# 10-701 Introduction to Machine Learning

# Perceptron &
# Neural Networks

Matt Gormley
Lecture 12
October 17, 2016

# Reminders

- Homework 3:
  - due 10/24/16

# Outline

- Discriminative vs. Generative

- Perceptron

- Neural Networks

- Backpropagation

# DISCRIMINATIVE AND GENERATIVE CLASSIFIERS

# Generative vs. Discriminative

- **Generative Classifiers:**
  - Example: Naïve Bayes
  - Define a joint model of the observations **x** and the labels y: $p(\boldsymbol{x}, y)$
  - Learning maximizes (joint) likelihood
  - Use Bayes' Rule to classify based on the posterior:
    $$p(y|\mathbf{x}) = p(\mathbf{x}|y)p(y)/p(\mathbf{x})$$
- **Discriminative Classifiers:**
  - Example: Logistic Regression
  - Directly model the conditional: $p(y|\mathbf{x})$
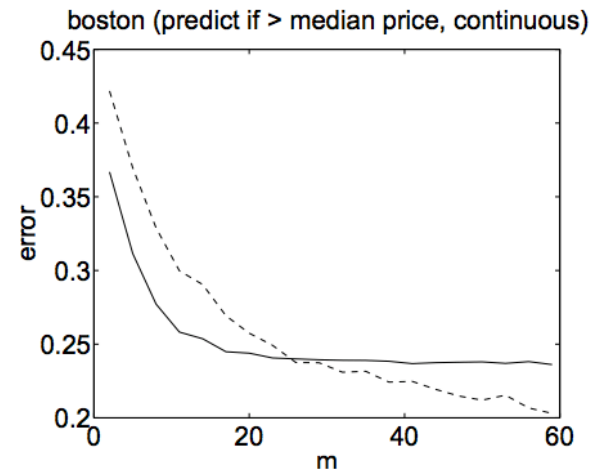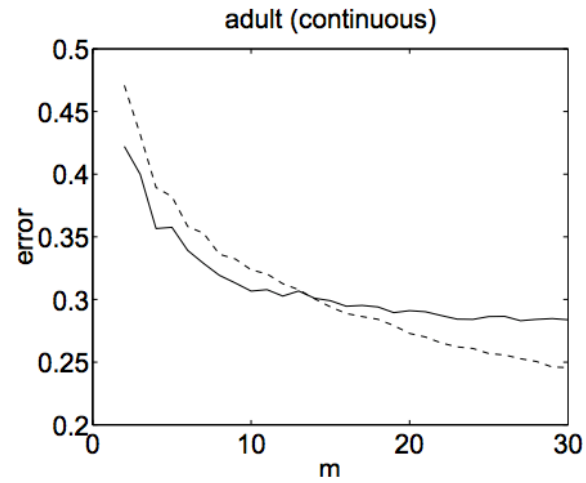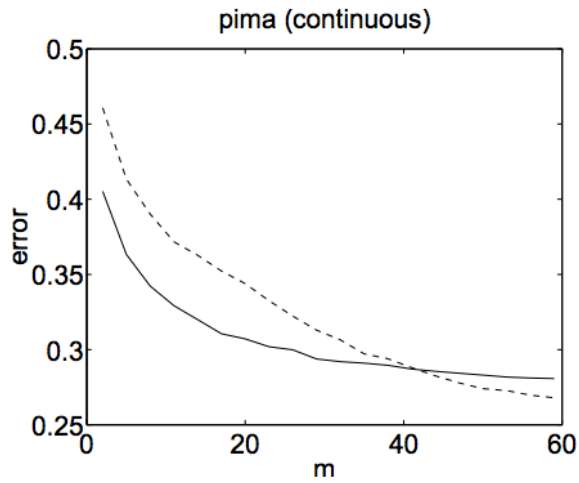  - Learning maximizes conditional likelihood

# Generative vs. Discriminative

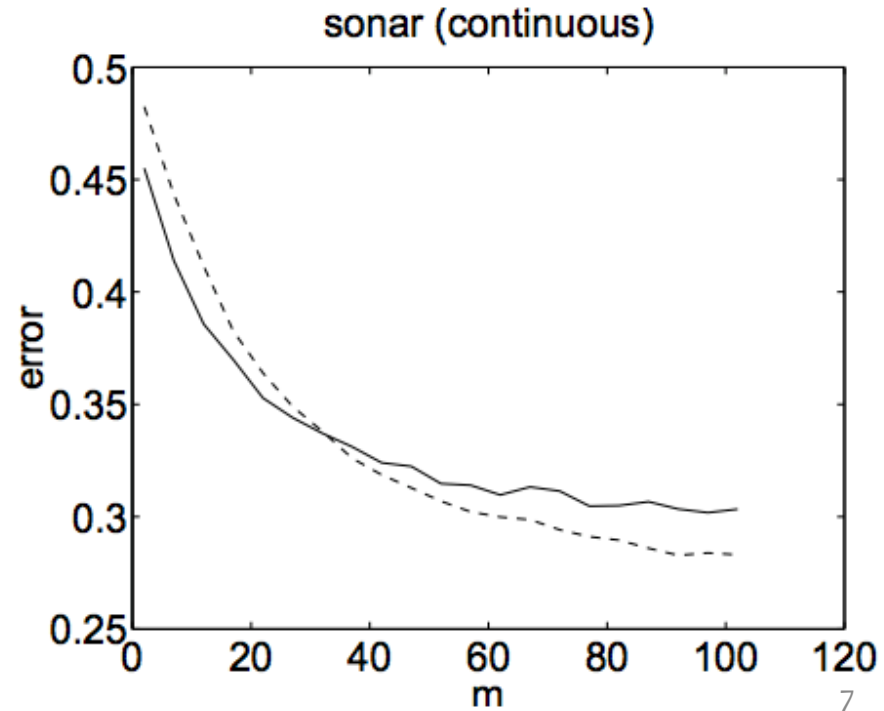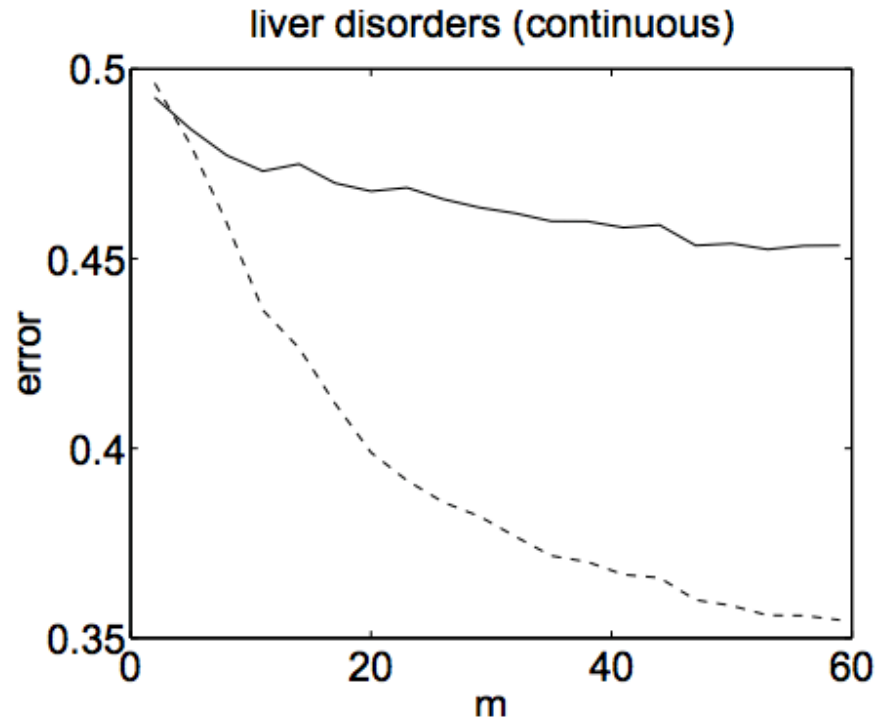**Finite Sample Analysis** (Ng & Jordan, 2002)

[Assume that we are learning from a finite training dataset]

**If model assumptions are correct:** Naive Bayes is a more efficient learner (requires fewer samples) than Logistic Regression
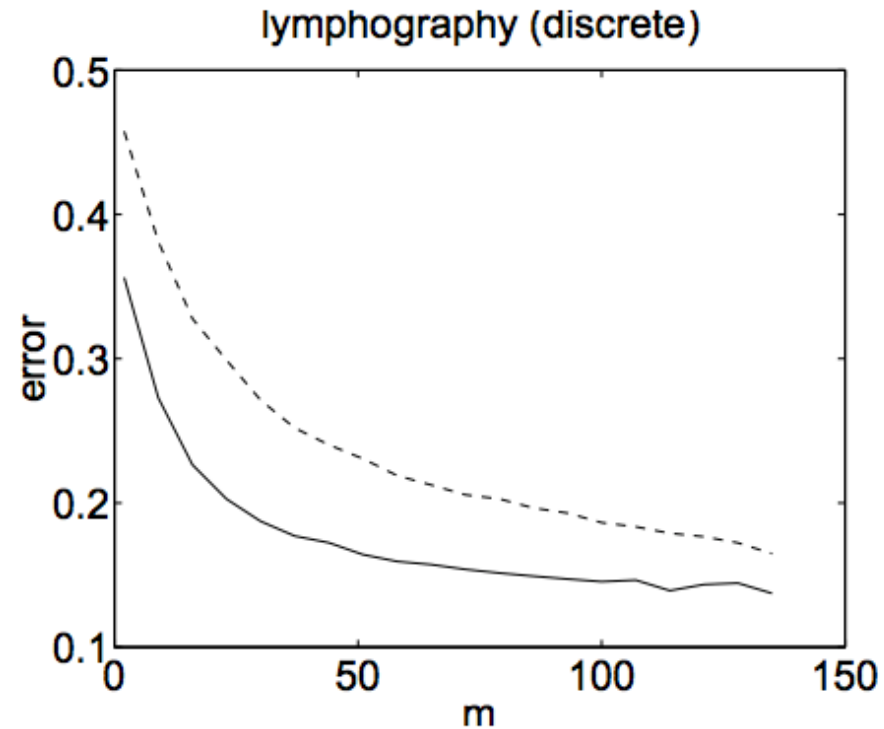
**If model assumptions are incorrect:** Logistic Regression has lower asymtotic error, and does better than Naïve Bayes

solid: NB dashed: LR

Slide courtesy of William Cohen

7

promoters (discrete)

lymphography (discrete)

Naïve Bayes makes stronger assumptions about the data but needs fewer examples to estimate the parameters

"On Discriminative vs Generative Classifiers: ….." Andrew Ng and Michael Jordan, NIPS 2001.

