

CS 360: Machine Learning

Prof. Sara Mathieson

Fall 2019



Admin

- **Midterm 2 due today at 6pm!**
- Office hours today: informal, I'm around this afternoon
- **Final project presentations:**
 - Wednesday Dec 18: 1-4pm (block out the entire time, but we may not need all of it)
 - Option to present last day of class (email me)

Outline for November 26

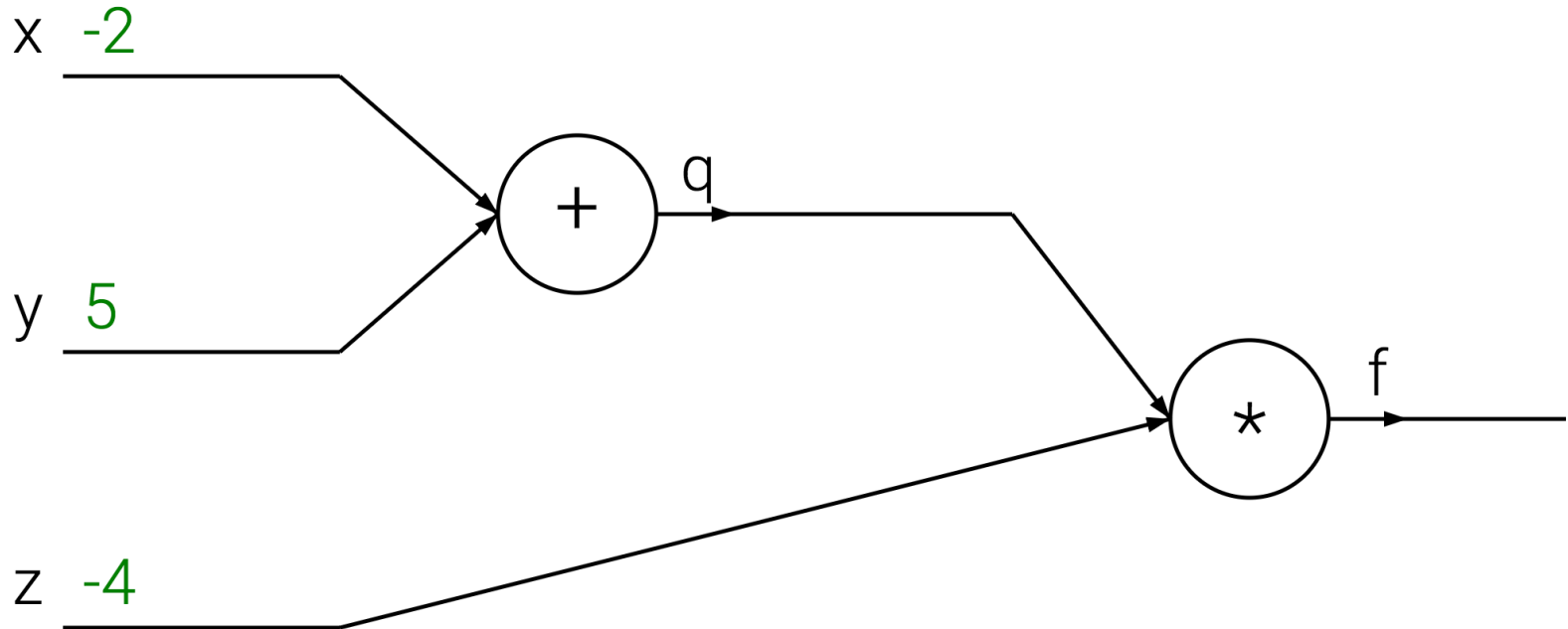
- Backpropagation
- Begin: unsupervised learning
- K-means
- Next week
 - Gaussian Mixture Models (GMM)
 - Principal Component Analysis (PCA)

Backpropagation

- *High-level goal:* we want to know how the output depends on the input
- *Issue:* network is very complicated and overall gradient may be difficult to compute
- *Idea:* use the chain rule to compute local gradients throughout the network
- *Takeaway:* nodes can know about their value and local gradient without knowing about the network they are imbedded in

