

CS 189: Introduction to Machine Learning - Discussion 7

1. Kernels

- a) Given a data point, \mathbf{x} , in \mathbb{R}^n , and the feature mapping ϕ corresponding to a quadratic kernel $k(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^T \mathbf{y} + c)^2$, what is the dimensionality of $\phi(\mathbf{x})$? What is the dimensionality of the feature mapping corresponding to a Gaussian/RBF kernel? How do we deal with this high-dimensional data?
- b) Why might we prefer to use kernels?

2. Feature Selection

Feature selection is the process of selecting a subset of a dataset's original features for use in learning tasks. Any algorithm for this task has two parts: a search technique for proposing new feature subsets, and an evaluation measure which scores the different feature subsets. Wrapper methods involve retraining a model with these different subsets. Filter methods don't interact with the model, but rather just the features; they model interactions between the features and eliminate highly correlated ones. Embedded methods embed feature selection directly into the loss function.

- a) Why might we want to perform feature selection?
- b) Describe the simplest wrapper method you can think of for performing feature selection. What is its runtime? Can we do better?
- c) Now, let's try the embedded method approach to feature selection.
 - i) What kind of penalty term corresponds directly to feature selection?
 - ii) What issues arise when optimizing a function with this penalty term?
 - iii) Describe two vectors such that their ℓ_1 norms are equal, but their ℓ_2 norms are not.
 - iv) Assume a method exists for optimizing non-differentiable functions, as long as they are convex.¹ How can we approximate the penalty term from part i, while still maintaining an "optimizable" loss function?

¹These are known as proximal gradient methods:

<http://www.eecs.berkeley.edu/~elghaoui/Teaching/EE227A/lecture18.pdf>

3. Leave One Out Cross-Validation²

K-Fold Cross Validation works pretty well, but imagine a case where our training data is really scarce. We could get even more use out of our data by training on all data points except for one, then attempt to classify the remaining point. This corresponds to setting $K = n$ in K-Fold Cross Validation, and is called Leave-One-Out-Cross-Validation (LOOCV). Although this would seem to require n model fits, we'll derive a method which saves us this expense in the **linear regression** setting.

Let's write the formal definition of LOOCV:

$LOOCV = \sum_{i=1}^n (y_i - \hat{y}_i^{-i})^2$, where \hat{y}_i^{-i} is the estimator of y_i with the i -th data point held-out when fitting the model.

- a) After fitting our regression model, we have $\hat{y} = X\hat{w}$. Express \hat{y} in terms of the actual labels, i.e. find H such that $\hat{y} = Hy$, where H is an $n \times n$ matrix³.
- b) By definition, \hat{y}^{-i} minimizes $\sum_{j \neq i} (y_j - \hat{y}_j^{-i})^2$. Prove \hat{y}^{-i} minimizes the squared error for z where

$$z_j = \begin{cases} y_j, & j \neq i \\ \hat{y}_i^{-i}, & j = i \end{cases}$$

- c) Prove $\hat{y}_i^{-i} = \hat{y}_i - H_{ii}y_i + H_{ii}\hat{y}_i^{-i}$
- d) Using this result, show that $LOOCV = \sum_{i=1}^n \left(\frac{y_i - \hat{y}_i}{1 - H_{ii}} \right)^2$. With this result, we can implement LOOCV while only fitting one regression model.
- e) Give an example of a dataset where LOOCV will give a poor estimate of the error.

²Problem Set-Up: Credit to Tom Mitchell and Andrew W. Moore, CMU

³This matrix is known as the hat matrix, and as a number of useful statistical properties:
https://en.wikipedia.org/wiki/Hat_matrix