# Logistic Regression

These slides were assembled by Eric Eaton, with grateful acknowledgement of the many others who made their course materials freely available online. Feel free to reuse or adapt these slides for your own academic purposes, provided that you include proper attribution. Please send comments and corrections to Eric.
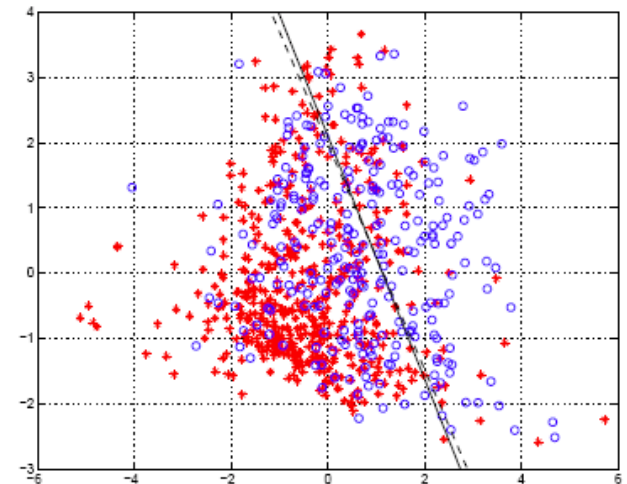
# Classification Based on Probability

- Instead of just predicting the class, give the probability of the instance being that class
  - i.e., learn $p(y \mid \boldsymbol{x})$

- Comparison to perceptron:
  - Perceptron doesn't produce probability estimate
  - Perceptron (and other discriminative classifiers) are only interested in producing a discriminative model

- Recall that:
$$0 \leq p(\text{event}) \leq 1$$
$$p(\text{event}) + p(\neg\text{event}) = 1$$

# Logistic Regression

- Takes a probabilistic approach to learning discriminative functions (i.e., a classifier)

- $h_{\boldsymbol{\theta}}(\boldsymbol{x})$ should give $p(y = 1 \mid \boldsymbol{x}; \boldsymbol{\theta})$

  – Want $0 \leq h_{\boldsymbol{\theta}}(\boldsymbol{x}) \leq 1$

Can't just use linear regression with a threshold

- Logistic regression model:

$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = g\left(\boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{x}\right)$$

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

$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \frac{1}{1 + e^{-\boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{x}}}$$

Logistic / Sigmoid Function

$g(z)$

# Interpretation of Hypothesis Output

$h_{\boldsymbol{\theta}}(\boldsymbol{x})$ = estimated $p(y = 1 \mid \boldsymbol{x}; \boldsymbol{\theta})$

Example: Cancer diagnosis from tumor size

$$\boldsymbol{x} = \left[ \begin{array}{c} x_0 \\ x_1 \end{array} \right] = \left[ \begin{array}{c} 1 \\ \mathrm{tumorSize} \end{array} \right]$$

$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = 0.7$

→ Tell patient that 70% chance of tumor being malignant

Note that: $p(y = 0 \mid \boldsymbol{x}; \boldsymbol{\theta}) + p(y = 1 \mid \boldsymbol{x}; \boldsymbol{\theta}) = 1$

Therefore, $p(y = 0 \mid \boldsymbol{x}; \boldsymbol{\theta}) = 1 - p(y = 1 \mid \boldsymbol{x}; \boldsymbol{\theta})$

# Another Interpretation

- Equivalently, logistic regression assumes that

$$\log \boxed{\frac{p(y = 1 \mid \boldsymbol{x}; \boldsymbol{\theta})}{p(y = 0 \mid \boldsymbol{x}; \boldsymbol{\theta})}} = \theta_0 + \theta_1 x_1 + \ldots + \theta_d x_d$$

odds of $y = 1$

> **Side Note**: the odds in favor of an event is the quantity $p\,/\,(1-p)$, where $p$ is the probability of the event
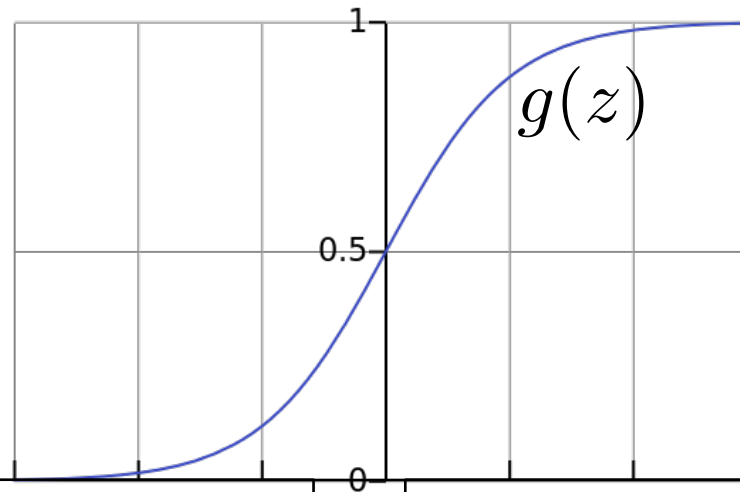>
> E.g., If I toss a fair dice, what are the odds that I will have a 6?