Backpropagation: Example 1

Forward pass: compute values

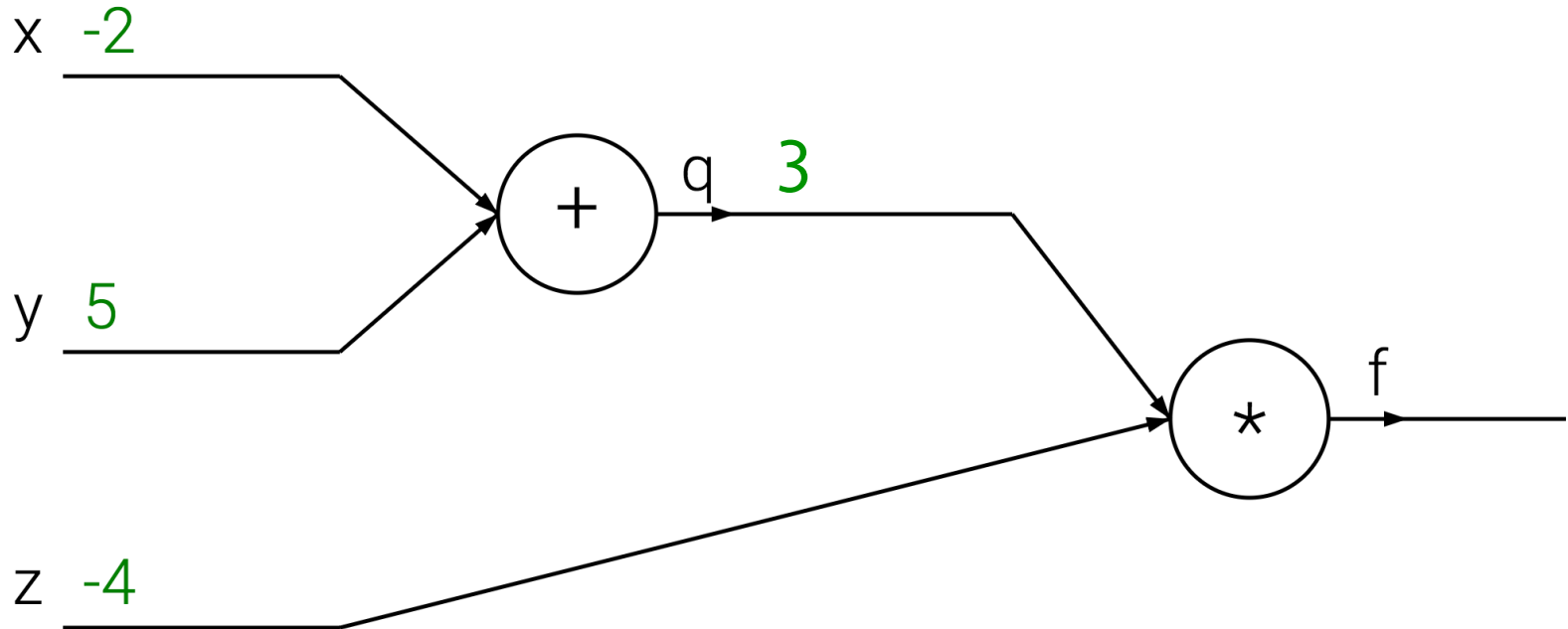


Example from:

<http://cs231n.github.io/optimization-2/>

Backpropagation: Example 1

Forward pass: compute values

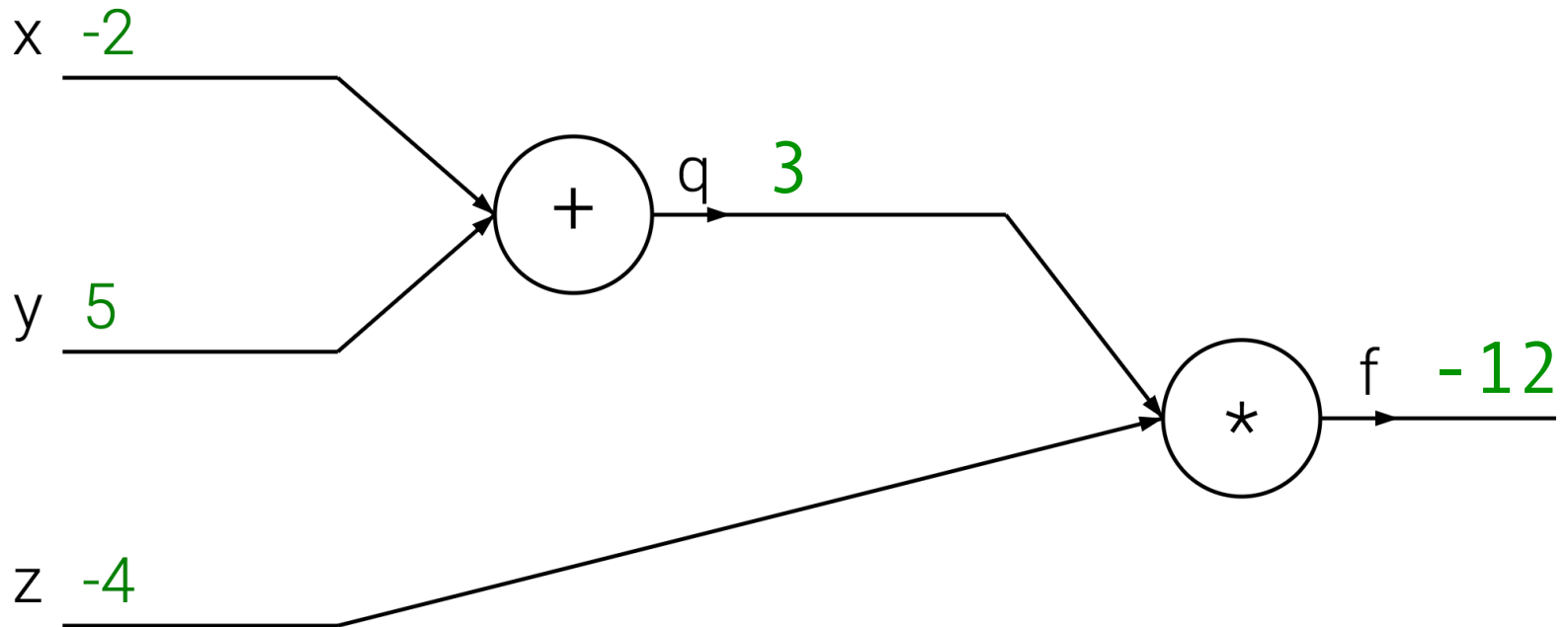


Example from:

<http://cs231n.github.io/optimization-2/>

Backpropagation: Example 1

Forward pass: compute values

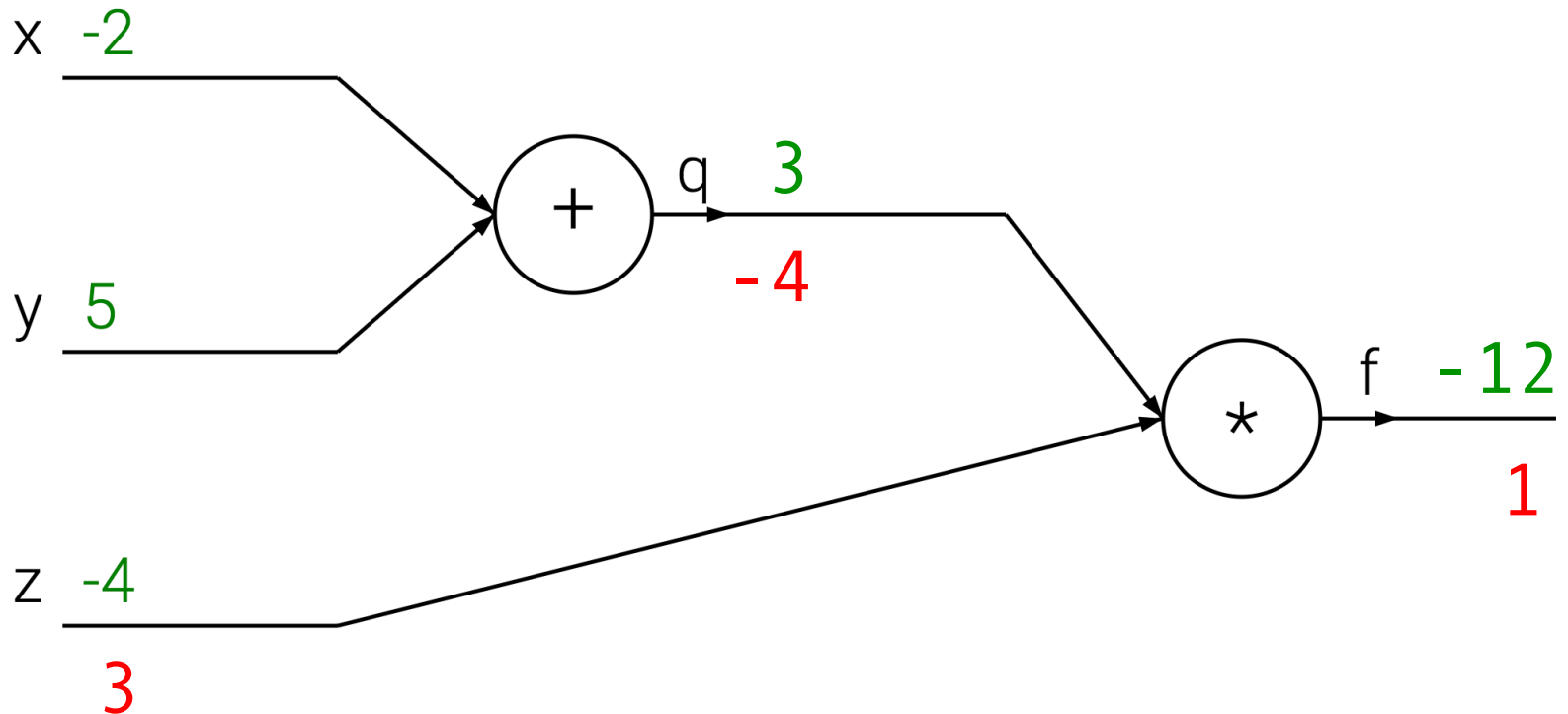


Example from:

<http://cs231n.github.io/optimization-2/>

Backpropagation: Example 1

Backward pass: compute local gradients

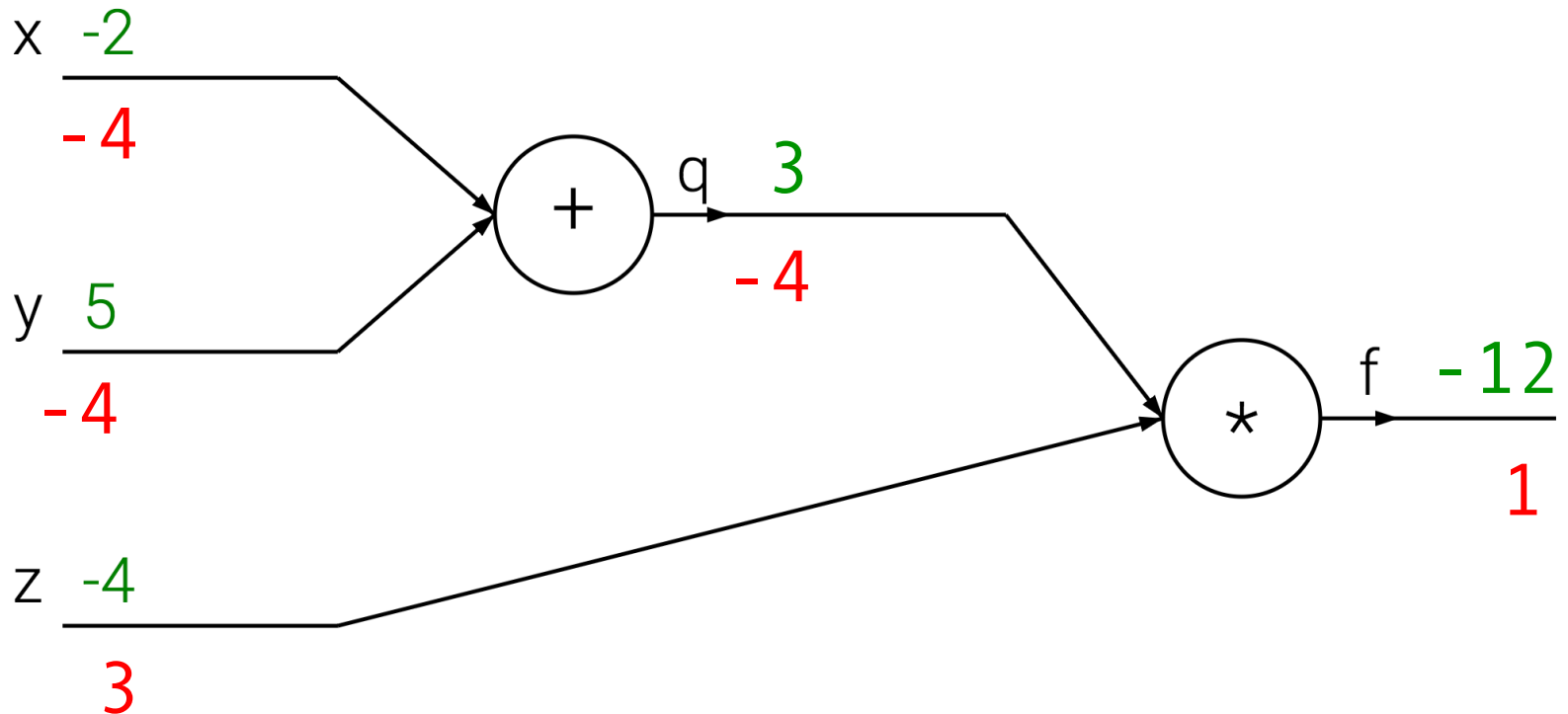


Example from:

<http://cs231n.github.io/optimization-2/>

Backpropagation: Example 1

Backward pass: compute local gradients



Example from:

<http://cs231n.github.io/optimization-2/>

Backpropagation

* Q3, Q5, Q6

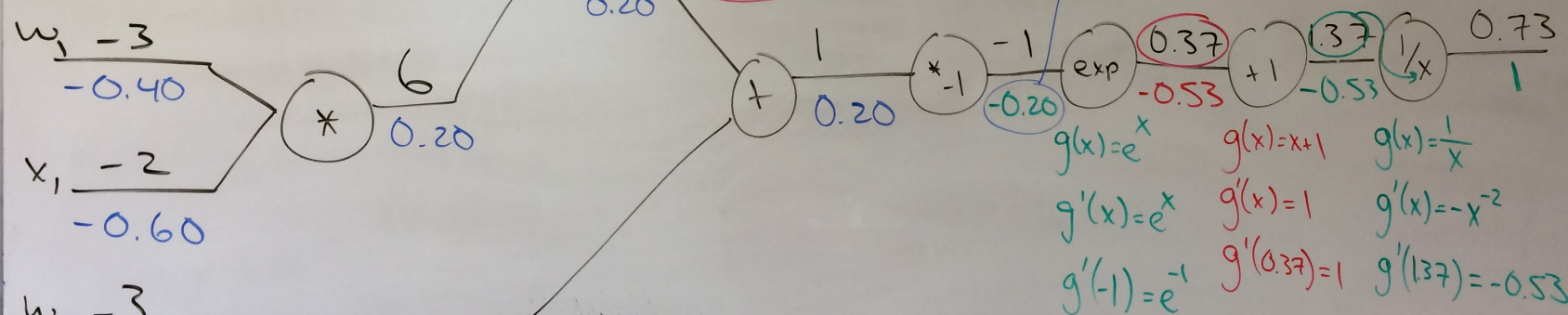
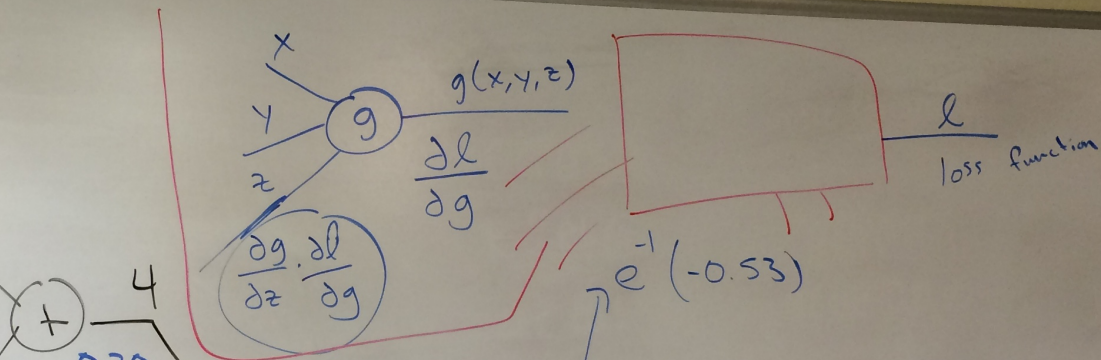
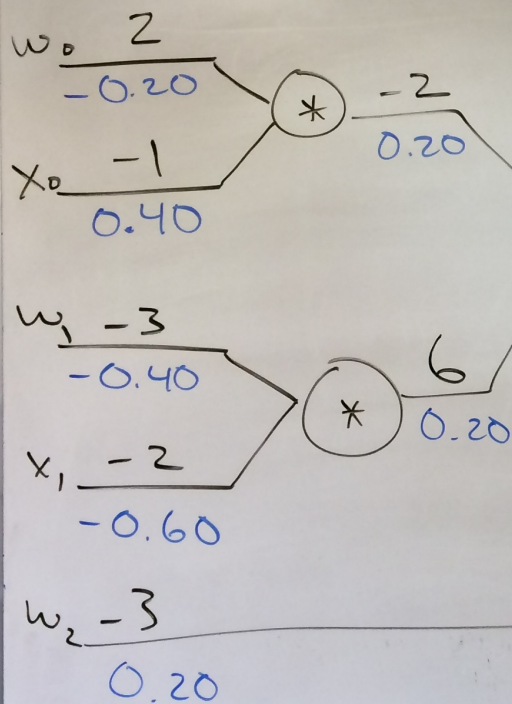
Idea: want to know how
output depends on input
(weights)

"loss" \rightarrow

$$f(x, y, z) = (\underbrace{x + y}_q) z = qz$$

chain rule

$$\frac{\partial f}{\partial x} = \underbrace{\frac{\partial f}{\partial q}}_z \cdot \underbrace{\frac{\partial q}{\partial x}}_1$$



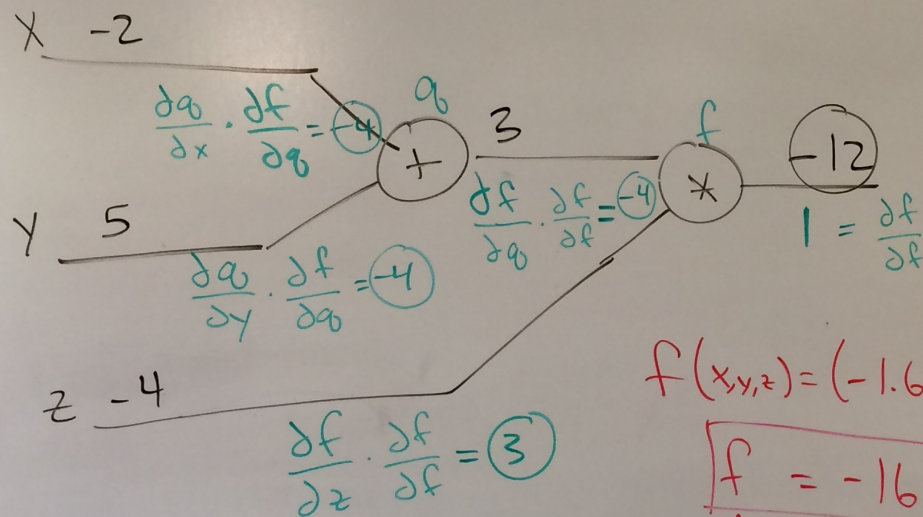
Derivatives of the activation functions:

- $g(x) = e^x$, $g'(x) = e^x$, $g'(-1) = e^{-1}$
- $g(x) = x+1$, $g'(x) = 1$, $g'(0.37) = 1$
- $g(x) = \frac{1}{x}$, $g'(x) = -\frac{1}{x^2}$, $g'(1.37) = -0.53$

