

# Warm-up as You Walk In

Bernouli distribution:

$$Y \sim \text{Bern}(z)$$

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



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

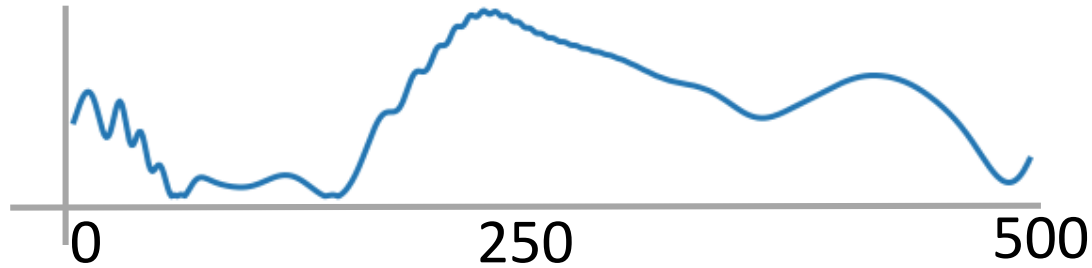
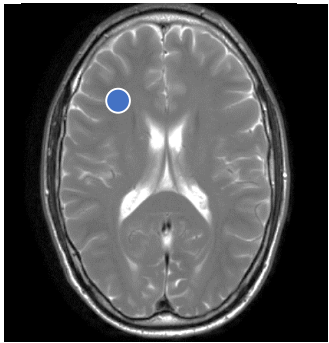
- Likelihood
- Density Estimation
- MLE for Density Estimation

## Today

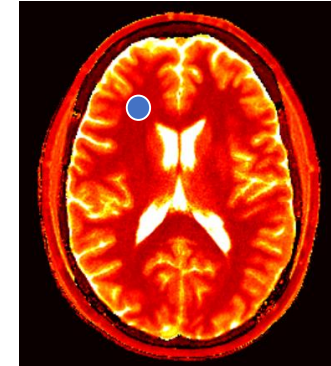
- Wrap up MLE for linear regression
- Classification models
- MLE for logistic regression

# MR Fingerprinting Assumptions

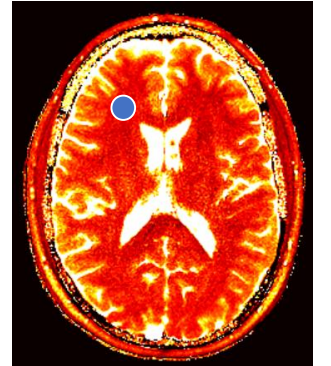
Forgot a really important assumption!!



