

Lecture 3: Linear Regression (Part 2)

Sep 10 2019

Lecturer: Steven Wu

Scribe: Steven Wu

Revisit SVD and Pseudoinverse

Given any matrix $M \in \mathbb{R}^{m \times n}$, we want to factorize the matrix as $M = USV^T$, where

- r is the rank of the matrix M ;
- $U \in \mathbb{R}^{m \times r}$ is orthonormal, that is $U^T U = I_r$;
- $V \in \mathbb{R}^{n \times r}$ is orthonormal, that is $V^T V = I_r$;
- $S \in \mathbb{R}^{r \times r}$ is a diagonal matrix $\text{diag}(s_1, \dots, s_r)$.

We could also express the factorization as a sum

$$M = \sum_{i=1}^r s_i u_i v_i^T$$

where each u_i is a column vector for U and each v_i is a column vector for V . Note that $\{u_i\}$ spans the column space of M and $\{v_i\}$ spans the row space of M . This allows us to define the (Moore-Penrose) pseudoinverse

$$M^+ = \sum_{i=1}^r \frac{1}{s_i} v_i u_i^T.$$

Basically, we take the inverse of the singular values and reverse the positions of the v_i and u_i within each term. Now let's return to the problem of least squares regression with We can define the design matrix and response vector respectively:

$$A = \begin{bmatrix} \leftarrow x_1^T \rightarrow \\ \vdots \\ \leftarrow x_n^T \rightarrow \end{bmatrix} \quad \mathbf{b} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

We learn that $\mathbf{w}^* = A^+ \mathbf{b}$ is a solution for the first-order condition $(A^T A) \mathbf{w} = A^T \mathbf{b}$, and thus is the solution to the least squares regression problem: $\arg \min_{\mathbf{w}} \|A \mathbf{w} - \mathbf{b}\|_2^2$. Note that if A is full rank, then $\mathbf{w}^* = A^+ \mathbf{b} = (A^T A)^{-1} A^T \mathbf{b}$, which is the unique minimizer of the least squares objective.

1 Ridge Regression

Here is another way to handle the case where $A^\top A$ is not invertible. The idea is to do *regularized regression*, which is also useful for preventing *overfitting*.

We know that if the covariance matrix $A^\top A$ is invertible, the least squares solution is given by $(A^\top A)^{-1}\mathbf{b}$. If the covariance matrix is not rank-deficient, then we can slightly change it into $(A^\top A + \lambda I)$ for some $\lambda > 0$, which is guaranteed to be full rank. There are many ways to show this matrix is full rank, and we show this via the *eigendecomposition* of the covariance matrix:

$$A^\top A = Q\Lambda Q^\top$$

where

- $Q \in \mathbb{R}^{d \times d}$ with orthonormal column vectors $\{q_i\}_{i=1}^d$ that are also eigenvectors of $A^\top A$.
- $\Lambda \in \mathbb{R}^{d \times d}$ is a diagonal matrix $\text{diag}(\lambda_1, \dots, \lambda_d)$ with diagonal entries that are eigenvalues of $A^\top A$.

(Of course, if you don't like matrices, you can also write the factorization as $A^\top A = \sum_{i=1}^d \lambda_i q_i q_i^\top$.)

Note that the set of eigenvalues may not be distinct. In fact, when the rank $r < d$, $(d - r)$ of the eigenvalues will be 0. The “augmented” matrix can similarly be factorized as:

$$A^\top A + \lambda I = A^\top A + \lambda Q I Q^\top = Q(\Lambda + \lambda I)Q^\top.$$

All of the eigenvalues of this matrix are positive, and the matrix is invertible.

Now let's replace $A^\top A$ by $(A^\top A + \lambda I)$, and consider the following solution

$$\hat{\mathbf{w}} = (A^\top A + \lambda I)^{-1} A \mathbf{b}. \quad (1)$$

Again, by first-order condition, we can show that $\hat{\mathbf{w}}$ is the the solution to the following regularized ERM problem:

$$\min_{\mathbf{w}} \|A\mathbf{w} - \mathbf{b}\|_2^2 + \lambda \|\mathbf{w}\|_2^2 \quad (2)$$

This is also called *ridge regression*, and the solution is always unique even if $n < d$. The *regularization* or *penalty* term $\lambda \|\mathbf{w}\|_2^2$ encourages “shorter” solutions \mathbf{w} with smaller ℓ_2 norm. The parameter λ manages the trade-off between fitting the data to minimize $\hat{\mathcal{R}}$ and shrinking the solution to minimize $\lambda \|\mathbf{w}\|_2^2$. Ridge regression can also be formulated as a *constrained optimization* problem:

$$\min_{\mathbf{w}} \|A\mathbf{w} - \mathbf{b}\|_2^2 \quad \text{such that} \quad \|\mathbf{w}\|_2 \leq \beta.$$

Why do we care to make the weights \mathbf{w} short or small? Intuitively, larger \mathbf{w} corresponds to higher *model complexity*. By bounding the model complexity, we can prevent *overfitting*—that is the model has small training error, but large test error. However, if we bound the norm of \mathbf{w} too aggressively (by setting λ to be very large), then we might run into the problem of *underfitting*—that is the model has large training error and test error.

