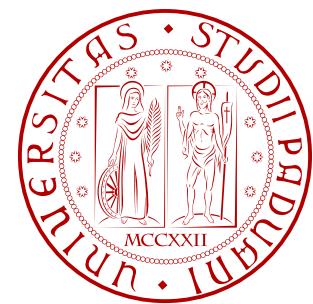


COMPETITIVE & UNSUPERVISED LEARNING FOR VECTOR QUANTIZATION

Michele Rossi
rossi@dei.unipd.it

Dept. of Information Engineering
University of Padova, IT



Outline (1/2)

- Vector quantization - definition
- Self Organizing Maps (SOM)
 - SOM & the human brain
 - SOM's goal
 - Input vectors & synaptic weights
 - 2D neural network example
- Learning steps
 - Competition
 - Cooperation
 - Synaptic adaptation
- Learning phases & experiments
- Pros & cons of SOM

Outline (2/2)

- Some useful concepts
 - Delaunay triangulation
 - Voronoi diagrams
- Topology preserving feature maps
 - Hebbian learning / Competitive Hebbian Learning (CHL)
 - Theorem “CHL preserves feature maps”
- Growing Neural Gas (GNG) networks
 - Algorithm
 - Remarks
 - Example results

Setup

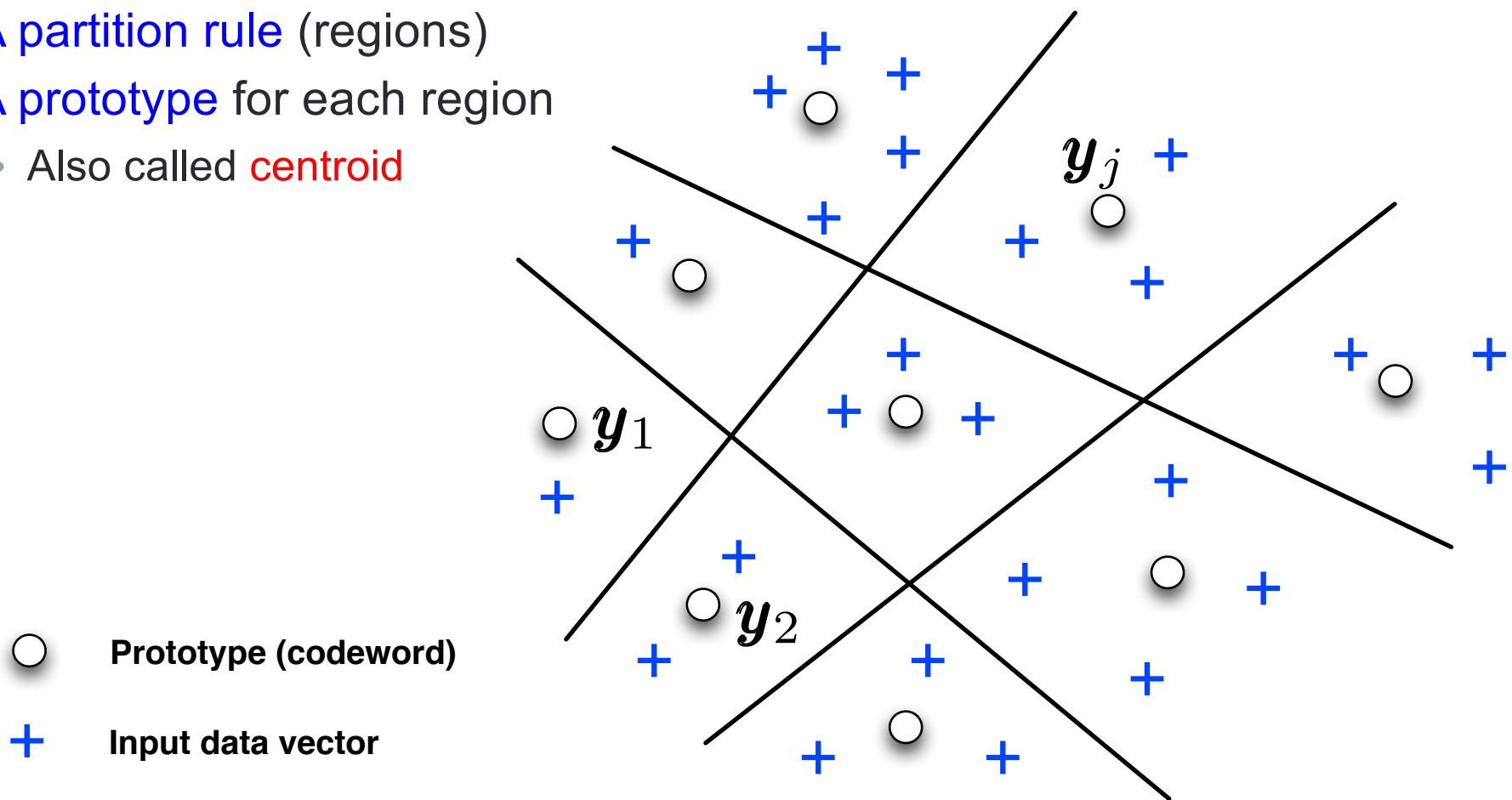
- Input space, vectors of m elements

$$\mathbf{x} = [x_1, x_2, \dots, x_m]^T$$

- i.i.d. sequentially sampled, one at a time,
 - from the same pdf $f(\mathbf{x})$
 - time is discrete $n=1, 2, \dots$
- A set of centroids is defined (e.g., cluster centers in K-means)
$$\{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_\ell\}$$
- with: $\mathbf{y}_j \in \mathbb{R}^m$

Vector Quantization (1/2)

- Given an input data distribution (multi-dimensional vectors)
- We need to find:
 - A partition rule (regions)
 - A prototype for each region
 - Also called **centroid**



Vector Quantization (2/2)

- **Distortion measure**

- Most common measure: **Euclidean distance**
- **Average distortion** is quantified as the mean square error

$$E[d(\mathbf{x}, \mathbf{y}_j)] = \sum_{j=1}^{\ell} \int_{I_j} \|\mathbf{a} - \mathbf{y}_j\| f(\mathbf{a}) d\mathbf{a} \quad (1)$$

Region j Centroid Input signal pdf
 $\|\mathbf{x}\| \rightarrow \text{norm-2}$

- **Optimal VQ**

- Find a set of **centroids** and a **partition rule** that **minimize** (1)
- **Nearest neighbor condition:** the optimal partition is the one returning the minimum distortion (region j)

$$I_j = \{\mathbf{x} : d(\mathbf{x}, \mathbf{y}_j) \leq d(\mathbf{x}, \mathbf{y}_h), j \neq h\}$$

What self organizing maps are (1/2)

- Neural networks based on **competitive learning**
- **Competitive learning**
 - Output neurons compete among themselves to be activated (or *fired*)
 - Result: *only one neuron is on at any one time*
 - The neuron that wins the competition is called
 - The “**winner-takes-all neuron**” or simply,
 - the “**winner**”
- **In a SOM**
 - Neurons are usually placed at the nodes of a lattice (usually 1D or 2D)
 - The neurons **become selectively tuned to the input patterns (stimuli)** in the course of a **competitive learning process**
 - Learning is **unsupervised**

What self organizing maps are (2/2)

- Each neuron in the SOM
 - Has a weight (or “feature vector”)
 - The weight corresponds to the position of the neuron in the lattice
 - Also referred to as “coordinates” of the neuron
 - The spatial locations (i.e., the neuron’s weights)
 - During the learning phase, **become ordered** with respect to each other in a way that a meaningful coordinate system is created over the lattice
 - The result is a **topographic map** of the input patterns, where the *coordinates of the neurons in the lattice are indicative of statistical features contained in the input patterns*
- Key properties
 - The SOM is inherently non-linear
 - Learning (adaptation of its weights) is unsupervised

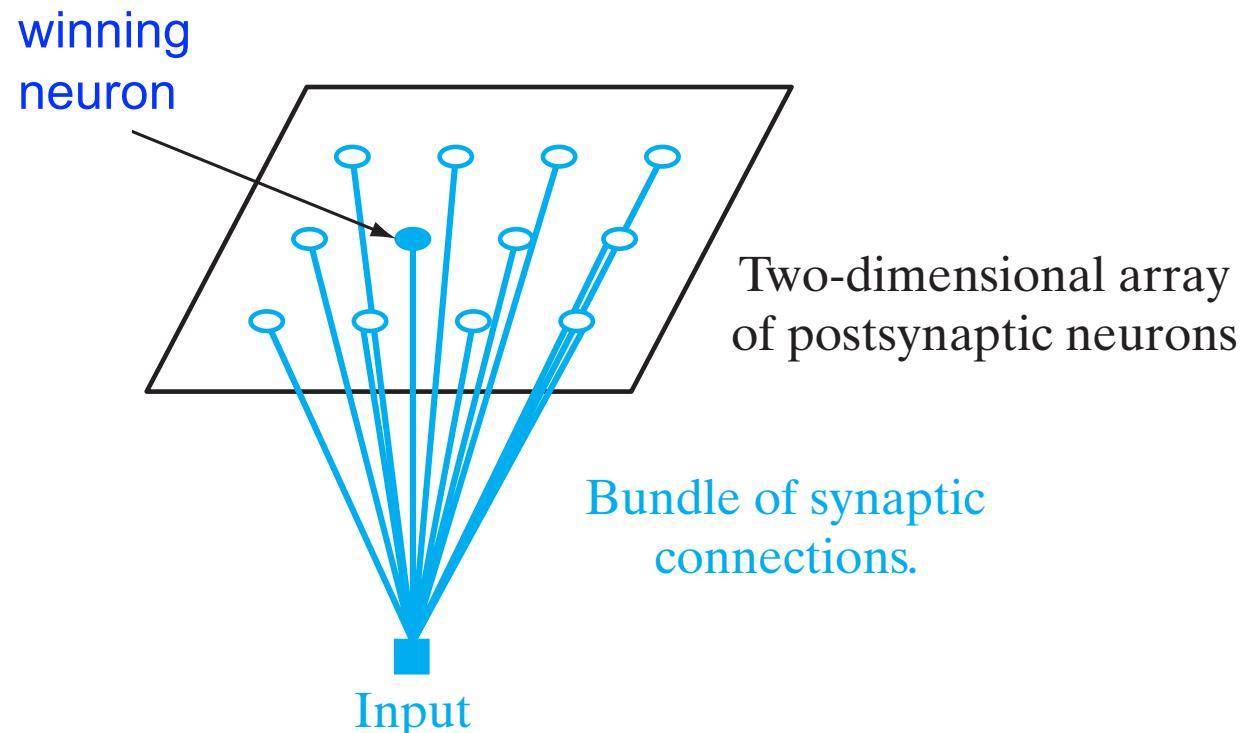
SOM & the human brain

- SOMs are motivated by features of the human brain
 - “The brain is organized in many places in a way that different sensory inputs are represented by topologically ordered computation maps”
- Key fact: sensory inputs such as
 - Tactile (Kaas, J.H., M.M. Merzenich, and H.P. Killackey, 1983. “The reorganization of somatosensory cortex following peripheral nerve damage in adult and developing mammals,” *Annual Review of Neurosciences*, vol. 6, pp. 325–356)
 - Visual (Hubel, D.H., and T.N. Wiesel, 1962. “Receptive fields, binocular interaction and functional architecture in the cat’s visual cortex,” *Journal of Physiology*, vol. 160, pp. 106–154, London)
 - Acoustic (Suga, N., 1985. “The extent to which bisonar information is represented in the bat auditory cortex,” in *Dynamic Aspects of Neocortical Function*, G.M. Edelman, W.E. Gall, and W.M. Cowan, eds. pp. 653–695, New York: Wiley)
- are mapped onto several areas of the cerebral cortex in a topologically ordered manner

SOM [1]



Prof. Teuvo Kohonen



- **Key observation:** the spatial location of an output neuron in a topographic map corresponds to a particular feature of the data drawn from the input space
- **SOM:** it is not meant to explain neurobiological details, but (i) *to capture the essential features* of computational maps in the brain and (ii) *to remain computationally tractable*
- [1] Teuvo Kohonen, “Self-organized formation of topologically correct feature maps,” *Biological Cybernetics*, vol. 43, pp. 59–69, 1982.

Self Organizing Map (SOM)

- Goal
 - Is to transform an incoming signal pattern of arbitrary dimension (m)
 - Into a one- or two-dimensional discrete map of neurons
- The algorithm in a nutshell
 - Initialization: the neuron's (synaptic) weights are initialized at random: no prior order is imposed on the map. Then, for each input vector DO:
 - 1. Competition: for each input pattern, each neuron computes a local value using a discriminant function. The neuron with the highest value wins the competition
 - 2. Cooperation: the winning neuron provides the location of a topological neighborhood of excited neurons. The neurons in this neighborhood will update their synaptic weights
 - 3. Synaptic adaptation: the winning neuron, as well as the neurons in its neighborhood update their synaptic weight, bringing it closer to the input pattern (vector)

Input vectors & synaptic weights

- Let \mathbf{m} denote the dimension of the input data space
- Let \mathbf{x} be a vector selected at random from the input space

$$\mathbf{x} = [x_1, x_2, \dots, x_m]^T$$

- The **synaptic-weight vector** of each neuron in the SOM has the same dimension of the input space. The weight vector associated with **neuron j** is:

$$\mathbf{w}_j = [w_{j1}, w_{j2}, \dots, w_{jm}]^T, j = 1, 2, \dots, \ell$$

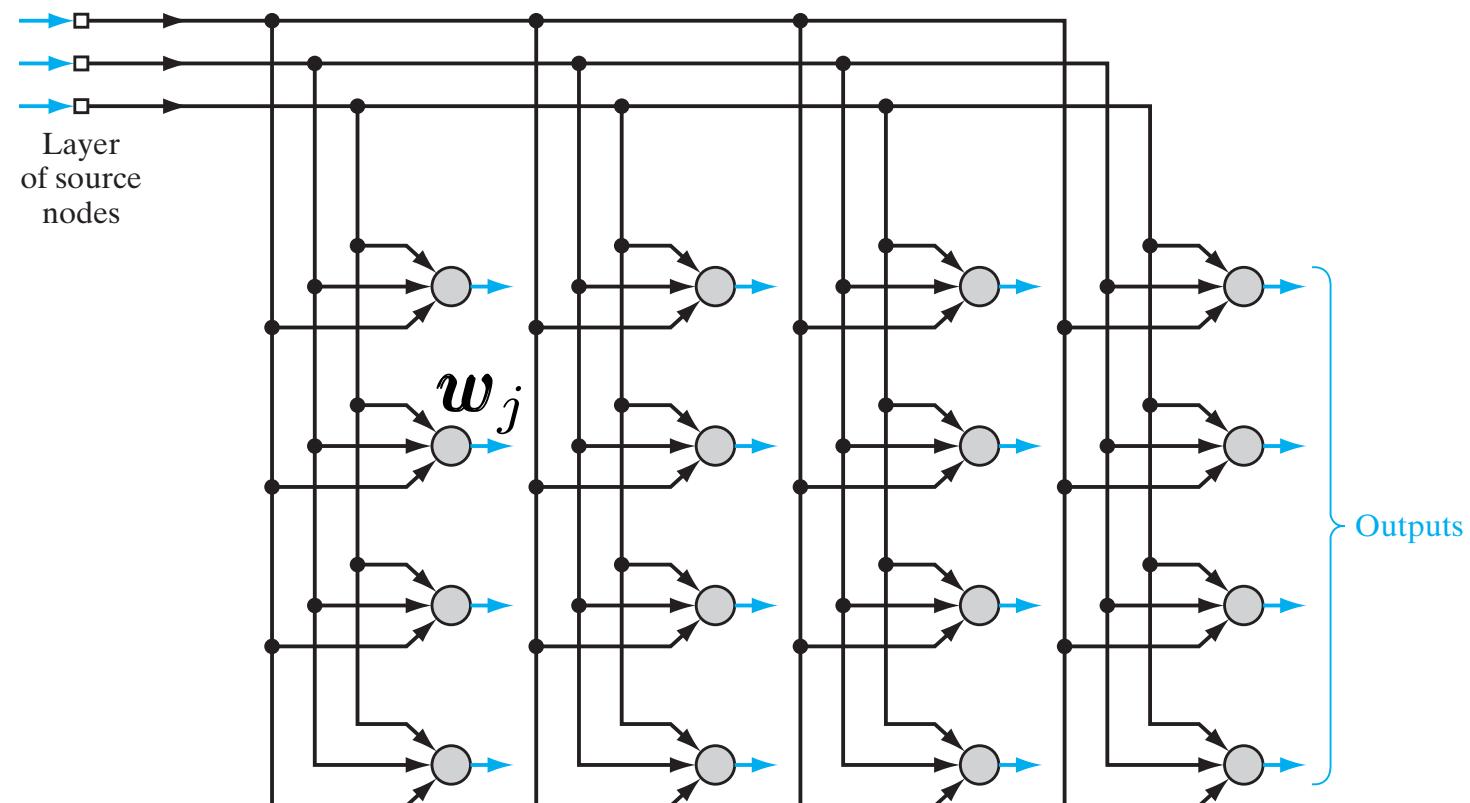
- ℓ represents the number of (output) neurons in the SOM

A 2D SOM example

2D lattice of $\ell = 4 \times 4 = 16$ neurons

Input patterns are vectors with three elements

$$\mathbf{x} = [x_1, x_2, x_3]^T$$



What SOM does

- Is a topological mapping algorithm
- “Optimally” places a fixed number of vectors (the neuron’s weights) into a higher dimensional space (of size m):
 - It is a dimensionality reduction algorithm (see below)
 - This means that it is suitable for clustering (data compression)
- More specifically, the SOM maps
 - Input (feature) vectors of m elements $\mathbf{x} \in \mathbb{R}^m$ into
 - ℓ vectors of m elements $\mathbf{w}_j \in \mathbb{R}^m$, $j = 1, 2, \dots, \ell$
 - with $\ell < m$

1. Competition

- Let \mathbf{x} be a vector from the input space

$$\mathbf{x} = [x_1, x_2, \dots, x_m]^T$$

- All the neurons in the SOM lattice **compete**
 - The best fit neuron will win the competition
 - Let $i(\mathbf{x})$ be the index associated with the **winning neuron**
 - This index $i(\mathbf{x})$ is computed as:

$$i(\mathbf{x}) = \arg \min_j \|\mathbf{x} - \mathbf{w}_j\|, j \in \mathcal{L}$$

- where \mathcal{L} indicates the SOM lattice (*output space*)
- and $\|\mathbf{x} - \mathbf{w}_j\|$ is the **Euclidean distance**
- Observation:** a *continuous* input space of activation patterns is mapped onto a *discrete output space* (lattice) of neurons by a process of competition among the neurons in the neural network

2. Cooperation (1/5)

- The winning neuron locates the center of a *topological neighborhood of cooperating neurons*
- The question now is: “*how do we define a neighborhood that is neuro-biologically correct?*”
- Answer: there is **neurobiological evidence** that
 - Lateral interaction occurs among excited neurons in the human brain
 - This means that: a neuron that is firing tends to excite *more the neurons that are located in its immediate neighborhood* rather than those that are located away from it. This is also intuitively satisfying
- Hence: this observation leads us to define a topological neighborhood around the winning neuron $i(x)$ and **make it decay smoothly with lateral distance**

2. Cooperation (2/5)

- Topological neighborhood $h_{i,j}$
 - Centered on **winning neuron i**
 - Encompassing a set of neighboring neurons, denoted by j
 - Let $d_{i,j}$ be the lateral distance between neuron i and neuron j
- We assume that:
 - **1.** the topological neighborhood $h_{i,j}$ is **symmetric** around the maximum point, which is defined by $d_{i,j} = 0$. In other words, it attains its maximum value at the winning neuron i
 - **2.** the amplitude of the lateral neighborhood $h_{i,j}$ **decreases monotonically** with the distance to the point that it goes to zero as the distance diverges to infinity (this condition is key for the convergence of the learning phase)

$$h_{i,j}(n) \rightarrow 0 \text{ as } d_{i,j} \rightarrow +\infty$$

2. Cooperation (3/5)

- A good choice for the neighborhood function is the **Gaussian function**

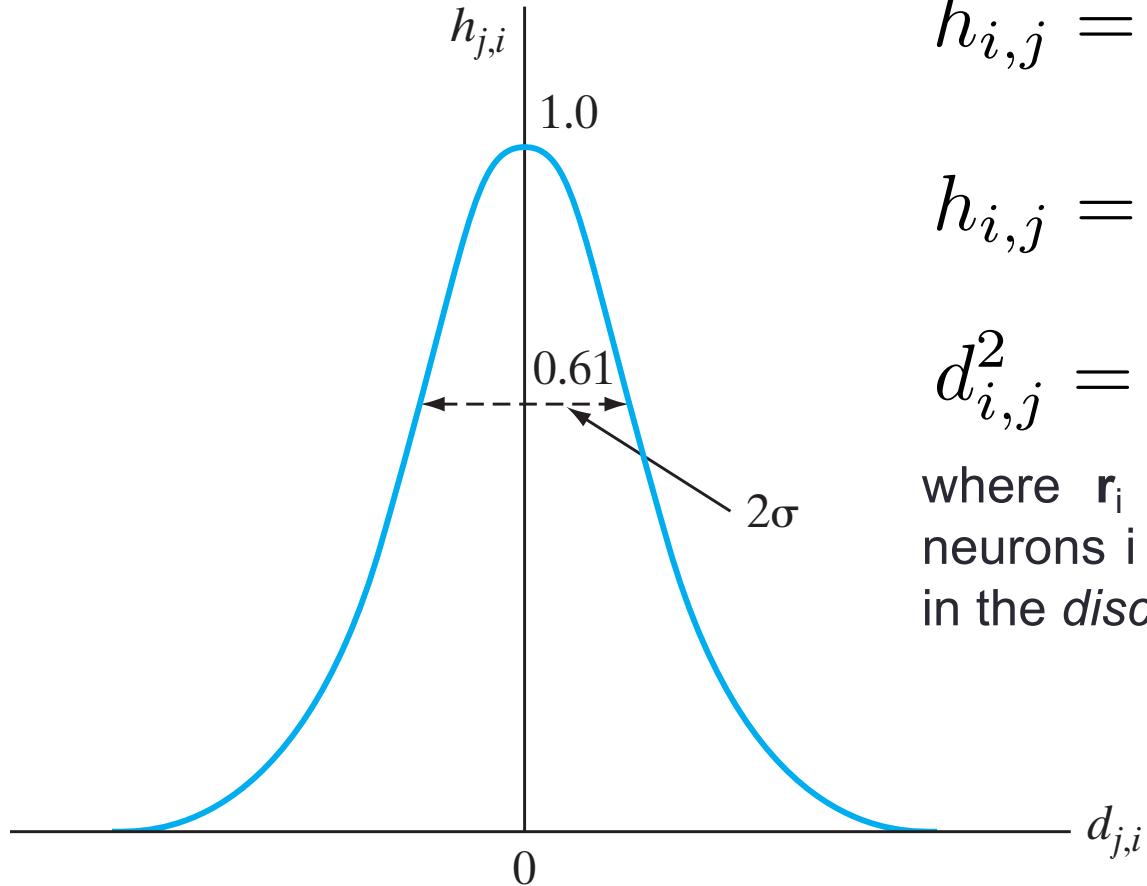
$$h_{i,j} = \exp\left(-\frac{d_{i,j}^2}{2\sigma^2}\right), j \in \mathcal{L}$$