T1

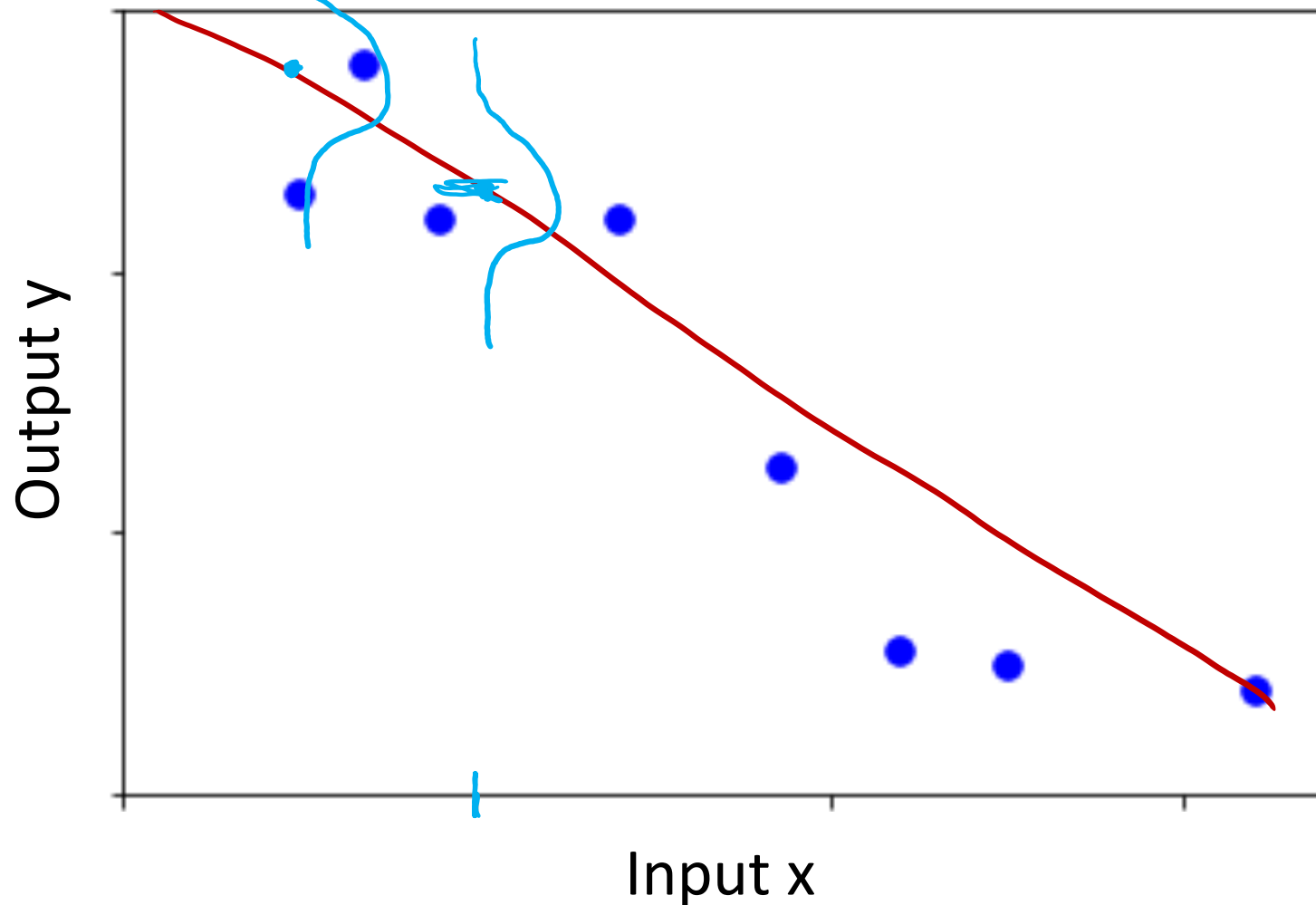


T2



# Assumptions

What assumptions do we make with this data?



$$\underline{y} = w^T x + \epsilon$$

$$\epsilon \sim \mathcal{N}(0, \sigma^2)$$

$$Y \sim \mathcal{N}(w^T x, \sigma^2)$$

Two purple arrows point from the  $w^T x$  term and the  $\sigma^2$  term in the equation above to the corresponding terms in the equation below.

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

$$f(D|\theta)$$

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

Conditional likelihood

$$f(y^{(n)}|x^{(n)}, \vec{w}, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{\left(\frac{-(y^{(n)} - \vec{w}^T x^{(n)})^2}{2\sigma^2}\right)}$$

$$f(\vec{y}|\vec{x}, \vec{w}, \sigma^2) = \prod_n f(y^{(n)}|x^{(n)}, \vec{w}, \sigma^2) \quad \leftarrow$$

Density Estimation

$$X \sim \mathcal{N}(\mu, \sigma^2)$$

$$f(\underline{x} | \underline{\mu}, \underline{\sigma^2})$$

$$f(\underline{D} | \underline{\theta})$$

# 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

$$e^{z_1} e^{z_2} = e^{z_1 + z_2}$$

$$L(\mathbf{w}, \sigma^2) = \frac{1}{\underbrace{(2\pi\sigma^2)^{N/2}}_{\text{red underline}}} e^{\left( \frac{-\sum_N (y^{(n)} - \mathbf{w}^T \mathbf{x}^{(n)})^2}{2\sigma^2} \right)}$$

$$l(\mathbf{w}, \sigma^2) = -\frac{N}{2} \log(2\pi) - \frac{N}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \underbrace{\sum (y^{(n)} - \mathbf{w}^T \mathbf{x}^{(n)})^2}_{\text{red underline}}$$

$$\frac{\partial l}{\partial \mathbf{w}} = 0$$

$$\hat{\mathbf{w}}_{ML} = ?$$

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

# M(C)LE for Linear Regression

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

$$l(w, \sigma^2) = -\frac{N}{2} \log(2\pi) - \frac{N}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \sum (y^{(n)} - w^T x^{(n)})^2$$