8

Slide courtesy of William Cohen

# Generative vs. Discriminative

## Learning (Parameter Estimation)

**Naïve Bayes:**
Parameters are decoupled → Closed form solution for MLE

**Logistic Regression:**
Parameters are coupled → No closed form solution – must use iterative optimization techniques instead

# Naïve Bayes vs. Logistic Reg.

## Learning (MAP Estimation of Parameters)

**Bernoulli Naïve Bayes:**
Parameters are probabilities → Beta prior (usually) pushes probabilities away from zero / one extremes

**Logistic Regression:**
Parameters are not probabilities → Gaussian prior encourages parameters to be close to zero

(effectively pushes the probabilities away from zero / one extremes)

# Naïve Bayes vs. Logistic Reg.

**Features**

**Naïve Bayes:**
Features $x$ are assumed to be conditionally independent given $y$. (i.e. Naïve Bayes Assumption)
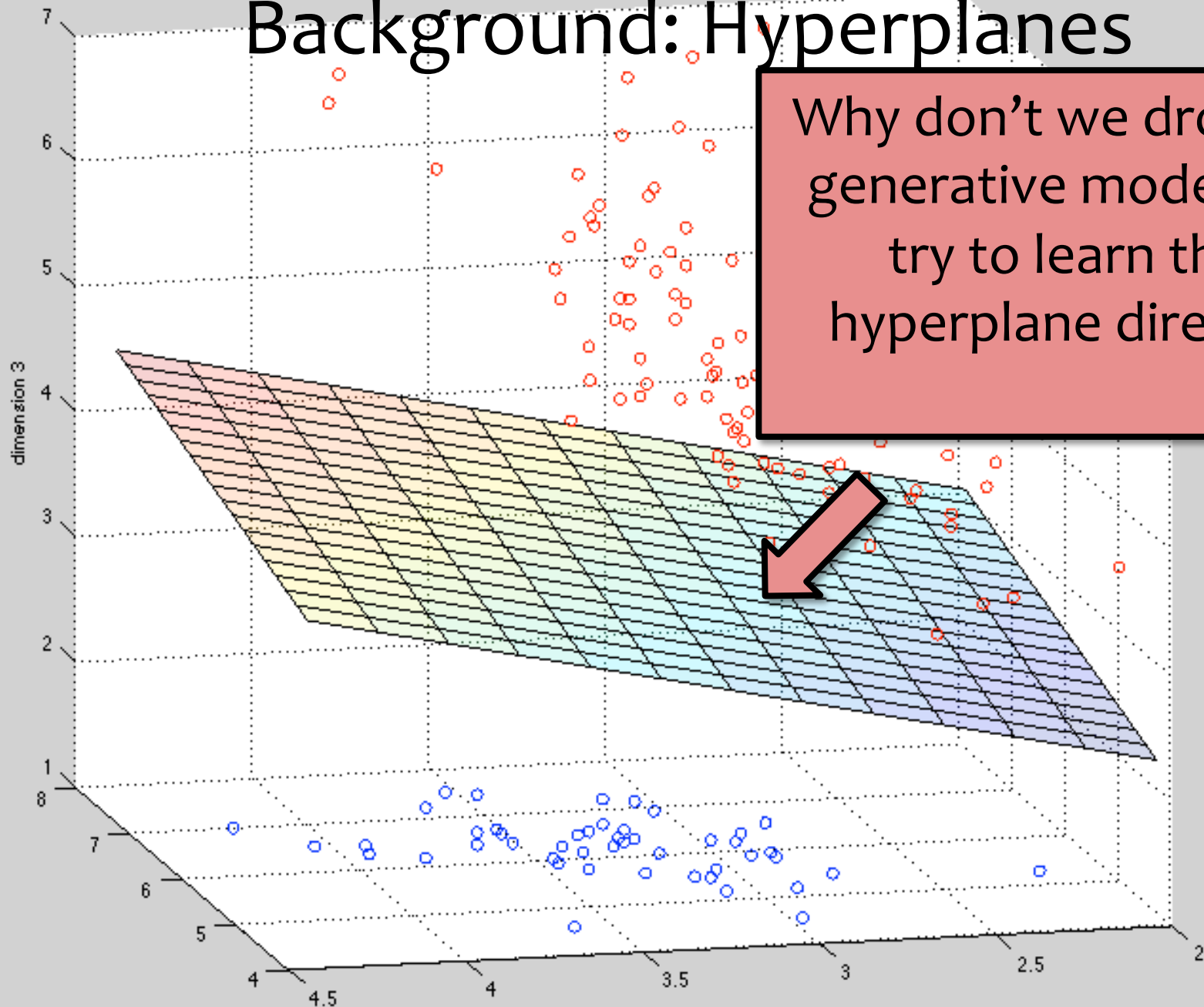
**Logistic Regression:**
No assumptions are made about the form of the features $x$. They can be dependent and correlated in any fashion.
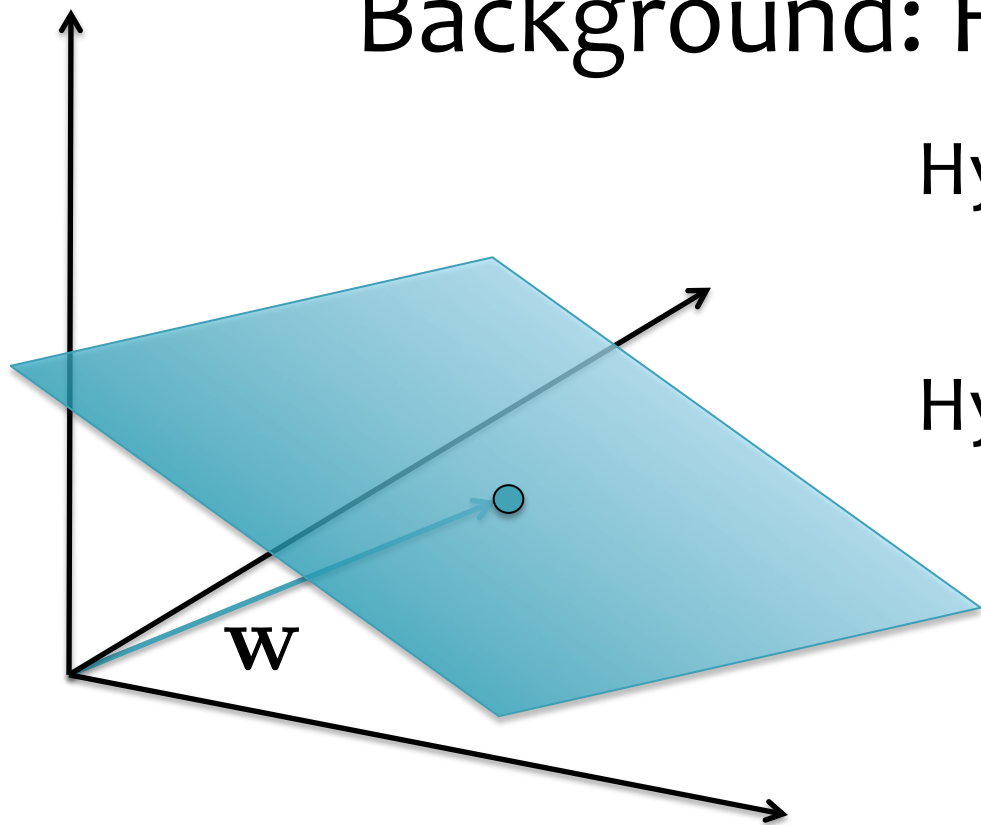
# THE PERCEPTRON ALGORITHM

# Background: Hyperplanes



Why don't we drop the generative model and try to learn this hyperplane directly?

# Background: Hyperplanes

Hyperplane (Definition 1):
$$\mathcal{H} = \{\mathbf{x} : \mathbf{w}^T\mathbf{x} = b\}$$

Hyperplane (Definition 2):
$$\mathcal{H} = \{\mathbf{x} : \mathbf{w}^T\mathbf{x} = 0$$
$$\text{and } x_1 = 1\}$$

$\mathbf{w}$

Half-spaces:
$$\mathcal{H}^+ = \{\mathbf{x} : \mathbf{w}^T\mathbf{x} > 0 \text{ and } x_1 = 1\}$$
$$\mathcal{H}^- = \{\mathbf{x} : \mathbf{w}^T\mathbf{x} < 0 \text{ and } x_1 = 1\}$$

d: Hyperplanes

Directly modeling the hyperplane would use a decision function:
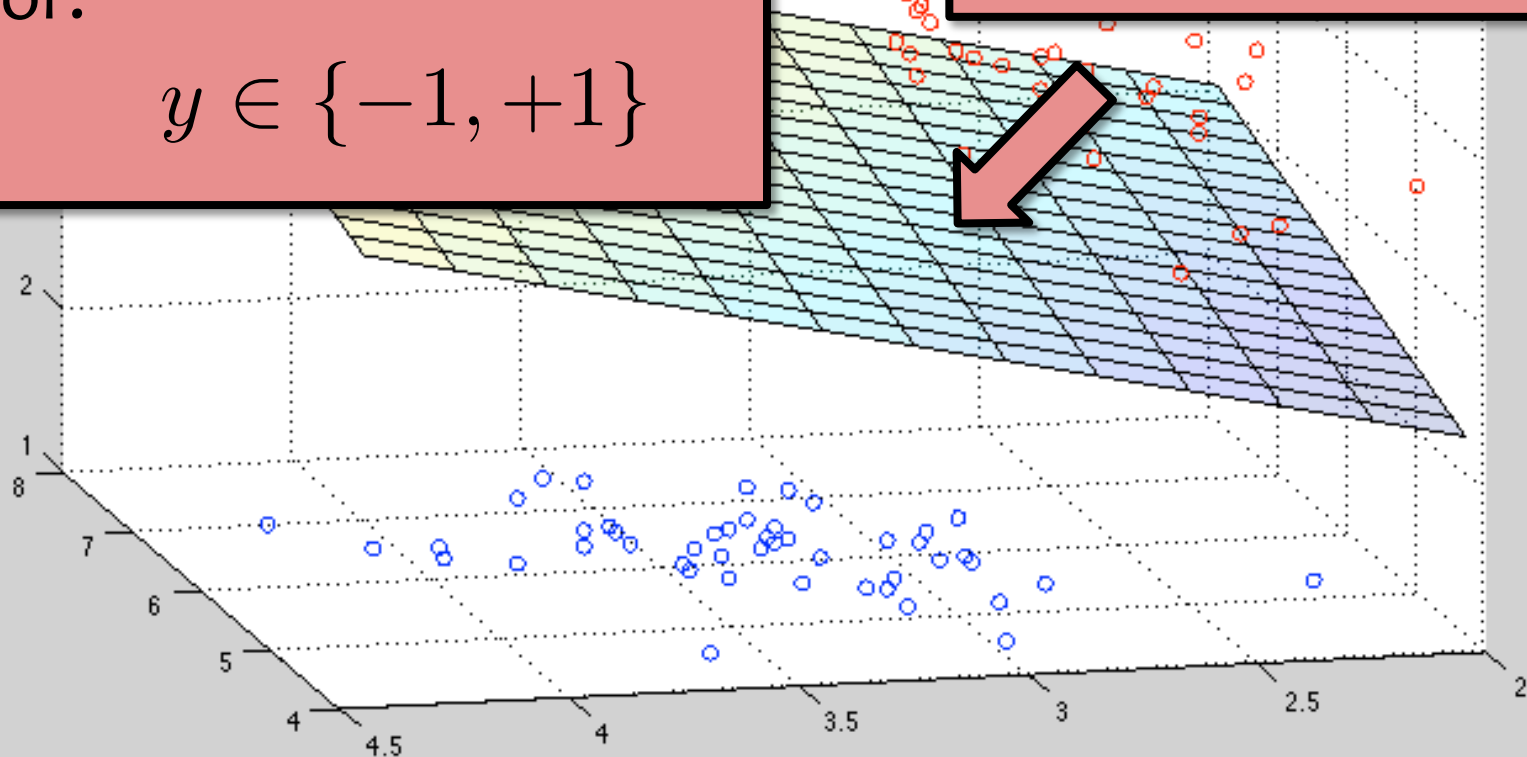
$$h(\mathbf{x}) = \text{sign}(\boldsymbol{\theta}^T \mathbf{x})$$

for:

$$y \in \{-1, +1\}$$

Why don't we drop the generative model and try to learn this hyperplane directly?

