

Big Picture for Today

Quick coverage of RNN

↳ recap from early in semester
(distance / similarity based classification)

Clustering

- hierarchical
- k -means
- mixture models / EM

Nearest Neighbor Classification

- idea: similar/close items should be likely to have the same label
- basically a "likeness" lookup table
 - ↳ for new example ~~find~~ **X** find closest and choose its class

K-nearest neighbor

- find k closest points to **X**
- choose class with most votes
 - OR
 - choose class with most moe than l
 - ↳ otherwise say unknown

Practical Considerations

- need distance metric
- Euclidean distance

$$d(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{\sum_k (x_i^k - x_j^k)^2} = \|\mathbf{x}_i - \mathbf{x}_j\|_2$$

- Manhattan distance

$$d(\mathbf{x}_i, \mathbf{x}_j) = \sum_k |x_i^k - x_j^k| = \|\mathbf{x}_i - \mathbf{x}_j\|_1$$

- Hamming Distance

Sum over all features:
0 if feature values are same
1 otherwise
(categorical)

- Different scales typically don't play well with distance

Ex: clustering houses

features = price

↳ $\approx 100,000$

beds

↳ ≈ 1

baths

↳ ≈ 1

↳ dominates distance

Solution: normalize/rescale

- Different scales and variances \rightarrow whitening
transform to unit mean and identity covariance

Compute

$$Z_i = \Lambda^{(-1/2)} U^T (x_i - \text{mean}(X))$$

where U = eigenvectors of $\text{covmat}(X)$

Λ = diagonal matrix of eigenvalues of $\text{covmat}(X)$

- Costly to find nearest neighbor for lots of points in high dimensions

↳ Solution: approximate nearest neighbor