$$J(\mathbf{w}) = \frac{1}{N} \|\vec{y} - X\mathbf{w}\|_2^2$$

## Piazza Poll 2:

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

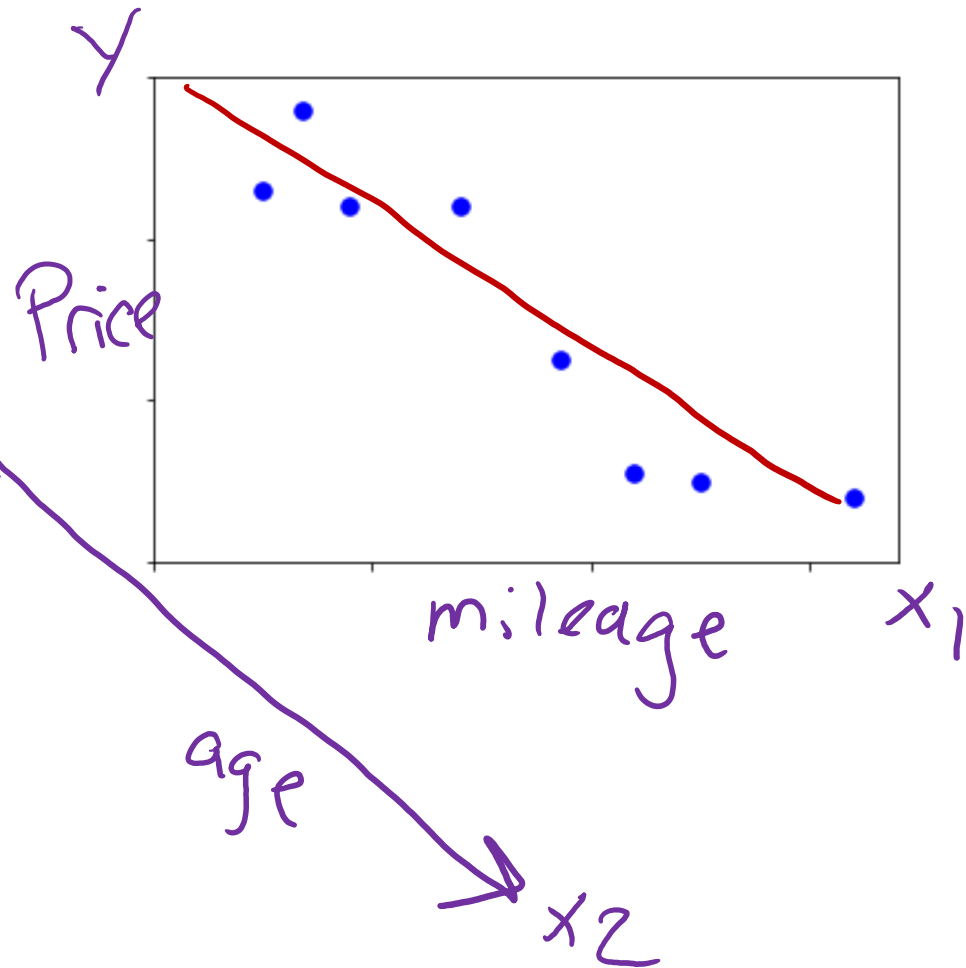
$$\ell(\vec{\mathbf{w}}, \sigma^2) = -\frac{N}{2} \log(2\pi) - \frac{N}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \sum (y^{(n)} - \vec{\mathbf{w}}^T \mathbf{x}^{(n)})^2$$

$$J(\vec{\mathbf{w}}) = \frac{1}{2} \|\vec{\mathbf{y}} - \mathbf{X} \vec{\mathbf{w}}\|_2^2$$

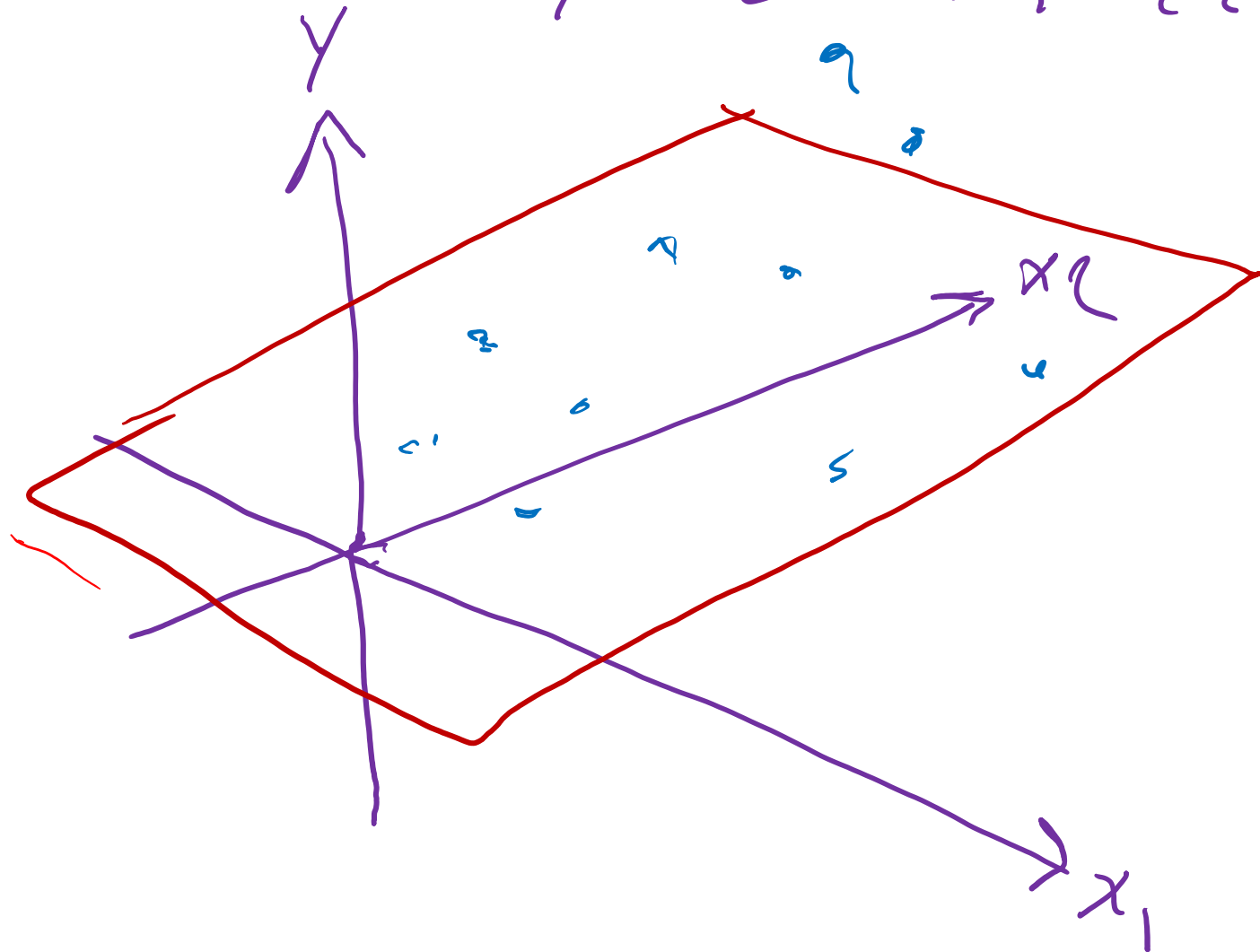
$$2 \sum (y^{(n)} - \mathbf{w}^T \mathbf{x}^{(n)}) \mathbf{x}^{(n)}$$