# Online Learning Model

**Setup:**

- We receive an example $(\mathbf{x}, y)$
- Make a prediction $h(\mathbf{x})$
- Check for correctness $h(\mathbf{x}) = y$?

**Goal:**

- **Minimize** the number of **mistakes**

# Margins

**Definition**: The margin of example $x$ w.r.t. a linear sep. $w$ is the distance from $x$ to the plane $w \cdot x = 0$ (or the negative if on wrong side)

**Definition**: The margin $\gamma_w$ of a set of examples $S$ wrt a linear separator $w$ is the smallest margin over points $x \in S$.

**Definition**: The margin $\gamma$ of a set of examples $S$ is the maximum $\gamma_w$ over all linear separators $w$.

# Perceptron Algorithm

**Data:** Inputs are continuous vectors of length K. Outputs are discrete.

$$\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^{N} \text{ where } \mathbf{x} \in \mathbb{R}^{K} \text{ and } y \in \{+1, -1\}$$

**Prediction:** Output determined by hyperplane.

$$\hat{y} = h_{\boldsymbol{\theta}}(\mathbf{x}) = \text{sign}(\boldsymbol{\theta}^{T} \mathbf{x})$$

$$\text{sign}(a) = \begin{cases} 1, & \text{if } a \geq 0 \\ -1, & \text{otherwise} \end{cases}$$

**Learning:** Iterative procedure:
- **while** not converged
    - **receive** next example (**x**, y)
    - **predict** y' = h(**x**)
    - **if** positive mistake: **add x** to parameters
    - **if** negative mistake: **subtract x** from parameters

# Perceptron Algorithm

**Learning:**

---

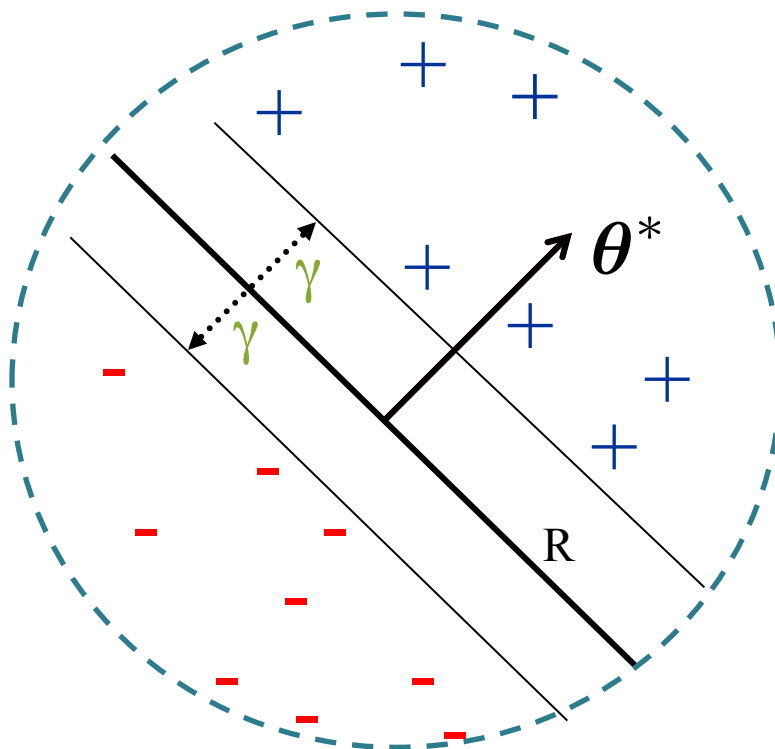**Algorithm 1** Perceptron Learning Algorithm (Batch)

---

1:    **procedure** PERCEPTRON($\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), \ldots, (\mathbf{x}^{(N)}, y^{(N)})\}$)
2:      $\boldsymbol{\theta} \leftarrow \mathbf{0}$                   ▷ Initialize parameters
3:      **while** not converged **do**
4:        **for** $i \in \{1, 2, \ldots, N\}$ **do**      ▷ For each example
5:          $\hat{y} \leftarrow \text{sign}(\boldsymbol{\theta}^T \mathbf{x}^{(i)})$          ▷ Predict
6:          **if** $\hat{y} \neq y^{(i)}$ **then**          ▷ If mistake
7:            $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + y^{(i)} \mathbf{x}^{(i)}$      ▷ Update parameters
8:      **return** $\boldsymbol{\theta}$

---

# Analysis: Perceptron

## Perceptron Mistake Bound

**Guarantee:** If data has margin $\gamma$ and all points inside a ball of radius $R$, then Perceptron makes $\leq (R/\gamma)^2$ mistakes.

(Normalized margin: multiplying all points by 100, or dividing all points by 100, doesn't change the number of mistakes; algo is invariant to scaling.)

Slide from Nina Balcan

# Analysis: Perceptron

**Perceptron Mistake Bound**
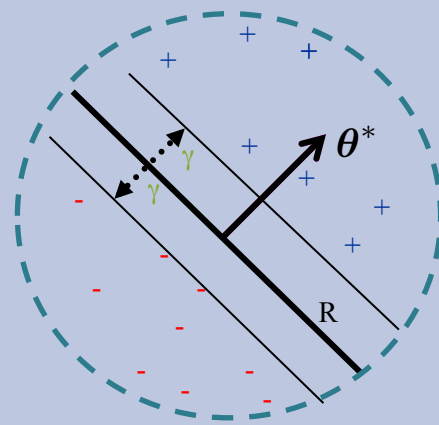
**Theorem 0.1** (Block (1962), Novikoff (1962))**.**
*Given dataset:* $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^{N}$.
*Suppose:*

1. *Finite size inputs:* $||x^{(i)}|| \leq R$
2. *Linearly separable data:* $\exists \boldsymbol{\theta}^*$ *s.t.* $||\boldsymbol{\theta}^*|| = 1$ *and* $y^{(i)}(\boldsymbol{\theta}^* \cdot \mathbf{x}^{(i)}) \geq \gamma, \forall i$

*Then: The number of mistakes made by the Perceptron algorithm on this dataset is*
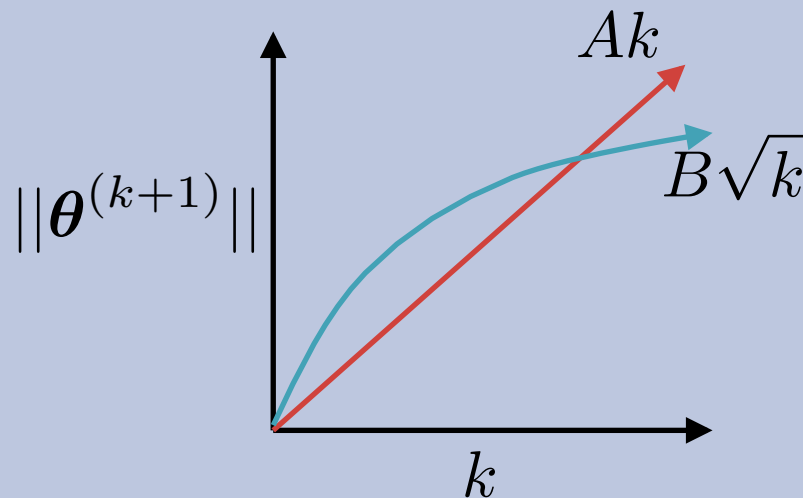
$$k \leq (R/\gamma)^2$$



Figure from Nina Balcan

# Analysis: Perceptron

**Proof of Perceptron Mistake Bound:**

We will show that there exist constants A and B s.t.

$$Ak \leq ||\boldsymbol{\theta}^{(k+1)}|| \leq B\sqrt{k}$$

# Analysis: Perceptron
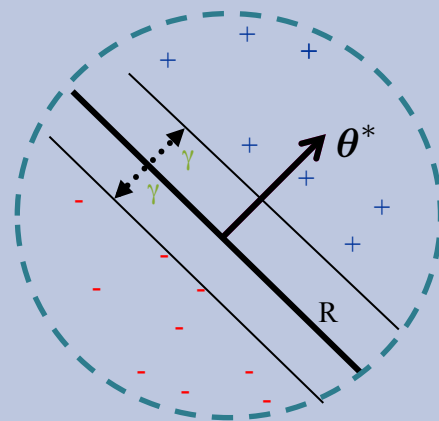
**Theorem 0.1** (Block (1962), Novikoff (1962)).

*Given dataset:* $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^{N}$.

*Suppose:*

1. *Finite size inputs:* $||x^{(i)}|| \leq R$
2. *Linearly separable data:* $\exists \boldsymbol{\theta}^*$ *s.t.* $||\boldsymbol{\theta}^*|| = 1$ *and*
   $y^{(i)}(\boldsymbol{\theta}^* \cdot \mathbf{x}^{(i)}) \geq \gamma, \forall i$

*Then: The number of mistakes made by the Perceptron algorithm on this dataset is*

$$k \leq (R/\gamma)^2$$



---

**Algorithm 1** Perceptron Learning Algorithm (Online)

---

1: **procedure** PERCEPTRON($\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \ldots\}$)
2:      $\boldsymbol{\theta} \leftarrow \mathbf{0}, k = 1$                           ▷ Initialize parameters
3:      **for** $i \in \{1, 2, \ldots\}$ **do**                 ▷ For each example
4:          **if** $y^{(i)}(\boldsymbol{\theta}^{(k)} \cdot \mathbf{x}^{(i)}) \leq 0$ **then**          ▷ If mistake
5:             $\boldsymbol{\theta}^{(k+1)} \leftarrow \boldsymbol{\theta}^{(k)} + y^{(i)}\mathbf{x}^{(i)}$     ▷ Update parameters
6:             $k \leftarrow k + 1$
7:      **return** $\boldsymbol{\theta}$

# Analysis: Perceptron

*Whiteboard*:
Proof of Perceptron Mistake Bound

# Analysis: Perceptron

**Proof of Perceptron Mistake Bound:**

Part 1: for some A, $Ak \leq ||\boldsymbol{\theta}^{(k+1)}||$

$\boldsymbol{\theta}^{(k+1)} \cdot \boldsymbol{\theta}^* = (\boldsymbol{\theta}^{(k)} + y^{(i)}\mathbf{x}^{(i)})\boldsymbol{\theta}^*$

       by Perceptron algorithm update

$= \boldsymbol{\theta}^{(k)} \cdot \boldsymbol{\theta}^* + y^{(i)}(\boldsymbol{\theta}^* \cdot \mathbf{x}^{(i)})$

$\geq \boldsymbol{\theta}^{(k)} \cdot \boldsymbol{\theta}^* + \gamma$

    by assumption

$\Rightarrow \boldsymbol{\theta}^{(k+1)} \cdot \boldsymbol{\theta}^* \geq k\gamma$

    by induction on $k$ since $\theta^{(1)} = \mathbf{0}$

$\Rightarrow ||\boldsymbol{\theta}^{(k+1)}|| \geq k\gamma$

    since $||\mathbf{w}|| \times ||\mathbf{u}|| \geq \mathbf{w} \cdot \mathbf{u}$ and $||\theta^*|| = 1$