- **Observations**

- This function is translation invariant: independent on the index (location) of the winning neuron
- parameter σ is the “effective width”: it measure the degree to which the neurons in the neighborhood of the winner participate in the learning process
- a Gaussian neighborhood function leads to faster convergence than, e.g., a rectangular neighborhood function
- in the definition of neighborhood there is no “wrapping around”

2. Cooperation (4/5)

- The Gaussian neighborhood function



$$h_{i,j} = \exp\left(-\frac{d_{i,j}^2}{2\sigma^2}\right), \quad j \in \mathcal{L}$$

$$h_{i,j} = h_{j,i}$$

$$d_{i,j}^2 = \|\mathbf{r}_i - \mathbf{r}_j\|^2$$

where \mathbf{r}_i and \mathbf{r}_j are the positions of neurons i and j and are both measured in the *discrete output space* (the lattice)

2. Cooperation (5/5)

- Another unique feature of a SOM algorithm is that the size of the topological neighborhood **is permitted to shrink with time n**, i.e.,

$$\sigma(n) = \sigma_0 \exp\left(-\frac{n}{\tau_1}\right), \quad n = 0, 1, 2, \dots$$

- Where τ_1 is a constant (empirically) picked by the designer. Hence, the topological neighborhood function assumes a time-varying form of its own:

$$h_{i,j}(n) = \exp\left(-\frac{d_{i,j}^2}{2\sigma(n)^2}\right), \quad j \in \mathcal{L}, \quad n = 0, 1, 2, \dots$$

3. Adaptation (1/3)

- Let n be the current time
- Time n is updated to $n+1$ at each new input
- Let $\mathbf{x} = [x_1, x_2, \dots, x_m]^T$ be the input vector at time n
- Let $i(\mathbf{x})$ be the index of the winning neuron
- Each neuron j (including the winning neuron i itself)
updates its synaptic weight vector using:

$$\Delta \mathbf{w}_j = \eta(n) h_{i,j}(n) (\mathbf{x} - \mathbf{w}_j), \quad \begin{cases} i : & \text{winning neuron} \\ j : & \text{excited neuron} \end{cases}$$

3. Adaptation (2/3)

- Using discrete-time formalism
- The update equation for the weight vectors is:

$$\mathbf{w}_j(n+1) = \mathbf{w}_j(n) + \eta(n)h_{i,j}(n)(\mathbf{x} - \mathbf{w}_j), \quad j \in \mathcal{L}$$

- $h_{i,j}(n)$: controls the size of the neighborhood vs space & time
- $\eta(n)$: is a learning rate parameter
- This update equation is of a **stochastic approximation** type
 - The new weight vector $\mathbf{w}(n+1)$ is equal to the old one $\mathbf{w}(n)$ plus a $\Delta\mathbf{w}$ (update) term, which depends on the distance between (i) the current input vector \mathbf{x} and (ii) the weight vector \mathbf{w}_j
 - The intensity of the update depends on (i) how much \mathbf{x} differs from \mathbf{w}_j , (ii) the learning rate parameter and (iii) the location of neuron j wrt to the winning neuron $i(\mathbf{x})$ (note that j and i can also coincide)

3. Adaptation (3/3)

- Learning rate parameter
- Usually an **exponentially decaying equation**, although may not be optimal, is adequate for the SOM algorithm to converge towards a correct topological map

$$\eta(n) = \eta_0 \exp\left(-\frac{n}{\tau_2}\right), \quad n = 0, 1, 2, \dots$$

- Parameter τ_2 is set empirically by the designer

Observations

- The neighborhood function **correlates neuron weights** in the output space
 - A wide function (large $h_{i,j}(n)$) will **correlate the directions of the weight (delta) updates** for the neurons around the winner $i(x)$, the closer the neuron j is, the more correlated its weight update will be with respect to that of the winner i
 - This, as learning evolves, creates **topological ordering**
 - This spatial correlation in the updates is the reason for the non-linearity of the SOM and **makes its mathematical analysis hard**
 - As time n goes by, the neighborhood function shrinks and, eventually, **it will be a spike equal to 1 only for $d_{i,j}=0$ and zero otherwise:**
 - At this point, learning is no longer correlated. **Each neuron independently updates its own synaptic weight**

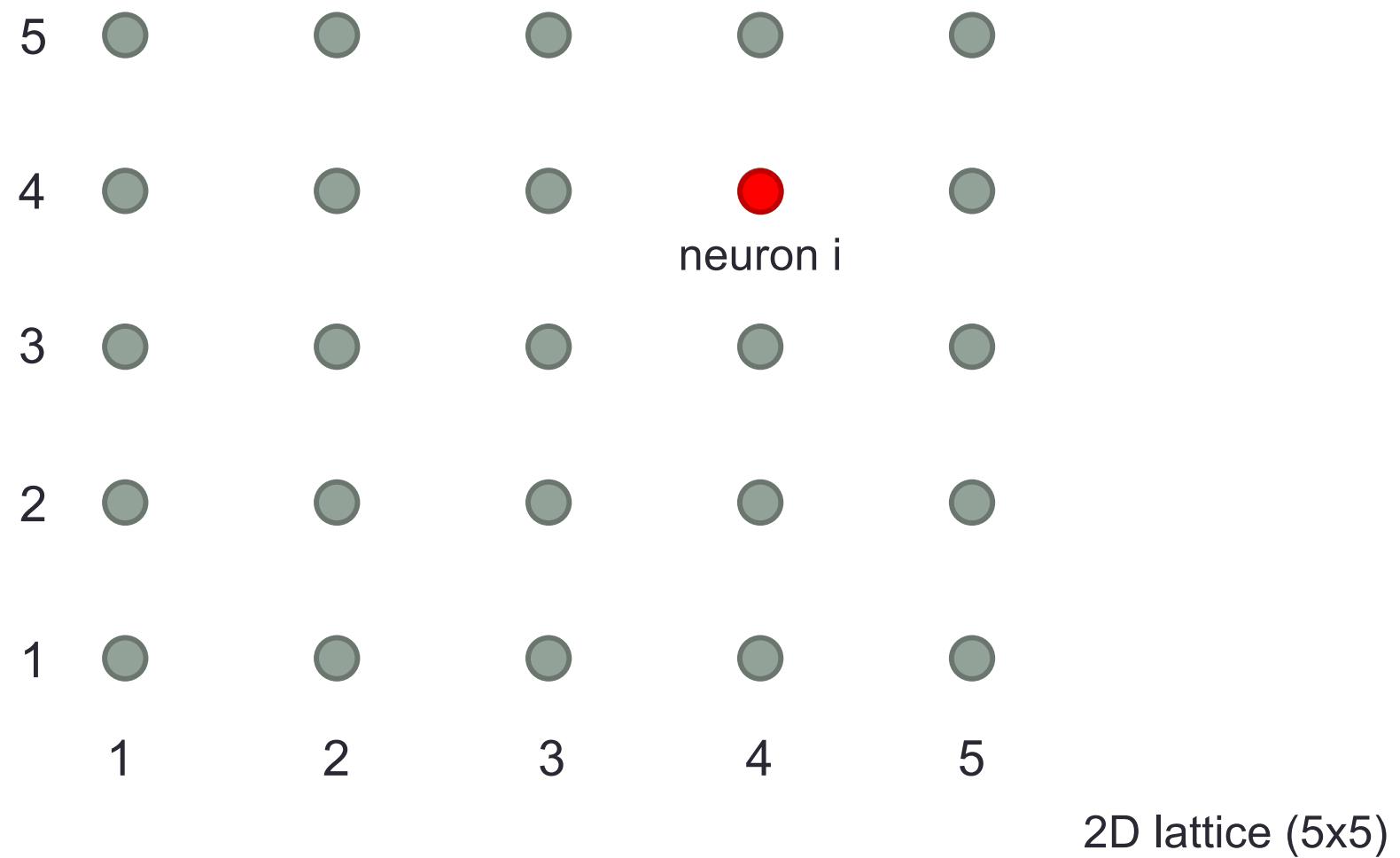
Review of the SOM algorithm

STEP 1: initialize SOM

$$w_i = [w_{i1} \ w_{i2} \dots \ w_{im}]^T$$

w_i picked at random

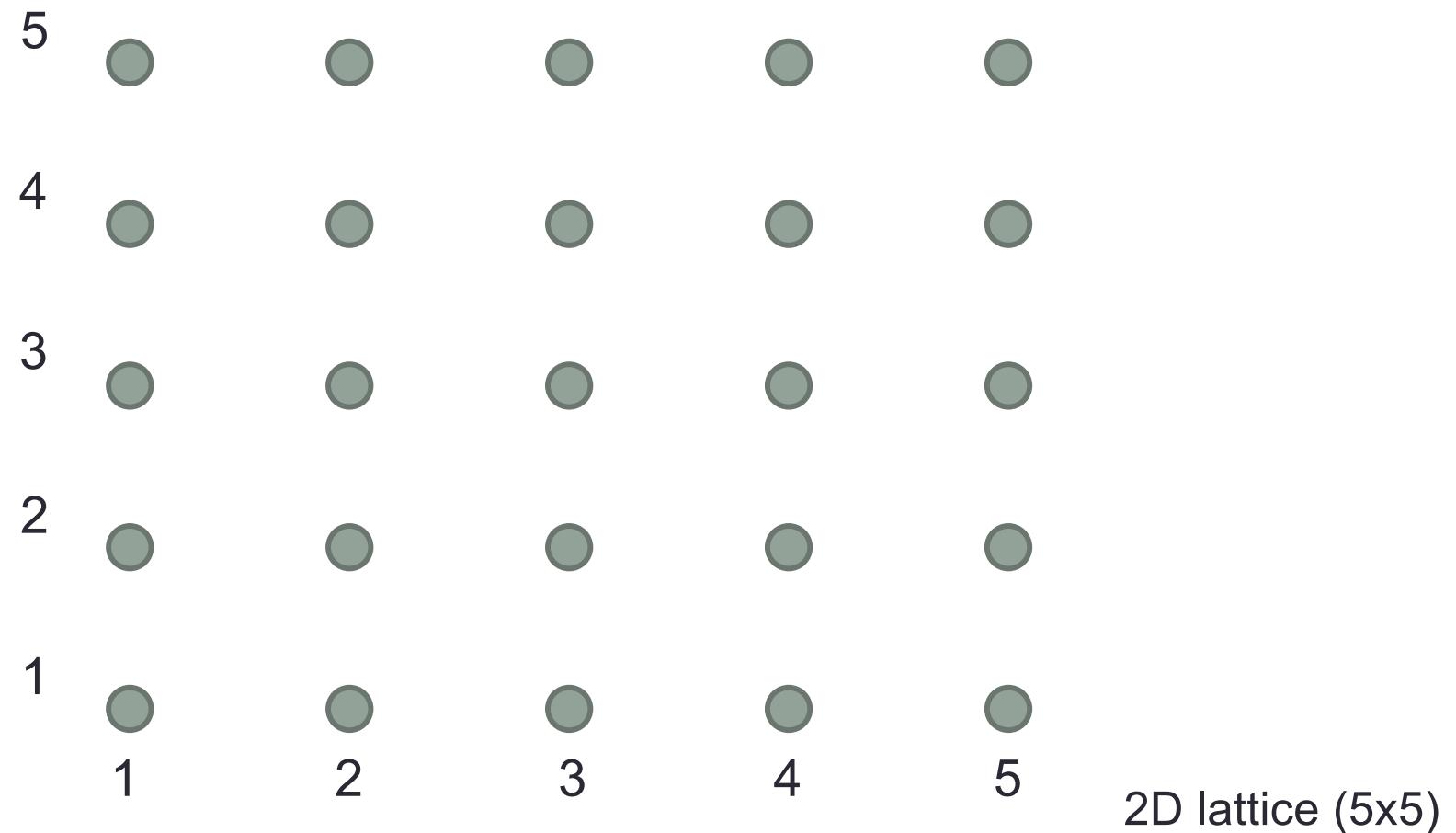
lattice
coordinates



STEP 2: sampling from input pdf

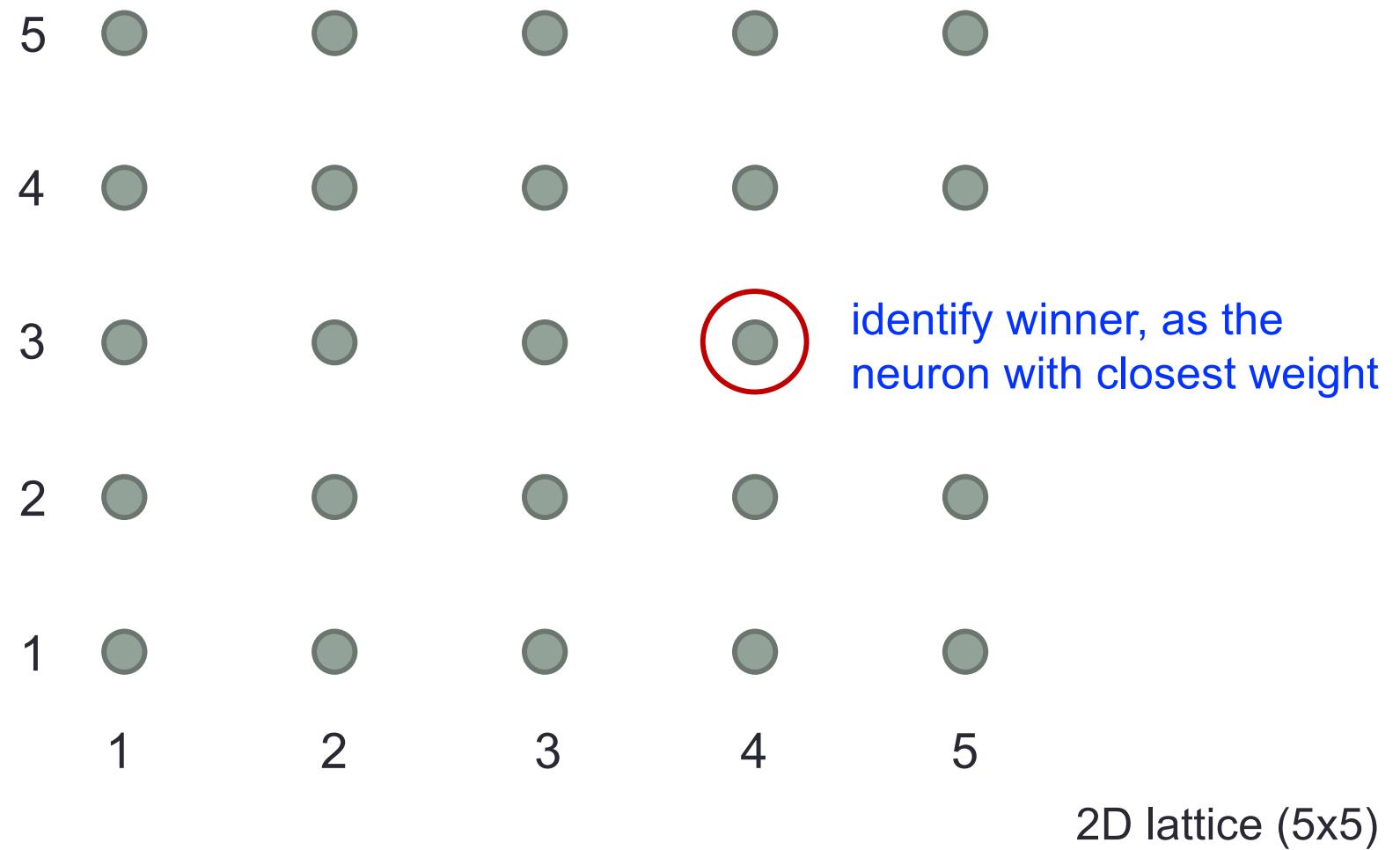
generic time $n+1$

sample from input distribution $\boldsymbol{x}_{n+1} \leftarrow f(\boldsymbol{x})$