Example 1

$$f(x, y, z) = (x + y)z$$

want to
minimize



$$f(x, y, z) = (-1.6 + 5.4)(-4.3)$$

$$f = -16.34$$

☆ lower than
before!

update using gradient descent

$$x \leftarrow x - \alpha \frac{\partial f}{\partial x}$$

$$x \leftarrow -2 - 0.1(-4)$$

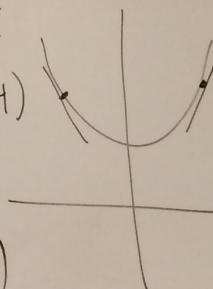
$$x \leftarrow -1.6$$

$$y \leftarrow 5 - 0.1(-4)$$

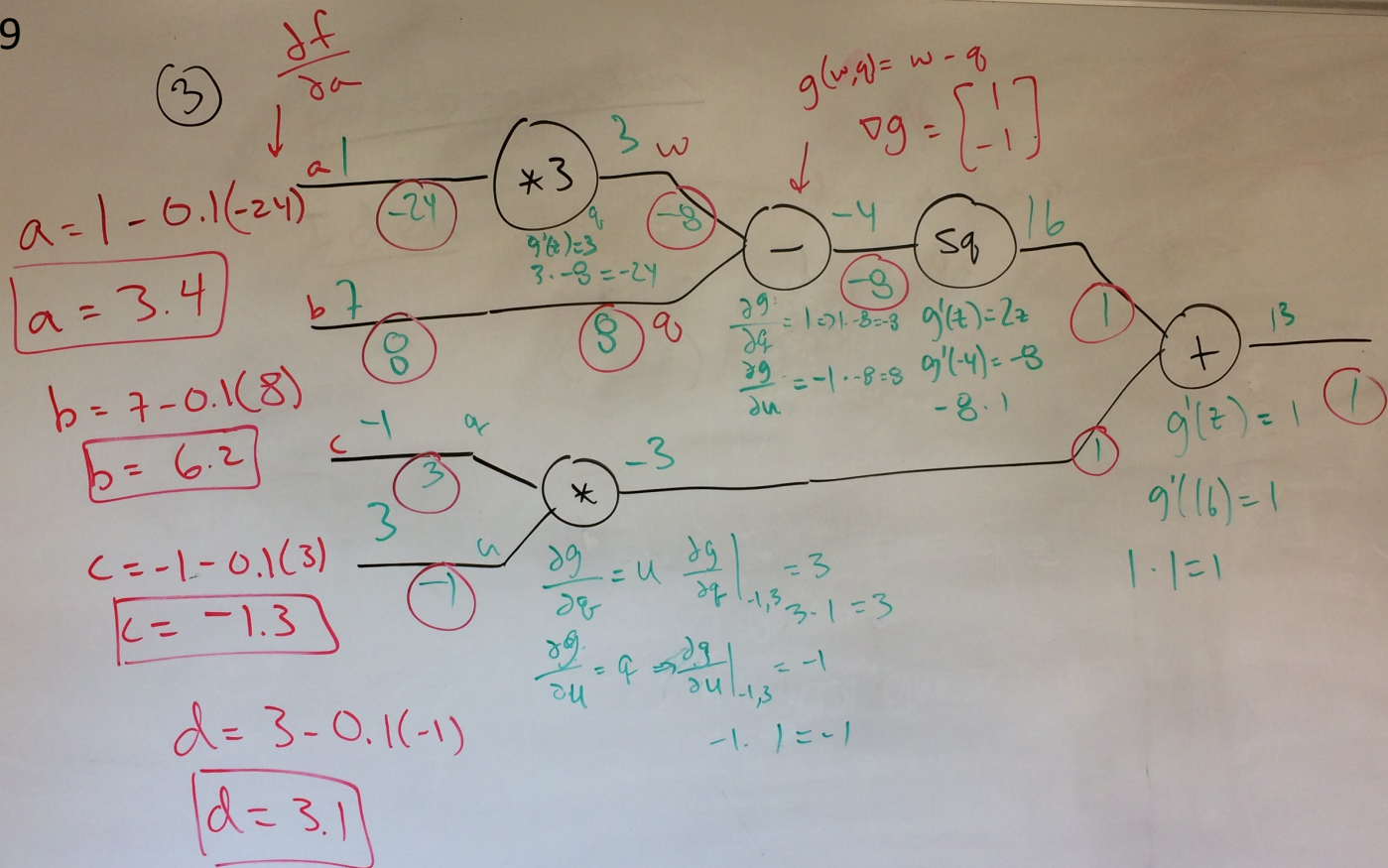
$$y \leftarrow 5.4$$

$$z \leftarrow -4 - 0.1(3)$$

$$z \leftarrow -4.3$$

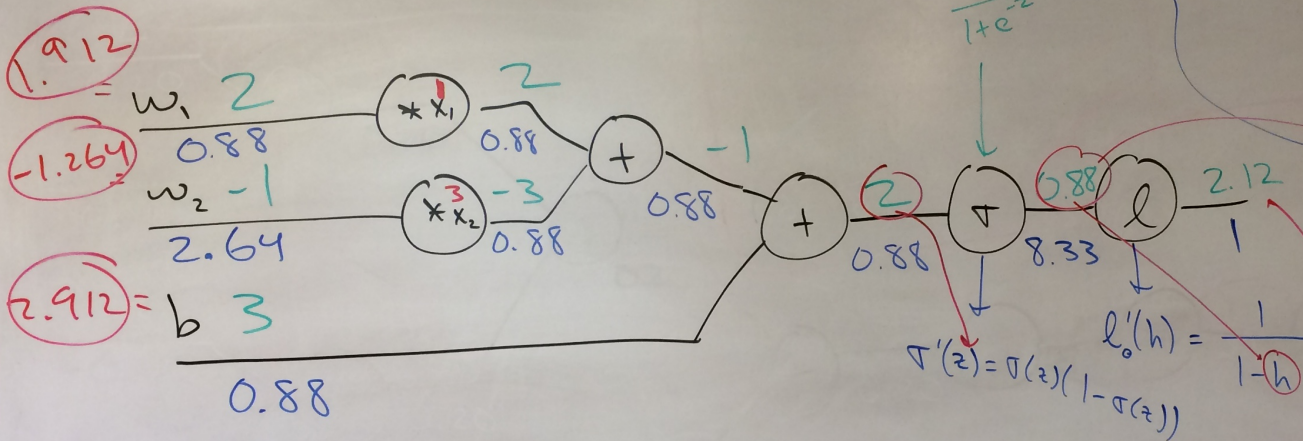


Handout 19



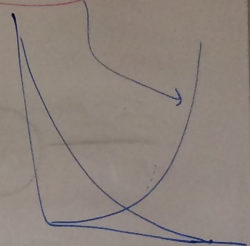
④

$$\vec{x} = (x_1, x_2) = (1, 3)$$



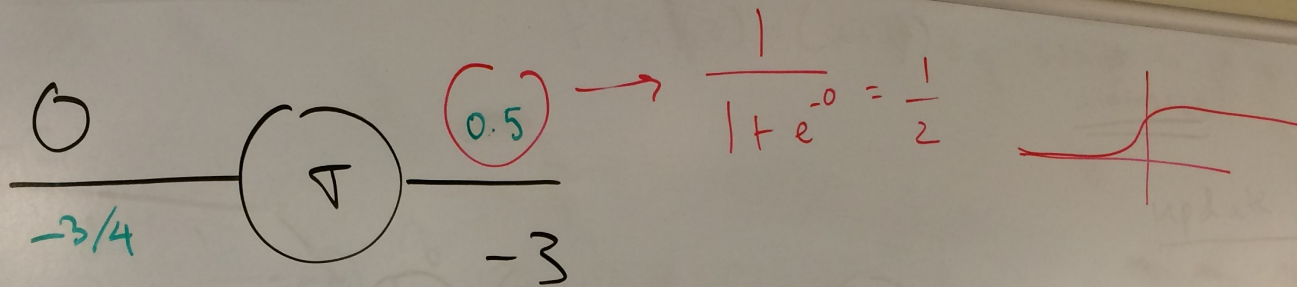
$$\ell_y(h) = -y \log h - (1-y) \log (1-h)$$

$$\ell_o(h) = -\log(1-h)$$



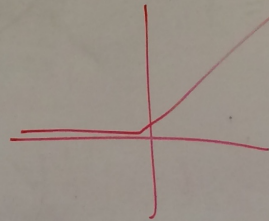
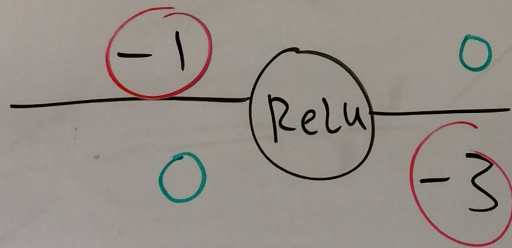
new loss
= 1.33

(5.)



$$\sigma'(0) = \underbrace{\sigma(0)}_{\frac{1}{2}} \left(1 - \underbrace{\sigma(0)}_{\frac{1}{2}} \right) = \frac{1}{4}$$

(6.)



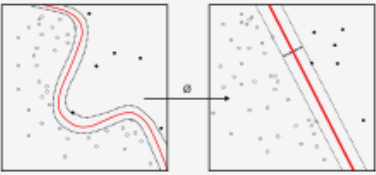
Outline for November 26

- Backpropagation
- **Begin: unsupervised learning**
- K-means
- Next week
 - Gaussian Mixture Models (GMM)
 - Principal Component Analysis (PCA)

Supervised Learning:

makes use of examples where we know the underlying “truth” (label/output)

Machine learning and data mining



Problems [show]

Supervised learning [hide]
(classification • regression)

Decision trees • Ensembles (Bagging, Boosting, Random forest) • *k*-NN • Linear regression • Naive Bayes • Neural networks • Logistic regression • Perceptron • Relevance vector machine (RVM) • Support vector machine (SVM)

Clustering [hide]
BIRCH • Hierarchical • *k*-means • Expectation-maximization (EM) • DBSCAN • OPTICS • Mean-shift

Dimensionality reduction [hide]
Factor analysis • CCA • ICA • LDA • NMF • PCA • t-SNE

Structured prediction [hide]
Graphical models (Bayes net, CRF, HMM)


Anomaly detection [hide]
k-NN • Local outlier factor

Neural nets [hide]
Autoencoder • Deep learning • Multilayer perceptron • RNN • Restricted Boltzmann machine • SOM • Convolutional neural network

Reinforcement Learning [hide]
Q-Learning • SARSA • Temporal Difference (TD)

Theory [show]

Machine learning venues [show]

 **Machine learning portal**

V • T • E

Unsupervised Learning:

Learn underlying structure or features without labeled training data

Unsupervised learning: 3 main areas

- 1) Clustering: group data points into clusters based on features only
- 2) Dimensionality reduction: remove feature correlation, compress data, visualize data
- 3) Structured prediction: model latent variables (example: Hidden Markov Models)

