

Bayesian Interpretation of

Regression

$$\underline{J} = \frac{1}{N} \sum_{i=1}^n (t_i - w^T \phi(x_i))^2 + \lambda \cdot \sum_{j=1}^m w_j^2$$

$$\arg \min_w J \Leftrightarrow \arg \max_w P(x|w) \cdot P(w)$$

Data likelihood prior

$$|| P(x|w) \sim G(t, 1)$$

$$|| P(w) \sim G(0, 1/\lambda)$$

$$\arg \max_w P(w|x)$$

= what is the value of w that
have higher probability for the data
 x .

$$\arg \min_w J = \arg \max_w -J$$

$$= \arg \max_w \exp(-J)$$

$$= \arg \max_w \exp \left(\frac{1}{N} \sum_{i=1}^N e_i^2 + \lambda \cdot \sum_{j=1}^M w_j^2 \right)$$

$$= \arg \max_w \underbrace{\exp \left(\frac{1}{N} \sum_{i=1}^N e_i^2 \right)}_{\sim G(e|0,1)} \cdot \underbrace{\exp \left(\lambda \cdot \sum_{j=1}^M w_j^2 \right)}_{\sim G(w|0, 1/\lambda)}$$

Bayesian interpretation says we

can model a given set of

data $\{(x_i, t_i)\}_{i=1}^n$ using ~~an~~

distributional form for the data likelihood

and for the prior distribution.

$$\arg \max_w P(x|w) \cdot \underbrace{P(w)}$$

$$\propto \arg \max_w \underbrace{P(w|x)}_{\text{(from Bayes' theorem)}}$$

Online setting

↳ New samples are arriving every day.

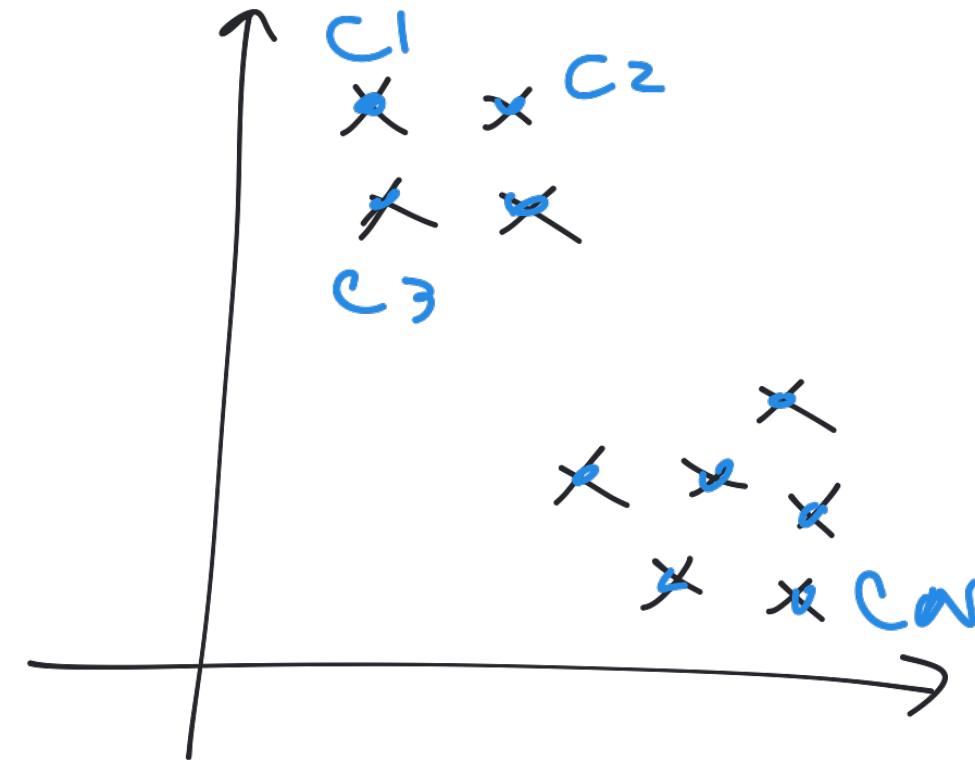
↳ If Posterior - Prior have a Conjugate Prior Relationship.