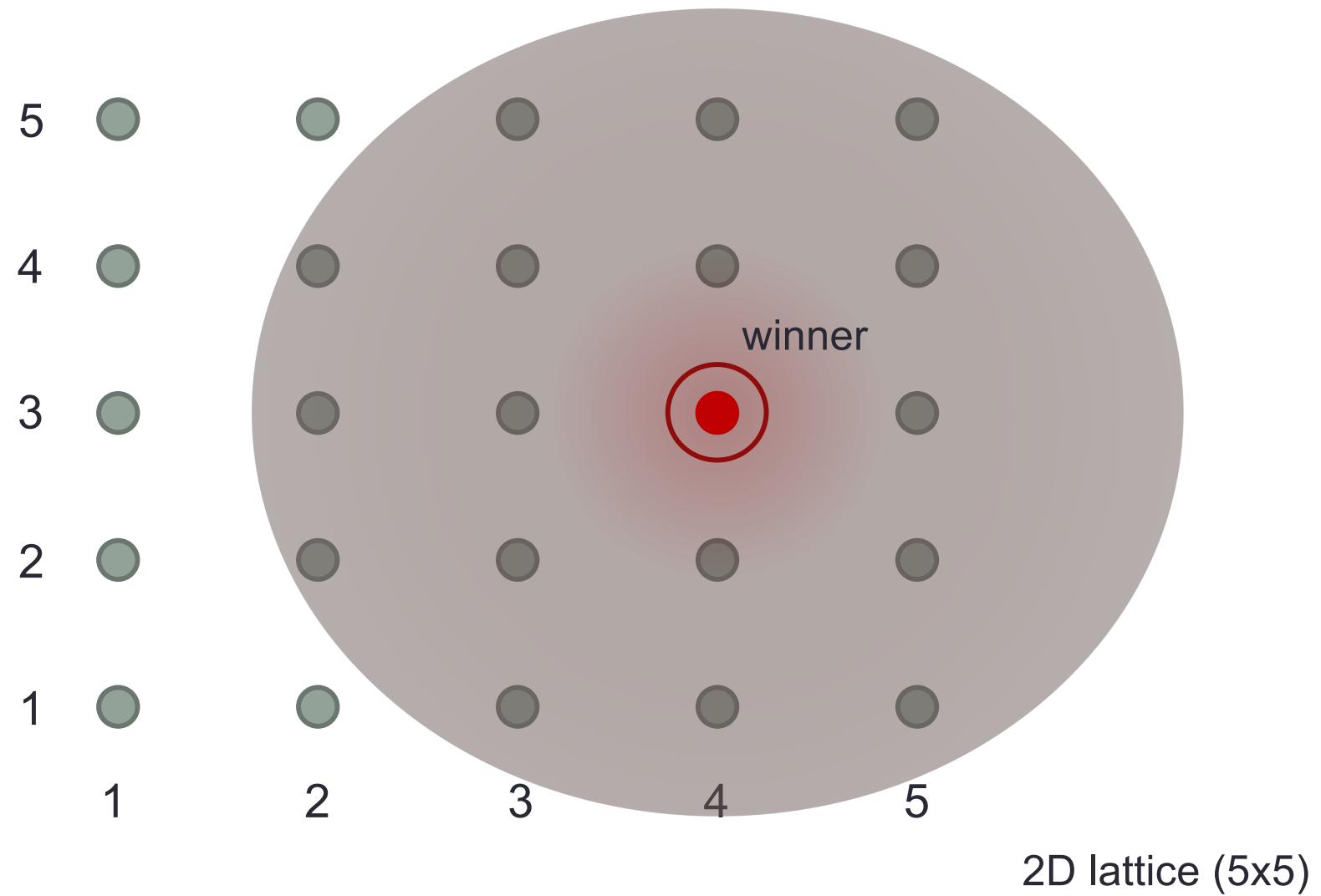
STEP 3: similarity matching

all neurons compute local value $\|\mathbf{x}_{n+1} - \mathbf{w}_i\|, \forall i$



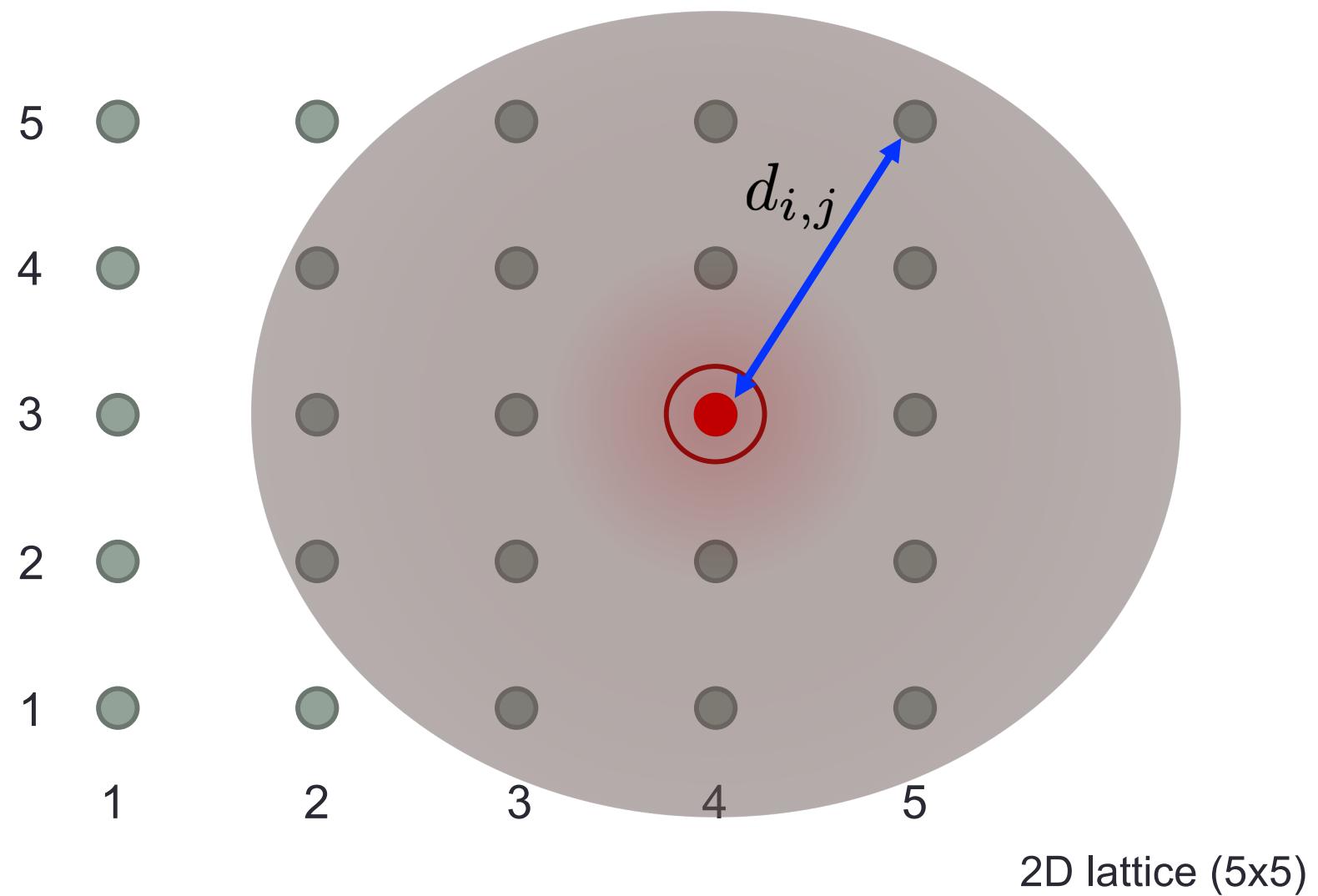
STEP 3: similarity matching

once winner identified → build Gaussian neighborhood around it



STEP 3: similarity matching

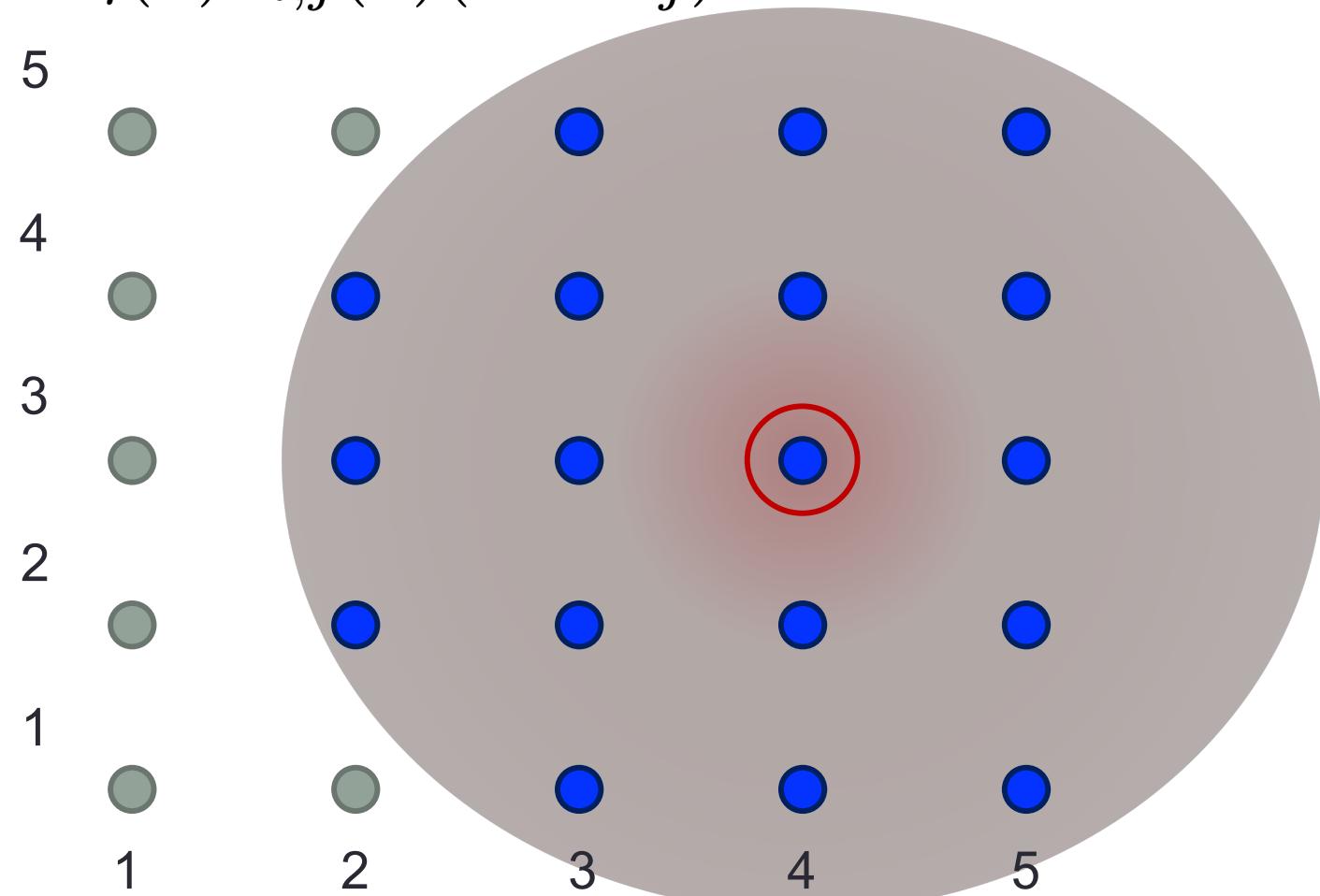
compute neighborhood distances → using lattice coordinates



STEP 4: update weights

update weights of all (blue) nodes in neighborhood

$$\mathbf{w}_j \leftarrow \mathbf{w}_j + \eta(n) h_{i,j}(n) (\mathbf{x} - \mathbf{w}_j)$$



STEP 5: update learning parameters

- Adjust neighborhood size $h_{i,j}(n + 1)$
- Adjust learning rate $\eta(n + 1)$
- Repeat from STEP 2

Review of the SOM algorithm

Algorithm 1 (SOM):

1. **Initialization:** set the initial time step to $n = 0$ and choose small random values for the initial synaptic-weight vectors $\mathbf{w}_j(0)$, $j \in \mathcal{L}$.
2. **Sampling:** set $n \leftarrow n + 1$ and sample the training input pattern \mathbf{x}_n , i.e., the feature vector
3. **Similarity matching:** find the winning neuron, whose index is $j^*(n)$ at time step n by using the minimum-distance criterion:

$$j^*(n) = \operatorname{argmin}_j \|\mathbf{x}_n - \mathbf{w}_j(n)\|, j \in \mathcal{L}$$

where with $\|\mathbf{a} - \mathbf{b}\|$, we mean the Euclidean distance between vectors \mathbf{a} and \mathbf{b} .

4. **Update:** adjust the synaptic-weight vectors of all neurons by using the update equation:

$$\mathbf{w}_j(n+1) = \mathbf{w}_j(n) + \eta(n)h_{ji}(n)(\mathbf{x}_n - \mathbf{w}_j(n)),$$

with $j \in \mathcal{L}$, where $\eta(n)$ is the learning rate parameter at iteration n and $h_{ji}(n)$ is the neighborhood function centered on $j^*(n)$ at iteration n .

5. **Adjust neighborhood:** adjust the neighborhood size ($h_{ji}(n)$), the learning rate ($\eta(n)$) and continue from Step 2 if $n < n_{\text{iter}}$; stop otherwise.
-

The two phases of the learning process: *ordering and convergence*

- Learning starts from a state of complete disorder
- Progressively, the SOM algorithm leads to an organized representation of patterns drawn from the input space (provided that good parameters are selected)
- The adaptation of the synaptic weights in the neural network *can be decomposed into two phases*:
 - 1. **Self-organizing or ordering phase**: during this phase the topological ordering of the weight vectors takes place
 - 2. **Convergence phase**: here the SOM map is fine tuned to provide an accurate statistical quantification of the input space

Self-organizing phase

- Rules of thumb:

- The learning rate parameter $\eta(n)$ should begin with a value close to 0.1 and should decrease gradually but **never getting smaller than 0.01**
- Considering that this phase takes approx. 1,000 iterations
- Suitable values for the parameters are:

$$\begin{cases} \eta_0 = 0.1 \\ \tau_2 = 1,000 \end{cases}$$

- The neighborhood function **should initially include all the neurons** in the network, centered on the winning neuron, and **shrink slowly with time**
- Assuming a 2D lattice, we may set σ_0 equal to the “radius” of the lattice and for τ_1 we may use:

$$\tau_1 = \frac{1,000}{\log \sigma_0}$$

The SOM feature map ϕ (1/2)

- Once the SOM is trained we have a non-linear **feature map**
- The feature map ϕ**
 - maps vectors \mathbf{x} from the input space
 - into a corresponding neuron in the lattice (a “prototype” for VQ)

$$\phi : \mathbb{R}^m \rightarrow \mathcal{L}$$

- SOM map is frozen** - for a new input vector \mathbf{x}
 - we compute the **best fit neuron** in the SOM lattice
 - whose synaptic weight vector is the closest to \mathbf{x}

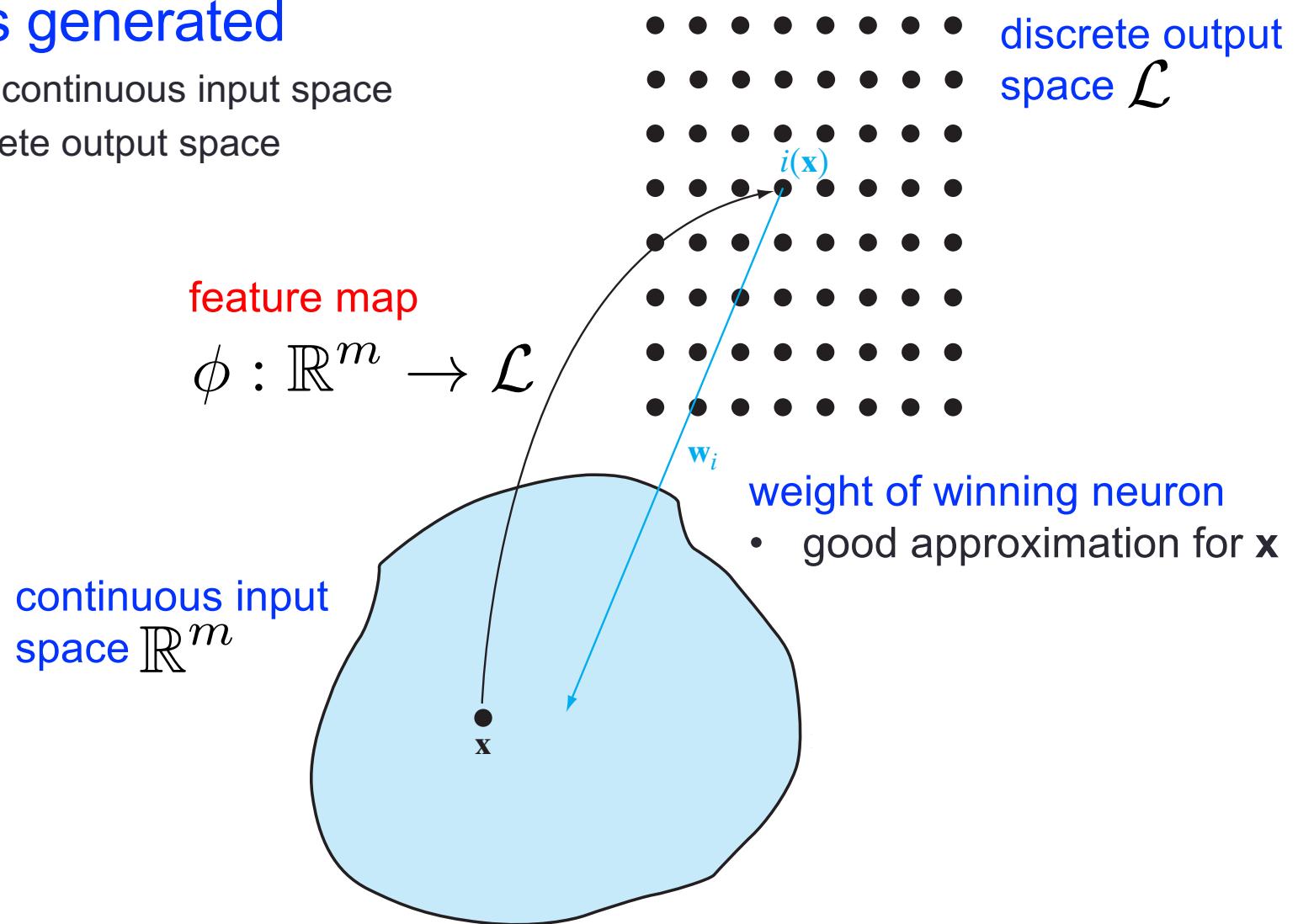
$$i(\mathbf{x}) = \arg \min_j \|\mathbf{x} - \mathbf{w}_j\|, j \in \mathcal{L}$$

- We have that:**

$$\phi(\mathbf{x}) \stackrel{\Delta}{=} i(\mathbf{x})$$

The SOM feature map ϕ (2/2)

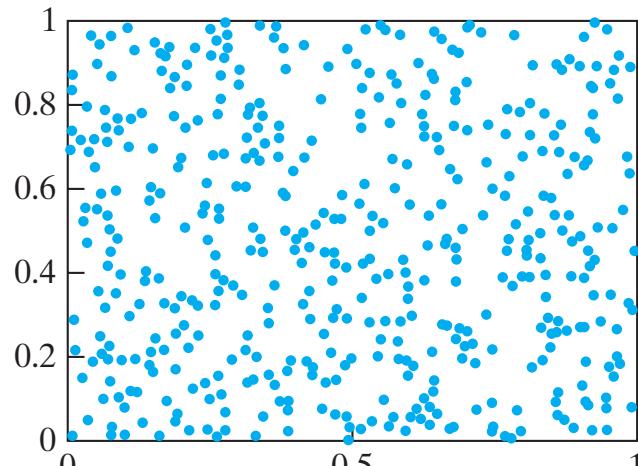
- A pointer is generated
 - Between the continuous input space
 - and the discrete output space



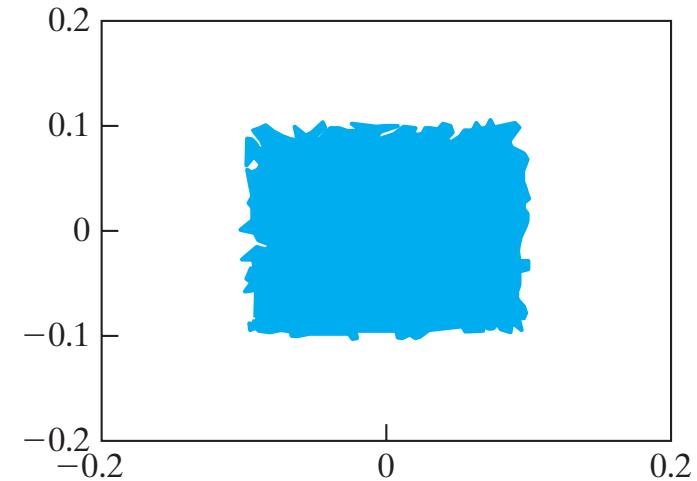
Simulation setup

- 576 neurons
 - Organized into a 2D lattice (24 rows & 24 columns)
 - Network is trained on 2D input vectors with **uniform distribution**
- **SOM learning** (initialization, ordering, convergence)
 - Initial weights are randomly assigned
 - Lines indicate neighboring neurons (synaptic weights are closest)
 - As learning progresses **synaptic weights become ordered**
 - In the refinement phase, **the weights specialize**
 - Mapping irregularities in the input distribution
 - Denser regions (more data points) tend to contain more neurons

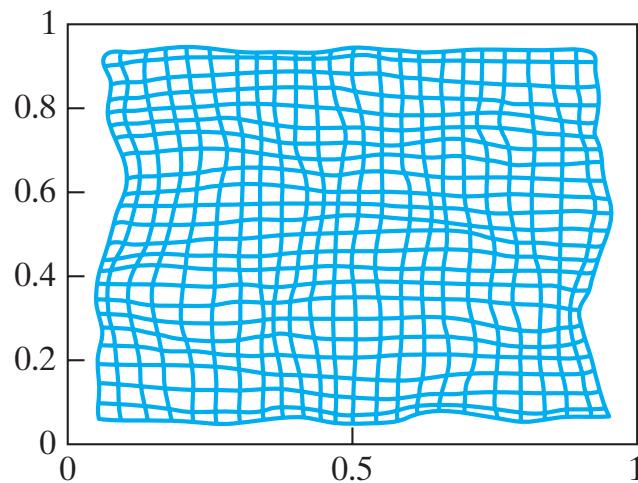
Simulation results (2D input)



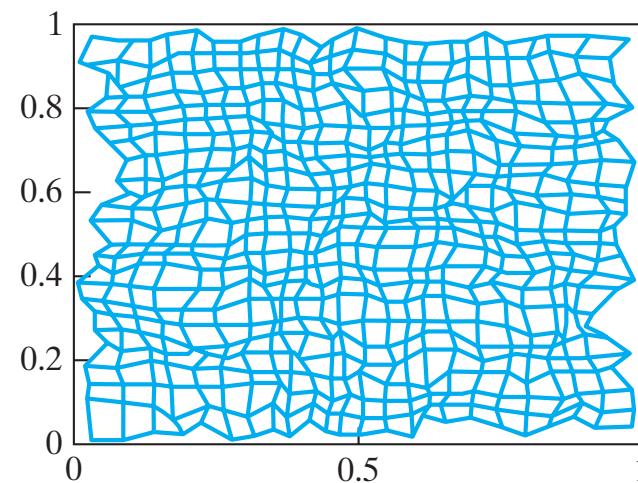
(a) Input distribution



Time = 0
(b) Initial weights



Time = 160 K
(c) Ordering phase



Time = 800 K
(d) Convergence phase

Intermission: variance unexplained (1/2)

- N input data points x_1, x_2, \dots, x_N
- A **regression function** returns for each x_i an estimate
- We define:

$$\hat{x}_i = f(x_i)$$

Average:

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i \quad \sigma_{\text{tot}}^2 \triangleq \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2$$

Total sample variance:

Variance of residuals:

$$\sigma_{\text{res}}^2 \triangleq \frac{1}{N} \sum_{i=1}^N (x_i - \hat{x}_i)^2$$

Intermission: variance unexplained (2/2)

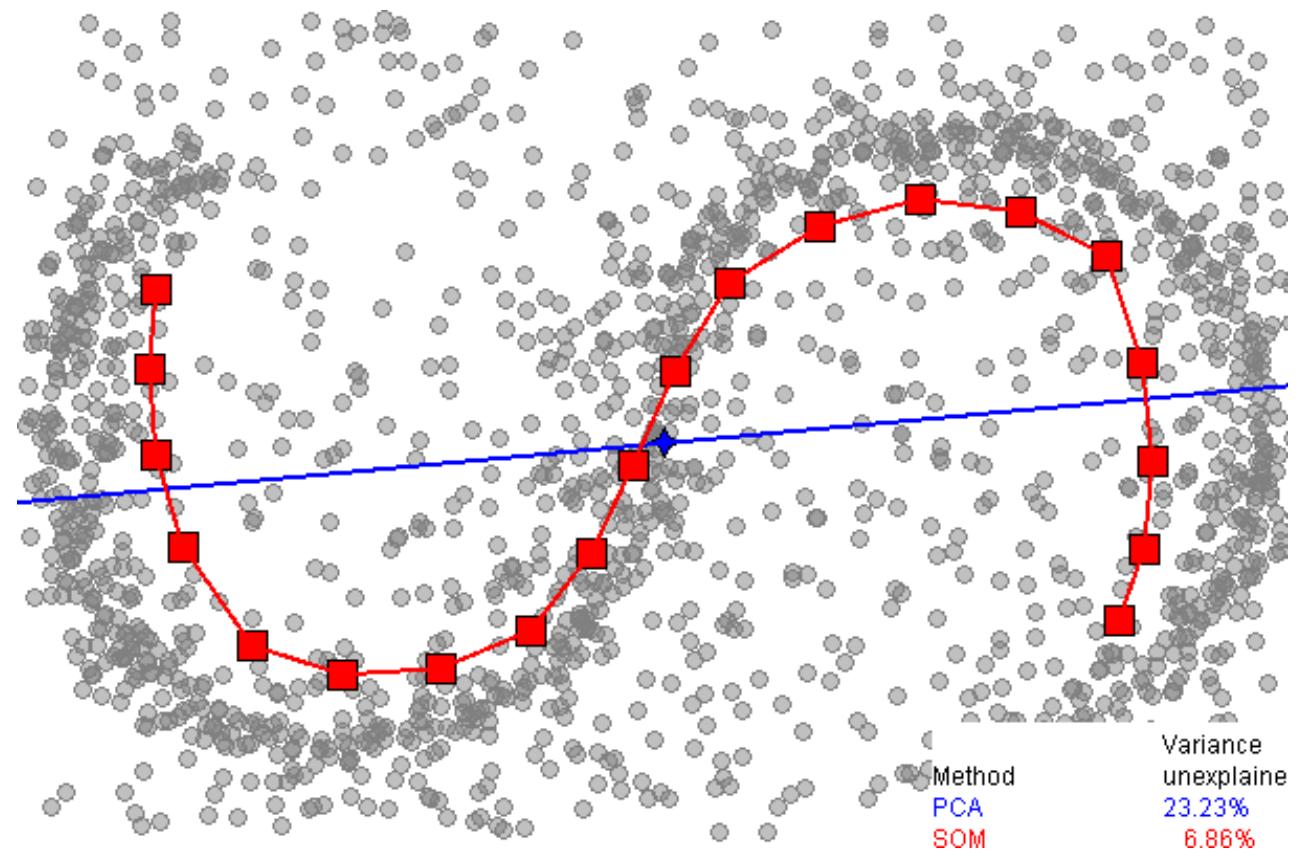
- Fraction of Variance Unexplained (FVU):

$$FVU \triangleq \frac{\sigma_{\text{res}}^2}{\sigma_{\text{tot}}^2} = \frac{\sum_{i=1}^N (x_i - \hat{x}_i)^2}{\sum_{i=1}^N (x_i - \bar{x})^2} = \frac{\sum_{i=1}^N (x_i - f(x_i))^2}{\sum_{i=1}^N (x_i - \bar{x})^2}$$

