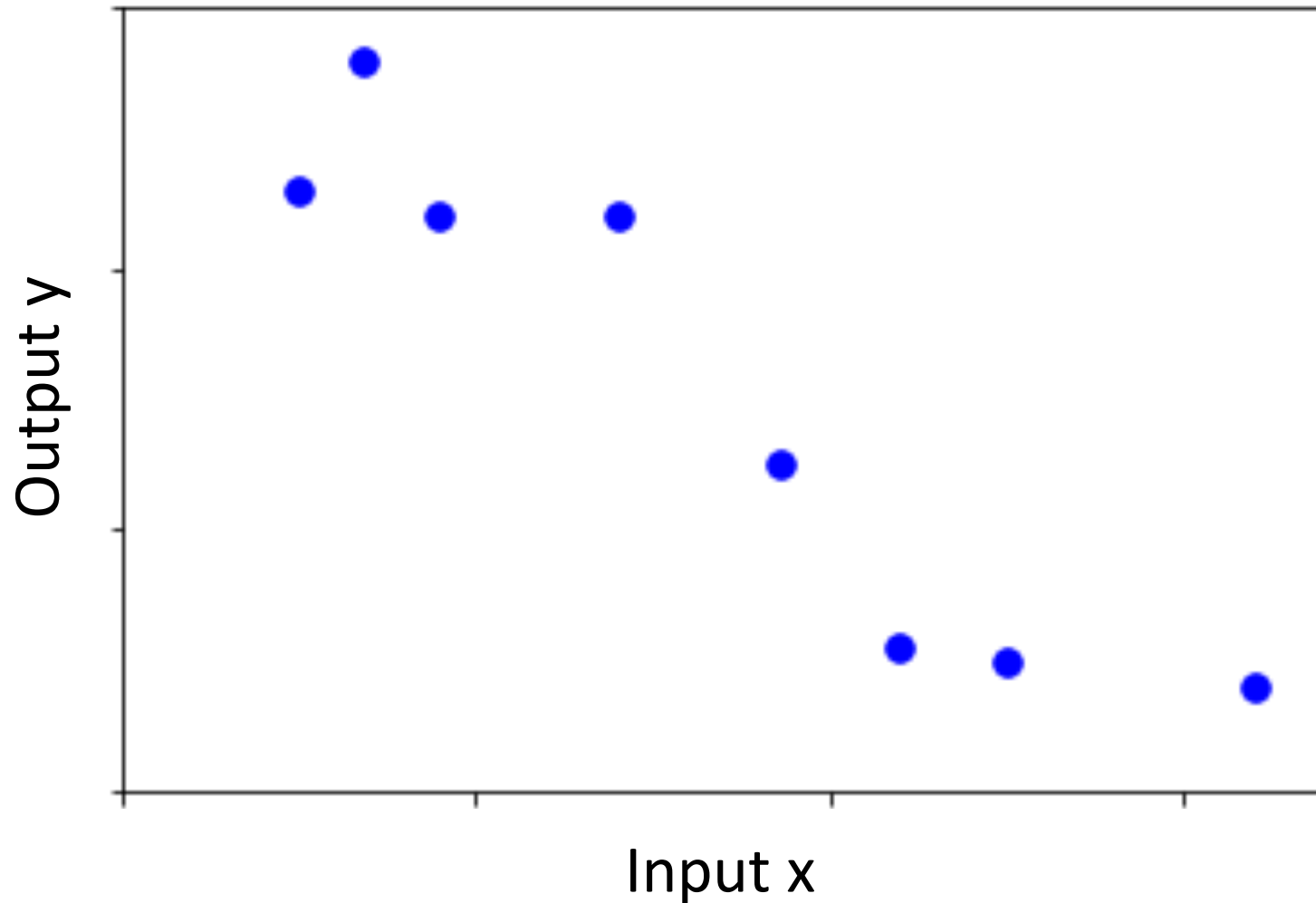
# Assumptions

What assumptions do we make with this data?