Cauchy-Schwartz inequality

# Analysis: Perceptron

**Proof of Perceptron Mistake Bound:**

Part 2: for some B, $||\boldsymbol{\theta}^{(k+1)}|| \leq B\sqrt{k}$

$||\boldsymbol{\theta}^{(k+1)}||^2 = ||\boldsymbol{\theta}^{(k)} + y^{(i)}\mathbf{x}^{(i)}||^2$

      by Perceptron algorithm update

$= ||\boldsymbol{\theta}^{(k)}||^2 + (y^{(i)})^2||\mathbf{x}^{(i)}||^2 + 2y^{(i)}(\boldsymbol{\theta}^{(k)} \cdot \mathbf{x}^{(i)})$

$\leq ||\boldsymbol{\theta}^{(k)}||^2 + (y^{(i)})^2||\mathbf{x}^{(i)}||^2$

    since $k$th mistake $\Rightarrow y^{(i)}(\boldsymbol{\theta}^{(k)} \cdot \mathbf{x}^{(i)}) \leq 0$

$= ||\boldsymbol{\theta}^{(k)}||^2 + R^2$

    since $(y^{(i)})^2||\mathbf{x}^{(i)}||^2 = ||\mathbf{x}^{(i)}||^2 = R^2$ by assumption and $(y^{(i)})^2 = 1$

$\Rightarrow ||\boldsymbol{\theta}^{(k+1)}||^2 \leq kR^2$

    by induction on $k$ since $(\theta^{(1)})^2 = 0$

$\Rightarrow ||\boldsymbol{\theta}^{(k+1)}|| \leq \sqrt{k}R$

# Analysis: Perceptron

**Proof of Perceptron Mistake Bound:**
Part 3: Combining the bounds finishes the proof.

$$k\gamma \leq ||\boldsymbol{\theta}^{(k+1)}|| \leq \sqrt{k}R$$

$$\Rightarrow k \leq (R/\gamma)^2$$

The total number of mistakes must be less than this

# Extensions of Perceptron

- **Kernel Perceptron**
  - Choose a kernel $K(x', x)$
  - Apply the **kernel trick** to Perceptron
  - Resulting algorithm is **still very simple**

- **Structured Perceptron**
  - Basic idea can also be applied when **y** ranges over an exponentially large set
  - Mistake bound **does not** depend on the size of that set

# Summary: Perceptron

- Perceptron is a simple **linear classifier**
- **Simple learning algorithm**: when a mistake is made, add / subtract the features
- For linearly separable and inseparable data, we can **bound the number of mistakes** (geometric argument)
- Extensions support nonlinear separators and structured prediction

# RECALL: LOGISTIC REGRESSION

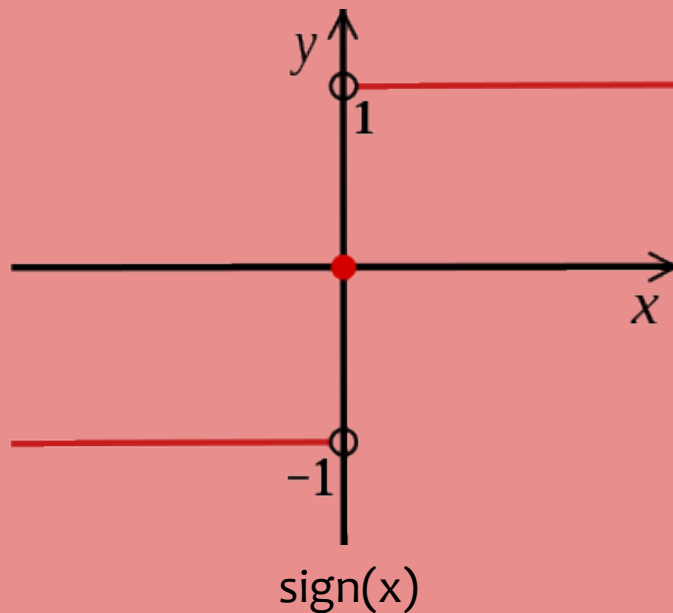# Using gradient ascent for linear classifiers

Key idea behind today's lecture:

1. Define a linear classifier (logistic regression)

2. Define an objective function (likelihood)

3. Optimize it with gradient descent to learn parameters

4. Predict the class with highest probability under the model
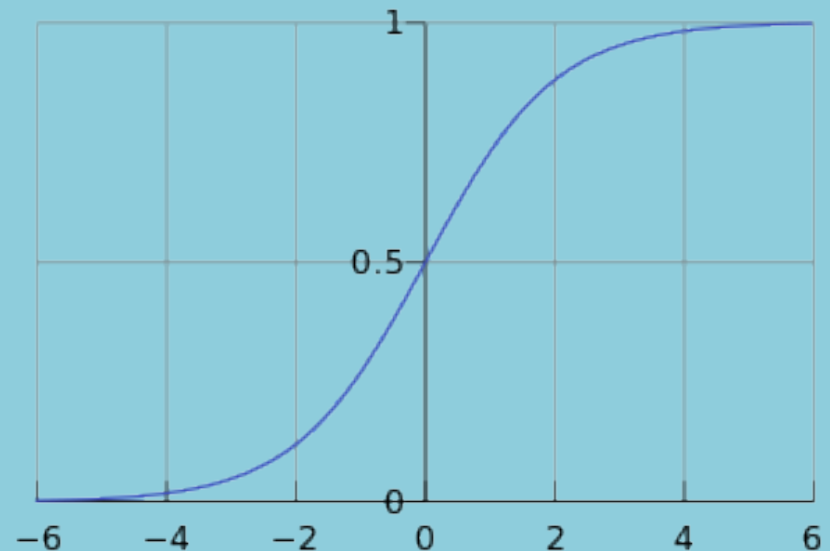
# Using gradient ascent for linear classifiers

## This decision function isn't differentiable:

$$h(\mathbf{x}) = \text{sign}(\boldsymbol{\theta}^T \mathbf{x})$$



sign(x)

## Use a differentiable function instead:

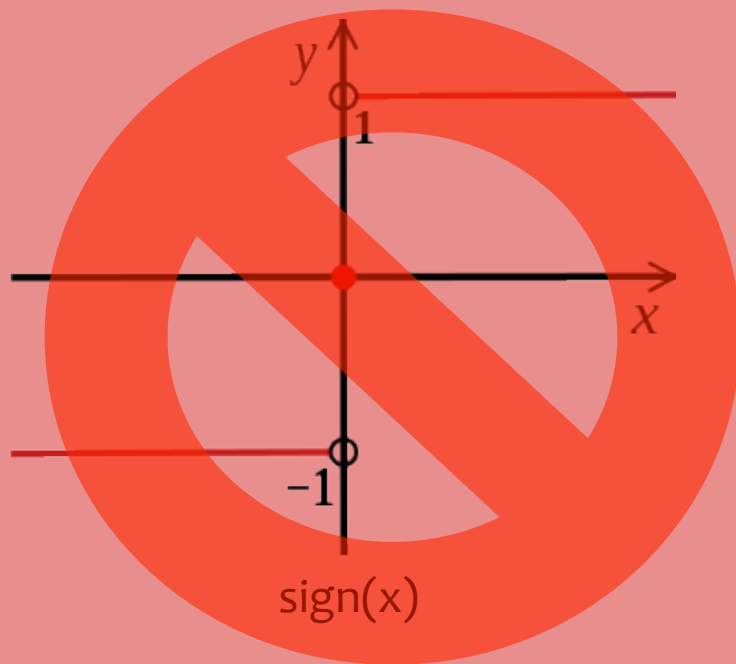$$p_{\boldsymbol{\theta}}(y = 1|\mathbf{x}) = \frac{1}{1 + \exp(-\boldsymbol{\theta}^T \mathbf{x})}$$



$$\text{logistic}(u) \equiv \frac{1}{1 + e^{-u}}$$
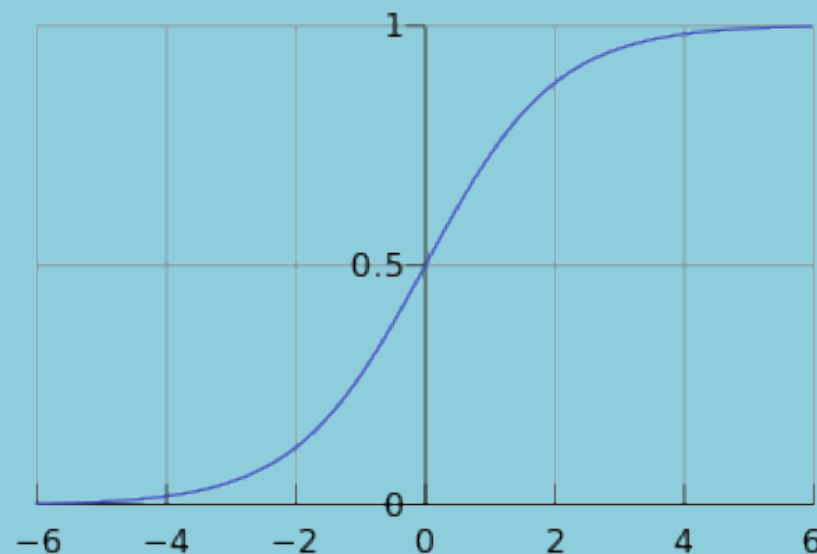
# Using gradient ascent for linear classifiers

This decision function isn't differentiable:

$$h(\mathbf{x}) = \text{sign}(\boldsymbol{\theta}^T \mathbf{x})$$



sign(x)

Use a differentiable function instead:

$$p_{\boldsymbol{\theta}}(y = 1|\mathbf{x}) = \frac{1}{1 + \exp(-\boldsymbol{\theta}^T \mathbf{x})}$$



$$\text{logistic}(u) \equiv \frac{1}{1 + e^{-u}}$$

# Logistic Regression

**Data:** Inputs are continuous vectors of length K. Outputs are discrete.
$$\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^{N} \text{ where } \mathbf{x} \in \mathbb{R}^K \text{ and } y \in \{0, 1\}$$

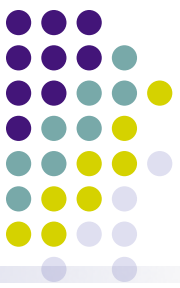**Model:** Logistic function applied to dot product of parameters with input vector.
$$p_{\boldsymbol{\theta}}(y = 1|\mathbf{x}) = \frac{1}{1 + \exp(-\boldsymbol{\theta}^T \mathbf{x})}$$

**Learning:** finds the parameters that minimize some objective function. $\boldsymbol{\theta}^* = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \, J(\boldsymbol{\theta})$

**Prediction:** Output is the most probable class.
$$\hat{y} = \underset{y \in \{0,1\}}{\operatorname{argmax}} \, p_{\boldsymbol{\theta}}(y|\mathbf{x})$$
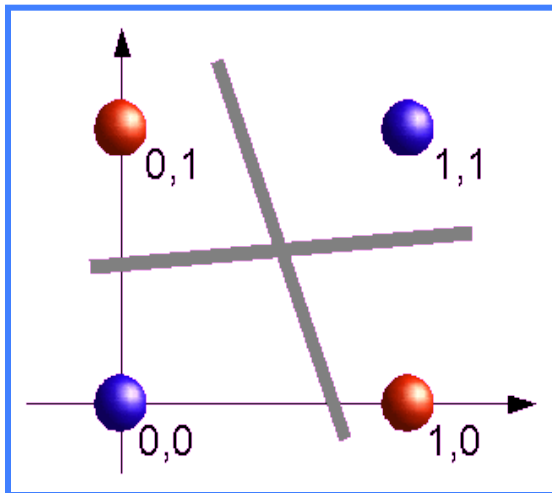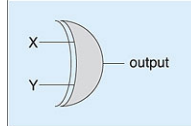
# NEURAL NETWORKS