- In other words, logistic regression assumes that the log odds is a linear function of $x$

# Logistic Regression

$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = g\left(\boldsymbol{\theta}^{\mathsf{T}}\boldsymbol{x}\right)$$

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



$g(z)$

| $\boldsymbol{\theta}^{\mathsf{T}}\boldsymbol{x}$ should be large <u>negative</u> values for negative instances | $\boldsymbol{\theta}^{\mathsf{T}}\boldsymbol{x}$ should be large <u>positive</u> values for positive instances |

- Assume a threshold and...
  - Predict $y$ = 1 if $h_{\boldsymbol{\theta}}(\boldsymbol{x}) \geq 0.5$
  - Predict $y$ = 0 if $h_{\boldsymbol{\theta}}(\boldsymbol{x}) < 0.5$



$y = 1$

$\theta$

$y = 0$

# Non-Linear Decision Boundary

- Can apply basis function expansion to features, same as with linear regression

$$\boldsymbol{x} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix} \rightarrow \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_1 x_2 \\ x_1^2 \\ x_2^2 \\ x_1^2 x_2 \\ x_1 x_2^2 \\ \vdots \end{bmatrix}$$

# Logistic Regression

- Given $\left\{ \left( \boldsymbol{x}^{(1)}, y^{(1)} \right), \left( \boldsymbol{x}^{(2)}, y^{(2)} \right), \ldots, \left( \boldsymbol{x}^{(n)}, y^{(n)} \right) \right\}$

  where $\boldsymbol{x}^{(i)} \in \mathbb{R}^d, \; y^{(i)} \in \{0, 1\}$

- Model: $h_{\boldsymbol{\theta}}(\boldsymbol{x}) = g\left( \boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{x} \right)$

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

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} \qquad \boldsymbol{x}^{\mathsf{T}} = \begin{bmatrix} 1 & x_1 & \ldots & x_d \end{bmatrix}$$

# Logistic Regression Objective Function

- Can't just use squared loss as in linear regression:

$$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^{n} \left( h_{\boldsymbol{\theta}} \left( \boldsymbol{x}^{(i)} \right) - y^{(i)} \right)^2$$

  – Using the logistic regression model

$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \frac{1}{1 + e^{-\boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{x}}}$$

  results in a non-convex optimization

# Deriving the Cost Function via Maximum Likelihood Estimation

- Likelihood of data is given by: $l(\boldsymbol{\theta}) = \prod_{i=1}^{n} p(y^{(i)} \mid \boldsymbol{x}^{(i)}; \boldsymbol{\theta})$

- So, looking for the $\boldsymbol{\theta}$ that maximizes the likelihood

$$\boldsymbol{\theta}_{\mathrm{MLE}} = \arg\max_{\boldsymbol{\theta}} l(\boldsymbol{\theta}) = \arg\max_{\boldsymbol{\theta}} \prod_{i=1}^{n} p(y^{(i)} \mid \boldsymbol{x}^{(i)}; \boldsymbol{\theta})$$

- Can take the log without changing the solution:

$$\boldsymbol{\theta}_{\mathrm{MLE}} = \arg\max_{\boldsymbol{\theta}} \log \prod_{i=1}^{n} p(y^{(i)} \mid \boldsymbol{x}^{(i)}; \boldsymbol{\theta})$$

$$= \arg\max_{\boldsymbol{\theta}} \sum_{i=1}^{n} \log p(y^{(i)} \mid \boldsymbol{x}^{(i)}; \boldsymbol{\theta})$$