# Modelling $f(Y|X, \theta)$

# MLE for Linear Regression

How does our model of $f(Y|X, \theta)$ with the likelihood function?

$L(\theta)$

Maximum (Conditional) Likelihood Estimate

# M(C)LE for Linear Regression

$$L(\boldsymbol{w}, \boldsymbol{\sigma^2}) = \frac{1}{(2\pi\sigma^2)^{N/2}} e^{\left(\frac{-\sum_N (y^{(n)} - \boldsymbol{w}^T \boldsymbol{x}^{(n)})^2}{2\sigma^2}\right)}$$

$$\ell(\mu) = -\frac{N}{2} \log(2\pi) - \frac{N}{2} \log(\sigma^2) - \frac{\sum_{n=1}^{N} (x^{(n)} - \mu)^2}{2\sigma^2}$$

# M(C)LE for Linear Regression

How does M(C)LE optimization relate to least squares optimization?

$\ell(\boldsymbol{w}) =$

$J(\boldsymbol{w}) =$

# Piazza Poll 2:

Does $\min_{\boldsymbol{w}} -\ell(\boldsymbol{w})$ equal $\min_{\boldsymbol{w}} J(w)$ ?

# Linear Regression with Multiple Input Features

# Poll 1: Which vector is the correct $\boldsymbol{\theta}$?

# Classification Models

Linear Regression

# Classification Models

Linear Regression with Decision Boundary

# Classification Models

Linear Regression with Probability

# Modelling $p(Y|X, \theta)$

Bernoulli distribution of logistic function of linear model

# MLE for Bernoulli

Bernoulli distribution:

$$Y \sim Bern(z)$$

$$p(y) = \begin{cases} z, & y = 1 \\ 1 - z, & y = 0 \end{cases}$$

What is the log likelihood for three i.i.d. samples, given parameter $z$?

$$\mathcal{D} = \{y^{(1)} = 1, y^{(2)} = 1, y^{(3)} = 0\}$$

$$L(z) =$$

$$\ell(z) =$$

# MLE for Bernoulli

Bernoulli distribution:

$$Y \sim Bern(z)$$

$$p(y) = \begin{cases} z, & y = 1 \\ 1 - z, & y = 0 \end{cases}$$

What is the log likelihood for three i.i.d. samples, given parameter $z$?

$$\mathcal{D} = \{y^{(1)} = 1, y^{(2)} = 1, y^{(3)} = 0\}$$

$$L(z) =$$

$$\ell(z) =$$

# MLE for Bernoulli

Bernoulli distribution:

$$Y \sim Bern(z)$$

$$p(y) = \begin{cases} z, & y = 1 \\ 1 - z, & y = 0 \end{cases}$$

What is the log likelihood for three i.i.d. samples, given parameter $z$?

$$\mathcal{D} = \{y^{(1)} = 1, y^{(2)} = 1, y^{(3)} = 0\}$$

$$L(z) = z \cdot z \cdot (1 - z) \qquad = \prod_n z^{y^{(n)}} (1 - z)^{(1 - y^{(n)})}$$

$$\ell(z) = \log z + \log z + \log(1 - z) \quad = \sum_n y^{(n)} \log z + \left(1 - y^{(n)}\right) \log(1 - z)$$

# M(C)LE for Logistic Regression

$p(Y \mid X, \theta)$

$p(Y \mid X, \boldsymbol{w}) = \prod_{n=1}^{N} p(y^{(n)} \mid \boldsymbol{x}^{(n)}, \boldsymbol{w})$

Model $Y$ as a Bernoulli distribution, but the temporary $z$ is now based on the logistic function of our linear model of input $\boldsymbol{x}$

$$Y \sim Bern(\mu), \qquad \mu = g(\boldsymbol{w}^T \boldsymbol{x}), \qquad g(z) = \frac{1}{1+e^{-z}}$$

What is the *conditional* log likelihood?