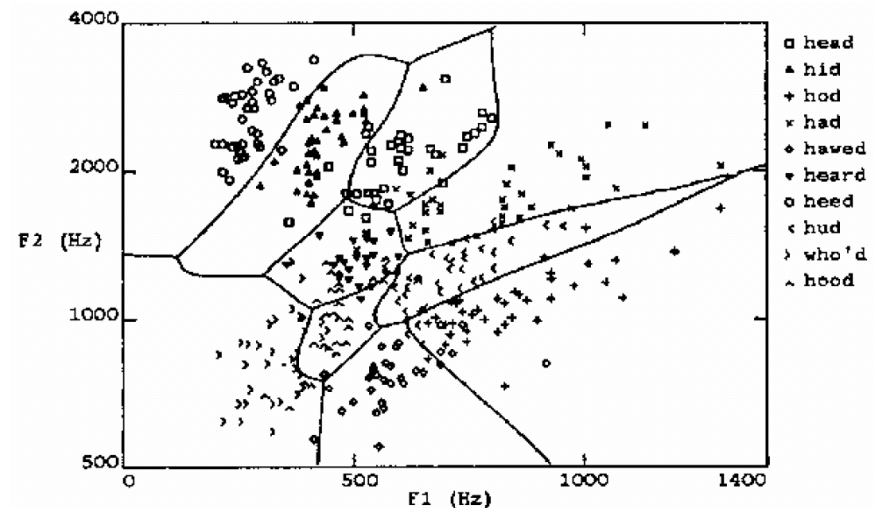# Learning highly non-linear functions

f: X $\rightarrow$ Y

- f might be non-linear function
- X (vector of) continuous and/or discrete vars
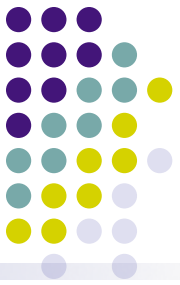- Y (vector of) continuous and/or discrete vars
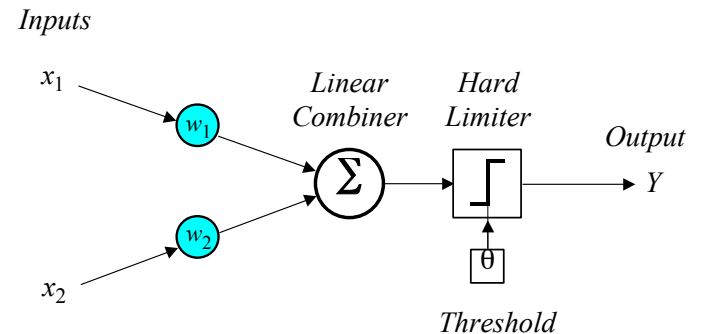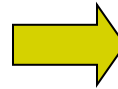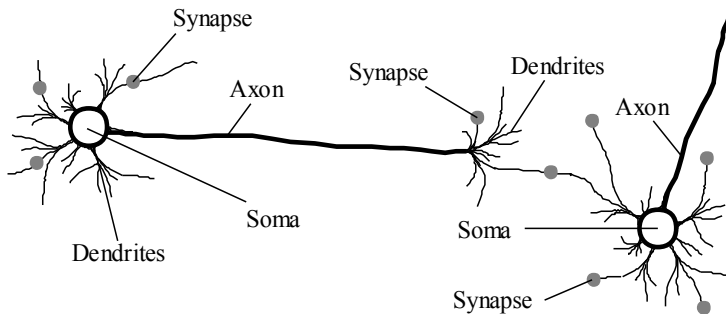
The XOR gate

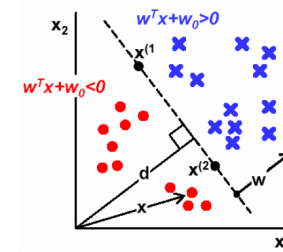Speech recognition

# Perceptron and Neural Nets

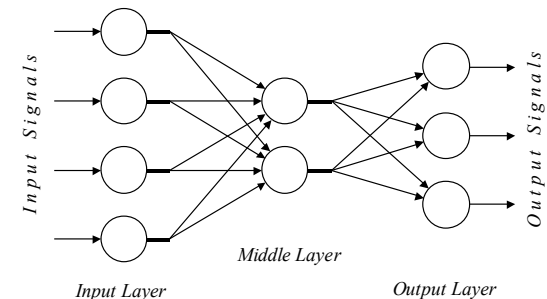- From biological neuron to artificial neuron (perceptron)



- Activation function

$$X = \sum_{i=1}^{n} x_i w_i \qquad Y = \begin{cases} +1, & \text{if } X \geq \omega_0 \\ -1, & \text{if } X < \omega_0 \end{cases}$$



- Artificial neuron networks
  - supervised learning
  - gradient descent

# Connectionist Models

- Consider humans:
  - Neuron switching time
    ~ 0.001 second
  - Number of neurons
    ~ $10^{10}$
  - Connections per neuron
    ~ $10^{4-5}$
  - Scene recognition time
    ~ 0.1 second
  - 100 inference steps doesn't seem like enough
    → much parallel computation

- Properties of artificial neural nets (ANN)
  - Many neuron-like threshold switching units
  - Many weighted interconnections among units
  - Highly parallel, distributed processes

# Why is everyone talking about Deep Learning?

- Because a lot of money is invested in it…
  - DeepMind:  Acquired by Google for **$400 million**
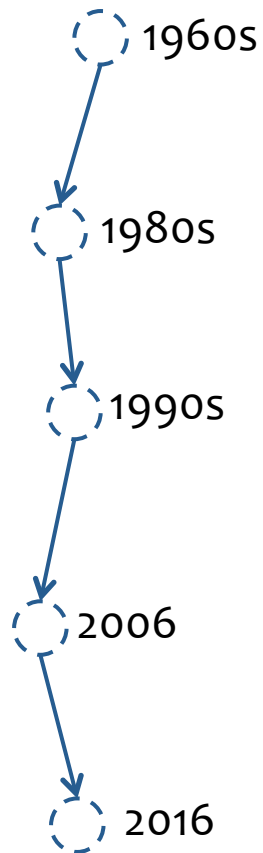
  - DNNResearch:  **Three person startup** (including Geoff Hinton) acquired by Google for unknown price tag

  - Enlitic, Ersatz, MetaMind, Nervana, Skylab: Deep Learning startups commanding **millions of VC dollars**

- Because it made the **front page** of the New York Times

# Why is everyone talking about Deep Learning?

1960s

1980s

1990s

2006

2016

## Deep learning:

- Has won numerous pattern recognition competitions
- Does so with minimal feature engineering

**This wasn't always the case!**

Since 1980s:  Form of models hasn't changed much, but lots of new tricks…

- More hidden units
- Better (online) optimization
- New nonlinear functions (ReLUs)
- Faster computers (CPUs and GPUs)

# A Recipe for
# Machine Learning

**1. Given training data:**

$$\{\boldsymbol{x}_i, \boldsymbol{y}_i\}_{i=1}^N$$

*Face*  *Face*  *Not a face*

**2. Choose each of these:**

– Decision function

$$\hat{\boldsymbol{y}} = f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)$$

– Loss function

$$\ell(\hat{\boldsymbol{y}}, \boldsymbol{y}_i) \in \mathbb{R}$$

**Examples:** Linear regression, Logistic regression, Neural Network

**Examples:** Mean-squared error, Cross Entropy

# A Recipe for Machine Learning

1. Given training data:

$$\{\boldsymbol{x}_i, \boldsymbol{y}_i\}_{i=1}^{N}$$

2. Choose each of these:
   – Decision function

$$\hat{\boldsymbol{y}} = f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)$$

   – Loss function

$$\ell(\hat{\boldsymbol{y}}, \boldsymbol{y}_i) \in \mathbb{R}$$

3. Define goal:

$$\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} \sum_{i=1}^{N} \ell(f_{\boldsymbol{\theta}}(\boldsymbol{x}_i), \boldsymbol{y}_i)$$

4. Train with SGD:

(take small steps opposite the gradient)

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta_t \nabla \ell(f_{\boldsymbol{\theta}}(\boldsymbol{x}_i), \boldsymbol{y}_i)$$

A Recipe for

## Gradients

1. Given training dat

$$\{\boldsymbol{x}_i, \boldsymbol{y}_i\}_{i=1}^N$$

2. Choose each of t

– Decision function

$$\hat{\boldsymbol{y}} = f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)$$

– Loss function

$$\ell(\hat{\boldsymbol{y}}, \boldsymbol{y}_i) \in \mathbb{R}$$

**Backpropagation** can compute this gradient!

And it's a **special case of a more general algorithm** called reverse-mode automatic differentiation that can compute the gradient of any differentiable function efficiently!

opposite the gradient)

$$\boldsymbol{\theta}^{(t} \qquad {}^{(t)} - \eta_t \nabla \ell(f_{\boldsymbol{\theta}}(\boldsymbol{x}_i), \boldsymbol{y}_i)$$

A Recipe for

B...

## Goals for Today's Lecture

1. **Explore a new class of decision functions** (Neural Networks)

2. Consider **variants of this recipe** for training

2. Choose each of these:

– Decision function

$$\hat{y} = f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)$$

– Loss function

$$\ell(\hat{\boldsymbol{y}}, \boldsymbol{y}_i) \in \mathbb{R}$$

Train with SGD:

take small steps opposite the gradient)

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta_t \nabla \ell(f_{\boldsymbol{\theta}}(\boldsymbol{x}_i), \boldsymbol{y}_i)$$

# Linear Regression

$$y = h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \sigma(\boldsymbol{\theta}^T \boldsymbol{x})$$

where $\sigma(a) = a$

Output

**y**

θ₁      θ₂      θ₃      θ_M

Input

$x_1$      $x_2$      $x_3$    ...    $x_M$

# Logistic Regression

$$y = h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \sigma(\boldsymbol{\theta}^T \boldsymbol{x})$$

Output

$$\text{where } \sigma(a) = \frac{1}{1 + \exp(-a)}$$

y

Input

$x_1$   $x_2$   $x_3$   $\cdots$   $x_M$

$\theta_1$   $\theta_2$   $\theta_3$   $\theta_M$

# Logistic Regression

$$y = h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \sigma(\boldsymbol{\theta}^T \boldsymbol{x})$$

Output

$$\text{where } \sigma(a) = \frac{1}{1 + \exp(-a)}$$

**y**

*Face*     *Face*     *Not a face*

$\theta_1$      $\theta_2$      $\theta_3$

Input

**x₁**     **x₂**     **x₃**

...