→ as new data comes in, we can update the prior with the posterior.

→ Posterior will represent the new data informative prior.

K-Means

$$\{x_i\}_{i=1}^N$$



$$J = \sum_{i=1}^N \sum_{j=1}^K u_{ij} \cdot d^2(x_i, c_j) = \sum_{i=1}^N \sum_{j=1}^K u_{ij} \|x_i - c_j\|_c^2$$

$$u_{ij} \in \{0, 1\}, \quad \sum_{j=1}^K u_{ij} = 1$$

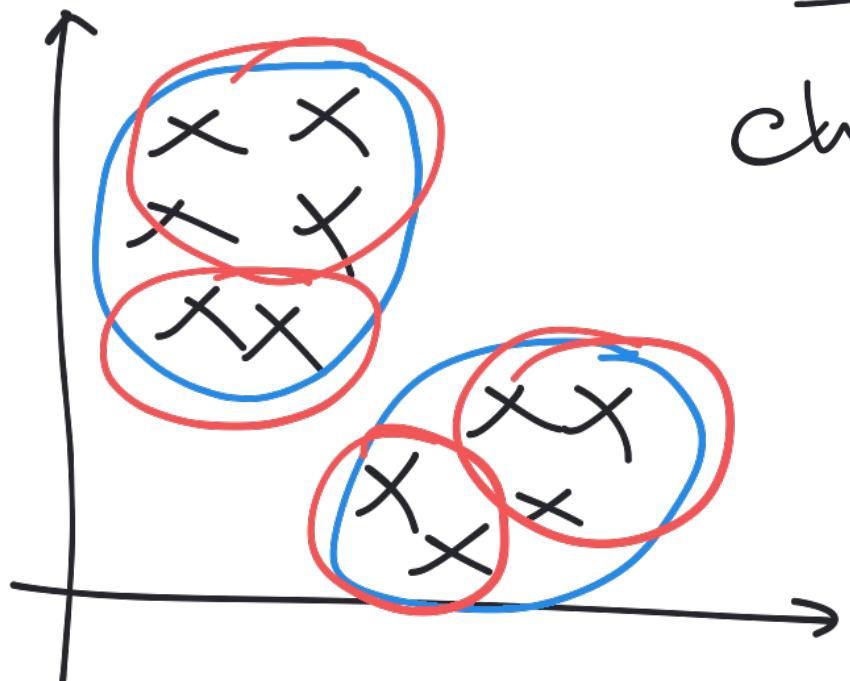
when $K = N$: $J = 0!$

what should we do + learn

clusters K?

