# CS/ECE 148 –
# **Data Science Fundamentals**

## NNs, Unsupervised & Semi-supervised Learning

**UCLA Computer Science**

# Outline

## Regularization of NN

- Norm Penalties
- Early Stopping
- Data Augmentation
- Sparse Representation
- **Dropout**

## Optimization

- **Challenges in Optimization**
- Momentum (next lectures)
- Adaptive Learning Rate (next lectures)
- Parameter Initialization (next lectures)
- Batch Normalization (next lectures)

# Outline

## Unsupervised Learning

- K-means
- Mean Shift
- Hierarchical Clustering
- DBSCAN
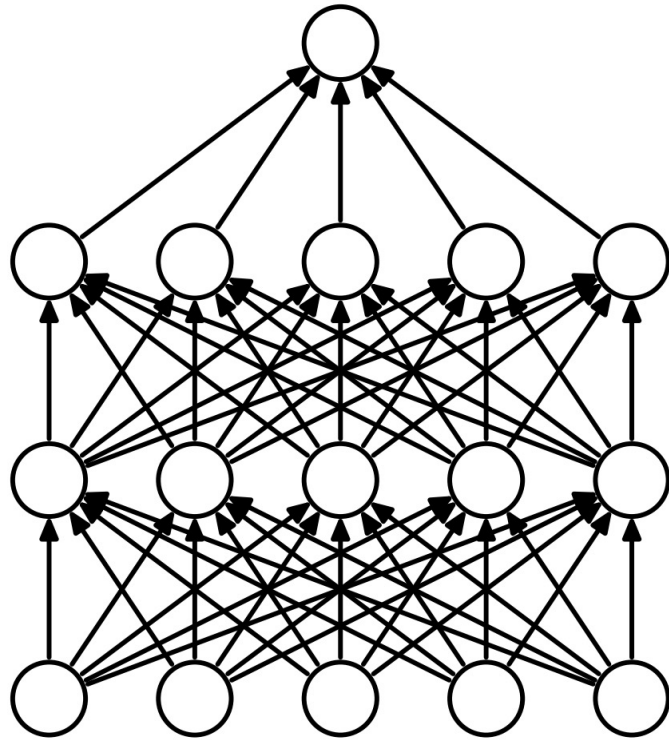
## Semi-supervised Learning

- Self Training
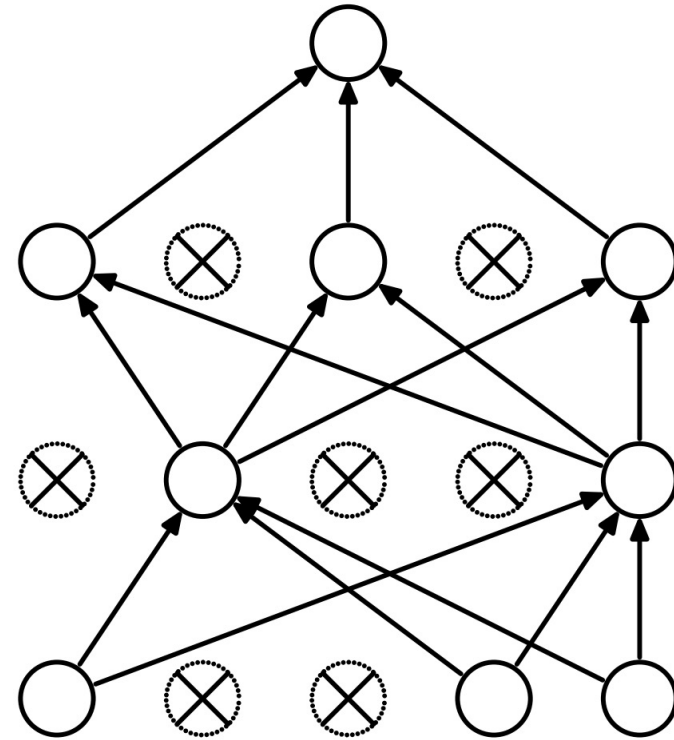
# Outline

**Regularization of NN**

- Norm Penalties
- Early Stopping
- Data Augmentation
- Sparse Representation
- **Dropout**

# Dropout

- Randomly set some neurons and their connections to zero (i.e. "dropped")
- Prevent overfitting by reducing co-adaptation of neurons
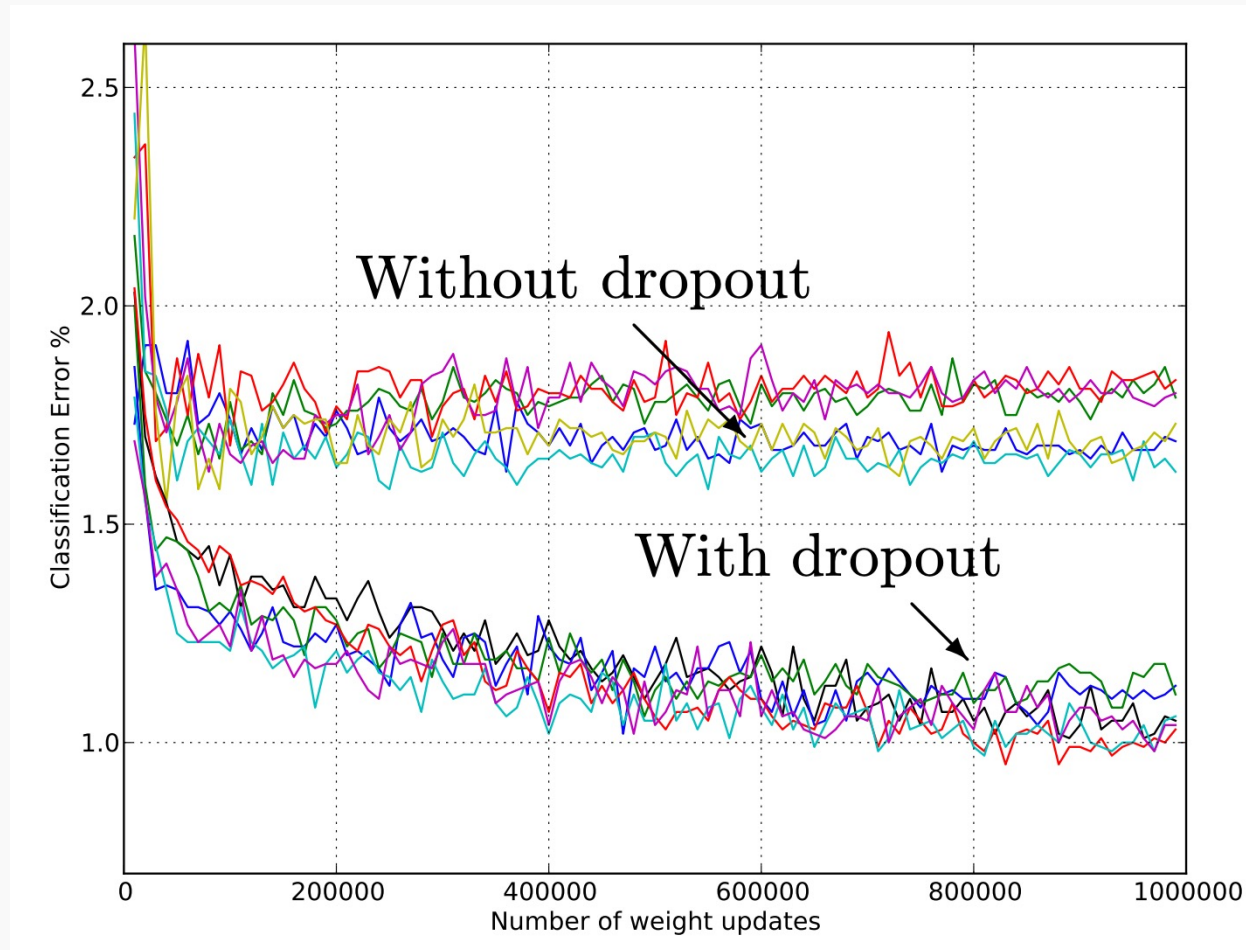- Like training many random sub-networks



(a) Standard Neural Net        (b) After applying dropout.