# Linear Regression with Multiple Input Features

$$y = w_0 + w_1 x$$



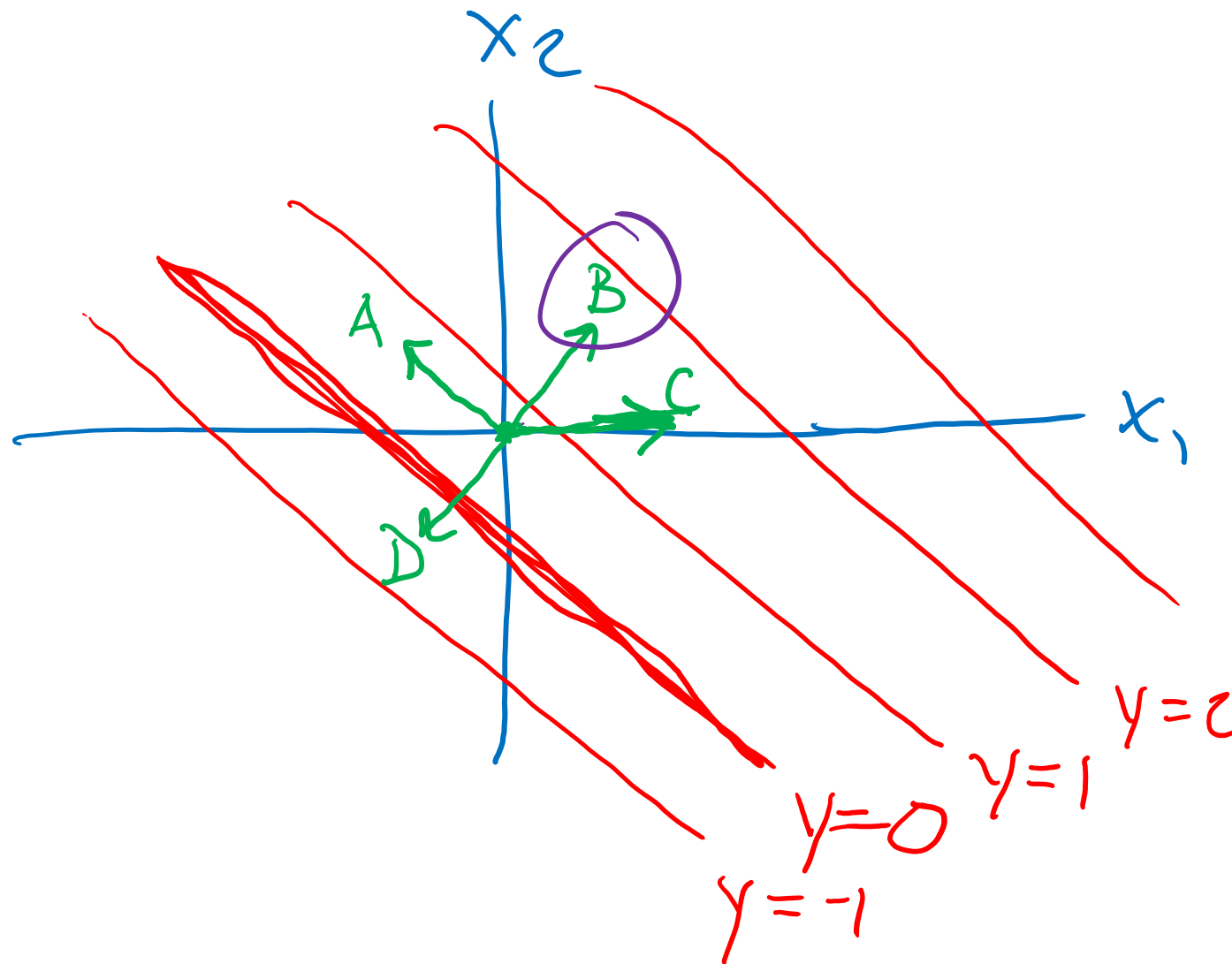
$$y = w_0 + w_1 x_1 + w_2 x_2$$