# Deriving the Cost Function via Maximum Likelihood Estimation

- Expand as follows:

$$\boldsymbol{\theta}_{\mathrm{MLE}} = \arg\max_{\boldsymbol{\theta}} \sum_{i=1}^{n} \log p(y^{(i)} \mid \boldsymbol{x}^{(i)}; \boldsymbol{\theta})$$

$$= \arg\max_{\boldsymbol{\theta}} \sum_{i=1}^{n} \left[ y^{(i)} \log p(y^{(i)} = 1 \mid \boldsymbol{x}^{(i)}; \boldsymbol{\theta}) + \left(1 - y^{(i)}\right) \log \left(1 - p(y^{(i)} = 1 \mid \boldsymbol{x}^{(i)}; \boldsymbol{\theta})\right) \right]$$

- Substitute in model, and take negative to yield

**Logistic regression objective**:

$$\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$$

$$J(\boldsymbol{\theta}) = - \sum_{i=1}^{n} \left[ y^{(i)} \log h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)}) + \left(1 - y^{(i)}\right) \log \left(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)})\right) \right]$$

# Intuition Behind the Objective

$$J(\boldsymbol{\theta}) = -\sum_{i=1}^{n} \left[ y^{(i)} \log h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)}) + \left( 1 - y^{(i)} \right) \log \left( 1 - h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)}) \right) \right]$$

- Cost of a single instance:

$$\text{cost}\left( h_{\boldsymbol{\theta}}(\boldsymbol{x}), y \right) = \begin{cases} -\log(h_{\boldsymbol{\theta}}(\boldsymbol{x})) & \text{if } y = 1 \\ -\log(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x})) & \text{if } y = 0 \end{cases}$$
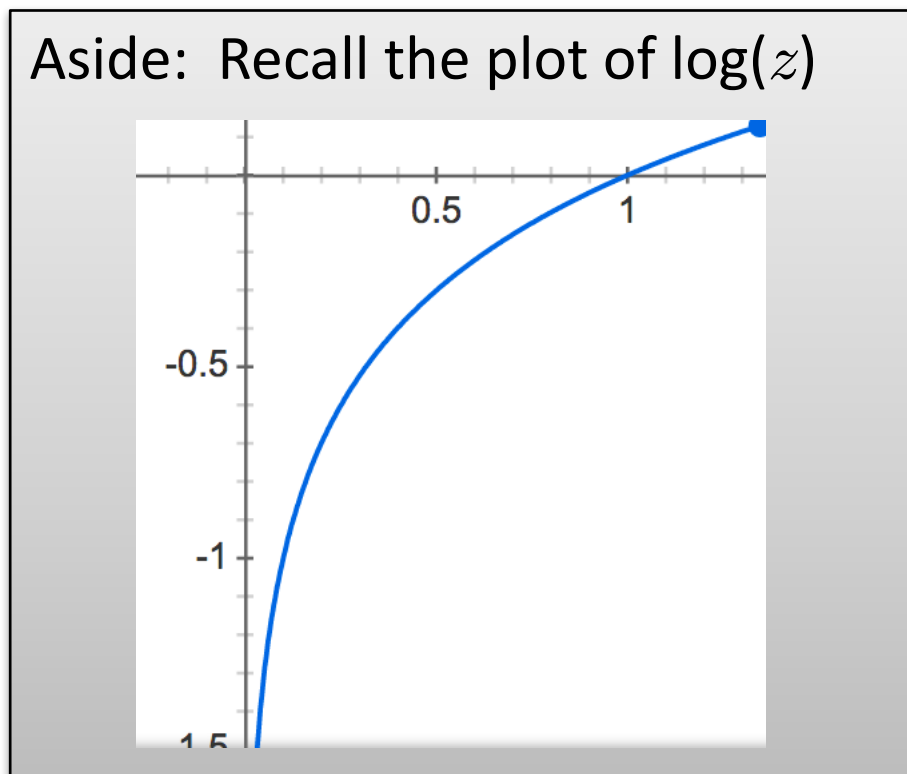
- Can re-write objective function as

$$J(\boldsymbol{\theta}) = \sum_{i=1}^{n} \text{cost}\left( h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)}), y^{(i)} \right)$$

Compare to linear regression: $J(\boldsymbol{\theta}) = \dfrac{1}{2n} \sum_{i=1}^{n} \left( h_{\boldsymbol{\theta}}\left( \boldsymbol{x}^{(i)} \right) - y^{(i)} \right)^2$

# Intuition Behind the Objective

$$\text{cost}\left(h_{\boldsymbol{\theta}}(\boldsymbol{x}), y\right) = \begin{cases} -\log(h_{\boldsymbol{\theta}}(\boldsymbol{x})) & \text{if } y = 1 \\ -\log(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x})) & \text{if } y = 0 \end{cases}$$