Unsupervised learning examples from biology: clustering

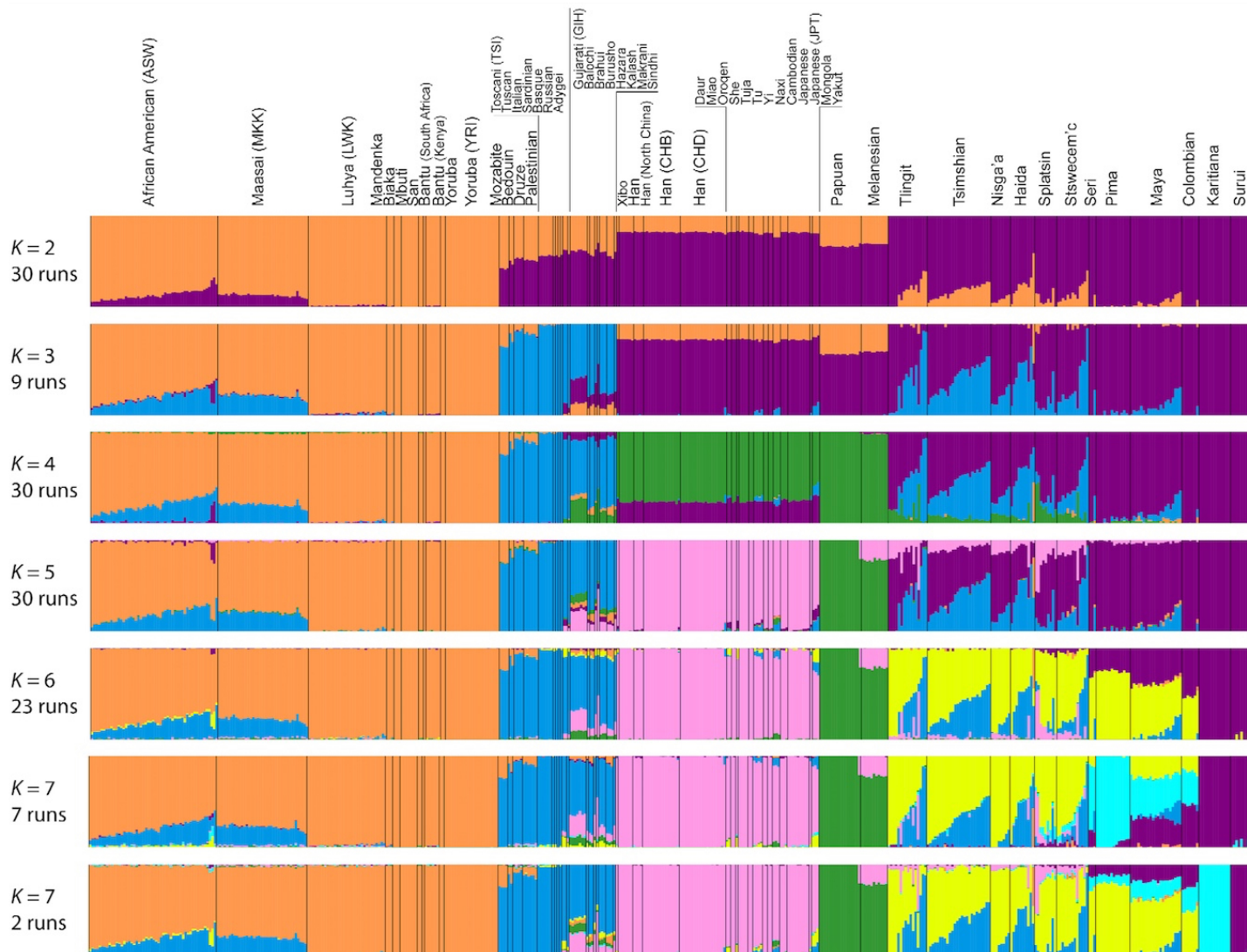