Poll 1: Which vector is the correct  $\theta$ ?

$$\theta = [w_1 \ w_2]^T$$

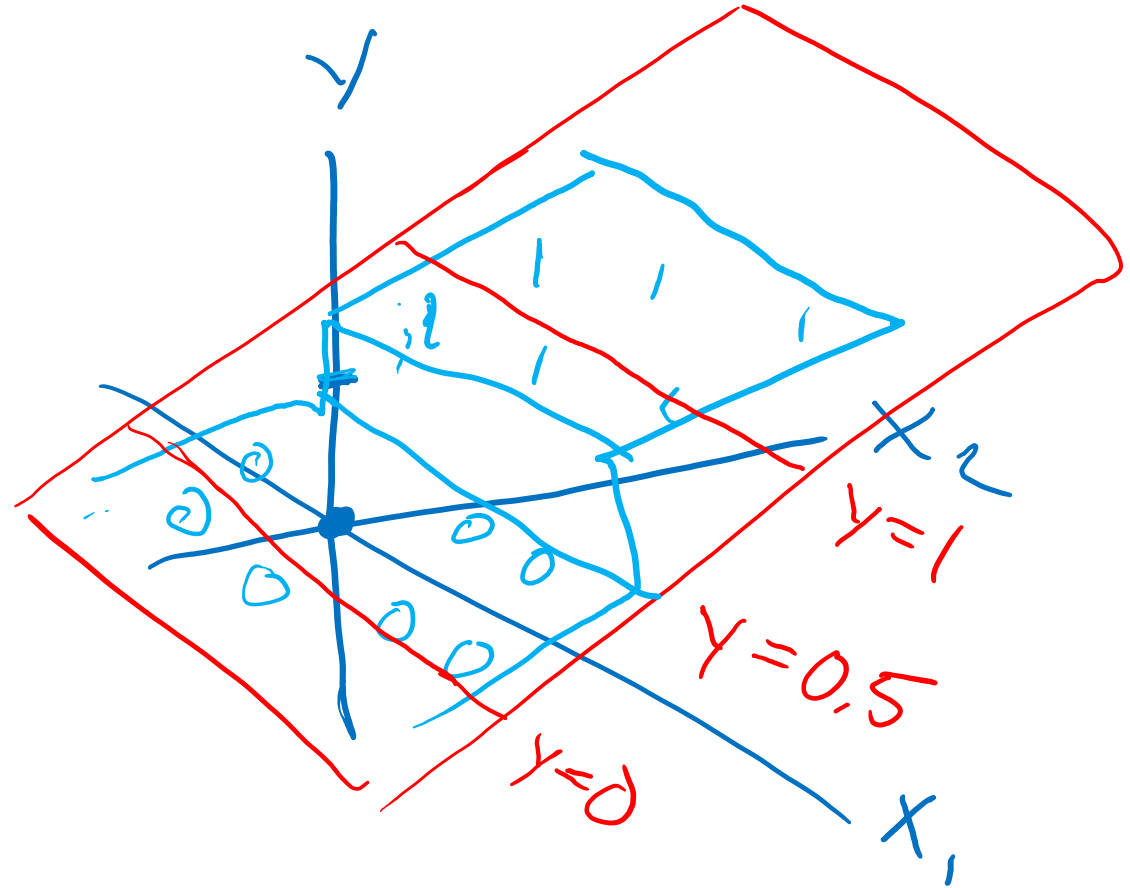
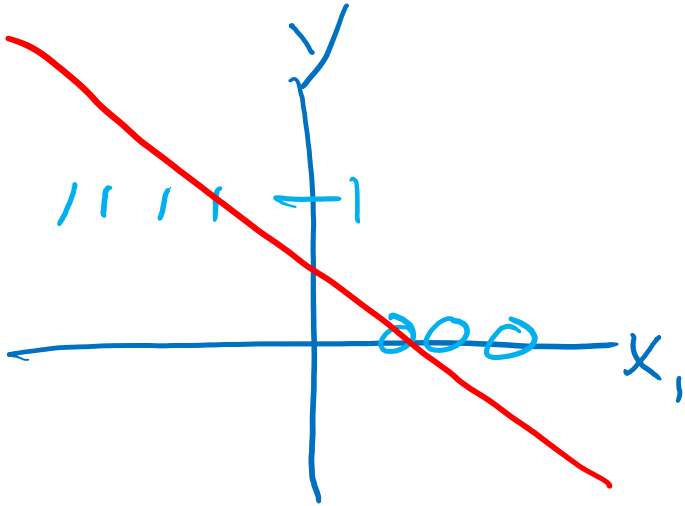
$$\vec{\theta}^T \vec{x}$$