- Is the fraction of the total variance that the regression model is not able to capture

1D SOM vs 1D PCA – FVU

- Data points rest on a non-linear manifold
- PCA performs poorly
- SOM does a better job



One-dimensional SOM versus principal component analysis (PCA) for data approximation. SOM is a red broken line with squares, 20 nodes. PCA is presented by a blue line. Data points are the small grey circles. For PCA, the fraction of variance unexplained in this example is 23.23%, for SOM it is 6.86%

SOM bibliography

Chapter 9 of Haikin's Book [1]

[1] Simon Haykin, "Neural Networks and Learning Machines," 3rd Edition, *Prentice Hall*, 2009. (+5k citations)

Further readings:

[2] Teuvo Kohonen, "Self Organizing Maps," *Springer-Verlag*, 3rd Edition, 2001. (+22kcitations)

[3] H. Shah-Hosseini, R. Safabakhsh, "TASOM: a new time adaptive self-organizing map," *IEEE Transactions on Systems, Man, and Cybernetics*, Part B, Vol. 33, No. 2, pp. 271-282, April 2003. (61 citations)

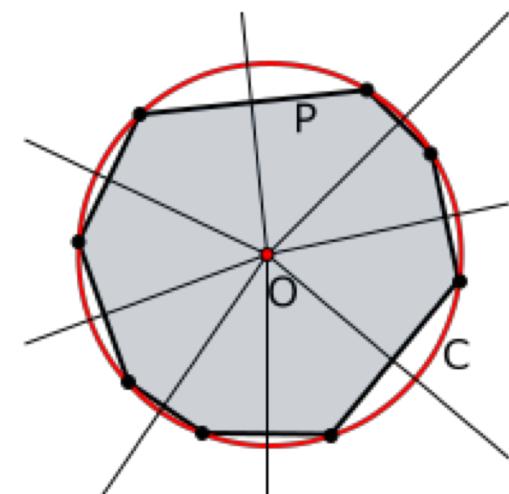
SOM pros and cons

- Pros
 - SOM learning leads to an **online** clustering algorithm
 - SOM learning is **unsupervised**, and precision can be tuned changing
 - training parameters,
 - number of neurons in the lattice
- Cons
 - Not clear how to “optimally” set the number of neurons
 - **It would be ideal** to automatically tune it
 - Add neurons until data pdf $f(\mathbf{x})$ is represented within some error tolerance
- Desideratum – adaptation
 - Once the SOM map stabilizes, **learning STOPS** :-)
 - We would like to **continuously adapt** (add, move & remove weights)
 - **Still, we want an unsupervised algorithm...**

GROWING NEURAL GAS (GNG) NETWORKS

Intermission: circumcircle

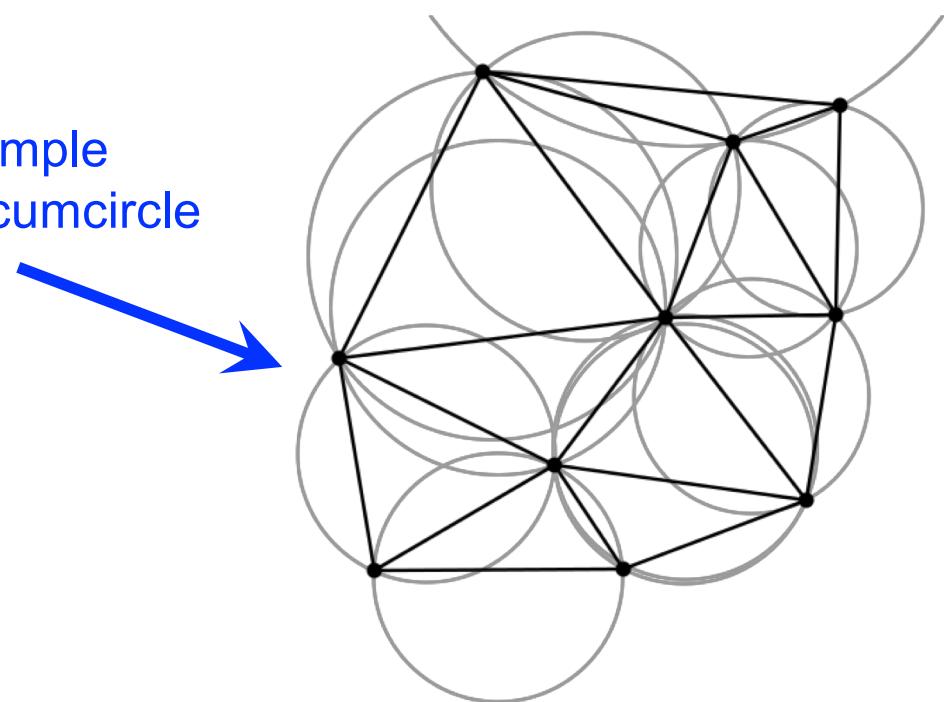
- **Circumscribed circle – “circumcircle”**
 - In geometry, the **circumscribed circle** or **circumcircle** of a polygon is a **circle which passes through all the vertices of the polygon**. The center of this circle is called the **circumcenter** and its radius is called the **circumradius**
 - A polygon which has a circumscribed circle is called a **cyclic polygon**. All regular simple polygons, all isosceles trapezoids, all triangles and all rectangles are **cyclic**



Intermission: Delaunay triangulation

- A **Delaunay triangulation** for a set P of discrete points in a plane is a triangulation $DT(P)$ such that
 - Points in P are connected through triangles
 - no point in P is inside the circumcircle of any triangle in $DT(P)$
 - **Delaunay triangulation:** maximizes the minimum angle of all the angles of the triangles in $DT(P)$

Delaunay triangulation example
no point in P is inside a circumcircle



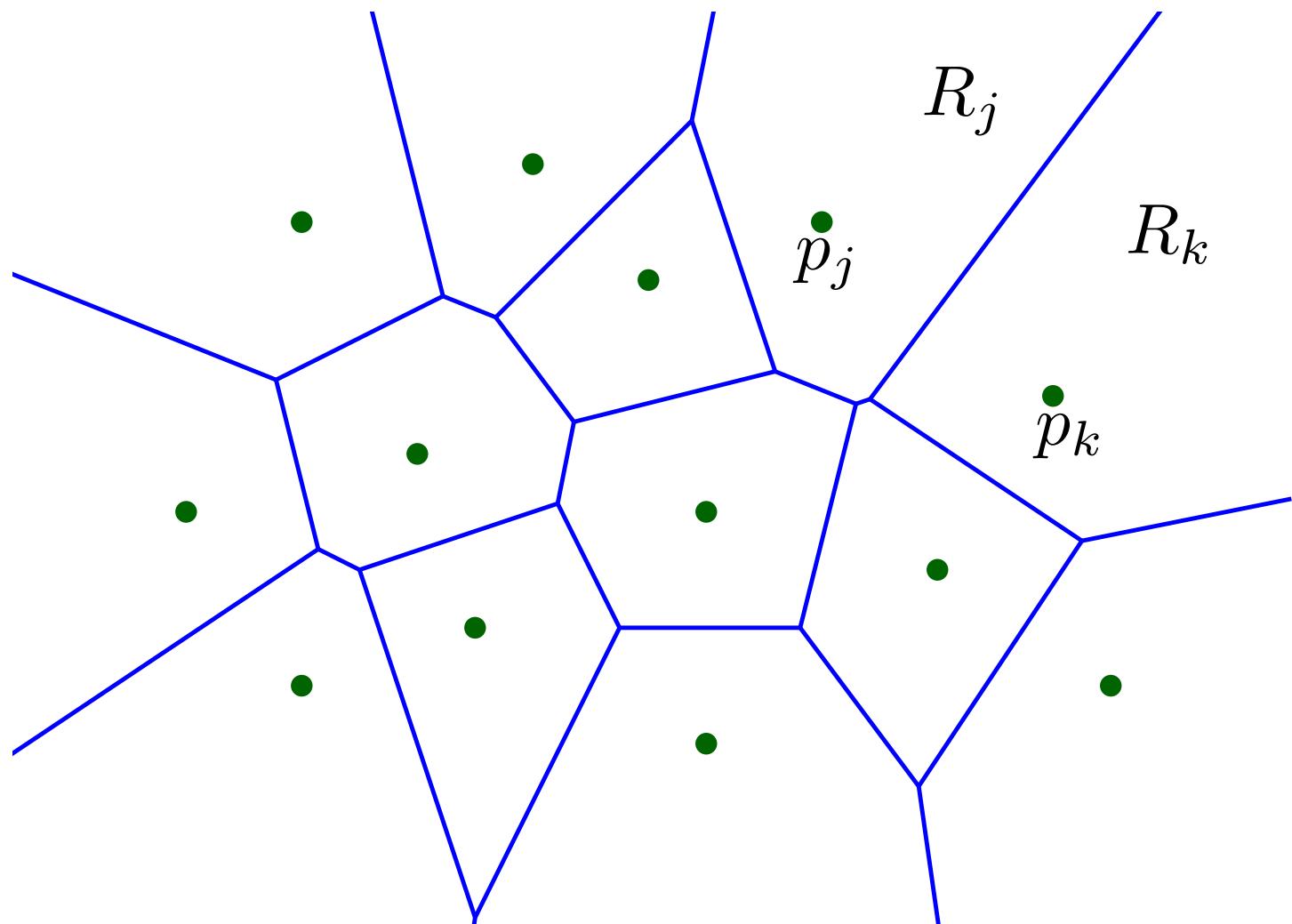
Voronoi diagram (1/6)

- Let X be a metric space
- Let $d(p_1, p_2)$ be a distance function between points p_1, p_2
- Let $P = \{p_k\}$ be a collection of K points in the space X
- The Voronoi region R_k associated with point p_k
 - Is the set of all points in X whose distance from p_k
 - Is no greater than their distance from any other point p_j
 - where j is any index different from k
- Formally:

$$R_k = \{x \in X \mid d(x, p_k) \leq d(x, p_j) \text{ for all } j \neq k\}$$

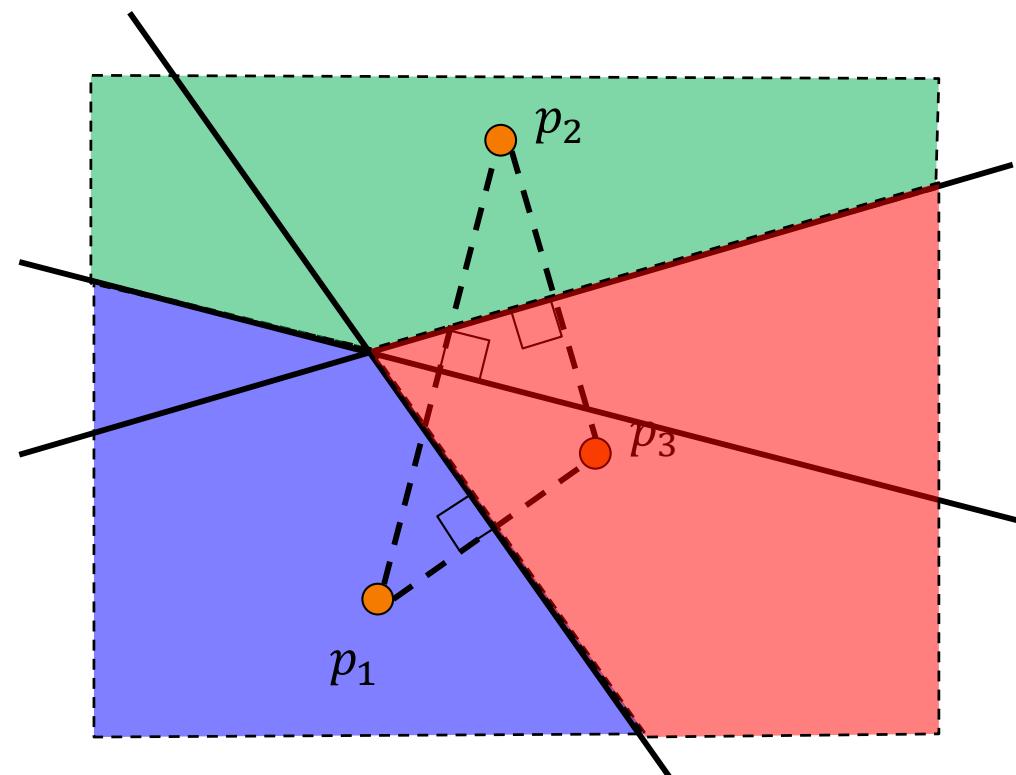
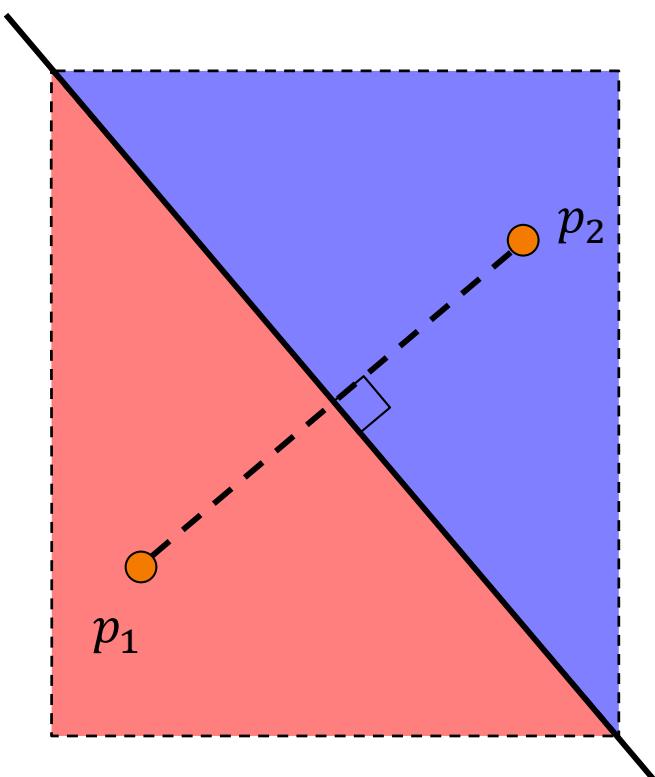
Voronoi diagram (2/6)

- Example: d is the Euclidean distance



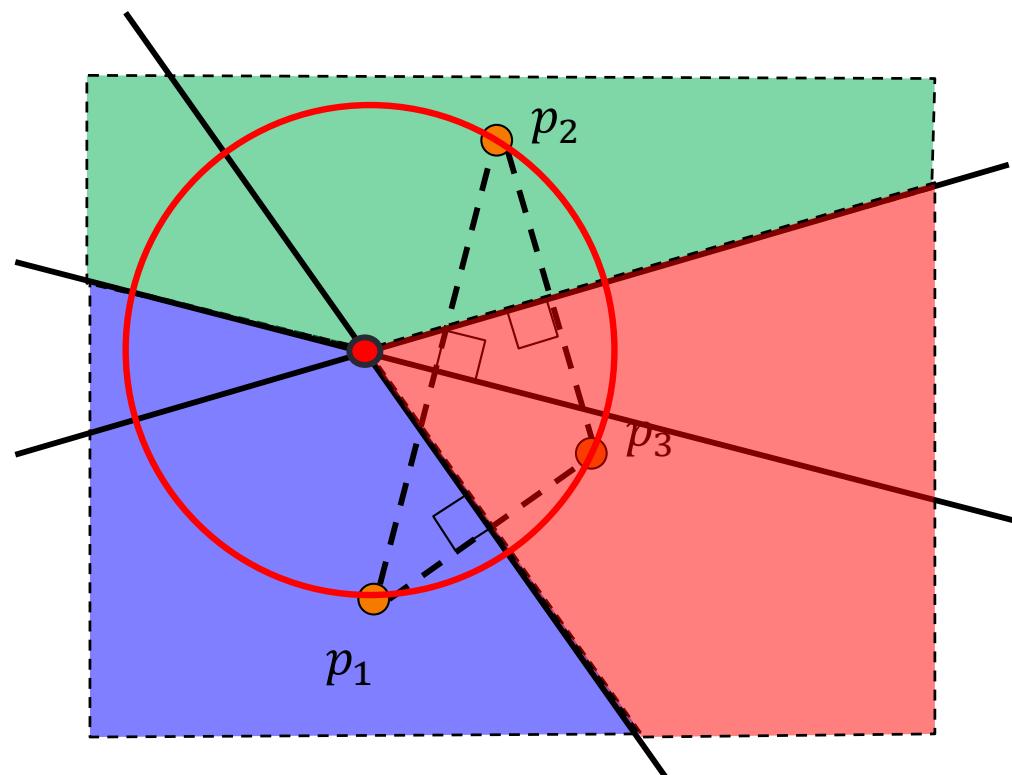
Voronoi diagram (3/6)

- 2 points: regions are defined by the perpendicular bisector
- 3 points: regions are **the intersection of the half spaces** $H(p_i, p_j)$ defined by the 3 perpendicular bisectors

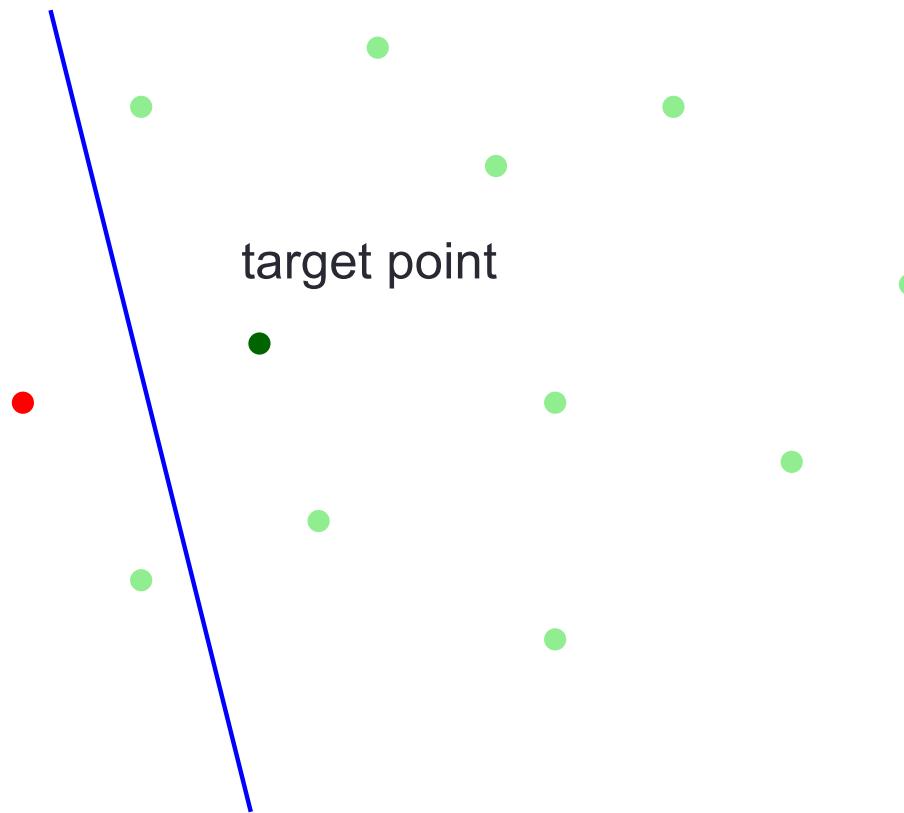


Voronoi diagram (4/6)

- 3 points: regions defined by the 3 perpendicular bisectors
 - The 3 bisectors intersect at a point (can be outside the triangle)
 - This point is the center of the circumcircle for the polygon (p_1 , p_2 , p_3)