# Dropout

- Widely used and highly effective
- Proposed as an alternative to ensembling, which is too expensive for neural nets



Test error for different architectures with and without dropout. The networks have 2 to 4 hidden layers each with 1024 to 2048 units.

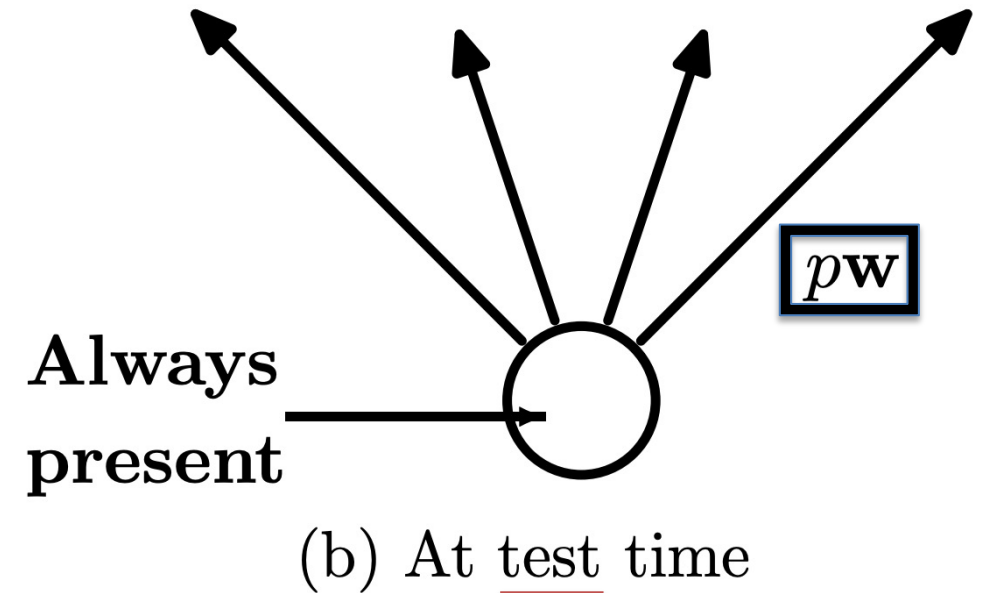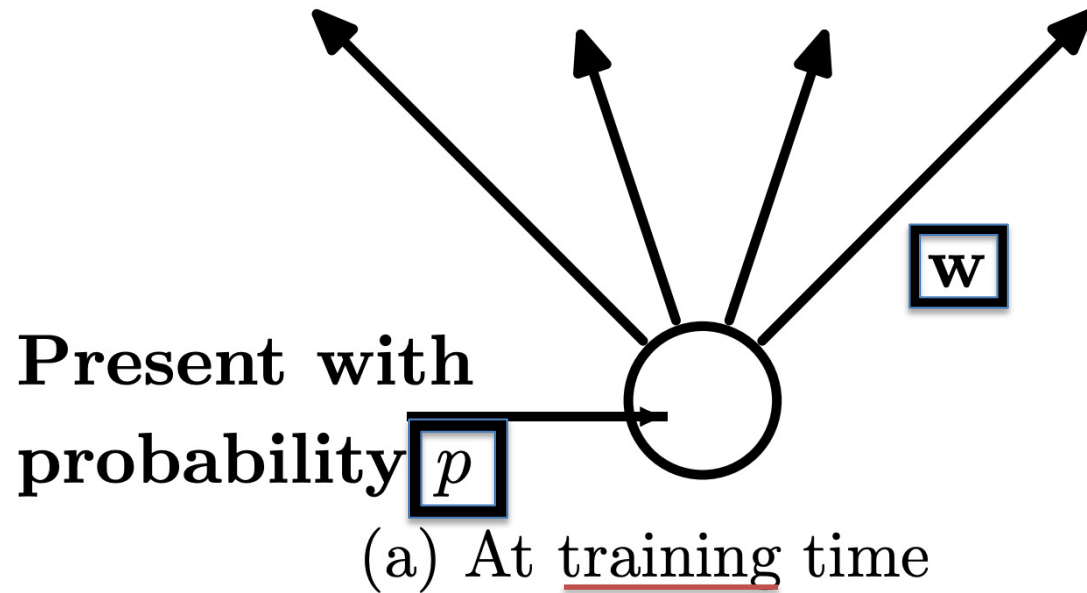http://jmlr.org/papers/volume15/srivastava14a/srivastava14a.pdf

# Dropout: Stochastic GD

For each new example/mini-batch:

- Randomly sample a binary mask $\mu$ independently, where $\mu_i$ indicates if input/hidden node $i$ is included

- Multiply output of node $i$ with $\mu_i$, and perform gradient update

Typically, an input node is **included** with **prob=0.8**, hidden node with **prob=0.5.**

# Dropout: Weight Scaling

- We can think of dropout as training many of sub-networks
- At **test time**, we can "aggregate" over these sub-networks by **reducing connection weights in proportion to dropout probability, _p_**



(a) At training time

(b) At test time

# Outline

## Optimization

- **Challenges in Optimization**
- Momentum (next lectures)
- Adaptive Learning Rate (next lectures)
- Parameter Initialization (next lectures)
- Batch Normalization (next lectures)

# Learning vs. Optimization

Goal of learning: minimize generalization error, or the loss function

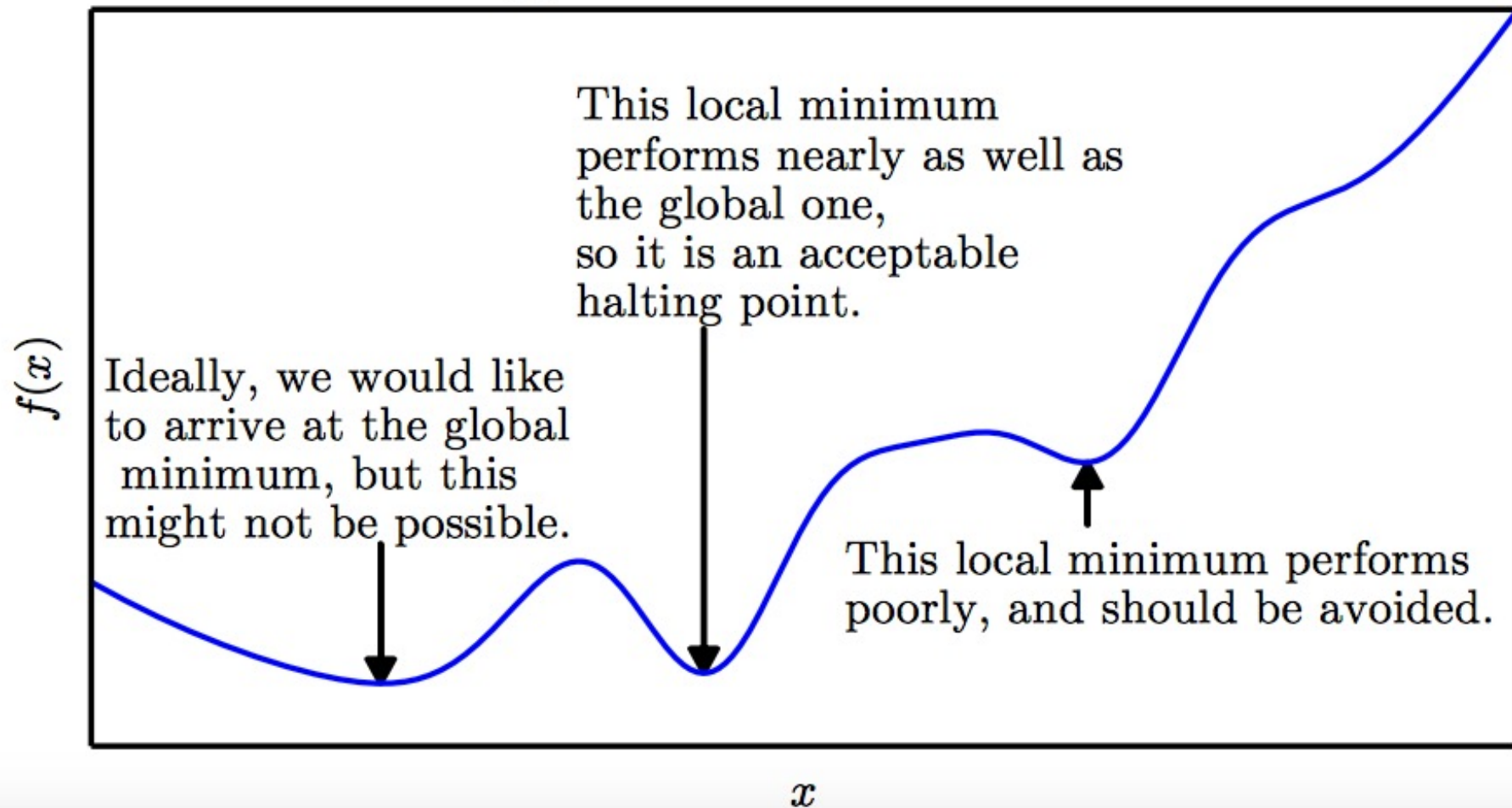$$\mathcal{L}(W) = \mathbb{E}_{(x,y) \sim p_{data}} \Big[ L(f(x, W), y) \Big]$$