Aside: Recall the plot of $\log(z)$

# Intuition Behind the Objective

$$\text{cost}\left(h_{\boldsymbol{\theta}}(\boldsymbol{x}), y\right) = \begin{cases} \boxed{-\log(h_{\boldsymbol{\theta}}(\boldsymbol{x})) \quad \text{if } y = 1} \\ -\log(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x})) \quad \text{if } y = 0 \end{cases}$$

If $y$ = 1

- Cost = 0 if prediction is correct
- As $h_{\boldsymbol{\theta}}(\boldsymbol{x}) \to 0, \text{cost} \to \infty$

- Captures intuition that larger mistakes should get larger penalties
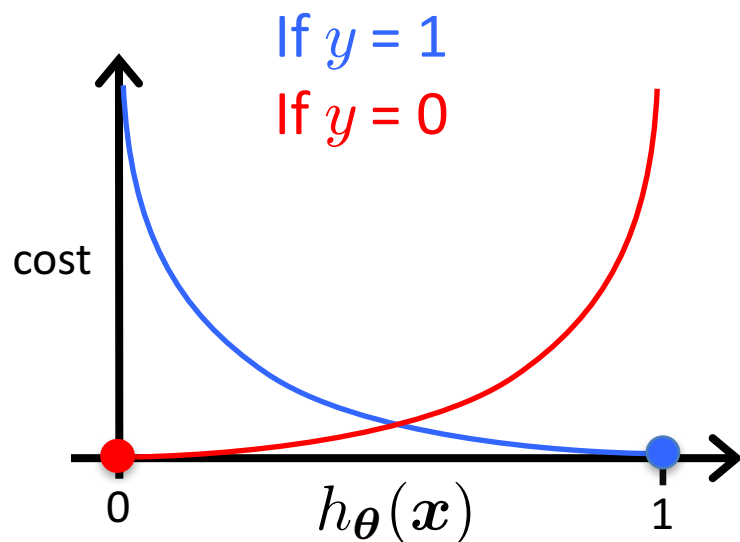  - e.g., predict $h_{\boldsymbol{\theta}}(\boldsymbol{x}) = 0$, but $y$ = 1

If $y$ = 1



cost

0     $h_{\boldsymbol{\theta}}(\boldsymbol{x})$     1

# Intuition Behind the Objective

$$\text{cost}(h_{\boldsymbol{\theta}}(\boldsymbol{x}), y) = \begin{cases} -\log(h_{\boldsymbol{\theta}}(\boldsymbol{x})) & \text{if } y = 1 \\ -\log(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x})) & \text{if } y = 0 \end{cases}$$

If $y = 0$

- Cost = 0 if prediction is correct
- As $(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x})) \to 0, \text{cost} \to \infty$

- Captures intuition that larger mistakes should get larger penalties

If $y = 1$
If $y = 0$

cost

0 $\qquad h_{\boldsymbol{\theta}}(\boldsymbol{x}) \qquad$ 1

# Regularized Logistic Regression

$$J(\boldsymbol{\theta}) = -\sum_{i=1}^{n} \left[ y^{(i)} \log h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)}) + \left(1 - y^{(i)}\right) \log \left(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)})\right) \right]$$

- We can regularize logistic regression exactly as before:

$$J_{\text{regularized}}(\boldsymbol{\theta}) = J(\boldsymbol{\theta}) + \frac{\lambda}{2} \sum_{j=1}^{d} \theta_j^2$$

$$= J(\boldsymbol{\theta}) + \frac{\lambda}{2} \|\boldsymbol{\theta}_{[1:d]}\|_2^2$$

# Gradient Descent for Logistic Regression

$$J_{\mathrm{reg}}(\boldsymbol{\theta}) = -\sum_{i=1}^{n}\left[y^{(i)}\log h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)}) + \left(1 - y^{(i)}\right)\log\left(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)})\right)\right] + \frac{\lambda}{2}\|\boldsymbol{\theta}_{[1:d]}\|_2^2$$

Want $\min\limits_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$

- Initialize $\boldsymbol{\theta}$
- Repeat until convergence

$$\theta_j \leftarrow \theta_j - \alpha\frac{\partial}{\partial\theta_j}J(\boldsymbol{\theta})$$

simultaneous update
for $j$ = 0 … d

Use the natural logarithm (ln = $\log_e$) to cancel with the exp() in $h_{\boldsymbol{\theta}}(\boldsymbol{x})$