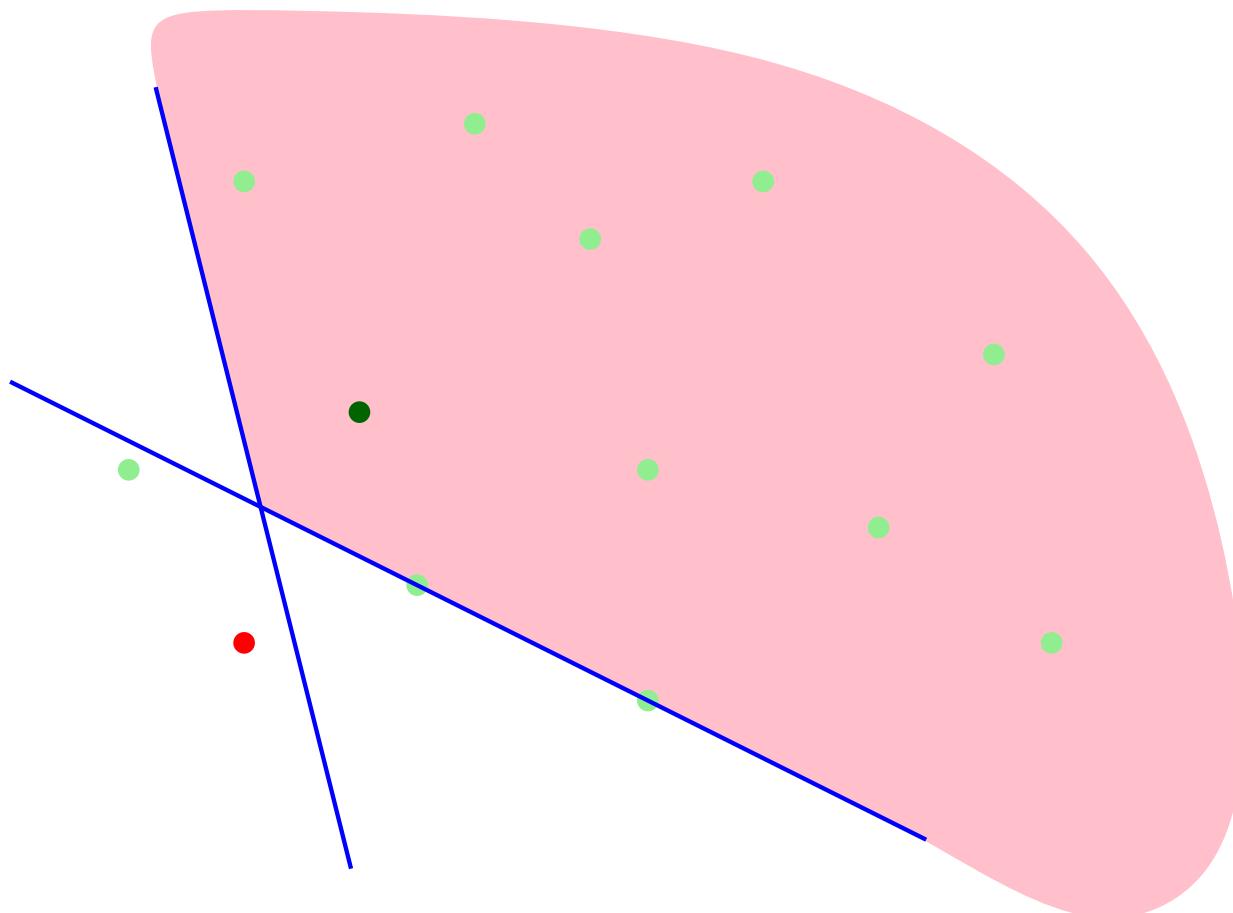


Faces of Voronoi diagram (1/5)



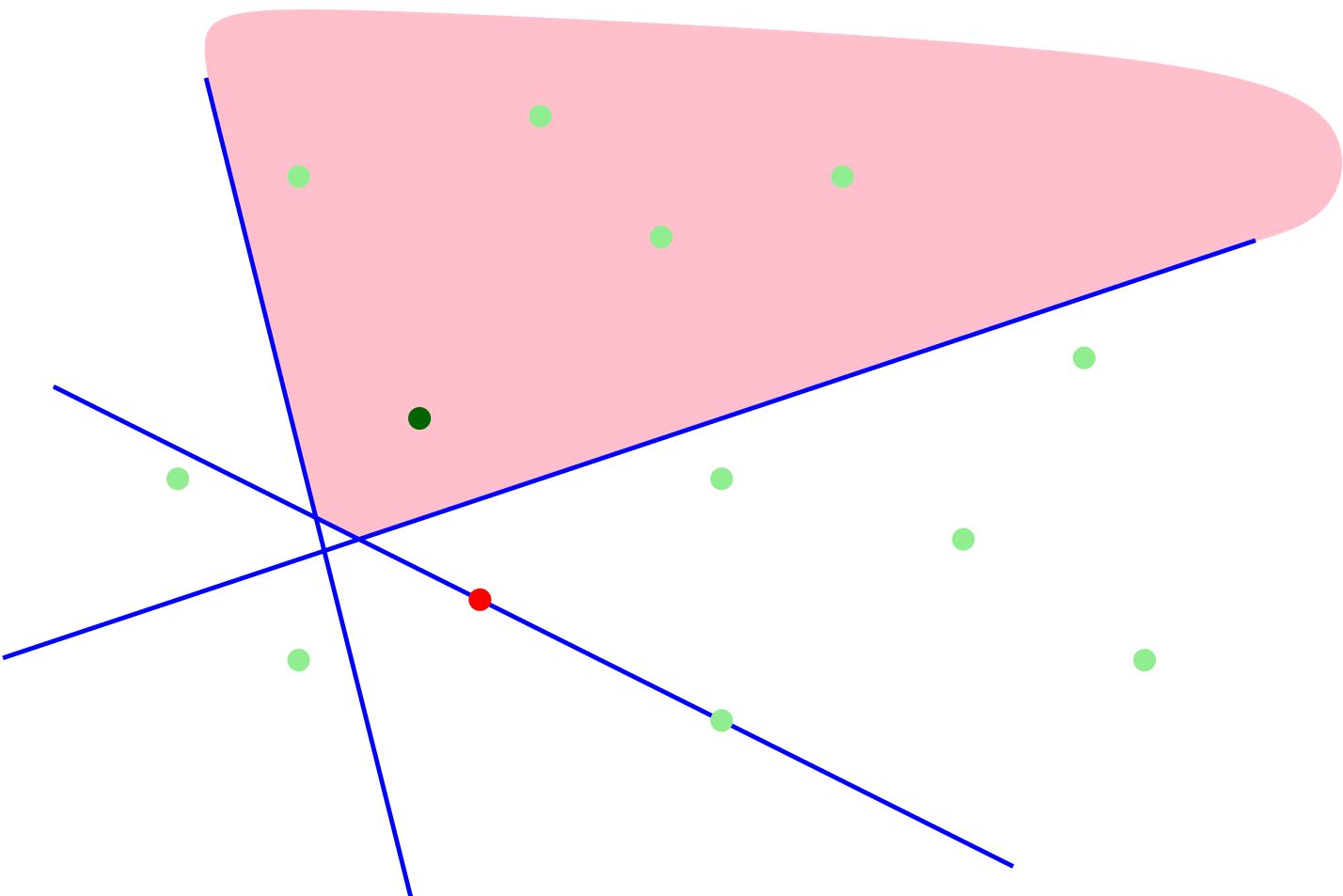
We consider all remaining points (differing from target), [one at a time](#)
Perpendicular bisector between target and red points → defines two half spaces

Faces of Voronoi diagram (2/5)



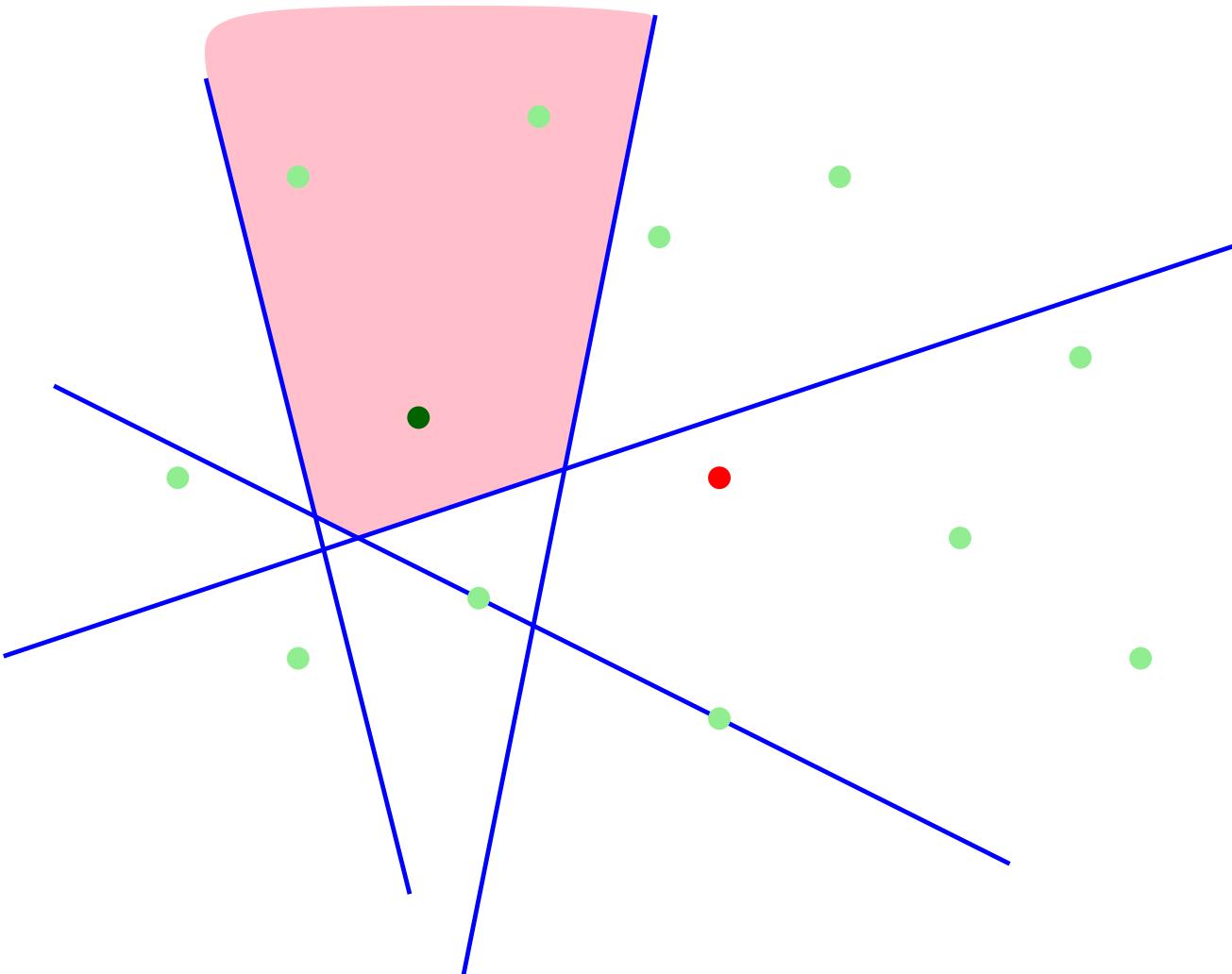
Perpendicular bisector between red and target points → defines two half spaces
Intersect those half spaces with previous one including the target point

Faces of Voronoi diagram (3/5)



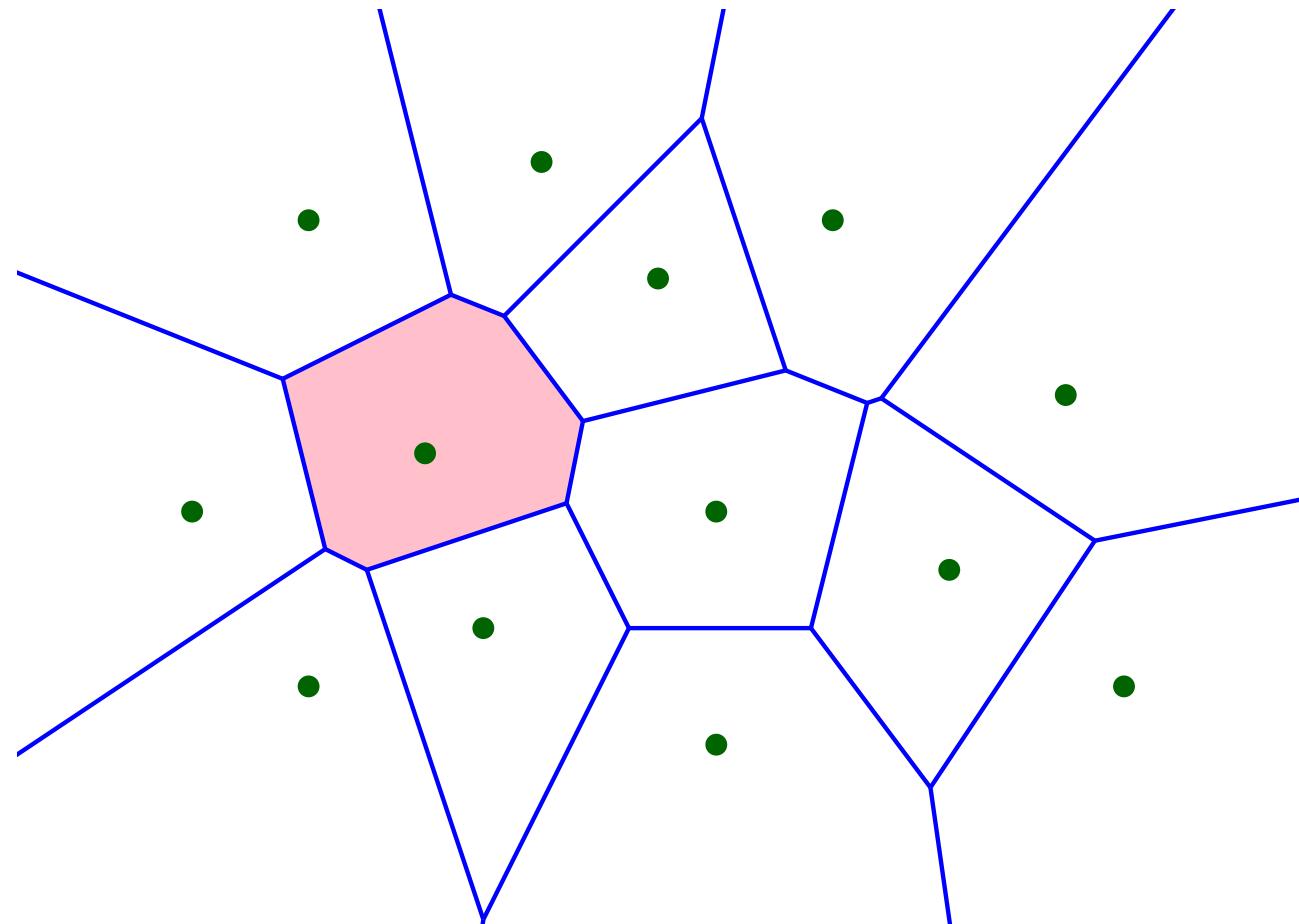
Perpendicular bisector between red and target points → define two half spaces
Intersect those half spaces with previous region including the target point

Faces of Voronoi diagram (4/5)



Perpendicular bisector between red and target points → define two half spaces
Intersect those half spaces with previous region including the target point

Faces of Voronoi diagram (5/5)



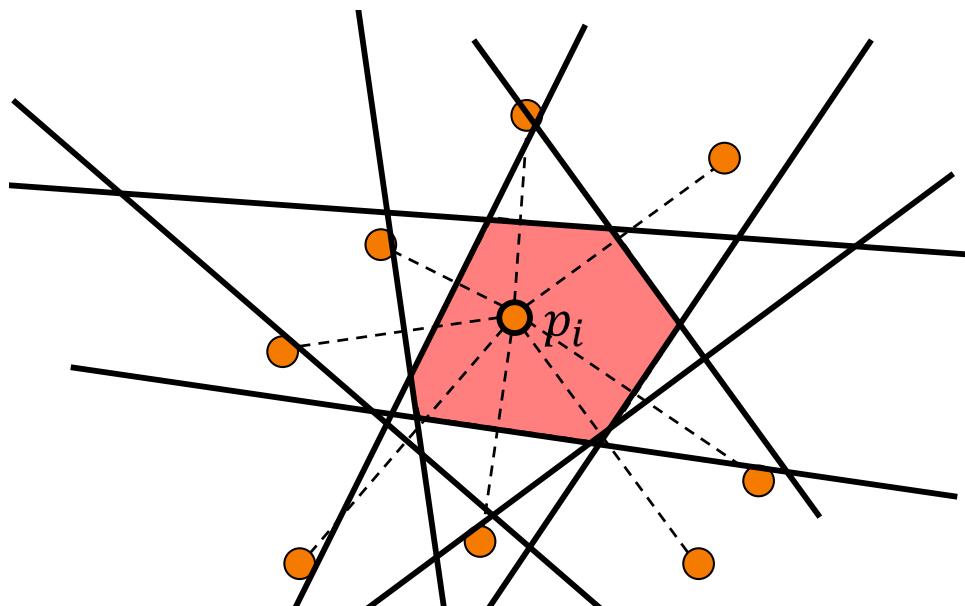
Upon cycling over all points differing from the target

The intersection of all the half spaces returns the Voronoi region for the target point

Voronoi diagram (5/6)

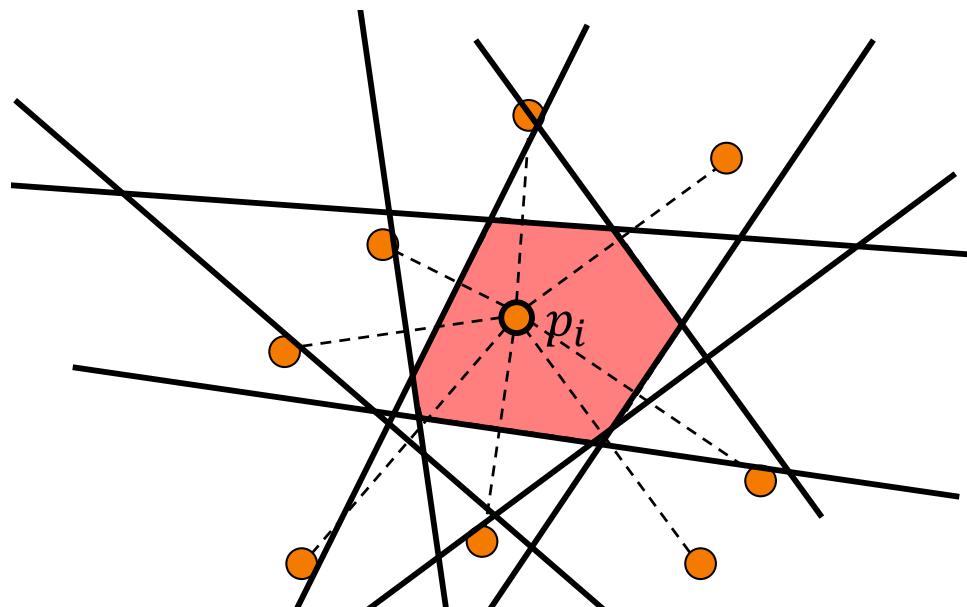
- **N points:** the Voronoi region R_i associated with point p_i corresponds to the intersection of the half spaces defined by the perpendicular bisectors

$$R_i = V(p_i) = \cap_{j \neq i} H(p_i, p_j)$$



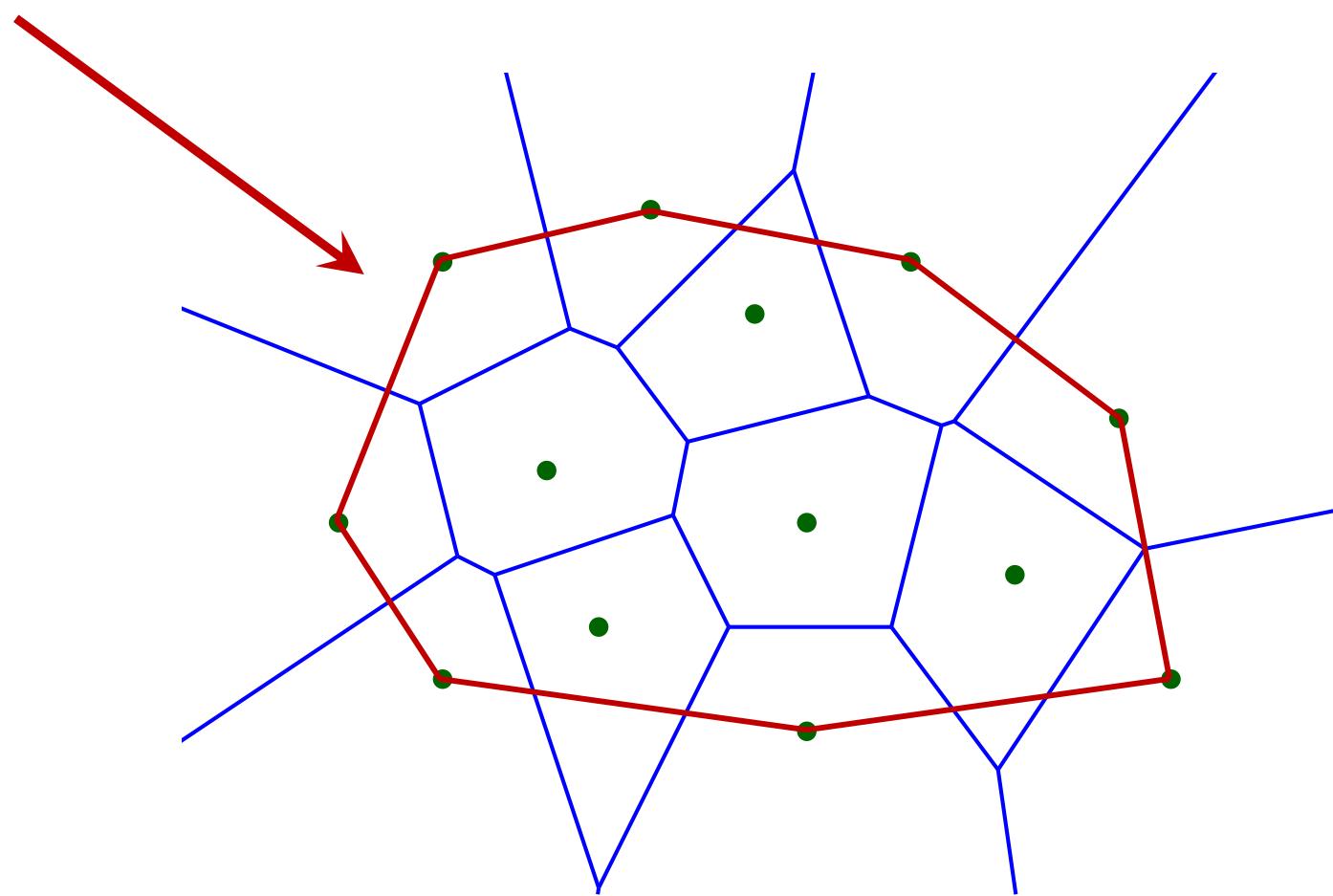
Voronoi diagram (6/6)

