

ECE/CS 559 Lecture 17/18 Th 10/24 / T 10/29

Last time: CNNs (max pooling, t-sne), PyTorch

① Unsupervised LearningBefore, data (x, y) . Now, data is just x .What is the task? Represent x "nicely".

- Clustering: Place x in one of a handful of "clusters"
 - Dimensionality reduction: Represent $x \in \mathbb{R}^n$ by $x \in \mathbb{R}^k$, $k \ll n$.
 - Density estimation: Find the distribution of x
- Relationship to generative models & feature extraction.
(represent \rightarrow create more) ("enough" to represent)

Lloyd's AlgorithmInput: Data, k

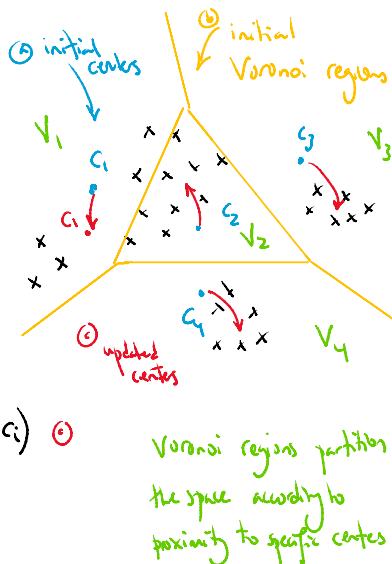
- Initialize $C = \{c_1, \dots, c_k\}$ (arbitrarily/randomly)

② Repeat:

- Cluster according to C (Voronoi regions)

$$\bullet c_i \leftarrow \text{Mean}(x : x \text{ closest to } c_i) \quad x \in V_i$$

③ Stop when convergence slow.

Output: Cluster centers C **② k-means Clustering**Input $x \in \mathbb{R}^n$, parameters $\overbrace{c_1, c_2, \dots, c_k}^C$, $c_i \in \mathbb{R}^n$.

$$\text{Output } f(x; C) = \underset{c \in C}{\arg \min} \|x - c\|_2^2$$

$$\text{Loss: } l(x, f) = \|x - f\|_2^2$$

$$\text{Risk: } R(C) = \frac{1}{|\text{Data}|} \sum_{x \in \text{Data}} l(x, f(x; C))$$

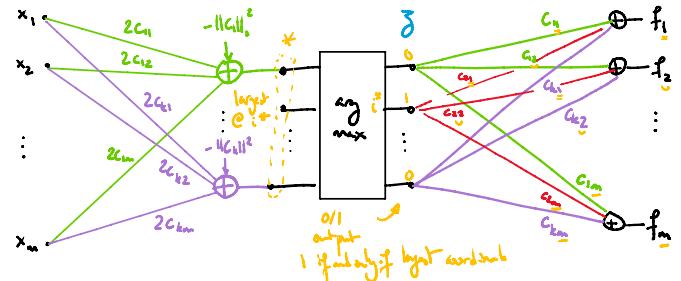
How do we minimize $R(C)$?

- It's possible to do gradient descent. It's non-convex!
- Or... use Lloyd's algorithm aka "the" k-means algorithm.

③ From k-means to Hebbian updates

$$\text{Since } \|x - c\|_2^2 = (x - c)^T (x - c) = \|x\|_2^2 + \|c\|_2^2 - 2x^T c$$

We can write the clustering as a neural network:

Note that f will be the closest center.

- Because $\arg \max$ is not differentiable (constant, except at jump), no gradient information backpropagates before it.
- We still have gradients after it, starting at: $\frac{\partial l}{\partial f_j} = -2(x_j - f_j)$

$$\text{Then: } \frac{\partial l}{\partial c_i} = \frac{\partial l}{\partial f_j} \cdot \frac{\partial f_j}{\partial c_i} = -2 \underbrace{z_i}_{\substack{\text{only ranges} \\ \text{at } i=i^* \text{ (argmax index)}}} (x_j - f_j) \quad \text{for those } f_j = c_{i^*}$$

$$\text{Thus: } \nabla_{c_i} R = \sum_{x \in V_i} -2(x_j - f_j) = -2 \sum_{x \in V_i} x + 2 \sum_{x \in V_i} c_i$$

Gradient descent becomes:

$$c_i \leftarrow c_i - \eta \nabla_{c_i} R = c_i + 2 \eta \sum_{x \in V_i} x - 2 \eta \sum_{x \in V_i} c_i$$

$$\text{If } \eta_i = \frac{1}{2 \sum_{x \in V_i} \|x\|^2} \text{ then } = \frac{1}{\sum_{x \in V_i} \|x\|^2} \sum_{x \in V_i} x \text{ same as Lloyd's algorithm!}$$

- Assume $\|x\| = \|w\|$ (or $\|w\| = 1$) for all weights.

- Reasonable by going one dimension higher \rightarrow \circlearrowright
- Distance \rightarrow (minus) inner-products! cosine(angle) \uparrow distance \downarrow

- That was full batch. With 1 point, we get instead

$$\text{H}: c_i \leftarrow c_i + \eta (x - c_i) \quad \underbrace{\{ \arg \max_i (c_i^T x) = i \}}_{\substack{\text{no bias needed} \\ \text{self-supervision}}} = z_i$$

- z is the neurons' output (decision). In general, it doesn't have to be restricted to 0/1.

- Definition: Competitive (winner-takes-all) learning rule:
 $w_i \leftarrow w_i + \eta z_i (x - w_i)$ only for $i = \arg \max_i z_i$

- If $\beta_i \in (0,1)$ then

$$w_i \leftarrow (1 - \eta_{\beta_i}) w_i + \eta_{\beta_i} x$$

→ weights that look not like x are made to look more like x
(they specialize at identifying things similar to x / clusters)

④ Oja's rule:

difference

Definition: $w_i \leftarrow w_i + \eta \beta_i (x - w_i \beta_i)$ doesn't affect it if off

- Claim, for small η , $w + \eta g(x - w^*) \in \frac{w + \eta x^*}{\|w + \eta x^*\|}$

vector scalar scalar vector
(for each i)

- Extracts principle component:
 $\rightarrow \underset{\|w\|=1}{\operatorname{argmax}} \sum_{x \in \text{data}} (w^T x)^2$
"minimally informative"
1-d linear feature

③ Hebbian rules

- Most of these rules are "Hebbian":
"Neurons that fire together, wire together"
 - If x is coming from another neuron and w is a neuron that mimics it, then it will strengthen.
 - Definition: Original Hebb's rule: $w_i \leftarrow (w_i + \eta x_j)$
(Compare to Perceptron with true supervision)
 - It can also cluster! But it's unstable:
Fixes: (Original) Rule, $w_i \leftarrow \alpha w_i + \eta x_j$
 \hookrightarrow damping < 1 .
 - Variants: $w_i \leftarrow w_i + \eta (x - \bar{x})(g_i - \bar{g}_i)$

Example: $z_i = \text{Sign}(w_i^T x)$

- Say $x=1$ vector, fixed $w_i(0)=1$ vector also initially. $\eta=1$.

$$w_i(1) = 1+1=2 \quad w_i(2) = 2+2=4 \quad w_i(3)=8 \dots$$

If we stabilize with $\alpha < 1$:

- In the notes: Examples of Hebb's rule and competitive