$L(\boldsymbol{w}) =$

$\ell(\boldsymbol{w}) =$

# M(C)LE for Logistic Regression

$p(Y \mid X, \theta)$

$p(Y \mid X, \boldsymbol{w}) = \prod_{n=1}^{N} p(y^{(n)} \mid \boldsymbol{x}^{(n)}, \boldsymbol{w})$

Model $Y$ as a Bernoulli distribution, but the temporary $z$ is now based on the logistic function of our linear model of input $\boldsymbol{x}$

$$Y \sim Bern(\mu), \qquad \mu = g(\boldsymbol{w}^T \boldsymbol{x}), \qquad g(z) = \frac{1}{1+e^{-z}}$$

What is the *conditional* log likelihood?

$$L(\boldsymbol{w}) = \prod_n g(\boldsymbol{w}^T \boldsymbol{x}^{(n)})^{y^{(n)}} \left(1 - g(\boldsymbol{w}^T \boldsymbol{x}^{(n)})\right)^{(1-y^{(n)})}$$

$$\ell(\boldsymbol{w}) = \sum_n \left(y^{(n)} \log g(\boldsymbol{w}^T \boldsymbol{x}^{(n)}) + (1 - y^{(n)}) \log\left(1 - g(\boldsymbol{w}^T \boldsymbol{x}^{(n)})\right)\right)$$

# M(C)LE for Logistic Regression

$$z = f(\boldsymbol{w}, \boldsymbol{x}) = \boldsymbol{w}^T \boldsymbol{x} \qquad \mu = g(z) = \frac{1}{1+e^{-z}}$$

$$\nabla_{\boldsymbol{w}} f(\boldsymbol{w}, \boldsymbol{x}) = \boldsymbol{x} \qquad \frac{dg}{dz} = g(z)\big(1 - g(z)\big) = \mu(1 - \mu)$$

$$\ell(\boldsymbol{w}) = \sum_n \Big(y^{(n)} \log \mu^{(n)} + \big(1 - y^{(n)}\big) \log\big(1 - \mu^{(n)}\big)\Big)$$

$$\frac{\partial \ell}{\partial \boldsymbol{w}} =$$

# M(C)LE for Logistic Regression

$$z = f(\boldsymbol{w}, \boldsymbol{x}) = \boldsymbol{w}^T \boldsymbol{x} \qquad \mu = g(z) = \frac{1}{1+e^{-z}}$$

$$\nabla_{\boldsymbol{w}} f(\boldsymbol{w}, \boldsymbol{x}) = \boldsymbol{x} \qquad \frac{dg}{dz} = g(z)\big(1 - g(z)\big) = \mu(1 - \mu)$$

$$\ell(\boldsymbol{w}) = \sum_n \left( y^{(n)} \log \mu^{(n)} + \left(1 - y^{(n)}\right) \log\left(1 - \mu^{(n)}\right) \right)$$

$$\frac{\partial \ell}{\partial \boldsymbol{w}} = \sum_n \left( \frac{y^{(n)}}{\mu^{(n)}} - \frac{1 - y^{(n)}}{1 - \mu^{(n)}} \right) \frac{\partial g}{\partial f} \frac{\partial f}{\partial \boldsymbol{w}}$$

$$= \sum_n \left( \frac{y^{(n)} - \mu^{(n)}}{\mu^{(n)}\left(1 - \mu^{(n)}\right)} \right) \mu^{(n)}\left(1 - \mu^{(n)}\right) \boldsymbol{x}^{(n)T}$$

$$= \sum_n \left( y^{(n)} - \mu^{(n)} \right) \boldsymbol{x}^{(n)T}$$

# M(C)LE for Logistic Regression

$$z = f(\boldsymbol{w}, \boldsymbol{x}) = \boldsymbol{w}^T \boldsymbol{x} \qquad \mu = g(z) = \frac{1}{1+e^{-z}}$$

$$\ell(\boldsymbol{w}) = \sum_n \left( y^{(n)} \log \mu^{(n)} + \left(1 - y^{(n)}\right) \log\left(1 - \mu^{(n)}\right) \right)$$

$$\nabla_{\boldsymbol{w}} \ell(\boldsymbol{w}) = \sum_n \left( y^{(n)} - \mu^{(n)} \right) \boldsymbol{x}^{(n)}$$

$$\nabla_{\boldsymbol{w}} \ell(\boldsymbol{w}) = 0?$$

No closed form solution ☹

Back to iterative methods. Solve with (stochastic) gradient descent, Newton's method, or Iteratively Reweighted Least Squares (IRLS)

# Logistic Function

Cool note: Logistic function is related the invers of logit function!