- Hence, Voronoi regions are *convex* polygons
- There is 1 Voronoi region for each point
- Voronoi faces can be *unbounded*



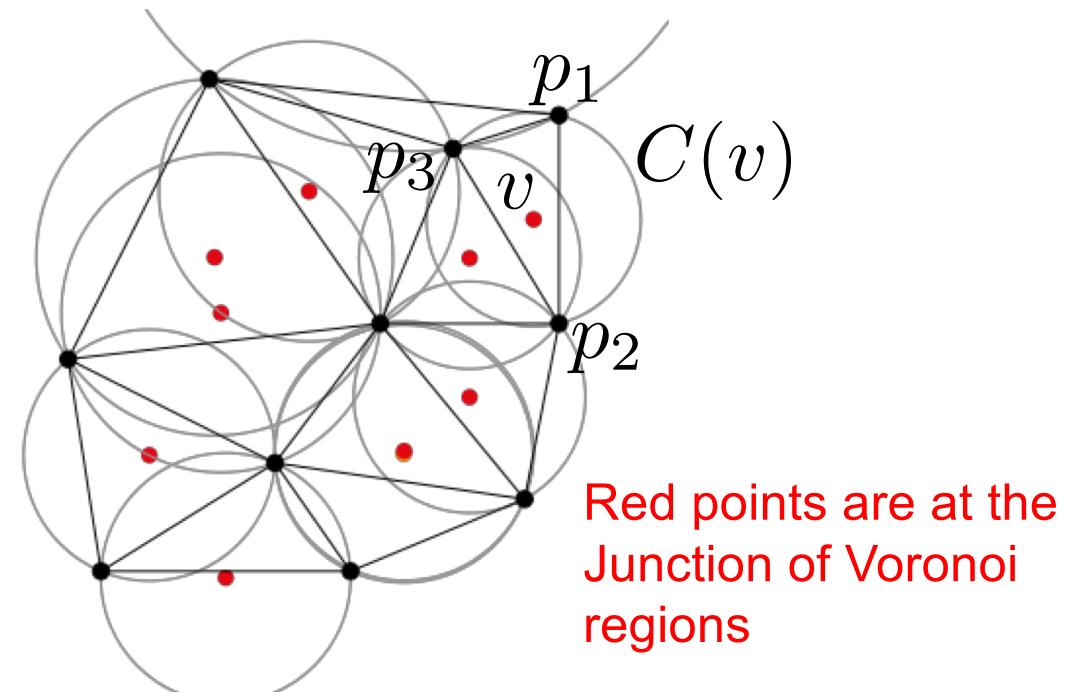
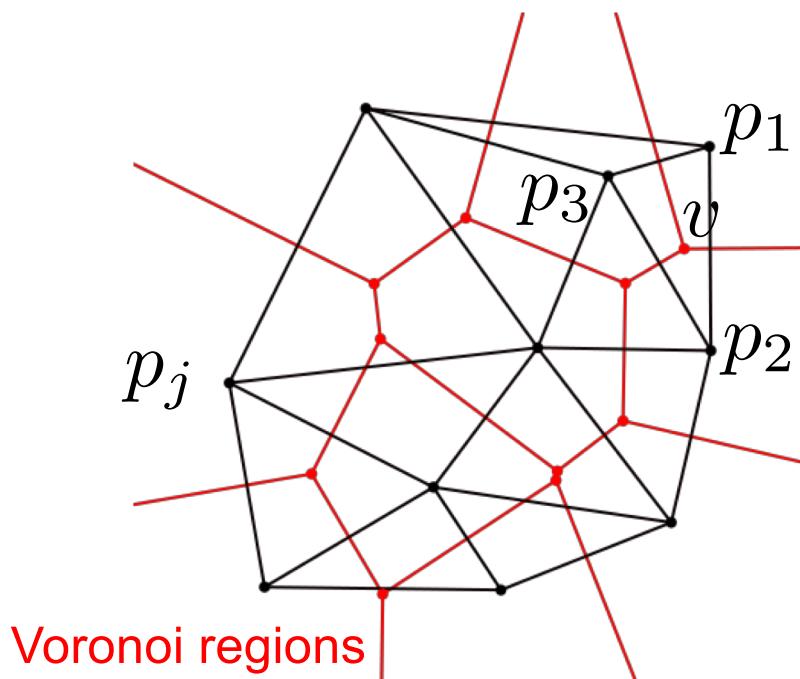
Properties of Voronoi diagrams (1/2)

- Given a set of data points \mathbf{P}
- **Convex hull of \mathbf{P} :** smallest convex set that contains \mathbf{P}



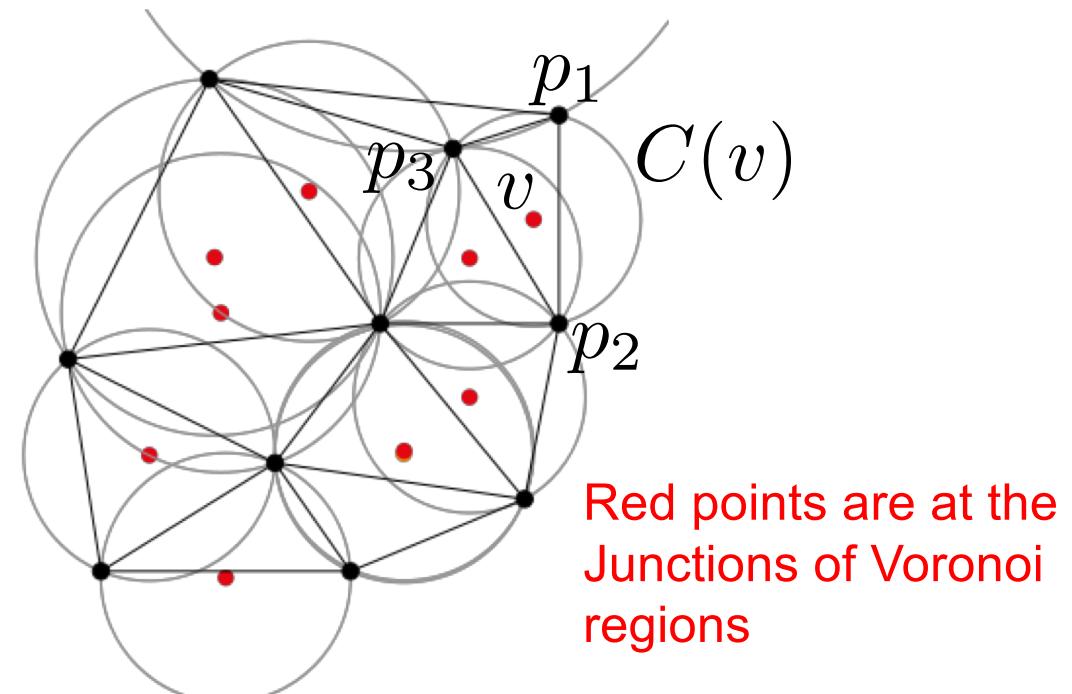
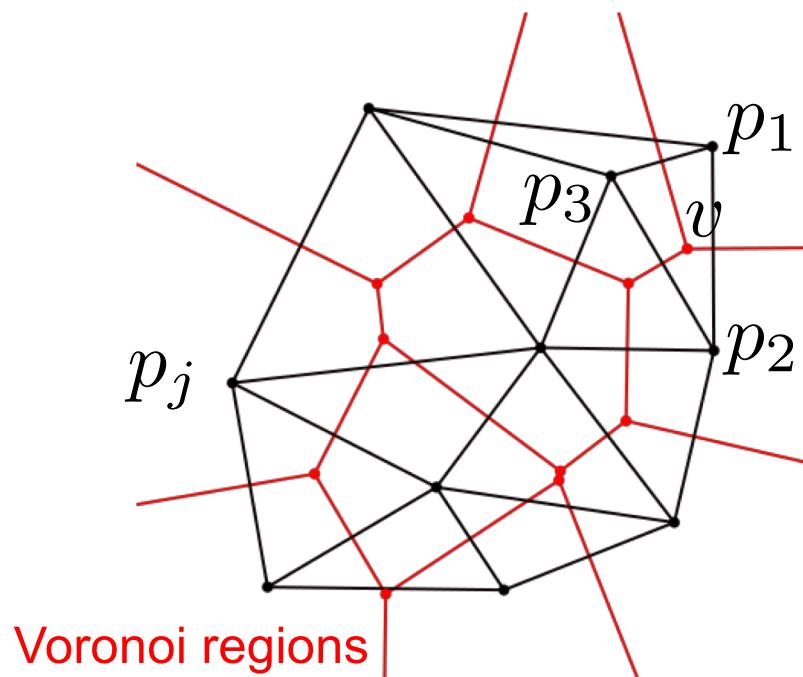
Properties of Voronoi diagrams (2/2)

- Each Voronoi region is **convex**
- $R_i = V(p_i)$ is *unbounded* IFF p_i is on the convex hull of P
- If v is at the junction of $V(p_1), V(p_2), \dots, V(p_k)$ with $k \geq 3$
 - v is the center of a circle $C(v)$ with p_1, p_2, \dots, p_k on the boundary
 - the interior of $C(v)$ **contains no points of P**



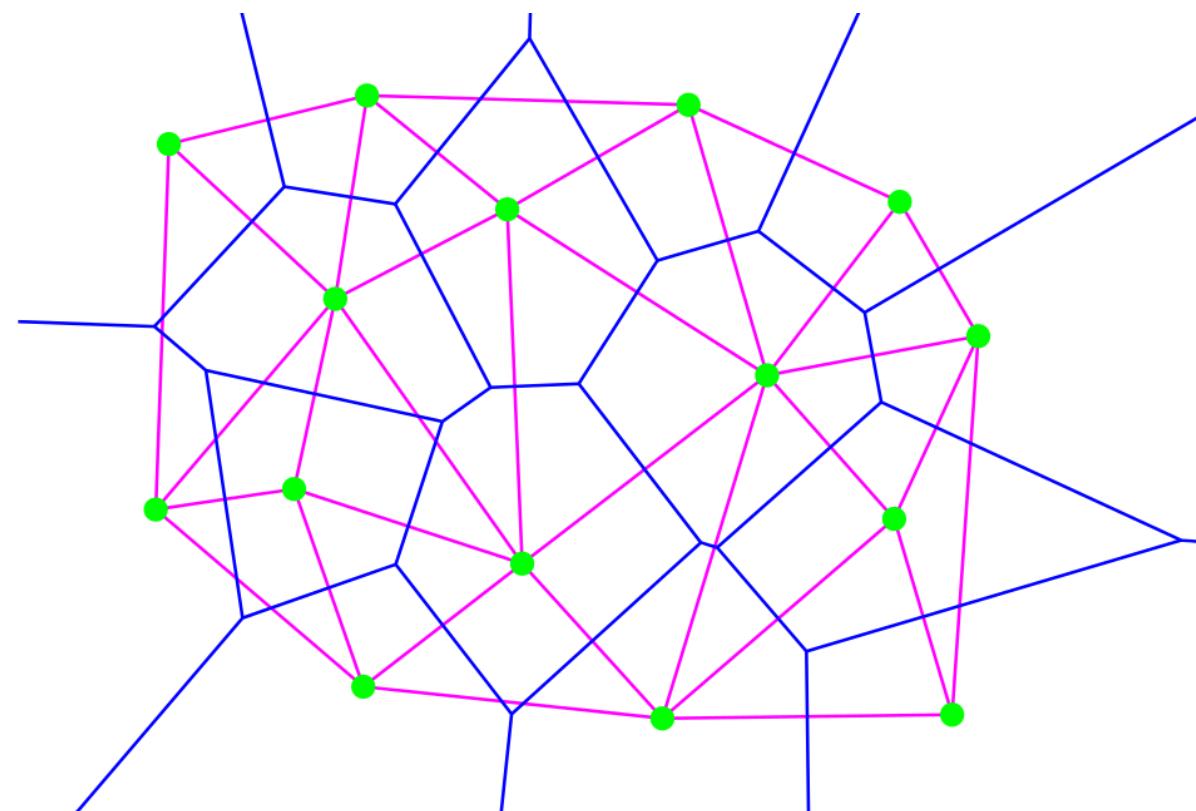
Delaunay triangulation

- It is the **straight line dual** of the Voronoi diagram
 - It is represented by the black straight lines in the **left plot** below
 - **Right plot:** we see that $\text{DT}(\mathbf{P})$ is a valid Delaunay triangulation
 - Starting from Voronoi diagram & connecting all the points *that share an edge* → we get the Delaunay triangulation



Voronoi diag. vs Delaunay triangulation

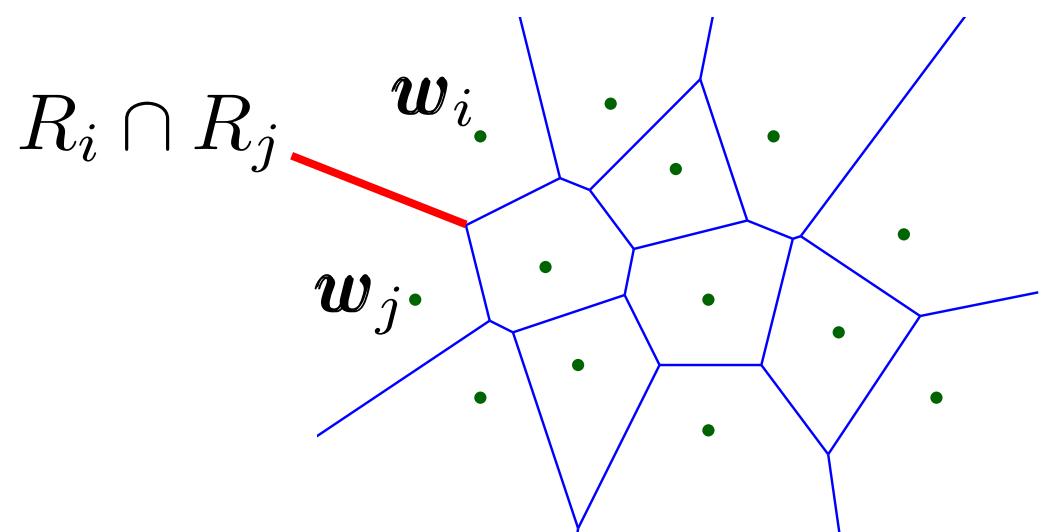
- The **Delaunay triangulation** of a discrete point set \mathbf{P} in general position corresponds to the **dual graph** of the **Voronoi diagram** for \mathbf{P}



Data points: **green**, Voronoi diagram: **blue**, Delaunay triangulation: **magenta**

Topology preserving structures

- We are given a **feature manifold** $M \subseteq \mathbb{R}^m$
- And a set of points $\mathbf{w}_i \in M$
- Any two points $\mathbf{w}_i, \mathbf{w}_j$ are adjacent if
 - The intersection of their Voronoi regions is non-empty, i.e.,
$$R_i \cap R_j \neq \emptyset$$
 - This means that they have a Voronoi edge in common



Adjacency in M

- **Definition.**
 - Let $M \subseteq \mathbb{R}^m$ be a given manifold
 - And let $S = \{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_N\}$ be a set of points $\mathbf{w}_i \in M$
 - The **Voronoi regions** of S are defined as $V(\mathbf{w}_i) \triangleq R_i$
 - Two points $\mathbf{w}_i, \mathbf{w}_j$ are adjacent **in M** if their masked Voronoi regions
 - Share (at least) an element $\mathbf{v} \in M$
 - Or equivalently if:

$$R_i^{(M)} = R_i \cap M \quad R_j^{(M)} = R_j \cap M$$
$$R_i^{(M)} \cap R_j^{(M)} \neq \emptyset$$

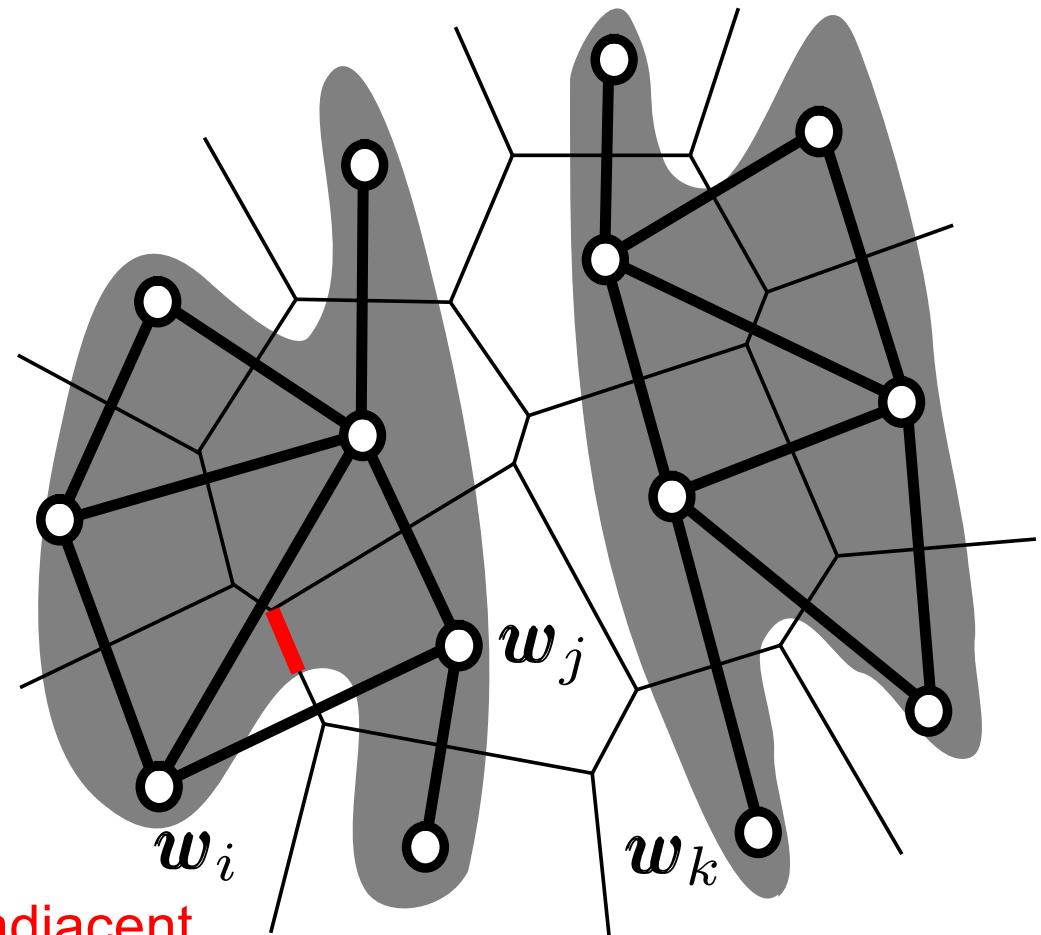
Adjacency in M

- **Example**

- M is the shaded area
- Is formed of two disjoint sets
- Here, w_i, w_j are adjacent

$$R_i^{(M)} \cap R_j^{(M)} \neq \emptyset$$

- However, w_j, w_k are NOT adjacent
 - Although their Voronoi regions have an edge in common
 - Their “masked” Voronoi regions (intersected with M) do not share any point



Topology preserving feature maps

- We are given a set of points

$$S = \{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_N\}$$

- These corresponds to centroids

- obtained, e.g., through SOM
- from an input distribution

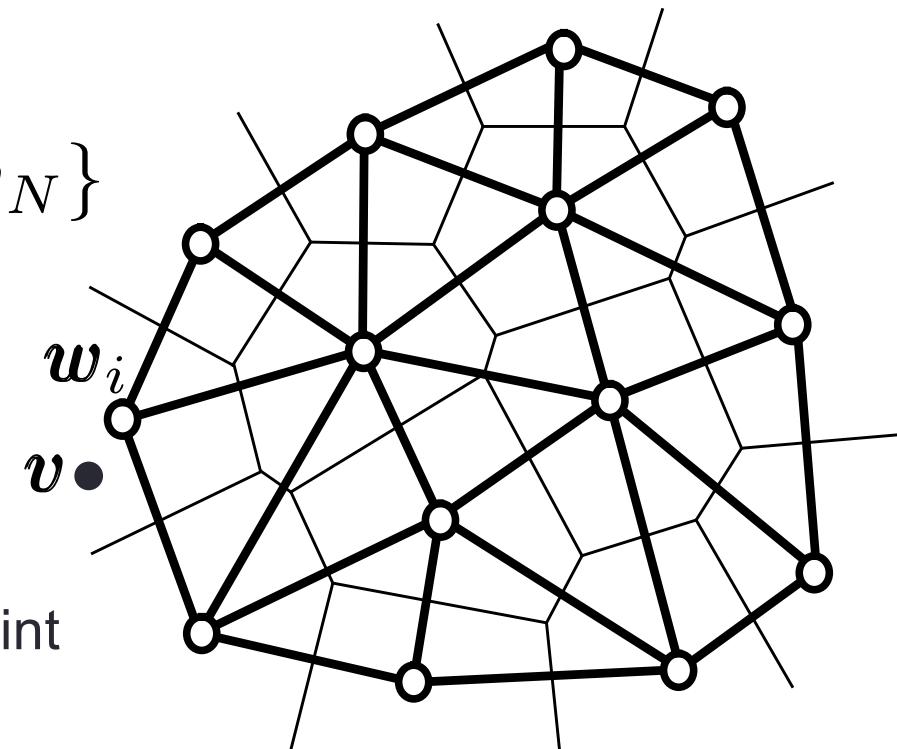
- If we sample \mathbf{v} from M

- this is associated with the closest point
- in the figure this closest point is \mathbf{w}_i