↳ Internal Criteria for our

cluster validity metric

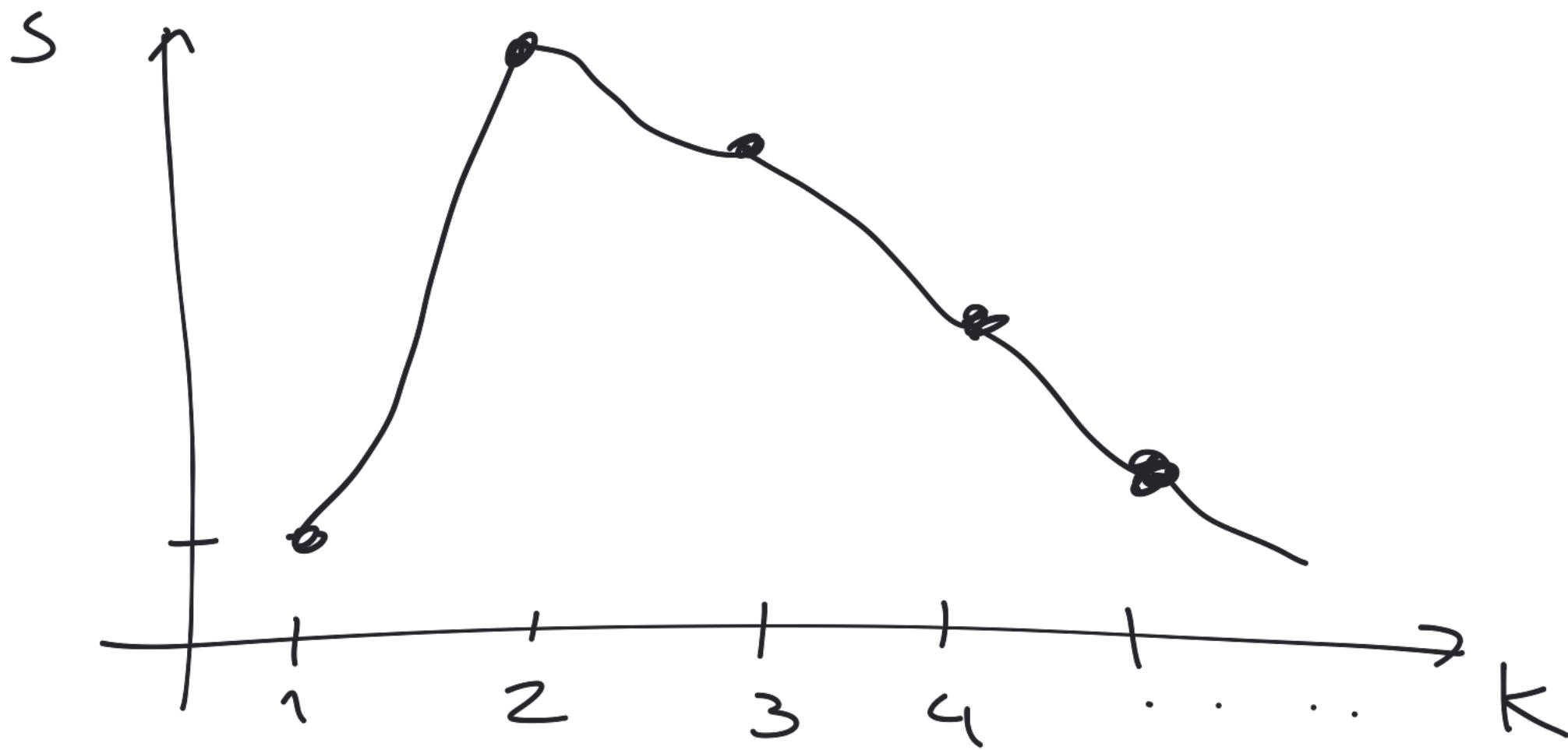


- ① Compact : all points assigned to the same cluster are near each other
- ② Distance between clusters to be large.