# Gradient Descent for Logistic Regression

$$J_{\text{reg}}(\boldsymbol{\theta}) = -\sum_{i=1}^{n}\left[y^{(i)}\log h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)}) + \left(1-y^{(i)}\right)\log\left(1-h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)})\right)\right] + \frac{\lambda}{2}\|\boldsymbol{\theta}_{[1:d]}\|_2^2$$

Want $\min\limits_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$

- Initialize $\boldsymbol{\theta}$

- Repeat until convergence    (simultaneous update for $j = 0 \dots d$)

$$\theta_0 \leftarrow \theta_0 - \alpha \sum_{i=1}^{n}\left(h_{\boldsymbol{\theta}}\left(\boldsymbol{x}^{(i)}\right) - y^{(i)}\right)$$

$$\theta_j \leftarrow \theta_j - \alpha\left[\sum_{i=1}^{n}\left(h_{\boldsymbol{\theta}}\left(\boldsymbol{x}^{(i)}\right) - y^{(i)}\right)x_j^{(i)} + \lambda\theta_j\right]$$

# Gradient Descent for Logistic Regression

- Initialize $\boldsymbol{\theta}$

- Repeat until convergence   (simultaneous update for $j = 0 \dots d$)

$$\theta_0 \leftarrow \theta_0 - \alpha \sum_{i=1}^{n} \left( h_{\boldsymbol{\theta}} \left( \boldsymbol{x}^{(i)} \right) - y^{(i)} \right)$$

$$\theta_j \leftarrow \theta_j - \alpha \left[ \sum_{i=1}^{n} \left( h_{\boldsymbol{\theta}} \left( \boldsymbol{x}^{(i)} \right) - y^{(i)} \right) x_j^{(i)} + \lambda \theta_j \right]$$

This looks IDENTICAL to linear regression!!!
- Ignoring the $1/n$ constant
- However, the form of the model is very different:

$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \frac{1}{1 + e^{-\boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{x}}}$$

# Stochastic Gradient Descent

# Consider Learning with Numerous Data

- Logistic regression objective:

$$J(\boldsymbol{\theta}) = -\frac{1}{n} \sum_{i=1}^{n} \underbrace{[y_i \log h_{\boldsymbol{\theta}}(\mathbf{x}_i) + (1 - y_i) \log (1 - h_{\boldsymbol{\theta}}(\mathbf{x}_i))]}_{\text{cost}_{\boldsymbol{\theta}}(\mathbf{x}_i, y_i)}$$

- Fit via gradient descent:

$$\theta_j \leftarrow \theta_j - \alpha \frac{1}{n} \sum_{i=1}^{n} (h_{\boldsymbol{\theta}}(\mathbf{x}_i) - y_i) x_{ij}$$

- What is the computational complexity in terms of $n$?

# Gradient Descent

## Batch Gradient Descent

Initialize θ

Repeat {

$$\theta_j \leftarrow \theta_j - \alpha \frac{1}{n} \underbrace{\sum_{i=1}^{n} \left( h_{\boldsymbol{\theta}} \left( \mathbf{x}_i \right) - y_i \right) x_{ij}}_{\frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta})} \qquad \text{for } j = 0...d$$