Odds: Ratio of two probabilities. For $Y \sim Bern(p)$, $\dfrac{p(Y=1)}{p(Y=0)} = \dfrac{p}{1-p}$

Logit function: Log odds. $\log \dfrac{p(Y=1)}{p(Y=0)} = \log \dfrac{p}{1-p}$

$$z = logit(p) = \log \frac{p}{1-p}$$

$$p = logit^{-1}(z) = \frac{1}{1+e^{-z}}$$

# Log Odds and Logistic Regression

Formulate log odds as linear model of X:

$$\log \frac{p(Y = 1 \mid X = \boldsymbol{x}, \boldsymbol{w})}{p(Y = 0 \mid X = \boldsymbol{x}, \boldsymbol{w})} = \boldsymbol{w}^T \boldsymbol{x}$$

Equivalent to logistic representation:

$$p(Y = 1 \mid X = \boldsymbol{x}, \boldsymbol{w}) = \frac{1}{1 + e^{-\boldsymbol{w}^T \boldsymbol{x}}}$$

# Log Odds and Logistic Regression (Multi-class!)

Formulate log odds as linear model of X:

$$\log \frac{p(Y = 1 \mid X = \boldsymbol{x}, \boldsymbol{W})}{p(Y = K \mid X = \boldsymbol{x}, \boldsymbol{W})} = \boldsymbol{w}_1^T \boldsymbol{x}$$

$$\log \frac{p(Y = 2 \mid X = \boldsymbol{x}, \boldsymbol{W})}{p(Y = K \mid X = \boldsymbol{x}, \boldsymbol{W})} = \boldsymbol{w}_2^T \boldsymbol{x}$$

$$\vdots$$

$$\log \frac{p(Y = K - 1 \mid X = \boldsymbol{x}, \boldsymbol{W})}{p(Y = K \mid X = \boldsymbol{x}, \boldsymbol{W})} = \boldsymbol{w}_{K-1}^T \boldsymbol{x}$$

Equivalent to softmax representation:

$$p(Y = k \mid X = \boldsymbol{x}, W) = \frac{e^{w_k^T x}}{1 + \sum_{j=1}^{K-1} e^{w_j^T x}}$$

$$p(Y = K \mid X = \boldsymbol{x}, W) = \frac{1}{1 + \sum_{j=1}^{K-1} e^{w_j^T x}}$$

OR

$$p(Y = k \mid X = \boldsymbol{x}, W) = \frac{e^{w_k^T x}}{\sum_{j=1}^{K} e^{w_j^T x}}$$

# Multi-class Logistic Regression

$$p(Y \mid X, \theta)$$

$$p(Y \mid X, \boldsymbol{W}) = \prod_{n=1}^{N} p(y^{(n)} \mid \boldsymbol{x}^{(n)}, \boldsymbol{W})$$

$$p\left(y^{(n)} = k \mid X = \boldsymbol{x}^{(n)}, W\right) = \frac{e^{w_k^T x^{(n)}}}{\sum_{j=1}^{K} e^{w_j^T x^{(n)}}}$$

What is the *conditional* likelihood?

$$L(\boldsymbol{w}) = \prod_n \frac{e^{w_k^T x^{(n)}}}{\sum_{j=1}^{K} e^{w_j^T x^{(n)}}}$$

What is the hypothesis function?

$$\hat{y} = h_{\boldsymbol{W}}(\boldsymbol{x}) =$$