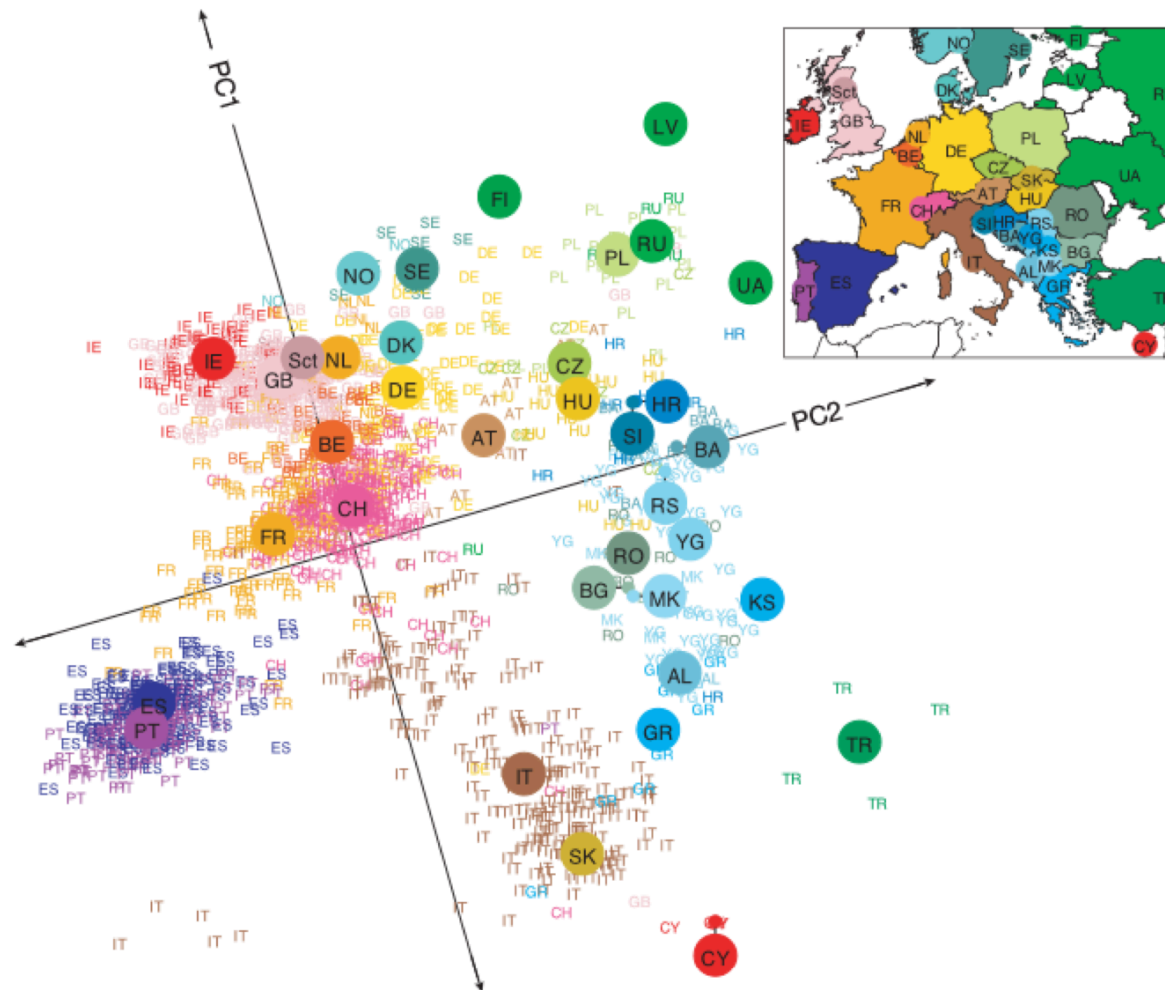
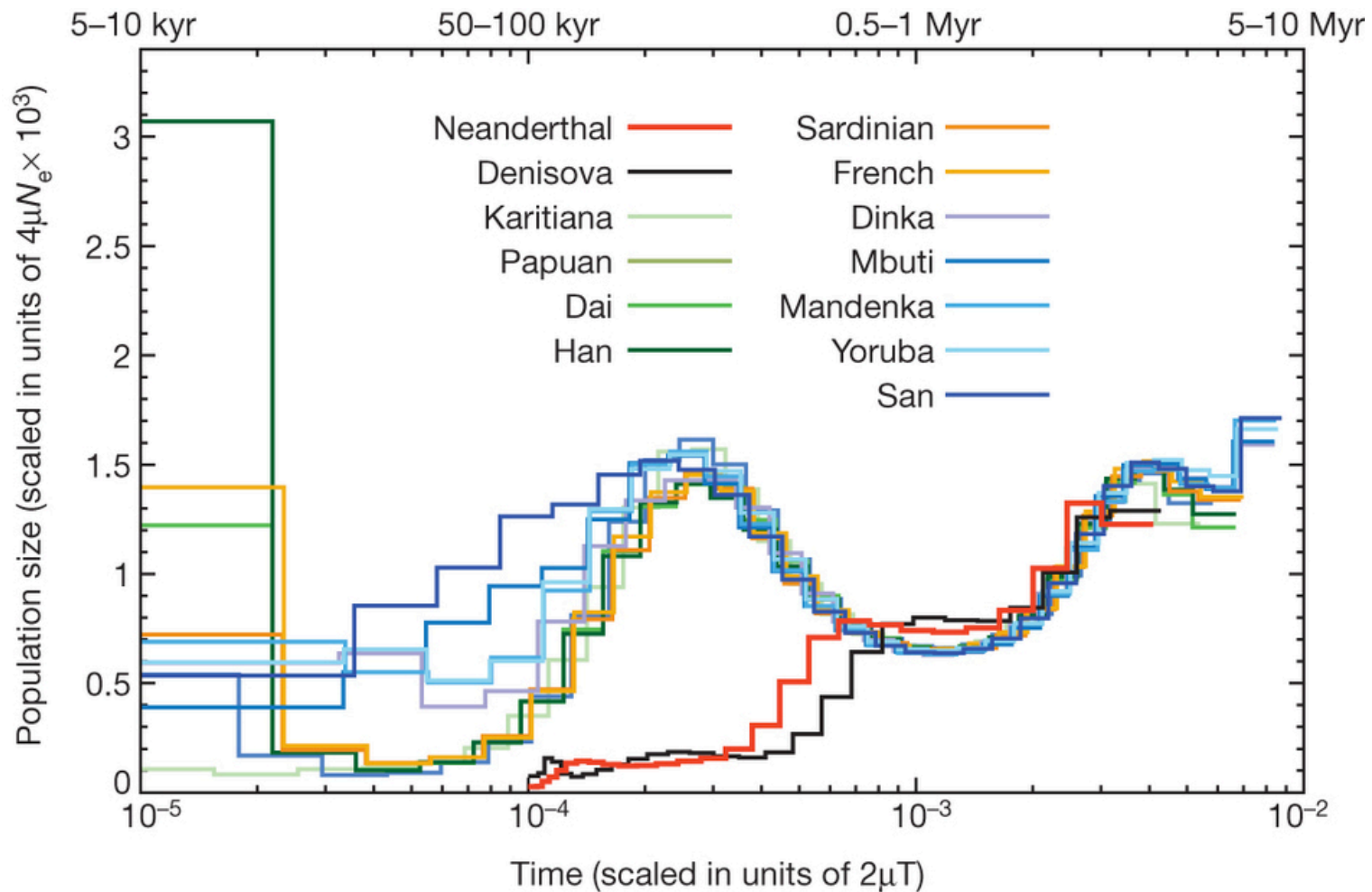