$f$ is the neural network

In practice, empirical risk minimization:

$$\mathcal{L}(W) = \sum_i \Big[ L(f(x_i; W), y_i) \Big]$$

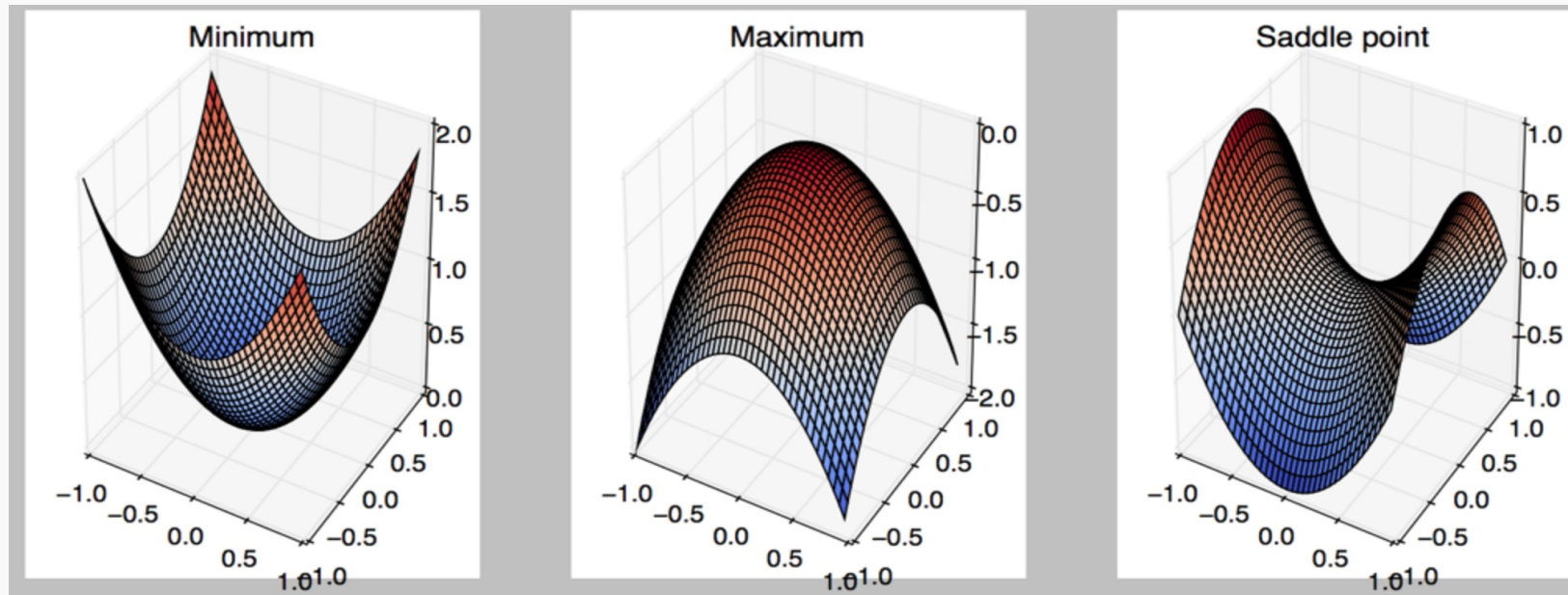Quantity optimized different from the quantity we care about

This local minimum
performs nearly as well as
the global one,
so it is an acceptable
halting point.

Ideally, we would like
to arrive at the global
minimum, but this
might not be possible.

This local minimum performs
poorly, and should be avoided.

$f(x)$

$x$

# Critical Points
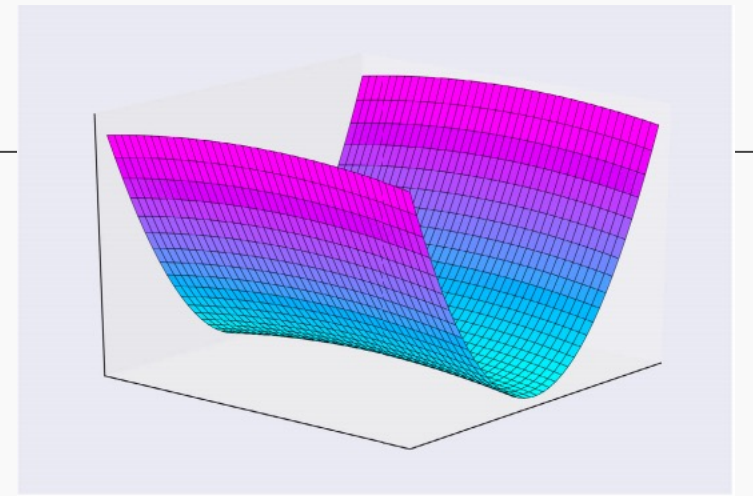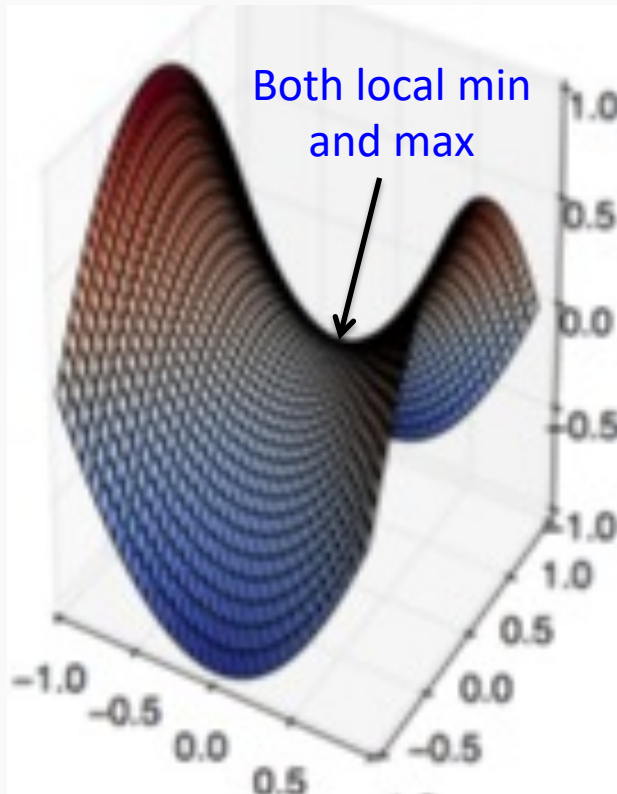
Points with zero gradient