E ?

# Classification Models

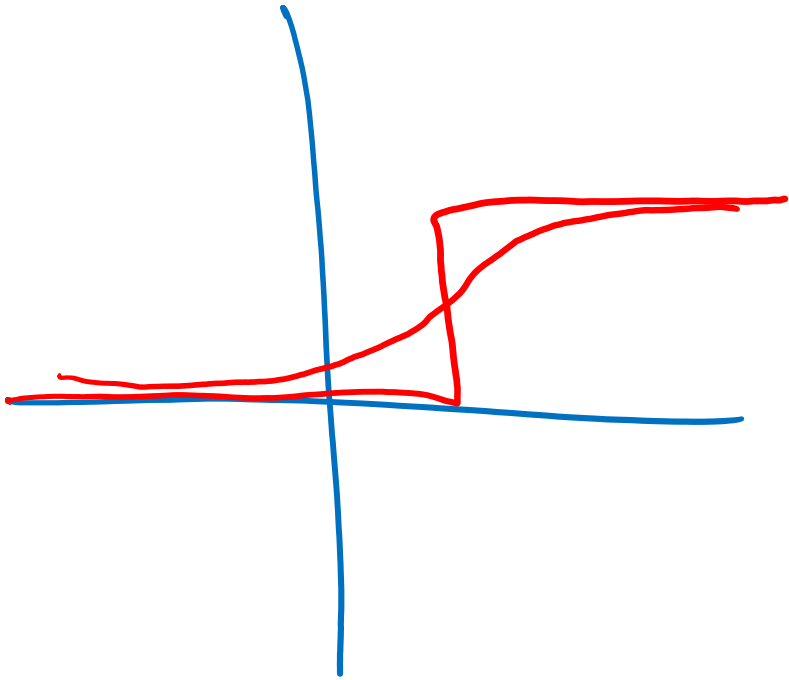
## Linear Regression



$$y = w_0 + x^T w$$

# Classification Models

## Linear Regression with Decision Boundary



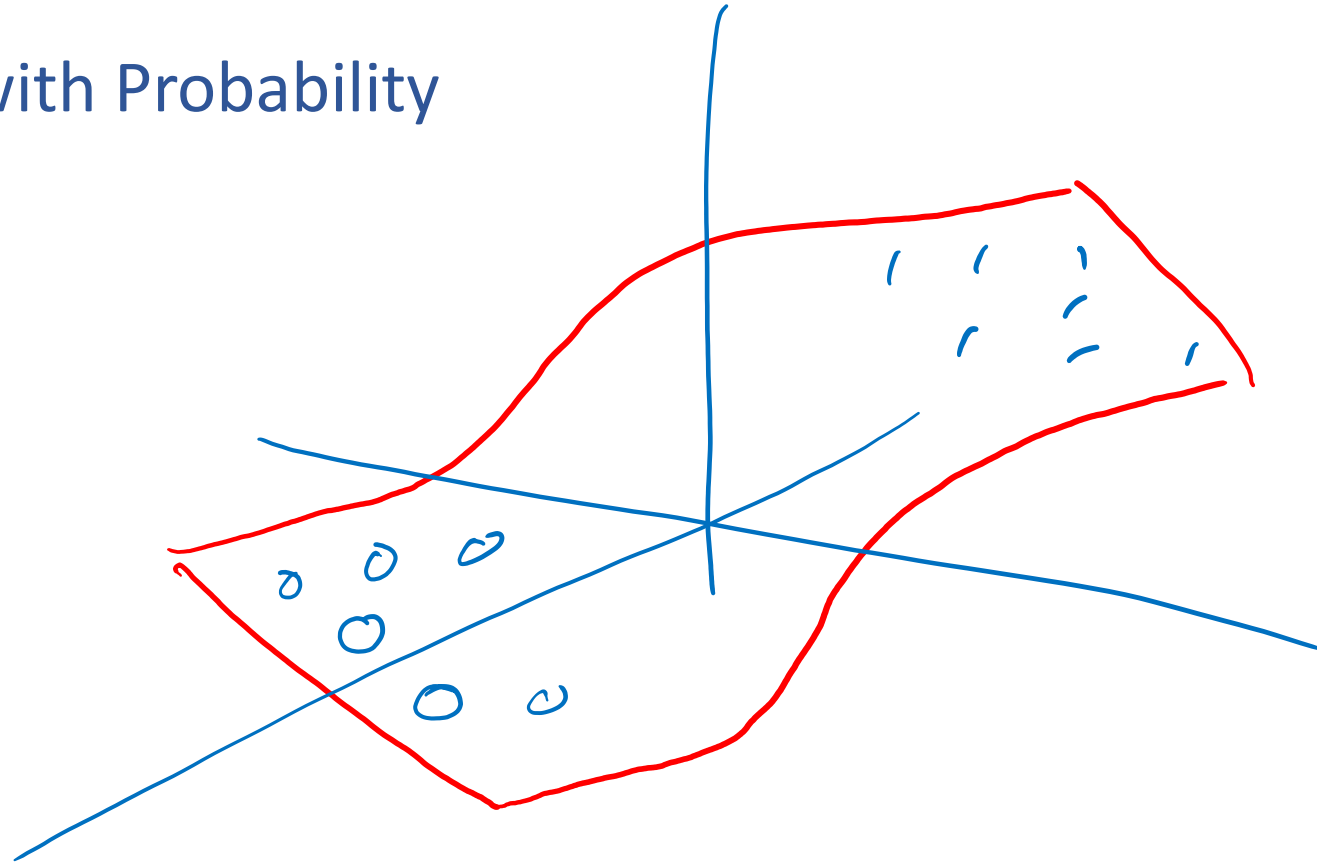
$$z = w^T x$$

$$y = \begin{cases} 1 & z \geq 0.5 \\ 0 & z < 0.5 \end{cases}$$

$$y = \text{step}(z, 0.5)$$

# Classification Models

## Linear Regression with Probability





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

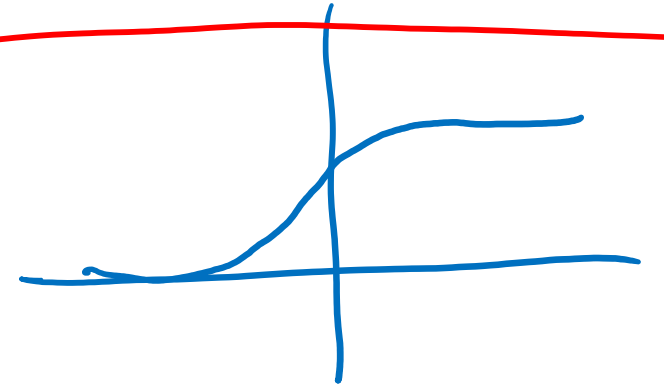
Bernoulli distribution of logistic function of linear model

$$z = x^T w$$

$$g = g(x^T w)$$

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

$$y \sim \text{Bern}(g) = \begin{cases} g & y=1 \\ 1-g & y=0 \end{cases}$$



# MLE for Bernoulli

Bernoulli distribution:

$$Y \sim \text{Bern}(z)$$

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

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

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

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

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

# MLE for Bernoulli

Bernoulli distribution:

$$Y \sim \text{Bern}(z)$$

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

$$y^{(3)2}$$

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

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

$$L(z) = z \cdot z \cdot (1-z)$$

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