Figure: German Dziebel

Unsupervised learning examples from biology: structured prediction



Unsupervised learning examples from biology: structured prediction



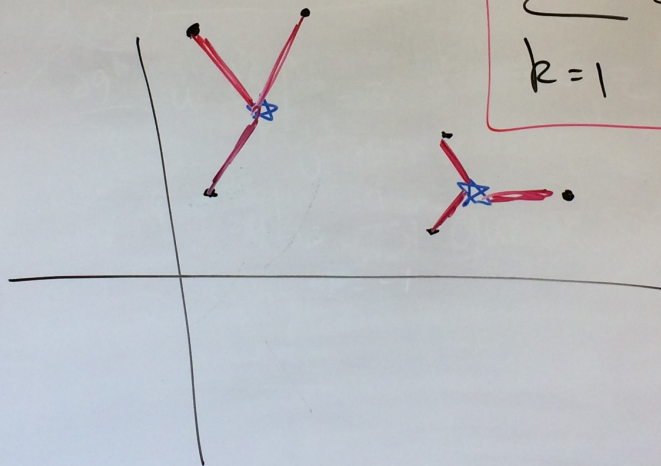
Outline for November 26

- Backpropagation
- Begin: unsupervised learning
- **K-means**
- Next week
 - Gaussian Mixture Models (GMM)
 - Principal Component Analysis (PCA)

Clustering Goals

* learn about the structure of the data.

* cluster new data ("test", prediction)



Goal find $\{\mathcal{C}_1, \mathcal{C}_2, \dots, \mathcal{C}_K\} = \mathcal{C}$
K clusters

s.t.

minimize

$$\sum_{k=1}^K \sum_{i \in \mathcal{C}_k} \|\bar{x}_i - \bar{\mu}_k\|^2 = J(\mathcal{C})$$

within cluster sum of square
WCSS

mean of cluster k

* cluster mean

WCSS

K-means

① initialization

choose means (centers)
from among the training data.

iterate $\vec{\mu}_1^{(1)}, \vec{\mu}_2^{(1)}, \dots, \vec{\mu}_K^{(1)}$

① E-step Assignment: assign
each datapoint to the cluster
with the closest mean.

② M-step Update: recompute each
mean as the average of all
cluster members Not data points