- Now, we would like to join the points in S with edges

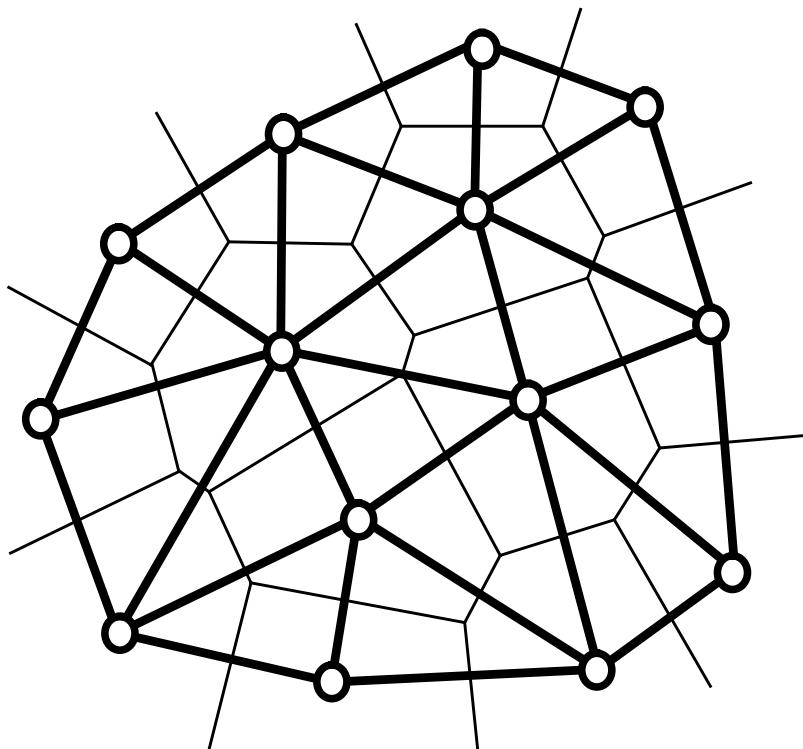
- to form a graph G (in the figure is obtained via Delaunay triangulation)

- **Definition.** Graph G forms a topology preserving map of M if points that are adjacent (edge) in G are also adjacent in M

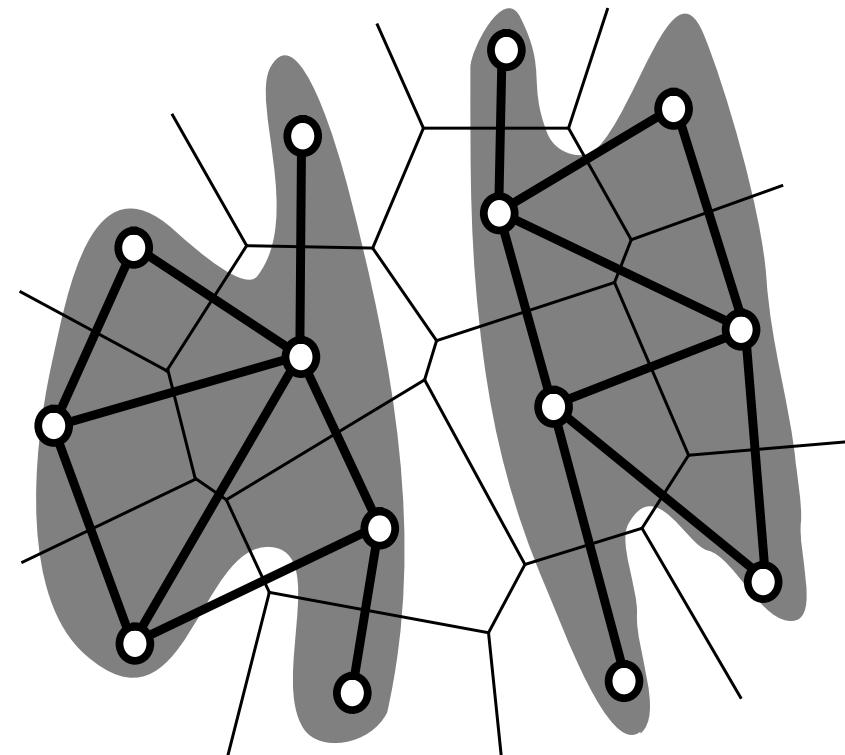


Topology preserving feature maps

- Hence, we see that for the given manifold M (shaded area)
 - Graph G on the left **is NOT** topology preserving
 - Graph G on the right **IS** topology preserving



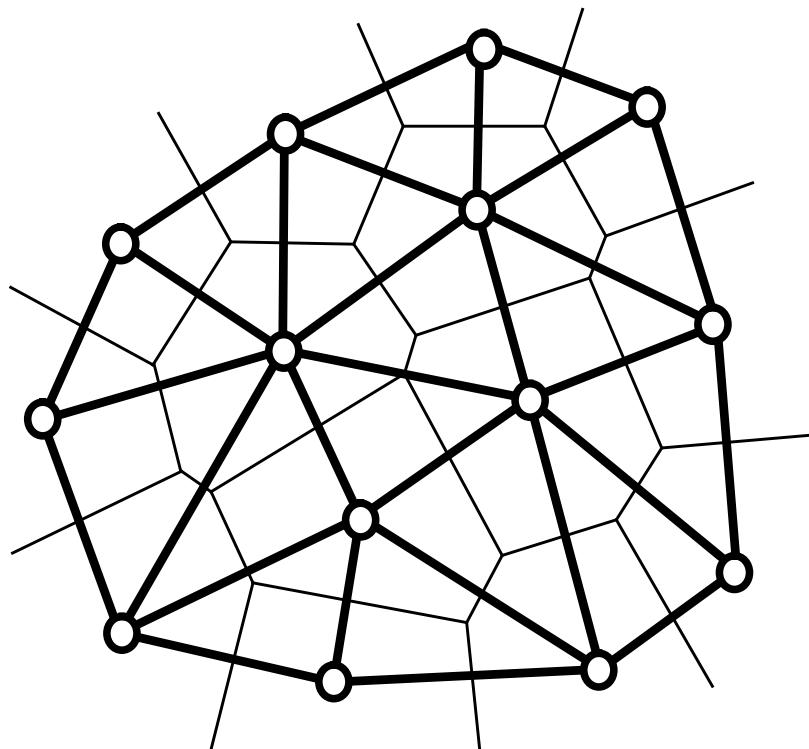
Delaunay triangulation



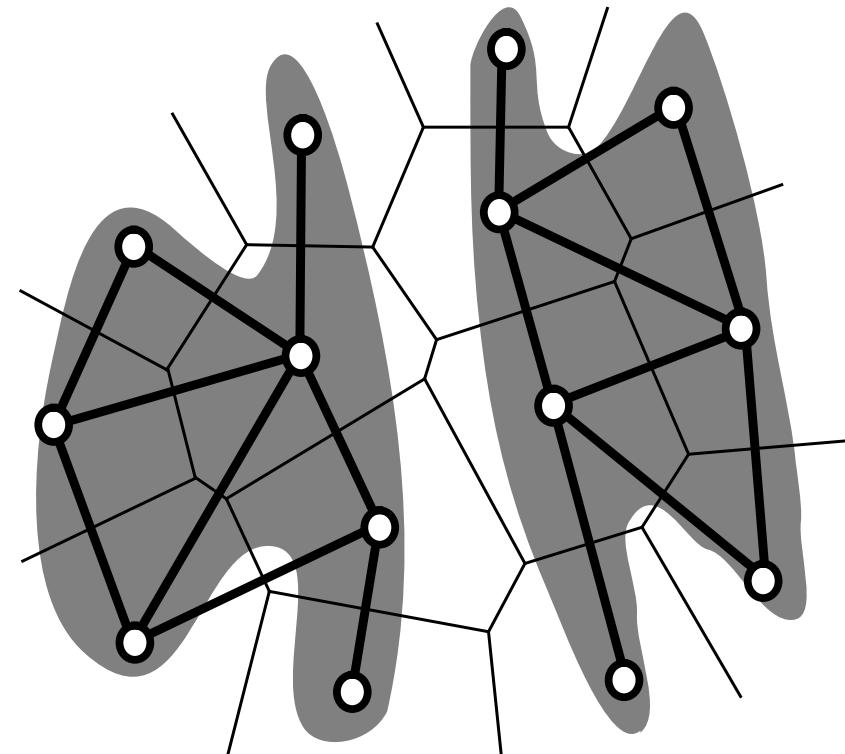
Masked Delaunay triangulation

Fundamental result

- The masked Delaunay triangulation
 - Leads to a topology preserving graph G
 - Preserves the topology relations in M



Delaunay triangulation

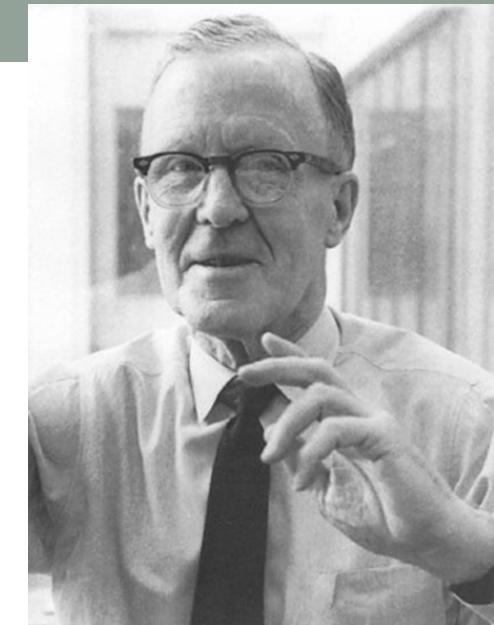


Masked Delaunay triangulation

Our objective

- The points in S are given $S = \{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_N\}$
- However, the edges connecting them **to form G are unknown**
- What we can do is sampling points \mathbf{v} from M, *one at a time*
- **From this procedure**
 - We would like to automatically understand which edges are valid
 - So as **to form a topology preserving map G**
- **From what we have seen this means that**
 - The edges that we form in G must also be valid in M
 - **That is: “adjacency” in G must correspond to “adjacency” in M and vice versa**

Hebbian learning (1/3)



- We assume a set of neural units $i = 1, 2, \dots, \ell$
- Neural units are
 - Initially isolated
 - Develop lateral connections (edges) with neighbors
 - A unit i develops a lateral connection to a unit j through a **synaptic link**
- **Lateral connections**
 - Represent *edges in the graph G*
 - Are represented through a **connection strength matrix \mathbf{C}** with elements $C_{ij} > 0$
 - If $C_{ij}=0$ then units i and j **are not connected**
- The dynamics by which nearby units develop neuronal connections
 - Were first studied by Hebb [Hebb49] – “*neurons that fire together wire together*”
- [Hebb49] D.O. Hebb: The Organization of Behavior, Wiley, New York, 1949.

Donald O. Hebb - psychologist
(1908 – 1985)

Hebbian learning (2/3)

'presynaptic' and 'postsynaptic' are used to indicate two neurons that are connected. Information flow in the nervous system basically goes *one way*. If one neuron fires (presynaptic cell) it can chemically activate another cell on which it "synapses" (the postsynaptic cell). A synapse is a neural pathway.

- Hebb's postulate:

- Unit i increases the strength of its synaptic link to unit j IF they are both active
 - "Neurons that fire together wire together"
 - This means that activities at units i and j are correlated
- It is a method to determine how to alter the weight between two neurons
- Hence, the change in the synaptic strength is proportional to:

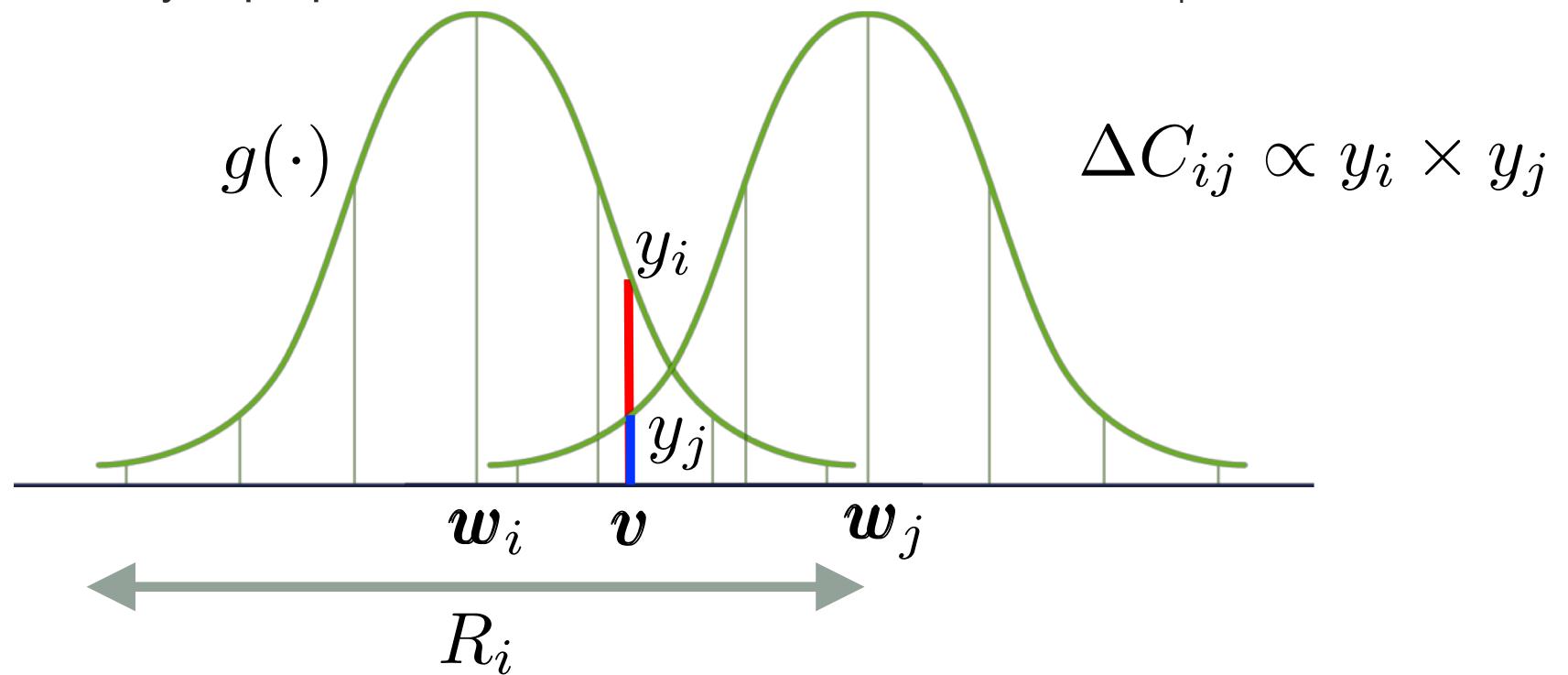
$$\Delta C_{ij} \propto y_i \times y_j$$

- y_i and y_j indicate the intensity of the synaptic activity at units i and j

$$C_{ij} \leftarrow C_{ij} + \eta y_i y_j$$

Intermission: receptive field

- Each neuron i , has an associated “receptive field” R_i
 - The center of the receptive field is the **neuron’s weight**
 - If any new input \mathbf{v} falls within that region \mathbf{w}_i
 - Neuron i fires with intensity: $y_i \propto g(\|\mathbf{w}_i - \mathbf{v}\|)$
 - The intensity is proportional to the distance between \mathbf{v} and \mathbf{w}_i



Hebbian learning (3/3)

- Verbatim application of Hebb's rule
 - Leads to edges among all units
 - With strength C_{ij} proportional to the overlap of receptive fields

$$C_{ij} \propto \int_{R_i \cup R_j} f(v)g(\|\mathbf{w}_i - \mathbf{v}\|)g(\|\mathbf{w}_j - \mathbf{v}\|)d\mathbf{v}$$

- This is not what we want
 - Activation of all units at every new input
 - Does not involve competition
 - Competition leads to meaningful structures !!!

Competitive Hebbian learning

- Intuition
 - Competition often leads to interesting structure (see, e.g., SOM)
- SOM
 - Competition is among the outputs of the units y_i (“winner-take-all” network)
- Here
 - Competition is introduced among synaptic links
 - The correlated output activities are computed as: $Y_{ij} = y_i y_j$
 - The link with the highest Y_{ij} wins and takes-all, the remaining ones take nothing
- Competitive Hebbian learning *among synaptic links*

$$\Delta C_{ij} \propto \begin{cases} y_i y_j & y_i y_j \geq y_k y_l, \forall k, l \\ 0 & \text{otherwise} \end{cases}$$

Competitive Hebbian Learning (CHL)

- Result no. 1 – “CHL leads to Delaunay triangulation”
 - If we sequentially present to the network input patterns \mathbf{v}
 - Having a PDF $f(\mathbf{v})$ which has support that is non-zero everywhere in \mathbb{R}^m

- Asymptotically, after a large number of input patterns have been collected
 - The weights converge to:

$$\lim_{n \rightarrow +\infty} \theta(C_{ij}(n)) = A_{ij}$$

- where: $\theta(\cdot)$ if the Heaviside step function
- Adjacency matrix (V_i is the Voronoi region for neuron i)

$$A_{ij} = \begin{cases} 1 & V_i \cap V_j \neq \emptyset \\ 0 & V_i \cap V_j = \emptyset \end{cases} \quad \text{This is the adjacency matrix of the Delaunay triangulation!!!}$$

Definition

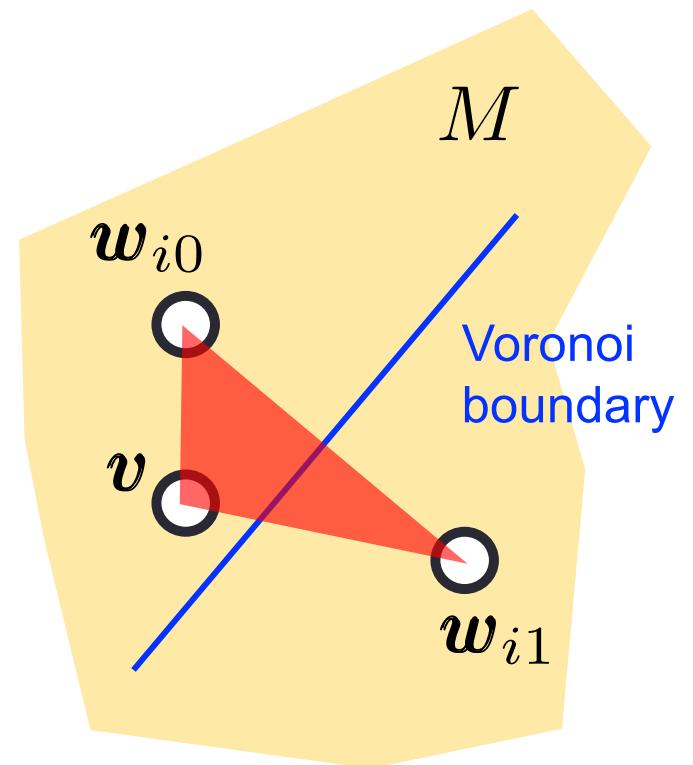
- Definition – S is dense on M

- Let v be an input pattern
- w_{i0} is the closest weight (neuron)
- w_{i1} is the second-closest weight (neuron)
- The triangle formed by
 - v , w_{i0} and w_{i1} called

$$\Delta(v, w_{i0}, w_{i1})$$

- Is entirely contained in M

$$\Delta(v, w_{i0}, w_{i1}) \subseteq M$$



Definition “second order polyhedron”

- Definition – Second order polyhedron V_{ij}
- Second order Voronoi polyhedron V_{ij} with $i, j = 1, 2, \dots, \ell$
- Is given by all the vectors $\mathbf{v} \in \mathbb{R}^m$
- For which $\mathbf{w}_i, \mathbf{w}_j$ are the two closest points to $\mathbf{v} \in \mathbb{R}^m$

$$V_{ij} = \{\mathbf{w} \in \mathbb{R}^m \mid \|\mathbf{v} - \mathbf{w}_i\| \leq \|\mathbf{v} - \mathbf{w}_k\| \cap \|\mathbf{v} - \mathbf{w}_j\| \leq \|\mathbf{v} - \mathbf{w}_k\| \forall k \neq i, j\}$$

Theorem 1 (without proof):

$$V_i \cap V_j \neq \emptyset \iff V_{ij} \neq \emptyset$$

$f(\mathbf{v})$ with support in M

- What happens when the input patterns PDF has support in M
(M is a subset of the input vector space)
- In this case
 - For some $V_{ij} \neq \emptyset$ (from Theorem 1, V_i and V_j have points in common)
 - We may have that:
$$V_{ij}^{(M)} \triangleq V_{ij} \cap M = \emptyset$$
(they do not have points in common in M)
- This means that
$$\int_{V_{ij}} f(\mathbf{v}) d\mathbf{v} \rightarrow 0$$
- That using CHL means that link (i,j) is not activated, i.e., $A_{ij} = 0$
 - This follows from Result 1, as none of the common points in the intersection between V_i and V_j has probability >0 of being generated

Theorem 2 (with proof)

- “CHL leads to perfect topology preserving maps”
- This can be rephrased as:

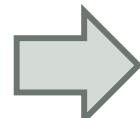
$$V_i^{(M)} \cap V_j^{(M)} \neq \emptyset \iff V_{ij}^{(M)} \neq \emptyset$$

- In words:
 - “masked Voronoi regions i and j have points in common (a link will be then created by CHL) iff the masked second order polyhedron is not empty”
- From previous slide: if the second order Voronoi polyhedron is $V_{ij}^{(M)} = \emptyset$
 - CHL will not create links between nodes i and j
 - But from Theorem 2 this implies that $V_i^{(M)} \cap V_j^{(M)} = \emptyset$
 - Which means that CHL is topology preserving (links will only be created between nodes that have Voronoi region with non-zero intersection in M, i.e., with common points in M)

Theorem 2

- We need to prove that: $V_i^{(M)} \cap V_j^{(M)} \neq \emptyset \iff V_{ij}^{(M)} \neq \emptyset$

Sufficient condition



If $V_i^{(M)} \cap V_j^{(M)} \neq \emptyset \Rightarrow$ there exists \mathbf{v} s.t.