# Logistic Regression

$$y = h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \sigma(\boldsymbol{\theta}^T \boldsymbol{x})$$

Output

$$\text{where } \sigma(a) = \frac{1}{1 + \exp(-a)}$$

y

θ₁ corresponds to $\theta_1$

θ₂ corresponds to $\theta_2$

θ₃ corresponds to $\theta_3$

Input

$x_1$

$x_2$

$x_3$

1     1     0

# Logistic Regression

$$y = h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \sigma(\boldsymbol{\theta}^T \boldsymbol{x})$$

Output

**y**

where $\sigma(a) = \dfrac{1}{1 + \exp(-a)}$

$\theta_1$  $\theta_2$  $\theta_3$  $\theta_M$

Input

$x_1$  $x_2$  $x_3$  $\cdots$  $x_M$

# Neural Network Model

Inputs

Output

*Age* 34

.6

.4

.2

.1

*Gender* 2

.5

0.6

.3

.2

.7

.8

"Probability of beingAlive"

*Stage* 4

.2

*Independent variables*

Weights

HiddenLayer

Weights

*Dependent variable*

*Prediction*

# "Combined logistic models"

Inputs

Output

*Age* 34

*Gender* 2

*Stage* 4

.2

.3

.2

.5

.8

Σ

0.6

"Probability of beingAlive"

*Independent variables*

Weights

HiddenLayer

Weights

*Dependent variable*

*Prediction*

Inputs

Output

*Age*  34

.6
.2
.1
.3
.7
.2

.5

.8

Σ

0.6

*Gender*  1

*Stage*  4

"Probability of beingAlive"

*Independent variables*

Weights

HiddenLayer

Weights

*Dependent variable*

*Prediction*

# Not really,
# no target for hidden units...



*Age* 34

*Gender* 2

*Stage* 4

.6  .2  .1  .3  .7  .2

Σ  .4

Σ  .2

.5  .8

Σ  0.6

"Probability of beingAlive"

*Independent variables*

Weights

HiddenL ayer

Weights

*Dependent variable*

*Prediction*

# Jargon Pseudo-Correspondence

- Independent variable = input variable

- Dependent variable = output variable

- Coefficients = "weights"

- Estimates = "targets"

## Logistic Regression Model (the sigmoid unit)

Inputs

Output

*Age* — 34 — 5

*Gender* — 1 — 4

*Stage* — 4 — 8

$\Sigma$

0.6

"Probability of beingAlive"

*Independent variables*

*x1, x2, x3*

Coefficients

*a, b, c*

*Dependent variable*

*p Prediction*

# Neural Network

Output

Hidden Layer

Input

y

$a_1$   $a_2$   ...   $a_D$

$x_1$   $x_2$   $x_3$   ...   $x_M$

# Neural Network

Output

y

Hidden Layer

$z_1$  $z_2$  ...  $z_D$

Input

$x_1$  $x_2$  $x_3$  ...  $x_M$

(E) **Output (sigmoid)**
$$y = \frac{1}{1+\exp(-b)}$$

(D) **Output (linear)**
$$b = \sum_{j=0}^{D} \beta_j z_j$$

(C) **Hidden (sigmoid)**
$$z_j = \frac{1}{1+\exp(-a_j)}, \ \forall j$$

(B) **Hidden (linear)**
$$a_j = \sum_{i=0}^{M} \alpha_{ji} x_i, \ \forall j$$

(A) **Input**
Given $x_i, \ \forall i$

# Building a Neural Net

Output

Features

$y$

$x_1$  $x_2$  $\cdots$  $x_M$

# Building a Neural Net

Output

$y$

Hidden Layer

$a_1$  $a_2$  $\cdots$  $a_D$

$D = M$

**1**  **1**  **1**

Input

$x_1$  $x_2$  $\cdots$  $x_M$

# Building a Neural Net

Output

Hidden Layer

Input

$D = M$

# Building a Neural Net



Output

Hidden Layer

$D = M$

Input

# Building a Neural Net



Output

Hidden Layer

Input

$D < M$

# Decision Boundary

- 0 hidden layers: linear classifier
  - Hyperplanes



Example from to Eric Postma via Jason Eisner

# Decision Boundary

- ## 1 hidden layer
  - Boundary of convex region (open or closed)

# Decision Boundary

- 2 hidden layers
  - Combinations of convex regions

Example from to Eric Postma via Jason Eisner

# Multi-Class Output

Output

$$y_1 \quad \cdots \quad y_K$$

Hidden Layer

$$a_1 \quad a_2 \quad \cdots \quad a_D$$

Input

$$x_1 \quad x_2 \quad x_3 \quad \cdots \quad x_M$$

66

Next lecture:

Output

Hidden Layer 1

Input

# Deeper Networks

## Next lecture:



Output: y

Hidden Layer 2: $b_1$, $b_2$, ..., $b_E$

Hidden Layer 1: $a_1$, $a_2$, ..., $a_D$

Input: $x_1$, $x_2$, $x_3$, ..., $x_M$

Next lecture: Making the neural networks deeper



Output

Hidden Layer 3

Hidden Layer 2

Hidden Layer 1

Input

# Decision Functions

# Different Levels of Abstraction

- We don't know the "right" levels of abstraction

- So let the model figure it out!

Feature representation



3rd layer "Objects"

2nd layer "Object parts"

1st layer "Edges"

Pixels

Example from Honglak Lee (NIPS 2010)

## Face Recognition:

– Deep Network can build up increasingly higher levels of abstraction

– Lines, parts, regions

Feature representation



3rd layer "Objects"

2nd layer "Object parts"

1st layer "Edges"

Pixels

Example from Honglak Lee (NIPS 2010)

# Different Levels of Abstraction



## Feature representation

Output — y

Hidden Layer 3 — $c_1$, $c_2$, …, $c_F$

Hidden Layer 2 — $b_1$, $b_2$, …, $b_E$

Hidden Layer 1 — $a_1$, $a_2$, …, $a_D$

Input — $x_1$, $x_2$, $x_3$, …, $x_M$

3rd layer "Objects"

2nd layer "Object parts"

1st layer "Edges"

Pixels

Example from Honglak Lee (NIPS 2010)

# ARCHITECTURES

# Neural Network Architectures

Even for a basic Neural Network, there are many design decisions to make:

1. # of hidden layers (depth)

2. # of units per hidden layer (width)

3. Type of activation function (nonlinearity)

4. Form of objective function

# Activation Functions

Neural Network with sigmoid activation functions



Output

Hidden Layer

Input

(F) **Loss**
$$J = \frac{1}{2}(y - y^*)^2$$

(E) **Output (sigmoid)**
$$y = \frac{1}{1+\exp(-b)}$$

(D) **Output (linear)**
$$b = \sum_{j=0}^{D} \beta_j z_j$$

(C) **Hidden (sigmoid)**
$$z_j = \frac{1}{1+\exp(-a_j)}, \ \forall j$$

(B) **Hidden (linear)**
$$a_j = \sum_{i=0}^{M} \alpha_{ji} x_i, \ \forall j$$

(A) **Input**
Given $x_i, \ \forall i$

# Activation Functions

Neural Network with arbitrary
nonlinear activation functions



Output

Hidden Layer

Input

(F) **Loss**
$$J = \tfrac{1}{2}(y - y^*)^2$$

(E) **Output (nonlinear)**
$$y = \sigma(b)$$

(D) **Output (linear)**
$$b = \sum_{j=0}^{D} \beta_j z_j$$

(C) **Hidden (nonlinear)**
$$z_j = \sigma(a_j), \ \forall j$$

(B) **Hidden (linear)**
$$a_j = \sum_{i=0}^{M} \alpha_{ji} x_i, \ \forall j$$

(A) **Input**
Given $x_i, \ \forall i$

# Activation Functions

## Sigmoid / Logistic Function

$$\text{logistic}(u) \equiv \frac{1}{1 + e^{-u}}$$



So far, we've assumed that the activation function (nonlinearity) is always the sigmoid function…



$net = \sum_{i=0}^{n} w_i x_i$ $\qquad o = \sigma(net) = \dfrac{1}{1 + e^{-net}}$

# Activation Functions

- A new change: modifying the nonlinearity
  - The logistic is not widely used in modern ANNs



Alternate 1:
tanh

Like logistic function but shifted to range [-1, +1]

$net = \sum_{i=0}^{n} w_i x_i$

$o = \sigma(net) = \dfrac{1}{1 + e^{-net}}$

# Understanding the difficulty of training deep feedforward neural networks

AI Stats 2010



Figure from Glorot & Bentio (2010)

# Activation Functions

- A new change: modifying the nonlinearity
  - reLU often used in vision tasks



$$\max(0, z)$$

Alternate 2: rectified linear unit

Linear with a cutoff at zero

(Implementation: clip the gradient when you pass zero)

$$\max(0, w \cdot x + b).$$



$$net = \sum_{i=0}^{n} w_i \, x_i \qquad o = \sigma(net) = \frac{1}{1 + e^{-net}}$$

# Activation Functions

- A new change: modifying the nonlinearity
  - reLU often used in vision tasks



Alternate 2: rectified linear unit

Soft version: log(exp(x)+1)

Doesn't saturate (at one end)
Sparsifies outputs
Helps with vanishing gradient

$net = \sum_{i=0}^{n} w_i x_i$

$o = \sigma(net) = \dfrac{1}{1+e^{-net}}$

# Objective Functions for NNs

- Regression:
  - Use the same objective as Linear Regression
  - Quadratic loss (i.e. mean squared error)
- Classification:
  - Use the same objective as Logistic Regression
  - Cross-entropy (i.e. negative log likelihood)
  - This requires probabilities, so we add an additional "softmax" layer at the end of our network

| | Forward | Backward |
|---|---|---|
| Quadratic | $J = \dfrac{1}{2}(y - y^*)^2$ | $\dfrac{dJ}{dy} = y - y^*$ |
| Cross Entropy | $J = y^* \log(y) + (1 - y^*) \log(1 - y)$ | $\dfrac{dJ}{dy} = y^* \dfrac{1}{y} + (1 - y^*) \dfrac{1}{y - 1}$ |

# Multi-Class Output

# Multi-Class Output

Softmax:

$$y_k = \frac{\exp(b_k)}{\sum_{l=1}^{K} \exp(b_l)}$$



Output

Hidden Layer

Input

(F) **Loss**
$J = \sum_{k=1}^{K} y_k^* \log(y_k)$

(E) **Output (softmax)**
$y_k = \frac{\exp(b_k)}{\sum_{l=1}^{K} \exp(b_l)}$

(D) **Output (linear)**
$b_k = \sum_{j=0}^{D} \beta_{kj} z_j \ \forall k$

(C) **Hidden (nonlinear)**
$z_j = \sigma(a_j), \ \forall j$

(B) **Hidden (linear)**
$a_j = \sum_{i=0}^{M} \alpha_{ji} x_i, \ \forall j$

(A) **Input**
Given $x_i, \ \forall i$

# Cross-entropy vs. Quadratic loss



Figure 5: Cross entropy (black, surface on top) and quadratic (red, bottom surface) cost as a function of two weights (one at each layer) of a network with two layers, $W_1$ respectively on the first layer and $W_2$ on the second, output layer.