}

## Stochastic Gradient Descent

Initialize θ

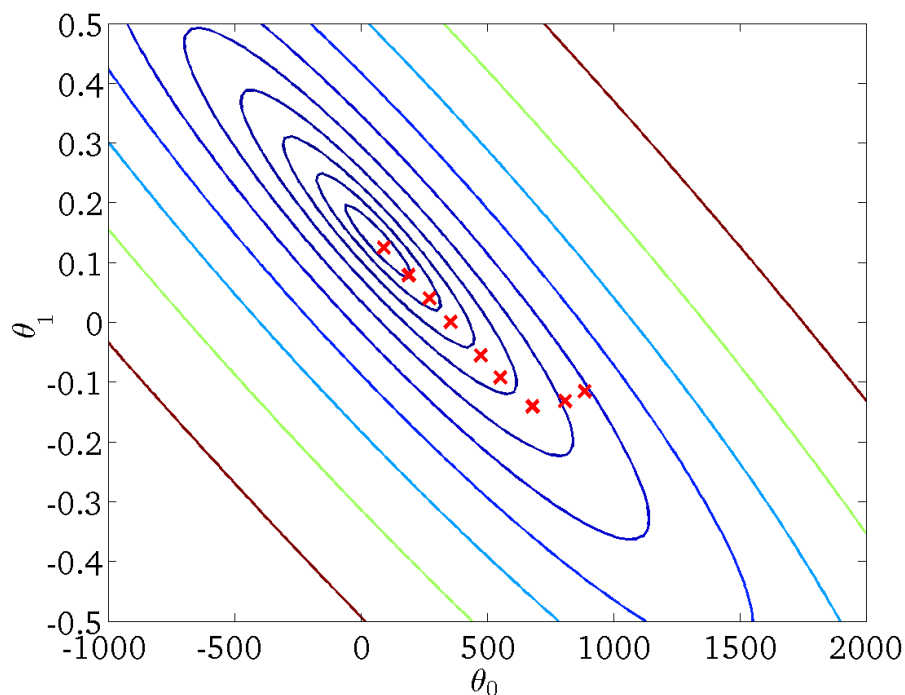Randomly shuffle dataset

Repeat {   (Typically $1 - 10x$)

For $i = 1...n$, *do*

$$\theta_j \leftarrow \theta_j - \alpha \underbrace{\left( h_{\boldsymbol{\theta}} \left( \mathbf{x}_i \right) - y_i \right) x_{ij}}_{\frac{\partial}{\partial \theta_j} \text{cost}_{\boldsymbol{\theta}} \left( \mathbf{x}_i, y_i \right)} \qquad \text{for } j = 0...d$$
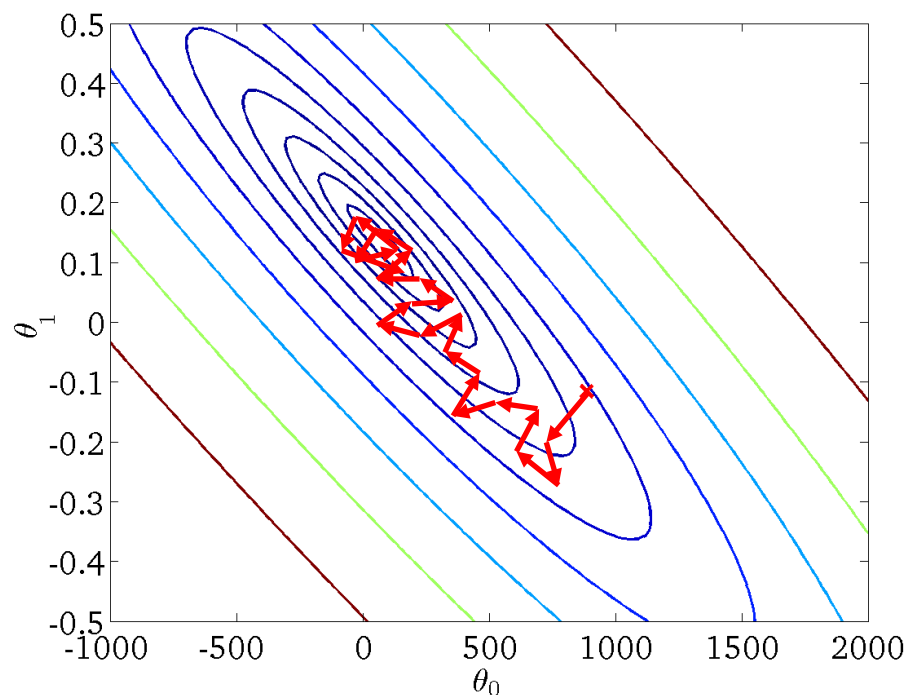
}

# Batch vs Stochastic GD

### Batch GD



### Stochastic GD



- Learning rate α is typically held constant
- Can slowly decrease α over time to force θ to converge:

$$\text{e.g.,} \quad \alpha_t = \frac{\text{constant1}}{\text{iterationNumber} + \text{constant2}}$$

# Adagrad

# New Stochastic Gradient Algorithms

- So far, we have considered:
  - a constant learning rate α
  - a time-dependent learning rate $α_t$ via a pre-set formula

- AdaGrad adjusts the learning rate based on historical information
  - Frequently occurring features in the gradients get small learning rates and infrequent features get higher ones
  - Key idea: "learn slowly" from frequent features but "pay attention" to rare but informative features

- Define a per-feature learning rate for feature $j$ as:

$$\alpha_{t,j} = \frac{\alpha}{\sqrt{G_{t,j}}} \quad \text{where} \quad G_{t,j} = \sum_{k=1}^{t} g_{k,j}^2 \underbrace{\qquad}\; \frac{\partial}{\partial \theta_j} \text{cost}_{\boldsymbol{\theta}}(\mathbf{x}_k, y_k)$$