$K=2 \rightarrow \underline{\underline{K=4}}$

Silhouette index

$$-1 \leq S \leq 1$$



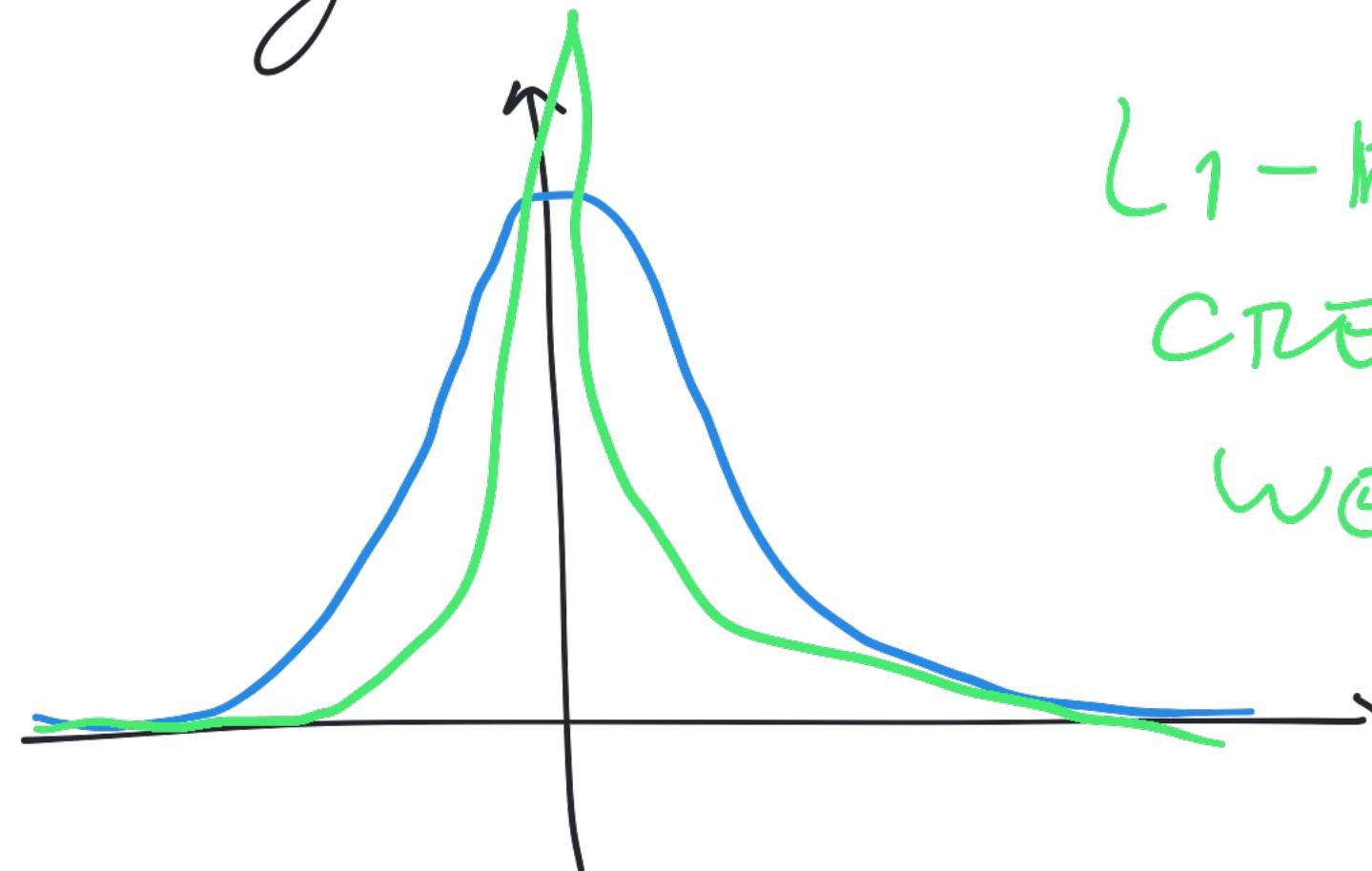
We choose K with maximum
Silhouette value.

→ in this case $K = 2$.