2nd-derivate (Hessian) determines curvature

# Local Minima

Old view: local minima is major problem in neural network training

Recent view:

- For sufficiently large neural networks, <span style="color:blue">most local minima incur low cost</span>

- Not important to find true global minimum

# Saddle Points





Both local min and max

Recent studies indicate that in high dim, saddle points are more likely than local min
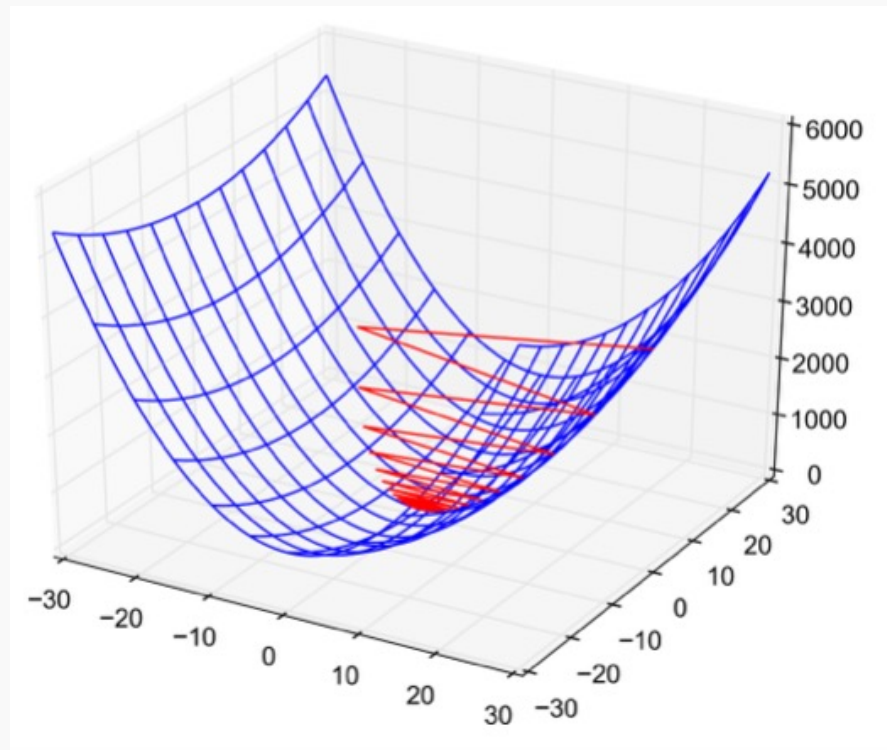
Gradient can be very small near saddle points

Goodfellow et al. (2016)

# Poor Conditioning

Poorly conditioned Hessian matrix

– High curvature: small steps leads to huge increase
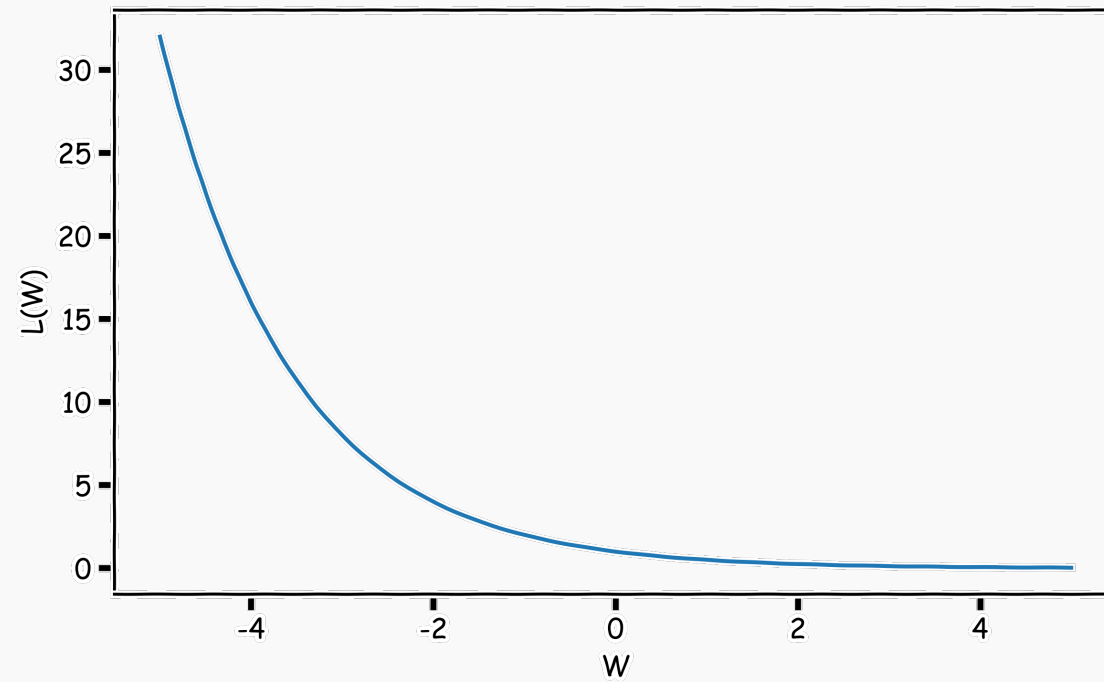
Learning is slow despite strong gradients

Oscillations slow down progress

# No Critical Points

Some cost functions do not have critical points. In particular classification.
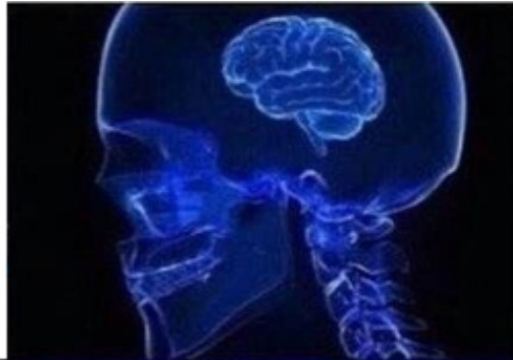
**WHY?**

# Optimization Challenges

We'll discuss some solutions in the next lectures

- Momentum (later)
- Adaptive Learning Rate (later)
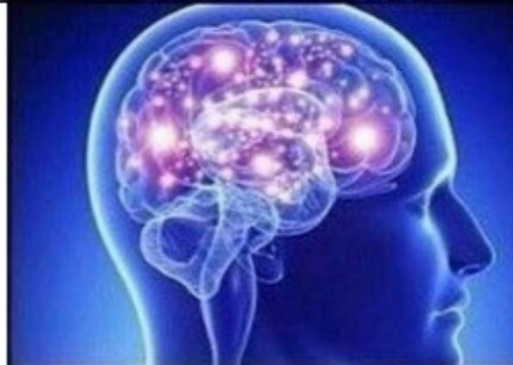- Parameter Initialization (later)
- Batch Normalization (later)