$$\mathbf{v} \in V_i^{(M)}, \mathbf{v} \in V_j^{(M)}$$

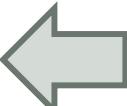
This means that:

$$\|\mathbf{v} - \mathbf{w}_i\| = \|\mathbf{v} - \mathbf{w}_j\| \leq \|\mathbf{v} - \mathbf{w}_k\|, \forall k \in S$$

Which exactly means that: $\mathbf{v} \in V_{ij}^{(M)}$

Theorem 2

- We need to prove that: $V_i^{(M)} \cap V_j^{(M)} \neq \emptyset \iff V_{ij}^{(M)} \neq \emptyset$

Necessary condition 

If $V_{ij}^{(M)} \neq \emptyset$ then there exists a vector \mathbf{v}^* in the input space such that

$$\mathbf{v}^* \in V_{ij}^{(M)} \text{ s.t. } \Delta(\mathbf{v}^*, \mathbf{w}_i, \mathbf{w}_j) \subseteq M$$

(this holds as we assume that S is dense in M)

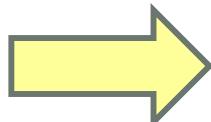
Note that, for each $\mathbf{v}^* \in V_{ij}^{(M)}$

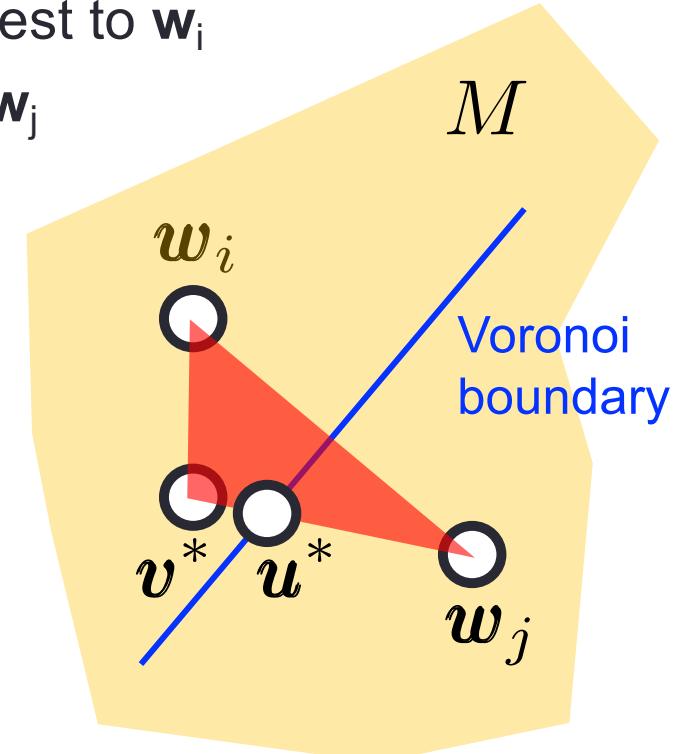
$\mathbf{w}_i, \mathbf{w}_j$ are the two closest neighbors

Theorem 2: continued

- Without loss of generality, assume that \mathbf{v}^* is closest to \mathbf{w}_i
- Let us take \mathbf{u} in the segment connecting \mathbf{v}^* and \mathbf{w}_j
- Since
 - For $\mathbf{u}=\mathbf{v}^* \rightarrow \mathbf{w}_i$ is closest to \mathbf{u}
 - For $\mathbf{u}=\mathbf{w}_j \rightarrow \mathbf{w}_j$ is closest to \mathbf{u}
- Then, there must exist an intermediate point \mathbf{u}^*
 - In the segment $\text{seg}(\mathbf{v}^*, \mathbf{w}_j)$
 - Such that:
$$\|\mathbf{u}^* - \mathbf{w}_i\| = \|\mathbf{u}^* - \mathbf{w}_j\|$$
- Hence, we have:

$$\mathbf{u}^* \in V_i, \mathbf{u}^* \in V_j, \mathbf{u}^* \in \Delta(\mathbf{v}^*, \mathbf{w}_i, \mathbf{w}_j)$$

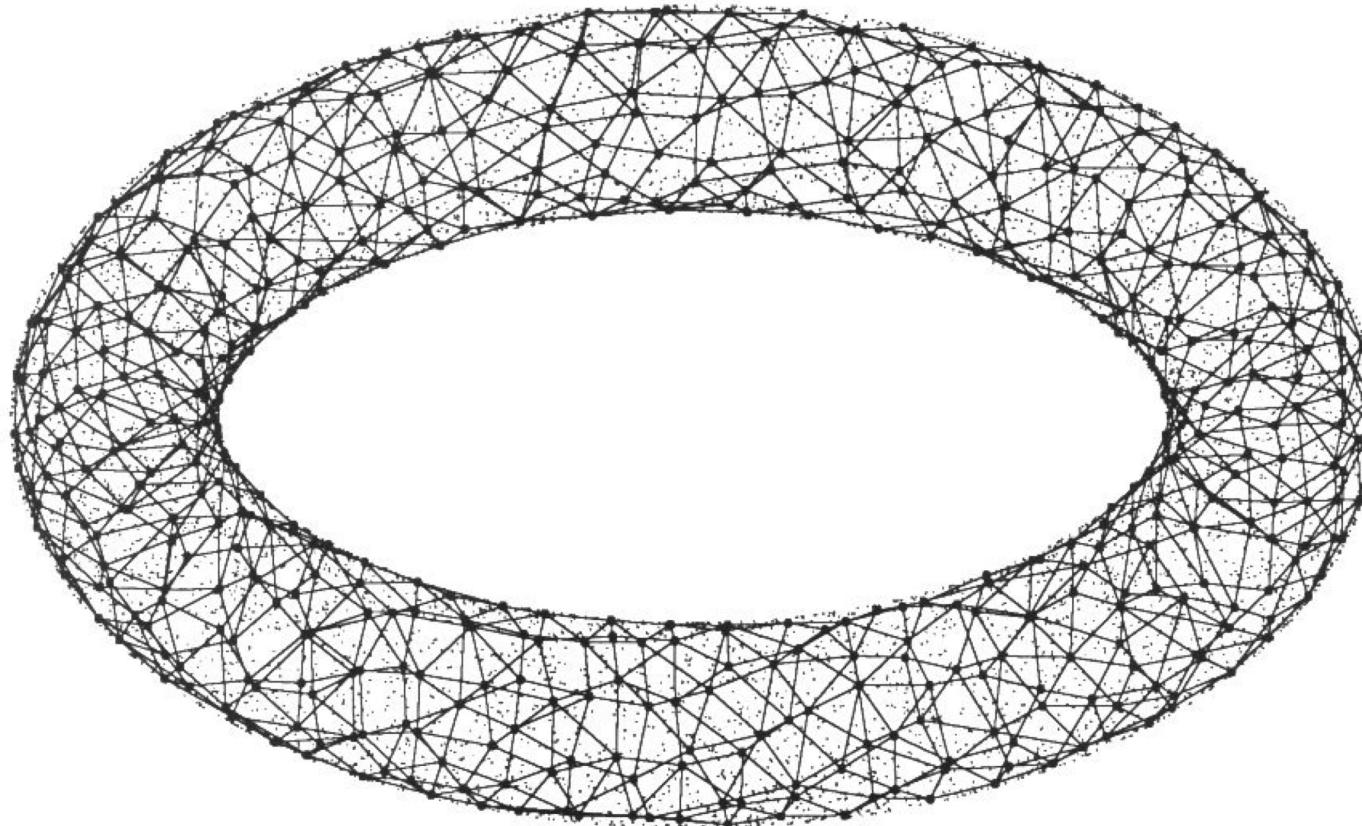
 $M \cap (V_i \cap V_j) \neq \emptyset \Rightarrow V_i^{(M)} \cap V_j^{(M)} \neq \emptyset$



QED

Example

- Points in S pre-obtained through a vector quantization algorithm
- Links among them obtained using CHL
- Structure of the 3D manifold is preserved



Reference paper

Result no. 1: Theorem 1 in [\[Martinetz93\]](#)

Result no. 2: Theorem 2 in [\[Martinetz93\]](#)

[\[Martinetz93\]](#) Thomas Martinetz, “Competitive Hebbian Learning Rule Forms Perfectly Topology Preserving Maps,” *Proceedings of the International Conference on Artificial Neural Networks (ICANN)*, 1993.
(375 citations)

Growing Neural Gas (GNG) networks

- Unsupervised dimensionality reduction technique
 - Performs unsupervised clustering (or *vector quantization*)
- Number of neurons (*centroids*)
 - Not fixed a priori but dynamically tuned based on input statistics
- Uses Competitive Hebbian Learning
 - To connect neighbors into a topology preserving manner
- Synaptic weights
 - Are not weighted
 - Their *sole purpose* is the definition of topological structure
 - An *age count* is maintained for each
 - Weights can be added and removed
- Neurons
 - *Addition and removal* are possible as learning evolves

GNG – basic construct

- Basic construct
 - A network of nodes (neurons)
 - Each node has a weight $\mathbf{w}_i \in \mathbb{R}^m$
 - Connections are *not weighted*
 - Their sole purpose is the definition of topological structure
- There is a (possibly infinite) number of m-dimensional input patterns
 $\mathbf{v} \in \mathbb{R}^m$
- Input patterns obey some unknown probability density function
 $f(\mathbf{v})$
- Main idea
 - Successively add new units to an initially small network by evaluating local statistical measures gathered during previous adaptation steps. Likewise, adapt the weight of existing units, possibly adding new or removing old ones.

GNG – steps of the algorithm (1/3)

- 0) Start with two nodes with random weights $\mathbf{w}_a, \mathbf{w}_b$
- 1) Generate an input pattern \mathbf{v} according to $f(\mathbf{v})$
- 2) Find the nearest and second-nearest nodes n_1, n_2
 - Closeness is measured as Euclidean distance between their weight and \mathbf{v}
- 3) Increment the age of all edges emanating from n_1
- 4) Add squared distance between the input pattern and the nearest node to a local variable

$$\text{error}(n_1) = \text{error}(n_1) + \|\mathbf{w}_{n_1} - \mathbf{v}\|^2$$

- 5) Move n_1 and its topological neighbors (all nodes that are linked with it) towards \mathbf{v} by fractions ϵ_b, ϵ_n of their “distance” from \mathbf{v}

$$\Delta \mathbf{w}_{n_1} = \epsilon_b (\mathbf{v} - \mathbf{w}_{n_1})$$

$$\Delta \mathbf{w}_{n_j} = \epsilon_n (\mathbf{v} - \mathbf{w}_{n_j}), \text{ for all neighbors } n_j \text{ of } n_1$$

GNG – steps of the algorithm (2/3)

- 6) if n_1, n_2 are already connected by an edge, set its age to zero, if this edge does not exist, create it (with zero age)
- 7) Remove edges with age larger than a_{\max} . If this results in nodes having no emanating edges, remove these nodes as well
- 8) Structure refinement: if number of input patterns that were generated is an integer multiple of a parameter λ , insert a new node r as follows
 - Determine the node q with maximum accumulated error
 - Insert the new node r halfway between q and its neighbor with the largest accumulated error variable

$$\mathbf{w}_r = 0.5(\mathbf{w}_q + \mathbf{w}_f)$$

- Insert edges connecting the new node r with nodes q and f, and remove the original edge between q and f
- Decrease the accumulated errors of q and f, by multiplying them by a constant α . Initialize the error variable of node r with the new value of the error variable of node q

GNG – steps of the algorithm (3/3)

- 9) decrease the error variables at all nodes, by multiplying them by a constant β
- 10) if some stopping criterion is not yet met (i.e., network size or some error measure, e.g., maximum error smaller than a threshold): GO TO 1

GNG – observations (1/2)

- The GNG algorithm **leads to a general movement** of all nodes (their weights) towards those areas of the input space where the input signal comes from (see **step 5**), i.e., where $f(\mathbf{v}) > 0$
- The **insertion of edges** (see **step 6**) between nearest and second-nearest node *generates a connection of the “induced Delaunay triangulation”* with respect to the current position of all nodes
- The **removal of edges** (**step 7**) is needed to **obtain adaptation of the topology, when a region is no longer excited, get rid of it !!!**
 - This is achieved through local edge aging around the nearest node (**step 3**), combined with age resetting of those edges that already exist between nearest and second-nearest nodes (**step 6**). This follows Hebbian learning, i.e., **only the winning edge is reinforced** (by resetting its age)

GNG – observations (2/2)

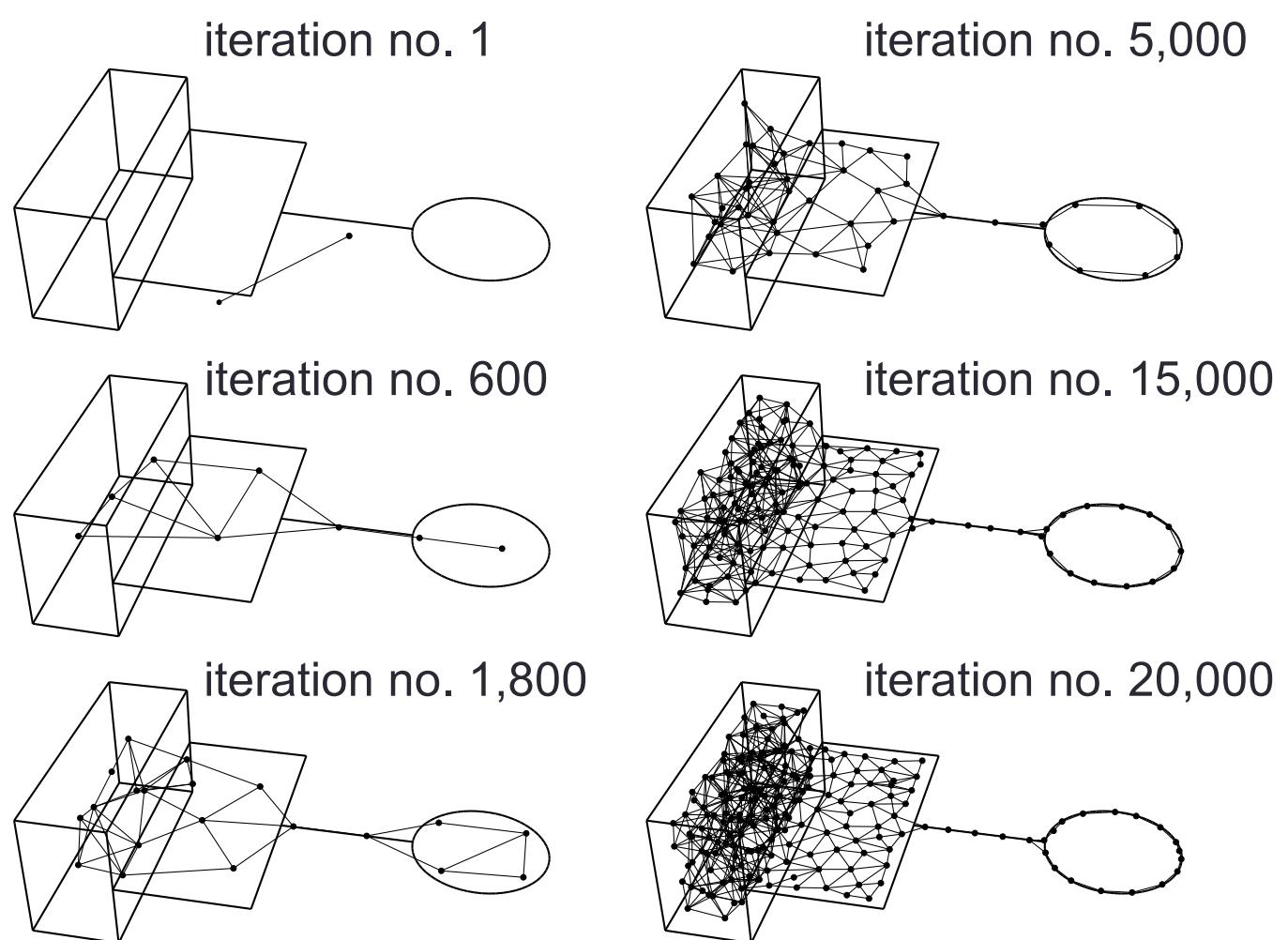
- The accumulation of squared error distances (step 4) during the adaptation helps identify those regions where the mapping causes a high error. To reduce this error, new nodes are inserted in such regions
- The maximum age a_{\max} determines when an edge has to be removed
 - Age tells how many steps have been elapsed since the last activation of one of its end nodes
 - In other words, this represents the last time that the corresponding region in space was hit by an input pattern
 - Hence, increasing a_{\max} corresponds to tracking smaller values of $f(\mathbf{v}) > 0$ (whose patterns occur less often)

Example results (1/2)

- Evolution of the GNG

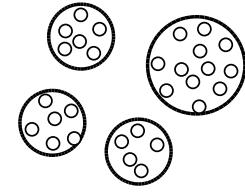
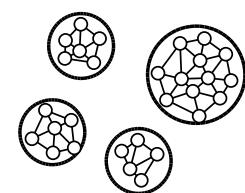
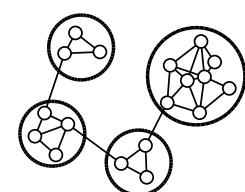
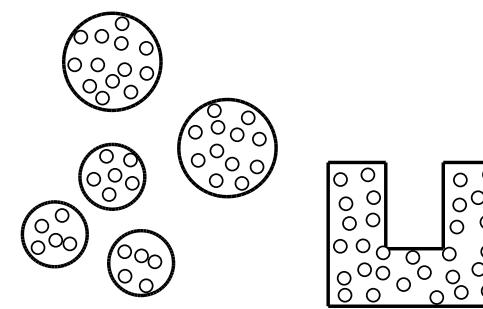
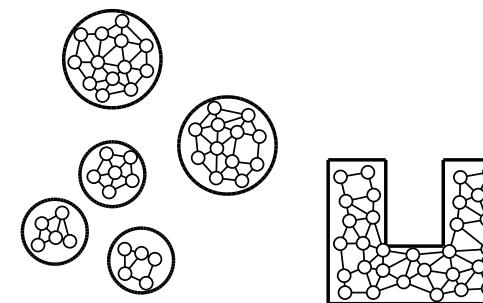
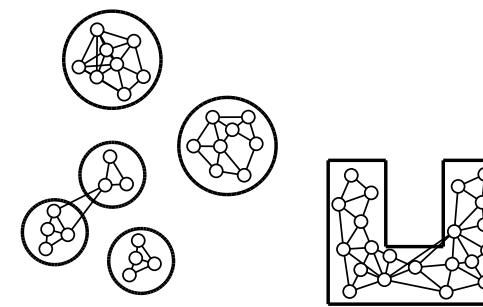
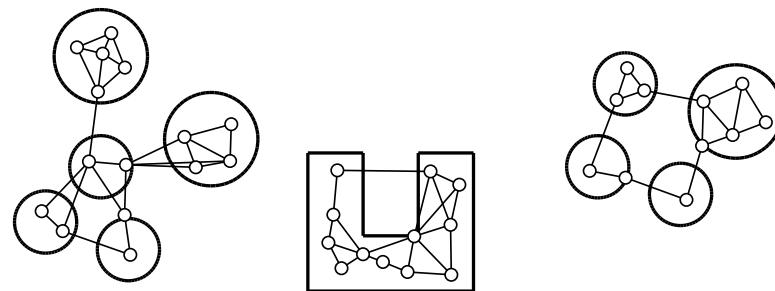
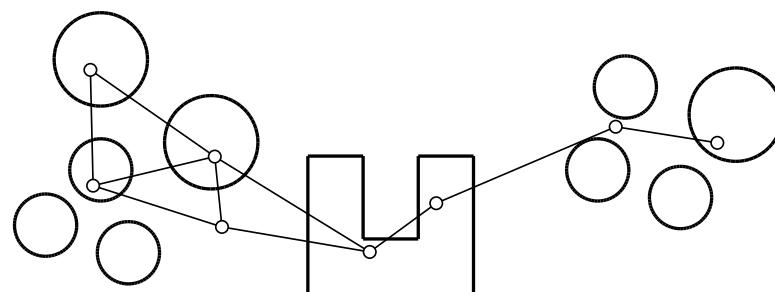
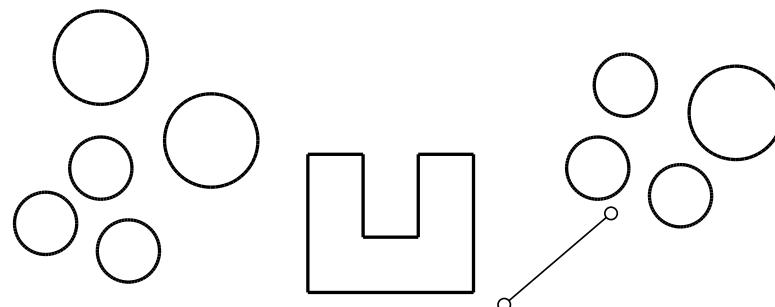
- Parameters are:

$$\left\{ \begin{array}{l} \lambda = 100 \\ \alpha = 0.5 \\ \beta = 0.995 \\ \epsilon_b = 0.2 \\ \epsilon_n = 0.006 \\ a_{\max} = 50 \end{array} \right.$$



Example results (2/2) (10,000 steps, 100 nodes at the end)

“growing neural gas”
(uses “competitive Hebbian learning”)



Gas neural networks bibliography

Slides are based on:

- [4] Thomas Martinetz, "Competitive Hebbian Learning Rule Forms Perfectly Topology Preserving Maps," *Proceedings of the International Conference on Artificial Neural Networks (ICANN)*, 1993. (375 citations)
- [5] Bernd Fritzke, "A Growing Neural Gas Network Learns Topologies," *Proceedings of the 7th International Conference on Neural Information Processing Systems (NIPS)*, 1994. (+2k citations)

Further readings:

- [6] Stephen Marsland, Jonathan Shapiro, Ulrich Nehmzow, "A self-organising network that grows when required," *Neural Networks Journal*, Vol. 15, No. 8-9, October 2002. (355 citations)

COMPETITIVE & UNSUPERVISED LEARNING FOR VECTOR QUANTIZATION

Michele Rossi
rossi@dei.unipd.it

Dept. of Information Engineering
University of Padova, IT