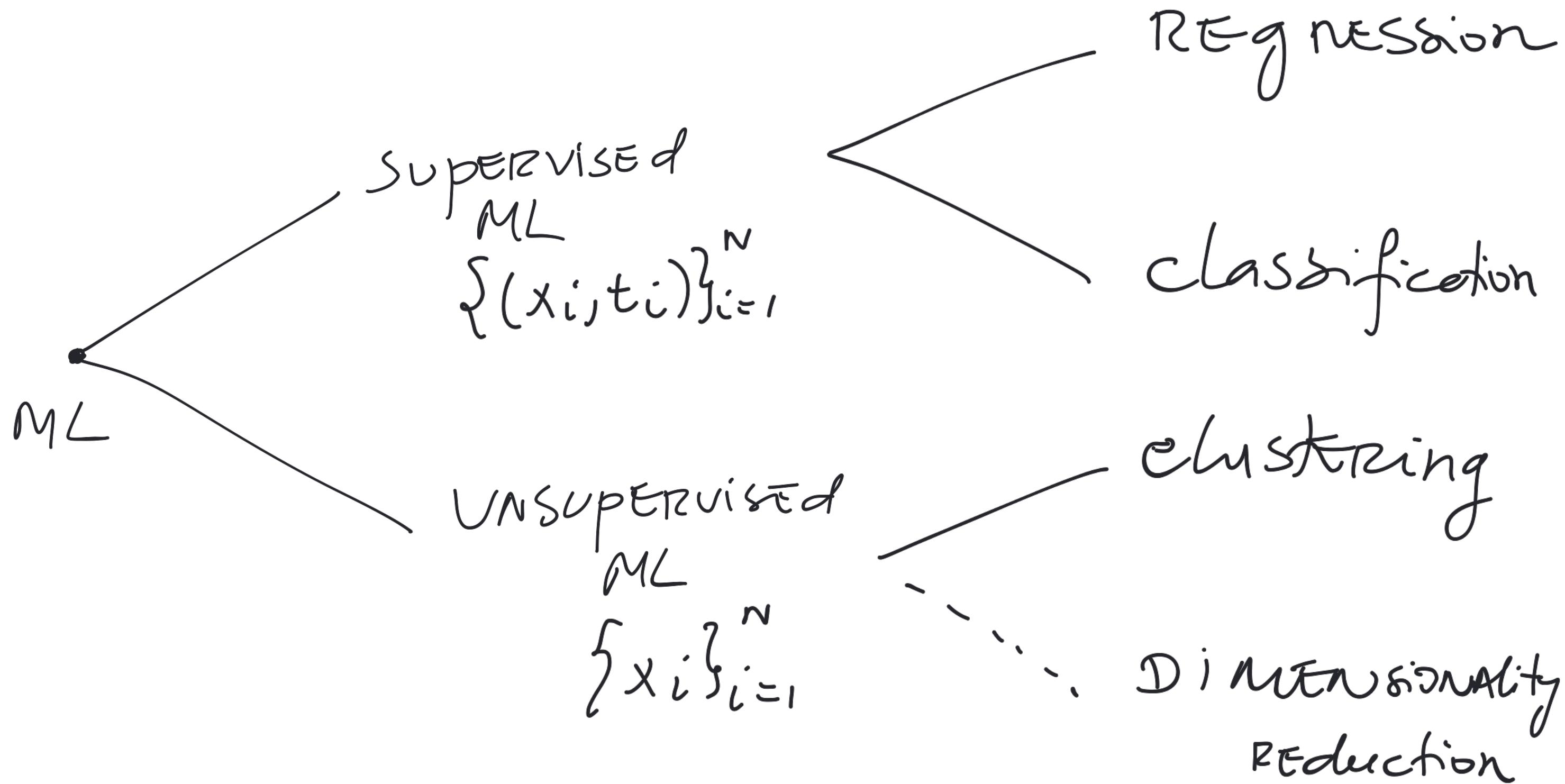
From Bayesian interpretation:

L2 - Regularizer is $G(0, \frac{1}{\lambda})$

L1 - Regularizer is $\mathcal{L}(0, \frac{1}{\lambda})$



L1-reg. will
CREATE SPARSE
WEIGHT VECTORS



For iteration t :

Data X

prior $\sim \text{Beta}(\mu | \alpha, \beta)$

$$\underline{P(\mu)}$$

data likelihood $\sim B(\mu)$

$$\underline{P(X|\mu)}$$

posterior = data likelihood \times prior

$$\underline{P(\mu|X)}$$

prior \leftarrow posterior

$$\alpha \leftarrow \alpha + N$$

$$\beta \leftarrow \beta + L$$