19

# Unsupervised Learning

# Unsupervised Learning

- K-means

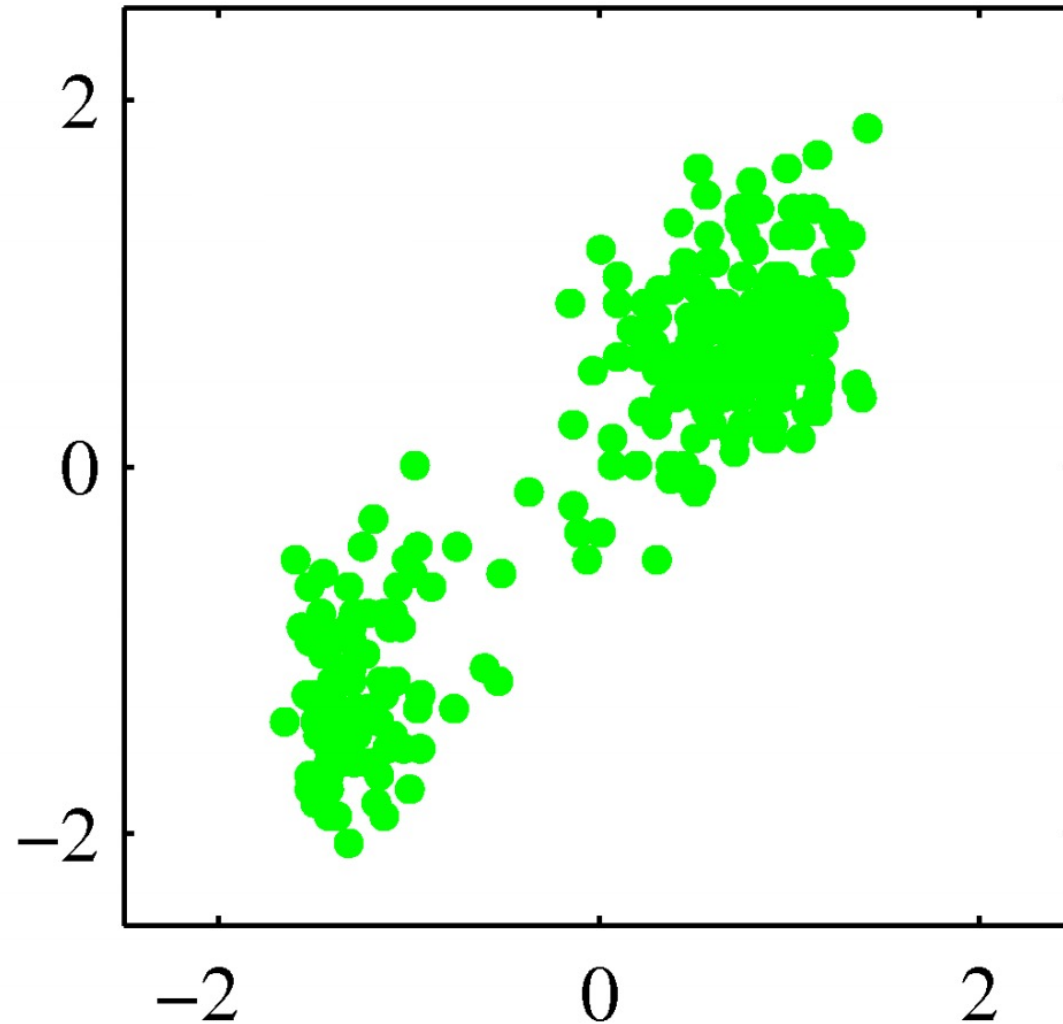- Mean-shift

- Hierarchical Clustering

- DBSCAN

- Applications

# Unsupervised Setting



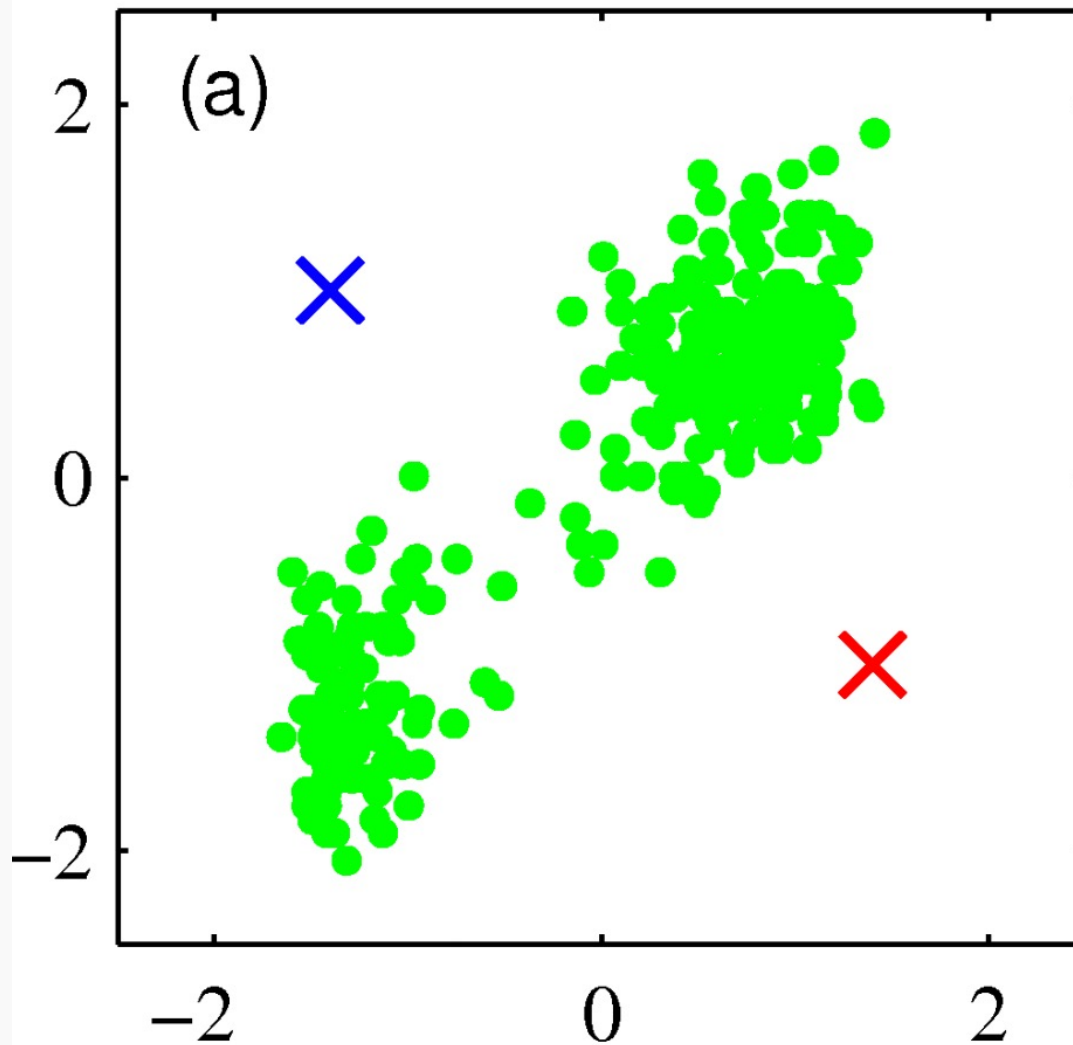Bishop, "Pattern Recognition and Machine Learning", Springer, 2006

# K-means – Algorithm

Initialization:

– choose K random positions

– assign cluster centers $\mu^j$ to these positions

# K-means



Bishop, "Pattern Recognition and Machine Learning", Springer, 2006
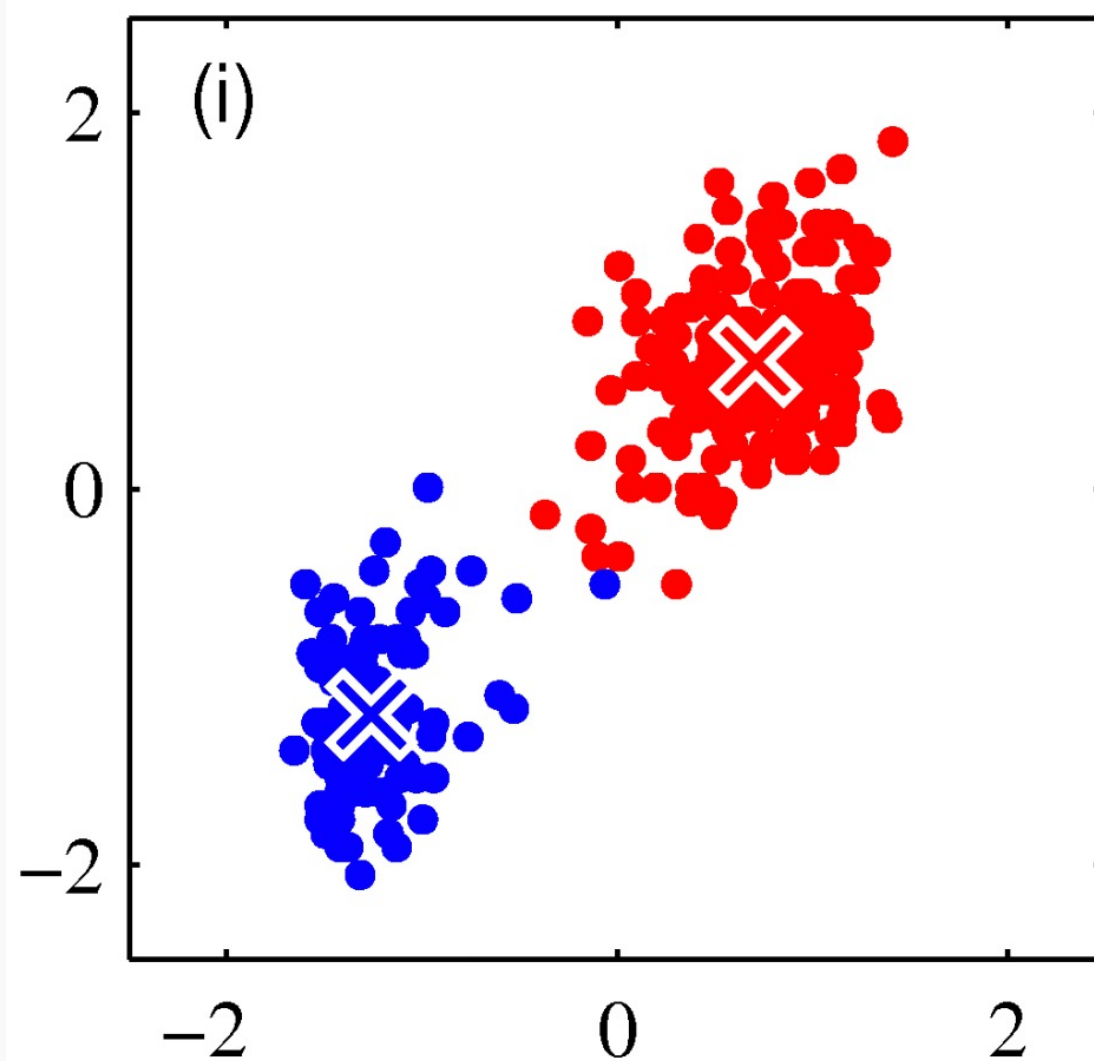
# K-means

Until Convergence:

– Compute distances $\left\|x^{(i)} - \mu^{(i)}\right\|$

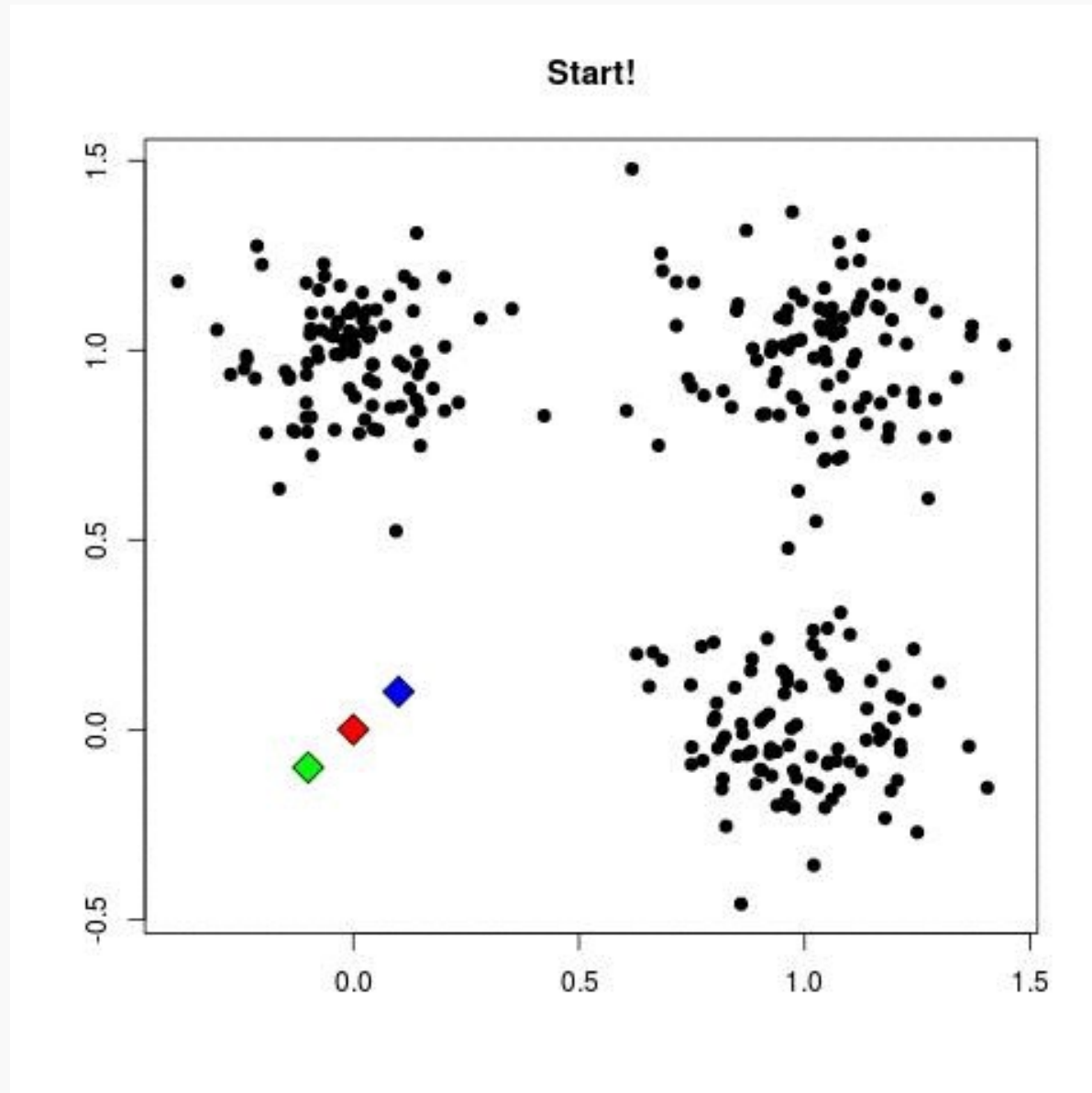– Assign points to nearest cluster center

– Update Cluster centers:

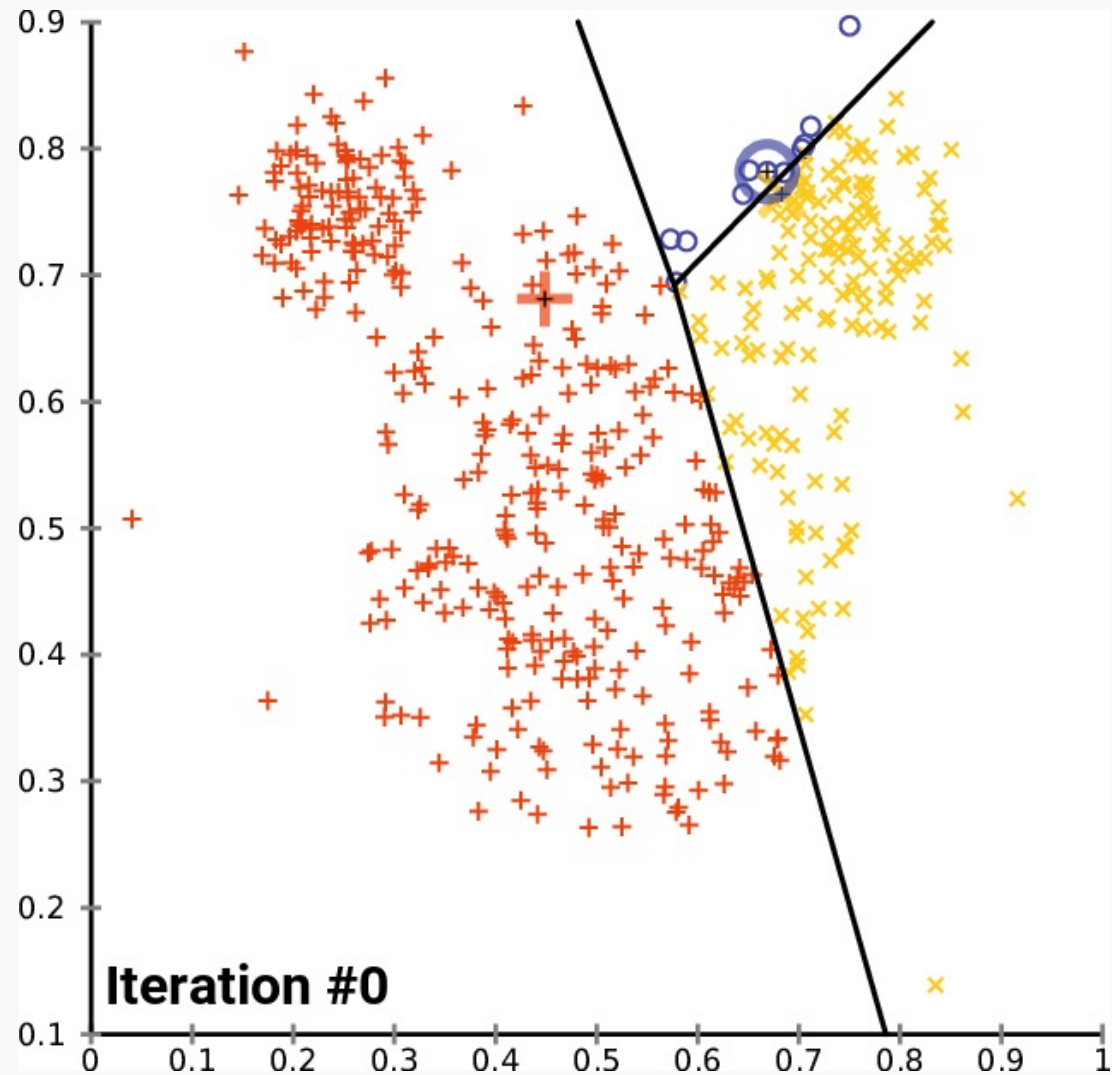$$\mu^{(j)} = \frac{1}{N_j} \sum_{x_i \in C_j} x_i$$

# K-means



(i)

Bishop, "Pattern Recognition and Machine Learning", Springer, 2006

# K-means



Start!

Iteration #0

# K-means Summary

- Guaranteed to converge

- Result depends on initialization

- Number of clusters is important

- Sensitive to outliers
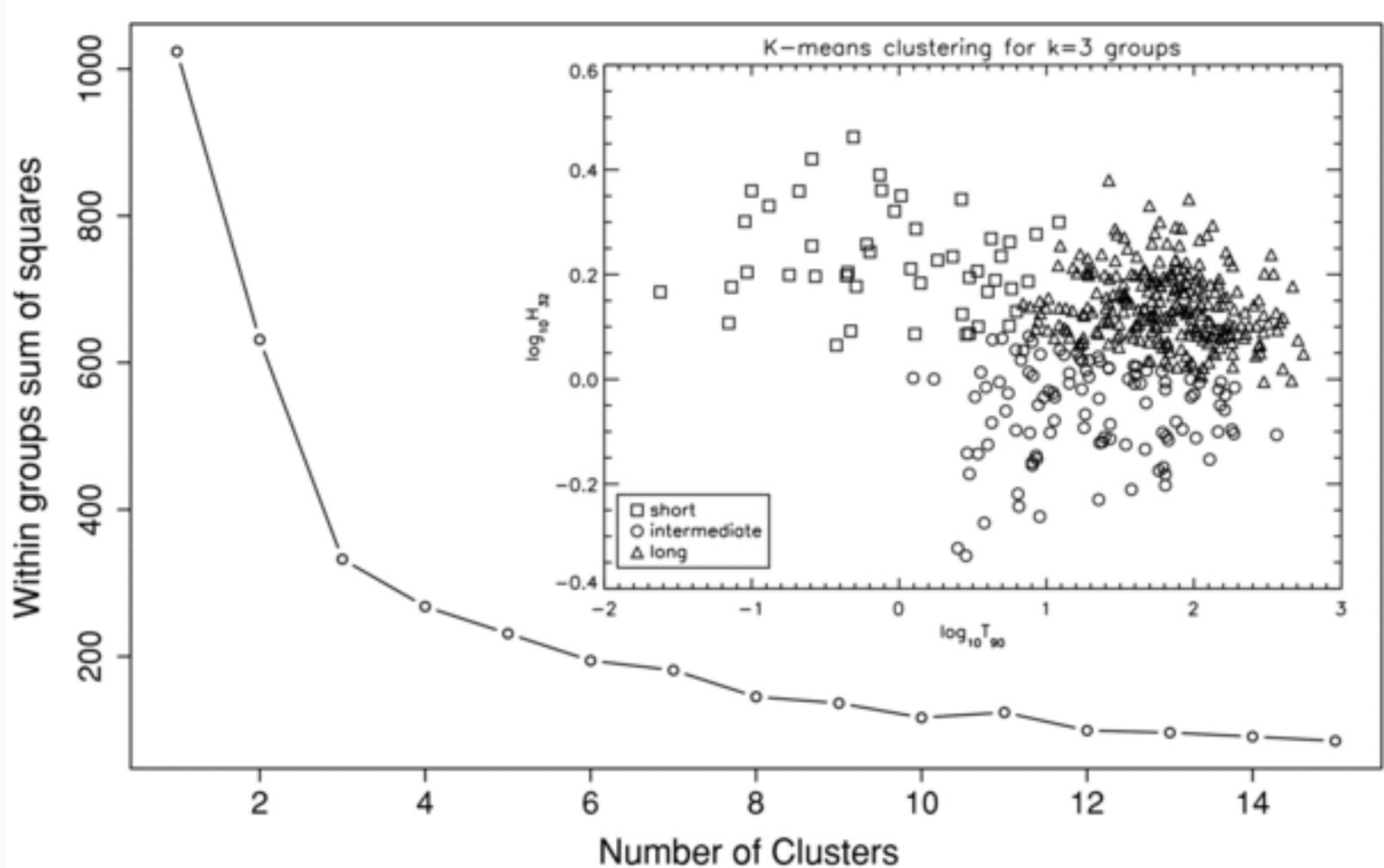  - Use median instead of mean for updates

# Initialization Methods

- Random Positions

- Random data points as Centers

- Random Cluster assignment to data points


- Start several times

# How to find K

- Extreme cases:
  - K=1
  - K=N

- Choose K such that increasing it does not model the data much better.