$$\ell(z) = \log z + \log z + \log(1-z) = \sum_n \underbrace{\log z^{y^{(n)}}}_{\log} + \underbrace{\log(1-z)^{1-y^{(n)}}}_{\log}$$

# MLE for Bernoulli

Bernoulli distribution:

$$Y \sim \text{Bern}(z)$$

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

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

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

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

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

# M(C)LE for Logistic Regression

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

$$p(Y | X, \underline{\mathbf{w}}) = \prod_{n=1}^N p(y^{(n)} | \mathbf{x}^{(n)}, \mathbf{w})$$

$$p(D | \theta)$$
$$p(y^{(1)}, y^{(2)}, y^{(3)} \dots | \mathbf{x}, \mathbf{w})$$

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

$$Y \sim \text{Bern}(\mu), \quad \mu = \underline{g(\mathbf{w}^T \mathbf{x})}, \quad \underline{g(z) = \frac{1}{1+e^{-z}}}$$

What is the *conditional* log likelihood?

$$L(\mathbf{w}) = \prod_n p(Y=y^{(n)} | X=\mathbf{x}^{(n)}, \mathbf{w}) = \prod_n \mu^{(n) y^{(n)}} (1-\mu^{(n)})^{1-y^{(n)}}$$

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

# M(C)LE for Logistic Regression

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

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

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

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

What is the *conditional* log likelihood?