# A Recipe for Machine Learning

1. Given training data:

$$\{\boldsymbol{x}_i, \boldsymbol{y}_i\}_{i=1}^{N}$$

2. Choose each of these:
   – Decision function

$$\hat{\boldsymbol{y}} = f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)$$

   – Loss function

$$\ell(\hat{\boldsymbol{y}}, \boldsymbol{y}_i) \in \mathbb{R}$$

3. Define goal:

$$\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} \sum_{i=1}^{N} \ell(f_{\boldsymbol{\theta}}(\boldsymbol{x}_i), \boldsymbol{y}_i)$$

4. Train with SGD:

(take small steps opposite the gradient)

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta_t \nabla \ell(f_{\boldsymbol{\theta}}(\boldsymbol{x}_i), \boldsymbol{y}_i)$$

# Objective Functions

**Matching Quiz:** Suppose you are given a neural net with a single output, y, and one hidden layer.

1) Minimizing sum of squared errors…

2) Minimizing sum of squared errors plus squared Euclidean norm of weights…

3) Minimizing cross-entropy…

4) Minimizing hinge loss…

…gives…

5) …MLE estimates of weights assuming target follows a Bernoulli with parameter given by the output value

6) …MAP estimates of weights assuming weight priors are zero mean Gaussian

7) …estimates with a large margin on the training data

8) …MLE estimates of weights assuming zero mean Gaussian noise on the output value

A. 1=5, 2=7, 3=6, 4=8
B. 1=5, 2=7, 3=8, 4=6
C. 1=7, 2=5, 3=5, 4=7
D. 1=7, 2=5, 3=6, 4=8
E. 1=8, 2=6, 3=5, 4=7
F. 1=8, 2=6, 3=8, 4=6

# BACKPROPAGATION

# A Recipe for Machine Learning

1. Given training data:
$$\{\boldsymbol{x}_i, \boldsymbol{y}_i\}_{i=1}^N$$

2. Choose each of these:
   – Decision function
   $$\hat{\boldsymbol{y}} = f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)$$

   – Loss function
   $$\ell(\hat{\boldsymbol{y}}, \boldsymbol{y}_i) \in \mathbb{R}$$

3. Define goal:
$$\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} \sum_{i=1}^N \ell(f_{\boldsymbol{\theta}}(\boldsymbol{x}_i), \boldsymbol{y}_i)$$

4. Train with SGD:

(take small steps opposite the gradient)

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta_t \nabla \ell(f_{\boldsymbol{\theta}}(\boldsymbol{x}_i), \boldsymbol{y}_i)$$

# Backpropagation

- **Question 1:**
  When can we compute the gradients of the parameters of an arbitrary neural network?

- **Question 2:**
  When can we make the gradient computation efficient?

# Chain Rule

**Given:** $\boldsymbol{y} = g(\boldsymbol{u})$ and $\boldsymbol{u} = h(\boldsymbol{x})$.

**Chain Rule:**

$$\frac{dy_i}{dx_k} = \sum_{j=1}^{J} \frac{dy_i}{du_j} \frac{du_j}{dx_k}, \quad \forall i, k$$

# Chain Rule

**Given:** $y = g(u)$ and $u = h(x)$.

**Chain Rule:**
$$\frac{dy_i}{dx_k} = \sum_{j=1}^{J} \frac{dy_i}{du_j} \frac{du_j}{dx_k}, \quad \forall i, k$$

**Backpropagation** is just repeated application of the **chain rule** from Calculus 101.

Given: $\boldsymbol{y} = g(\boldsymbol{u})$ and $\boldsymbol{u} = h(\boldsymbol{x})$.

Chain Rule:

$$\frac{dy_i}{dx_k} = \sum_{j=1}^{J} \frac{dy_i}{du_j} \frac{du_j}{dx_k}, \quad \forall i, k$$

**Backpropagation:**
1. **Instantiate the computation as a directed acyclic graph**, where each intermediate quantity is a node
2. At each node, store (a) the quantity computed in the forward pass and (b) the **partial derivative** of the goal with respect to that node's intermediate quantity.
3. **Initialize** all partial derivatives to 0.
4. Visit each node in **reverse topological order**. At each node, add its contribution to the partial derivatives of its parents

This algorithm is also called **automatic differentiation in the reverse-mode**

93

# Backpropagation

**Simple Example:** The goal is to compute $J = \cos(\sin(x^2) + 3x^2)$ on the forward pass and the derivative $\frac{dJ}{dx}$ on the backward pass.

Forward

$J = cos(u)$

$u = u_1 + u_2$

$u_1 = sin(t)$

$u_2 = 3t$

$t = x^2$

# Backpropagation

**Simple Example:** The goal is to compute $J = \cos(\sin(x^2) + 3x^2)$ on the forward pass and the derivative $\frac{dJ}{dx}$ on the backward pass.

| Forward | Backward | | |
|---|---|---|---|
| $J = cos(u)$ | $\dfrac{dJ}{du} \mathrel{+}= -sin(u)$ | | |
| $u = u_1 + u_2$ | $\dfrac{dJ}{du_1} \mathrel{+}= \dfrac{dJ}{du}\dfrac{du}{du_1}, \quad \dfrac{du}{du_1} = 1$ | | $\dfrac{dJ}{du_2} \mathrel{+}= \dfrac{dJ}{du}\dfrac{du}{du_2}, \quad \dfrac{du}{du_2} = 1$ |
| $u_1 = sin(t)$ | $\dfrac{dJ}{dt} \mathrel{+}= \dfrac{dJ}{du_1}\dfrac{du_1}{dt}, \quad \dfrac{du_1}{dt} = \cos(t)$ | | |
| $u_2 = 3t$ | $\dfrac{dJ}{dt} \mathrel{+}= \dfrac{dJ}{du_2}\dfrac{du_2}{dt}, \quad \dfrac{du_2}{dt} = 3$ | | |
| $t = x^2$ | $\dfrac{dJ}{dx} \mathrel{+}= \dfrac{dJ}{dt}\dfrac{dt}{dx}, \quad \dfrac{dt}{dx} = 2x$ | | |

# Backpropagation

Output **y**

**Case 1:
Logistic
Regression**

$\theta_1$   $\theta_2$   $\theta_3$   $\theta_M$

Input   **x₁**   **x₂**   **x₃**   ...   **x_M**

## Forward

$$J = y^* \log y + (1 - y^*) \log(1 - y)$$

$$y = \frac{1}{1 + \exp(-a)}$$

$$a = \sum_{j=0}^{D} \theta_j x_j$$

## Backward

$$\frac{dJ}{dy} = \frac{y^*}{y} + \frac{(1 - y^*)}{y - 1}$$

$$\frac{dJ}{da} = \frac{dJ}{dy}\frac{dy}{da}, \quad \frac{dy}{da} = \frac{\exp(-a)}{(\exp(-a) + 1)^2}$$

$$\frac{dJ}{d\theta_j} = \frac{dJ}{da}\frac{da}{d\theta_j}, \quad \frac{da}{d\theta_j} = x_j$$

$$\frac{dJ}{dx_j} = \frac{dJ}{da}\frac{da}{dx_j}, \quad \frac{da}{dx_j} = \theta_j$$

# Backpropagation

Output

Hidden Layer

Input

(E) **Output (sigmoid)**
$$y = \frac{1}{1+\exp(-b)}$$

(D) **Output (linear)**
$$b = \sum_{j=0}^{D} \beta_j z_j$$

(C) **Hidden (sigmoid)**
$$z_j = \frac{1}{1+\exp(-a_j)}, \ \forall j$$

(B) **Hidden (linear)**
$$a_j = \sum_{i=0}^{M} \alpha_{ji} x_i, \ \forall j$$

(A) **Input**
Given $x_i, \ \forall i$

97

# Backpropagation



(F) **Loss**
$$J = \frac{1}{2}(y - y^*)^2$$

(E) **Output (sigmoid)**
$$y = \frac{1}{1+\exp(-b)}$$

(D) **Output (linear)**
$$b = \sum_{j=0}^{D} \beta_j z_j$$

(C) **Hidden (sigmoid)**
$$z_j = \frac{1}{1+\exp(-a_j)}, \ \forall j$$

(B) **Hidden (linear)**
$$a_j = \sum_{i=0}^{M} \alpha_{ji} x_i, \ \forall j$$

(A) **Input**
Given $x_i, \ \forall i$

Output

Hidden Layer

Input

98

# Backpropagation

**Case 2:
Neural
Network**



Forward

$$J = y^* \log y + (1 - y^*) \log(1 - y)$$

$$y = \frac{1}{1 + \exp(-b)}$$

$$b = \sum_{j=0}^{D} \beta_j z_j$$

$$z_j = \frac{1}{1 + \exp(-a_j)}$$

$$a_j = \sum_{i=0}^{M} \alpha_{ji} x_i$$

Backward

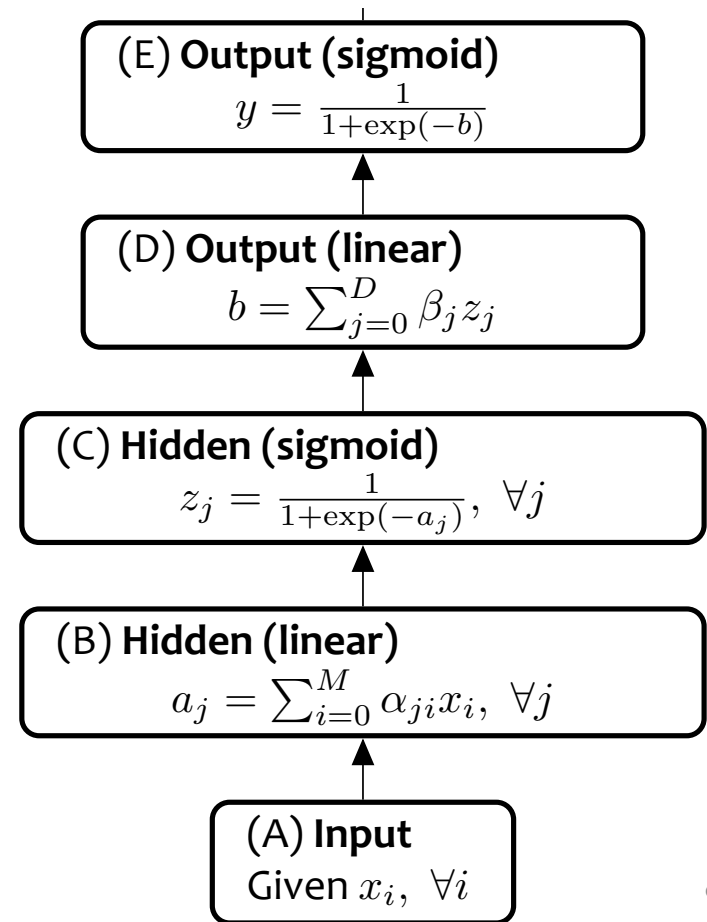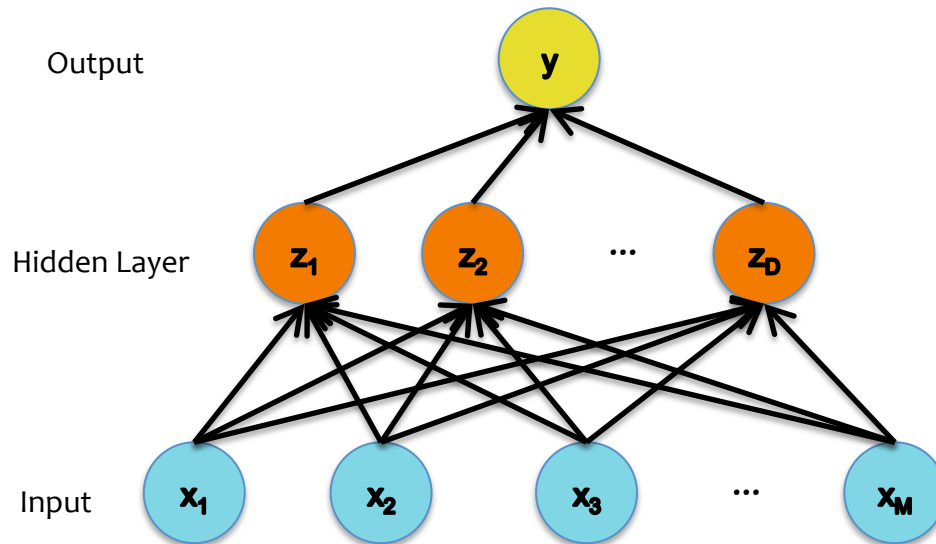$$\frac{dJ}{dy} = \frac{y^*}{y} + \frac{(1 - y^*)}{y - 1}$$

$$\frac{dJ}{db} = \frac{dJ}{dy} \frac{dy}{db}, \quad \frac{dy}{db} = \frac{\exp(-b)}{(\exp(-b) + 1)^2}$$

$$\frac{dJ}{d\beta_j} = \frac{dJ}{db} \frac{db}{d\beta_j}, \quad \frac{db}{d\beta_j} = z_j$$