# "Knee" or "Elbow" method

# Cross Validation

- Use this if you want to apply your clustering solution to new unseen data


- Partition data into n folds

- Cluster on n-1 folds

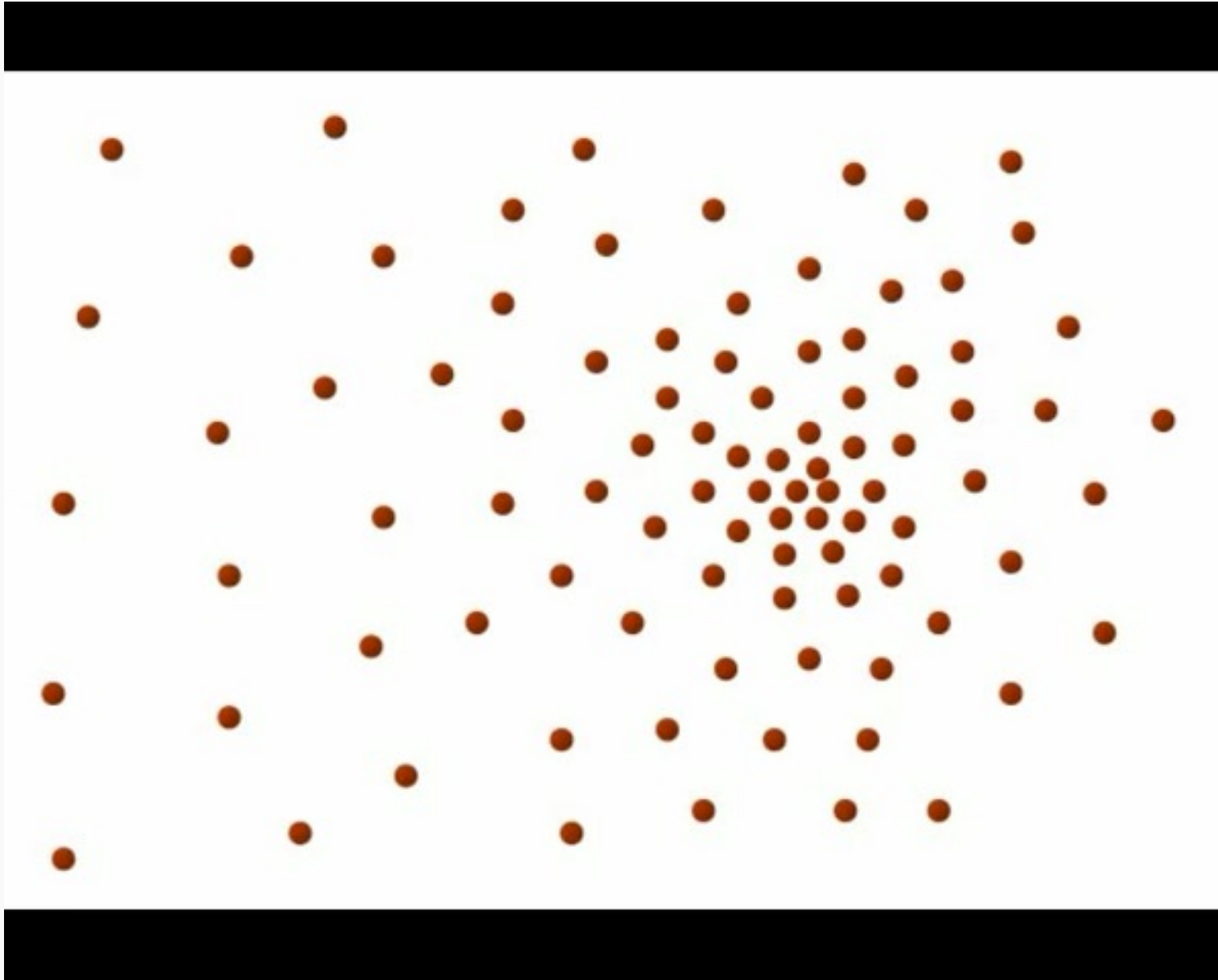- Compute sum of squared distances to centroids for validation set

# Getting Rid of K

- Having to specify K is annoying
- Can we do without?

# Mean Shift

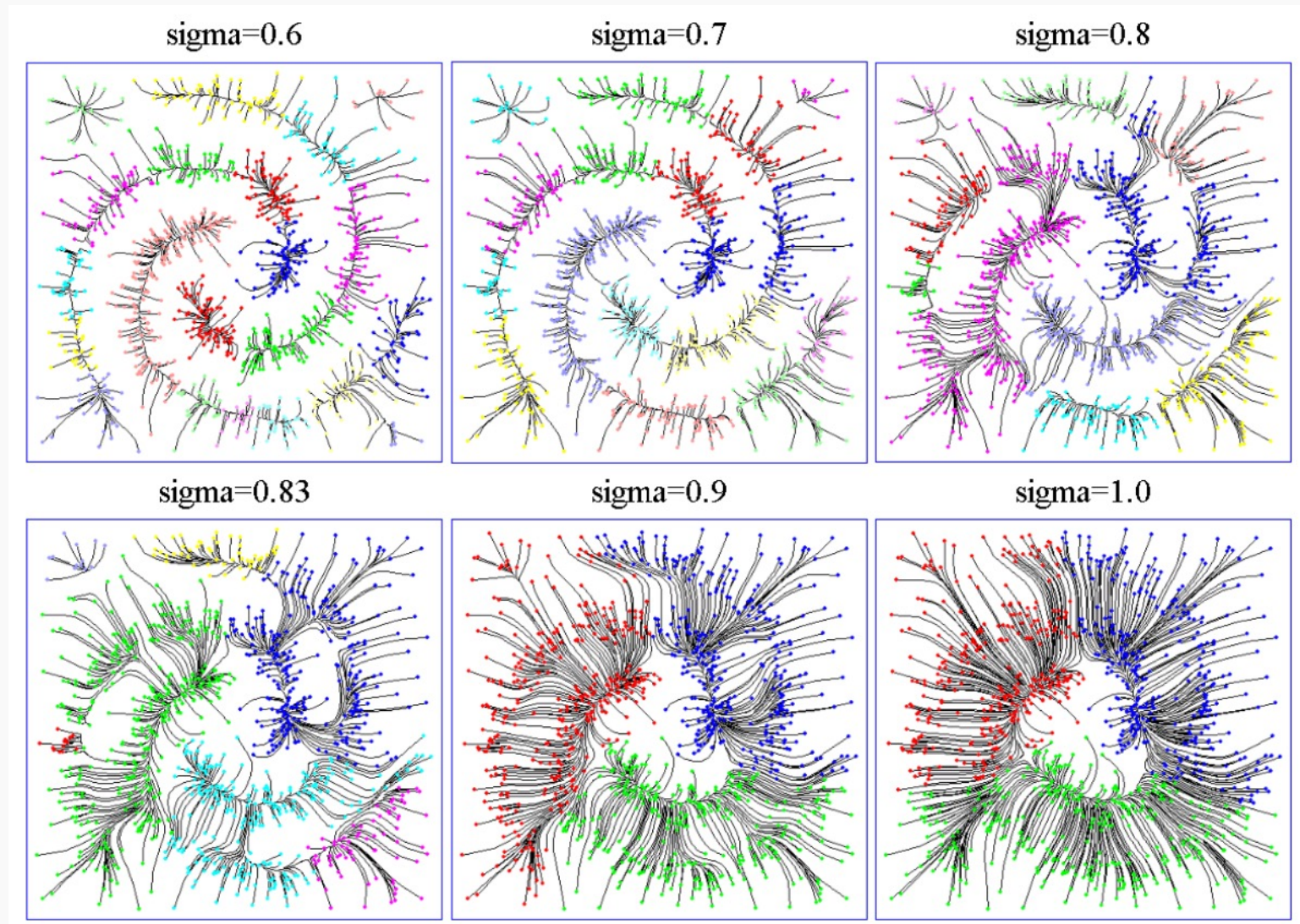1. Put a window around each point

2. Compute mean of points in the frame.

3. Shift the window to the mean

4. Repeat until convergence

https://www.youtube.com/watch?v=kmaQAsotT9s

# Mean Shift



Fischer et al., "Clustering with the Connectivity Kernel", NIPS (2003)

# Mean Shift Summary

- Does not need to know number of clusters

- Can handle arbitrary shaped clusters

- Robust to initialization

- Needs bandwidth parameter (window size)

- Computationally expensive

- Very good article:

http://saravananthirumuruganathan.wordpress.com/2010/04/01/introduction-to-mean-shift-algorithm/