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

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

# M(C)LE for Logistic Regression

$$\underline{z} = \underline{f}(\underline{w}, \underline{x}) = \underline{w}^T \underline{x}$$

$$\nabla_{\underline{w}} f(\underline{w}, \underline{x}) = \underline{x}$$

$$\underline{\mu} = g(z) = \frac{1}{1+e^{-z}}$$

$$\frac{dg}{dz} = g(z)(1 - g(z)) = \mu(1 - \mu)$$

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

$$\begin{aligned} \frac{\partial \ell}{\partial \underline{w}} &= \sum_n \left( \frac{y^{(n)}}{\mu^{(n)}} - \frac{1 - y^{(n)}}{1 - \mu^{(n)}} \right) \frac{\partial g}{\partial z} \frac{df}{d\vec{w}} \\ &= \sum_n \left( \frac{y^{(n)} - \mu^{(n)}}{\mu^{(n)}(1 - \mu^{(n)})} \right) \mu^{(n)}(1 - \mu^{(n)}) \underline{x}^{(n)T} \end{aligned}$$

# M(C)LE for Logistic Regression

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

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

$$\ell(\mathbf{w}) = \sum_n (y^{(n)} \log \mu^{(n)} + (1 - y^{(n)}) \log(1 - \mu^{(n)})) \quad \leftarrow$$

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

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

$$= \sum_n (y^{(n)} - \mu^{(n)}) \mathbf{x}^{(n)T} \quad \leftarrow$$



# M(C)LE for Logistic Regression

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

$$\ell(\mathbf{w}) = \sum_n (y^{(n)} \log \mu^{(n)} + (1 - y^{(n)}) \log(1 - \mu^{(n)})) \quad \leftarrow$$

$$\nabla_{\mathbf{w}} \ell(\mathbf{w}) = \sum_n (y^{(n)} - \underline{\mu}^{(n)}) \mathbf{x}^{(n)}$$

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

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  $Y \sim \text{Bern}(p)$ ,  $\frac{p(Y=1)}{p(Y=0)} = \frac{p}{1-p}$

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

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

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

# Log Odds and Logistic Regression

Formulate log odds as linear model of  $X$ :

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

Equivalent to logistic representation:

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

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

Formulate log odds as linear model of  $X$ :

$$\begin{aligned}\log \frac{p(Y = 1 \mid X = \mathbf{x}, \mathbf{W})}{p(Y = K \mid X = \mathbf{x}, \mathbf{W})} &= \mathbf{w}_1^T \mathbf{x} \\ \log \frac{p(Y = 2 \mid X = \mathbf{x}, \mathbf{W})}{p(Y = K \mid X = \mathbf{x}, \mathbf{W})} &= \mathbf{w}_2^T \mathbf{x} \\ &\vdots \\ \log \frac{p(Y = K - 1 \mid X = \mathbf{x}, \mathbf{W})}{p(Y = K \mid X = \mathbf{x}, \mathbf{W})} &= \mathbf{w}_{K-1}^T \mathbf{x}\end{aligned}$$

$\mathbf{W}$

Equivalent to softmax representation:

$$p(Y = k \mid X = \mathbf{x}, \mathbf{W}) = \frac{e^{\mathbf{w}_k^T \mathbf{x}}}{1 + \sum_{j=1}^{K-1} e^{\mathbf{w}_j^T \mathbf{x}}}$$

$k \in \{1, \dots, K-1\}$

OR

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

$$p(Y = k \mid X = \mathbf{x}, \mathbf{W}) = \frac{e^{\mathbf{w}_k^T \mathbf{x}}}{\sum_{j=1}^K e^{\mathbf{w}_j^T \mathbf{x}}}$$

# Multi-class Logistic Regression

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

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

$$\underline{p(y^{(n)} = k | X = \mathbf{x}^{(n)}, \mathbf{W})} = \frac{e^{\mathbf{w}_k^T \mathbf{x}^{(n)}}}{\sum_{j=1}^K e^{\mathbf{w}_j^T \mathbf{x}^{(n)}}}$$

What is the *conditional* likelihood?

$$L(\mathbf{w}) = \prod_n \frac{e^{\mathbf{w}_k^T \mathbf{x}^{(n)}}}{\sum_{j=1}^K e^{\mathbf{w}_j^T \mathbf{x}^{(n)}}}$$

What is the hypothesis function?

$$\hat{y} = h_{\mathbf{w}}(\mathbf{x}) = \underset{k}{\operatorname{argmax}} \operatorname{softmax}(\mathbf{x}, \mathbf{w})$$