  - $G_{t,j}$ is the sum of squares of gradients of feature $j$ through time $t$

# New Stochastic Gradient Algorithms

Adagrad per-feature learning rate

$$\alpha_{t,j} = \frac{\alpha}{\sqrt{G_{t,j}}} \qquad \text{where} \qquad G_{t,j} = \sum_{k=1}^{t} g_{k,j}^2$$

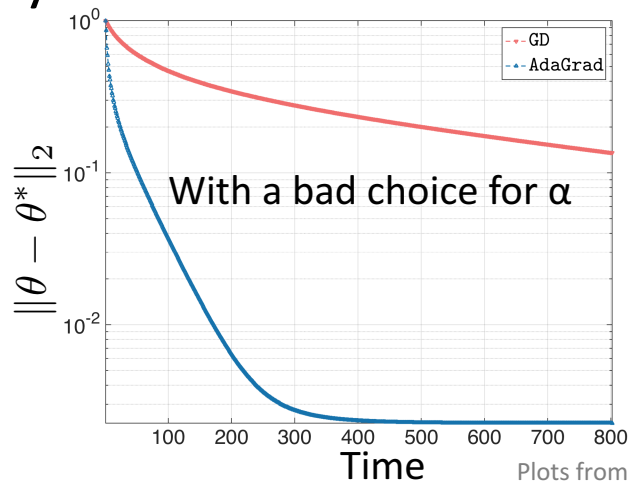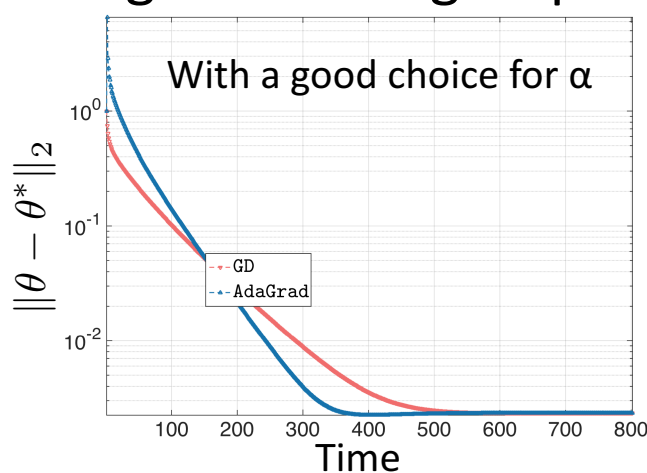- Adagrad changes the update rule for SGD at time $t$ from

$$\theta_j \leftarrow \theta_j - \alpha g_{t,j}$$

to

$$\theta_j \leftarrow \theta_j - \frac{\alpha}{\sqrt{G_{t,j}} + \zeta} g_{t,j}$$

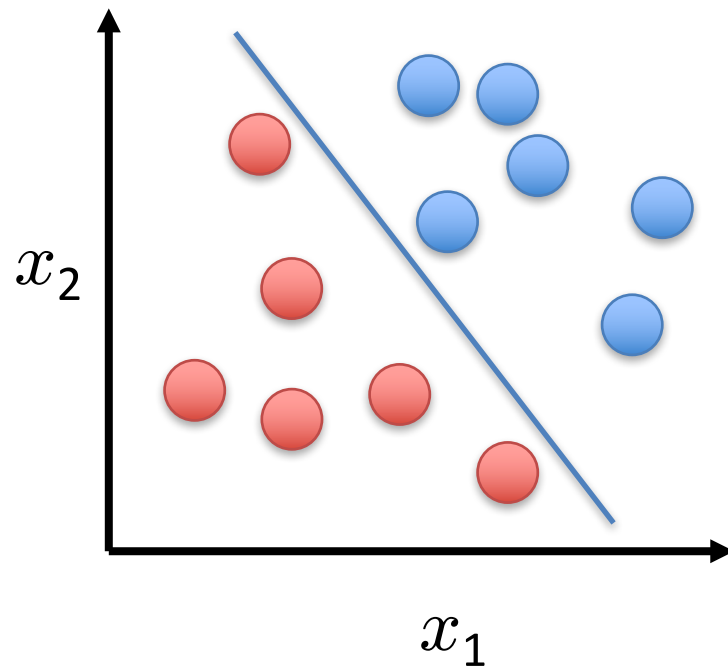In practice, we add a small constant $\zeta > 0$ to prevent dividing by zero errors
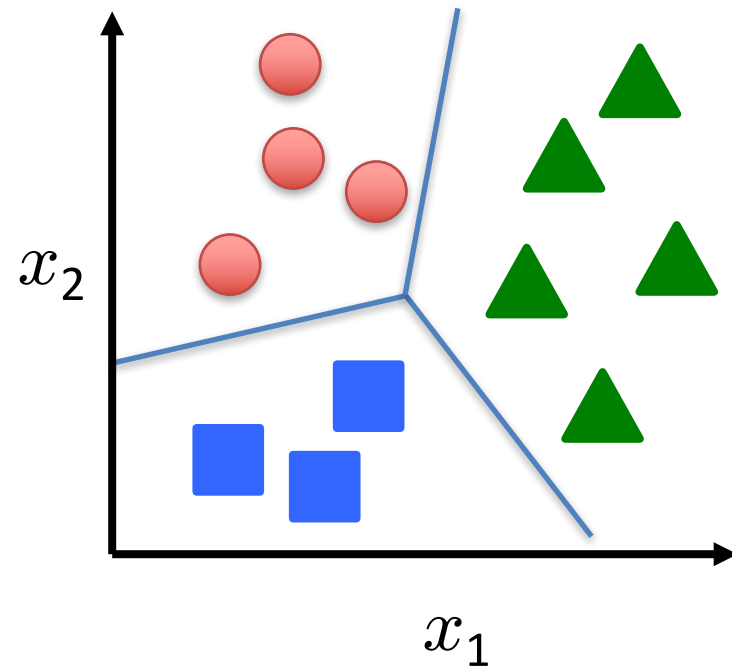
- Adagrad converges quickly:



With a good choice for α



With a bad choice for α

Plots from http://akyrillidis.github.io/notes/AdaGrad

# Multi-Class Classification

# Multi-Class Classification

Binary classification:
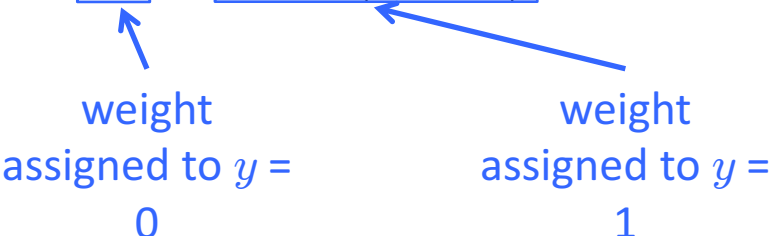
Multi-class classification:



Disease diagnosis:   healthy / cold / flu / pneumonia

Object classification:  desk / chair / monitor / bookcase

# Multi-Class Logistic Regression

- For 2 classes:

$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \frac{1}{1 + \exp(-\boldsymbol{\theta}^\mathsf{T}\boldsymbol{x})} = \frac{\exp(\boldsymbol{\theta}^\mathsf{T}\boldsymbol{x})}{\boxed{1} + \boxed{\exp(\boldsymbol{\theta}^\mathsf{T}\boldsymbol{x})}}$$

weight assigned to $y = 0$

weight assigned to $y = 1$

- For $C$ classes $\{1, ..., C\}$:

$$p(y = c \mid \boldsymbol{x}; \boldsymbol{\theta}_1, \ldots, \boldsymbol{\theta}_C) = \frac{\exp(\boldsymbol{\theta}_c^\mathsf{T}\boldsymbol{x})}{\sum_{c=1}^{C} \exp(\boldsymbol{\theta}_c^\mathsf{T}\boldsymbol{x})}$$

  – Called the **softmax** function

# Multi-Class Logistic Regression

Split into One vs Rest:



- Train a logistic regression classifier for each class $i$ to predict the probability that $y = i$ with

$$h_c(\boldsymbol{x}) = \frac{\exp(\boldsymbol{\theta}_c^\mathsf{T} \boldsymbol{x})}{\sum_{c=1}^{C} \exp(\boldsymbol{\theta}_c^\mathsf{T} \boldsymbol{x})}$$

# Implementing Multi-Class Logistic Regression

- Use  $h_c(\boldsymbol{x}) = \dfrac{\exp(\boldsymbol{\theta}_c^\mathsf{T} \boldsymbol{x})}{\sum_{c=1}^{C} \exp(\boldsymbol{\theta}_c^\mathsf{T} \boldsymbol{x})}$  as the model for class $c$

- Gradient descent simultaneously updates all parameters for all models
  - Same derivative as before, just with the above $h_c(\boldsymbol{x})$

- Predict class label as the most probable label

$$\max_c h_c(\boldsymbol{x})$$