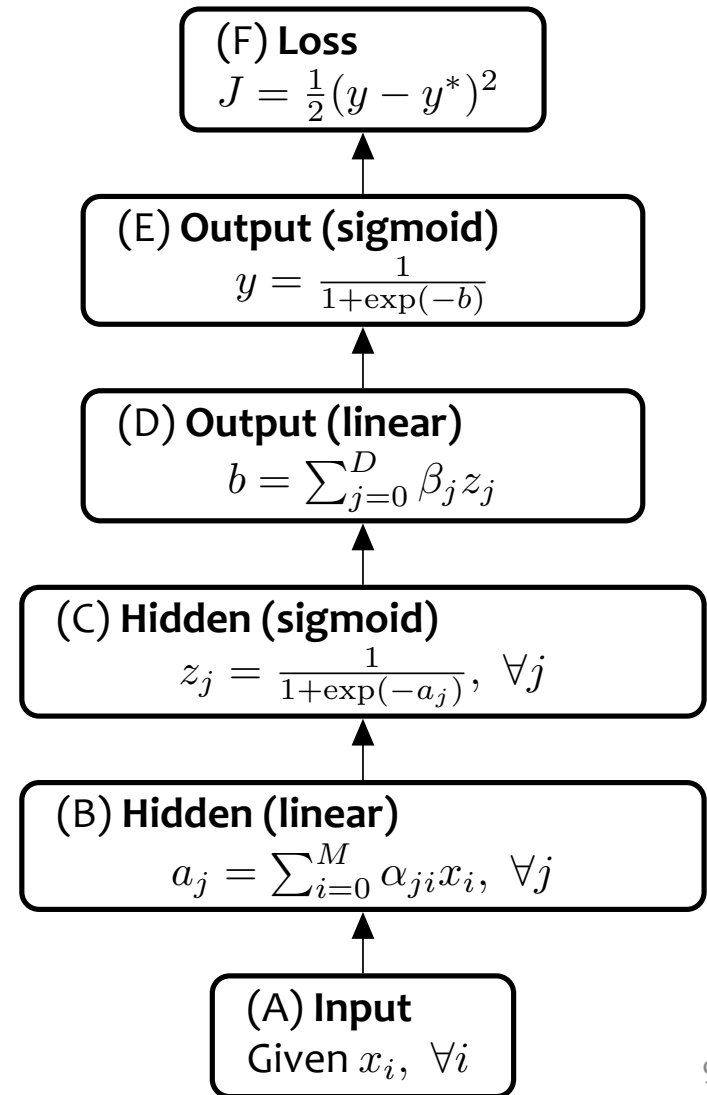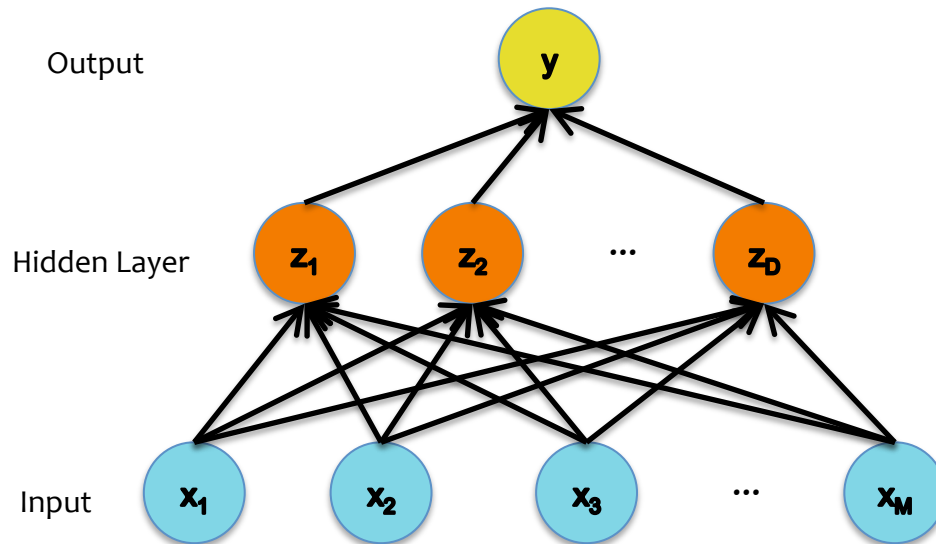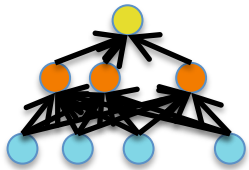$$\frac{dJ}{dz_j} = \frac{dJ}{db} \frac{db}{dz_j}, \quad \frac{db}{dz_j} = \beta_j$$

$$\frac{dJ}{da_j} = \frac{dJ}{dz_j} \frac{dz_j}{da_j}, \quad \frac{dz_j}{da_j} = \frac{\exp(-a_j)}{(\exp(-a_j) + 1)^2}$$

$$\frac{dJ}{d\alpha_{ji}} = \frac{dJ}{da_j} \frac{da_j}{d\alpha_{ji}}, \quad \frac{da_j}{d\alpha_{ji}} = x_i$$

$$\frac{dJ}{dx_i} = \frac{dJ}{da_j} \frac{da_j}{dx_i}, \quad \frac{da_j}{dx_i} = \sum_{j=0}^{D} \alpha_{ji}$$
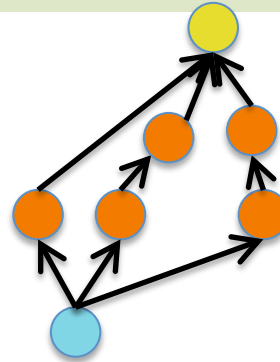
# Training

# Chain Rule

Given: $\boldsymbol{y} = g(\boldsymbol{u})$ and $\boldsymbol{u} = h(\boldsymbol{x})$.

Chain Rule:

$$\frac{dy_i}{dx_k} = \sum_{j=1}^{J} \frac{dy_i}{du_j} \frac{du_j}{dx_k}, \quad \forall i, k$$

**Backpropagation:**
1. **Instantiate the computation as a directed acyclic graph**, where each intermediate quantity is a node
2. At each node, store (a) the quantity computed in the forward pass and (b) the **partial derivative** of the goal with respect to that node's intermediate quantity.
3. **Initialize** all partial derivatives to 0.
4. Visit each node in **reverse topological order**. At each node, add its contribution to the partial derivatives of its parents

This algorithm is also called **automatic differentiation in the reverse-mode**

# Training

# Chain Rule

**Given:** $\boldsymbol{y} = g(\boldsymbol{u})$ and $\boldsymbol{u} = h(\boldsymbol{x})$.

**Chain Rule:**

$$\frac{dy_i}{dx_k} = \sum_{j=1}^{J} \frac{dy_i}{du_j} \frac{du_j}{dx_k}, \quad \forall i, k$$

**Backpropagation:**

1. **Instantiate the computation as a directed acyclic graph**, where each node represents a Tensor.
2. At each node, store (a) the quantity computed in the forward pass and (b) the **partial derivatives** of the goal with respect to that node's Tensor.
3. **Initialize** all partial derivatives to 0.
4. Visit each node in **reverse topological order**. At each node, add its contribution to the partial derivatives of its parents

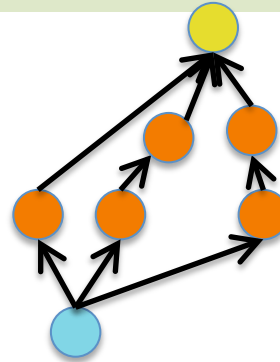This algorithm is also called **automatic differentiation in the reverse-mode**

**Case 2:**

| | Forward | Backward |
|---|---|---|
| Module 5 | $J = y^* \log y + (1 - y^*) \log(1 - y)$ | $\dfrac{dJ}{dy} = \dfrac{y^*}{y} + \dfrac{(1 - y^*)}{y - 1}$ |
| Module 4 | $y = \dfrac{1}{1 + \exp(-b)}$ | $\dfrac{dJ}{db} = \dfrac{dJ}{dy}\dfrac{dy}{db}, \; \dfrac{dy}{db} = \dfrac{\exp(-b)}{(\exp(-b) + 1)^2}$ |
| Module 3 | $b = \displaystyle\sum_{j=0}^{D} \beta_j z_j$ | $\dfrac{dJ}{d\beta_j} = \dfrac{dJ}{db}\dfrac{db}{d\beta_j}, \; \dfrac{db}{d\beta_j} = z_j$ <br><br> $\dfrac{dJ}{dz_j} = \dfrac{dJ}{db}\dfrac{db}{dz_j}, \; \dfrac{db}{dz_j} = \beta_j$ |
| Module 2 | $z_j = \dfrac{1}{1 + \exp(-a_j)}$ | $\dfrac{dJ}{da_j} = \dfrac{dJ}{dz_j}\dfrac{dz_j}{da_j}, \; \dfrac{dz_j}{da_j} = \dfrac{\exp(-a_j)}{(\exp(-a_j) + 1)^2}$ |
| Module 1 | $a_j = \displaystyle\sum_{i=0}^{M} \alpha_{ji} x_i$ | $\dfrac{dJ}{d\alpha_{ji}} = \dfrac{dJ}{da_j}\dfrac{da_j}{d\alpha_{ji}}, \; \dfrac{da_j}{d\alpha_{ji}} = x_i$ <br><br> $\dfrac{dJ}{dx_i} = \dfrac{dJ}{da_j}\dfrac{da_j}{dx_i}, \; \dfrac{da_j}{dx_i} = \displaystyle\sum_{j=0}^{D} \alpha_{ji}$ |

A Recipe for

## Gradients

1. Given training dat

$$\{\boldsymbol{x}_i, \boldsymbol{y}_i\}_{i=1}^{N}$$

**Backpropagation** can compute this gradient!

And it's a **special case of a more general algorithm** called reverse-mode automatic differentiation that can compute the gradient of any differentiable function efficiently!

2. Choose each of th

– Decision functio

$$\hat{y} = f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)$$

opposite the gradient)

– Loss function

$$\ell(\hat{\boldsymbol{y}}, \boldsymbol{y}_i) \in \mathbb{R}$$

$$\boldsymbol{\theta}^{(t)} \qquad {}^{(t)} - \eta_t \nabla \ell(f_{\boldsymbol{\theta}}(\boldsymbol{x}_i), \boldsymbol{y}_i)$$

# Summary

1. **Neural Networks…**
   - provide a way of learning features
   - are highly nonlinear prediction functions
   - (can be) a highly parallel network of logistic regression classifiers
   - discover useful hidden representations of the input

2. **Backpropagation…**
   - provides an efficient way to compute gradients
   - is a special case of reverse-mode automatic differentiation