Lasso regression. Another common regularization is the *Lasso regression* that uses ℓ_1 penalty:

$$\arg \min_{\mathbf{w}} \|A\mathbf{w} - \mathbf{b}\|_2^2 + \lambda \|\mathbf{w}\|_1$$

Lasso encourages sparse solutions, and is commonly used when d is much greater than the number of observations n . However, it does not admit a closed-form solution.

2 Feature Transformation

We can enrich linear regression models by transforming the features: first transform each feature vector x into $\phi(x)$, and then predict by using linear function over the transformed features, that is $\hat{f}(x) = \mathbf{w}^\top \phi(x)$. Consider the following examples of feature transformation:

- for $x \in \mathbb{R}$, $\phi(x) = \ln(1 + x)$
- for $x \in \{0, 1\}^d$, we can apply boolean functions such as

$$\phi(x) = (x_1 \wedge x_2) \vee (x_3 \vee x_4)$$

- for $x \in \mathbb{R}^d$, we can also apply polynomial expansion:

$$\phi(x) = (1, x_1, \dots, x_d, x_1^2, \dots, x_d^2, x_1x_2, \dots, x_{d-1}x_d)$$

- for $x \in \mathbb{R}$, we can also apply trigonometry expansion:

$$\phi(x) = (1, \sin(x), \cos(x), \sin(2x), \cos(2x), \dots)$$

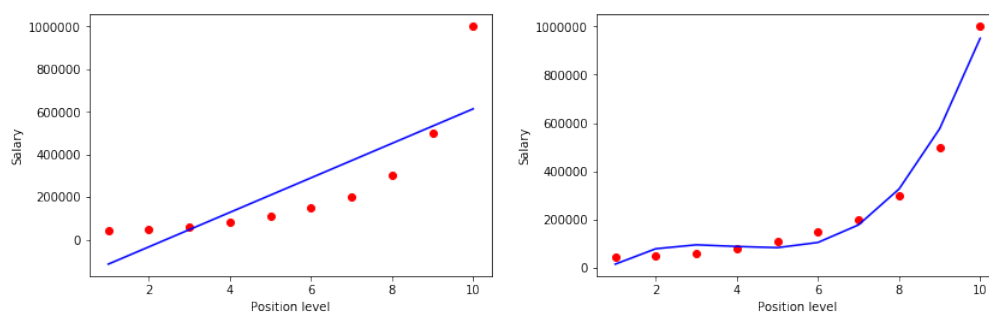


Figure 1: Examples shown in class. Fitting a linear function versus fitting a degree-3 polynomial. (More details here.)

Can we just use complicated linear mapping though? No, we won't gain anything: $\mathbf{w}^\top \phi(x)$ is just another linear function of x , when ϕ is also a linear mapping of x .

Feature engineering can get messy, and often requires a lot of domain knowledge. For example, we probably should not use polynomial expansion for periodic data.

3 Hyperparameters, Validation Set, and Test Set

The parameter λ in ridge regression and Lasso regression, and the order of polynomials in polynomial expansion, and also the parameter k in k -nearest neighbor are often called *hyperparameters* for the machine learning algorithms, which requires tuning. How do we optimize these parameters? A standard way is to perform the following three-way data splits:

- Training set: learn the predictor \hat{f} (e.g. weight vector \mathbf{w}) by “fitting” this dataset.
- Validation set: a set of examples to tune the hyperparameters. We use the loss on this dataset to find the “best” hyperparameter.
- Test set: we use this data to assess the *risk* of the final model:

$$\mathcal{R}(f) = \mathbf{E}_{(X,Y) \sim P}[\ell(Y, f(X))]$$

In the case of squared loss, this is

$$\mathcal{R}(f) = \mathbf{E}_{(X,Y) \sim P}[(f(X) - Y)^2]$$

In general, we want to predict well on future instances, so the goal is formulated as finding a predictor \hat{f} that minimizes the risk (instead of empirical risk on the training set).

What if we did not start with a validation set? We can always create a validation set from the training set. One standard method is *cross validation*.

k -fold cross validation We split the training set into k parts or folds of roughly equal size: F_1, \dots, F_k . (Typically, $k = 5$ or 10 , but it also depends on the size of your dataset.)

1. For $j = 1, \dots, k$:
 - We will train on the union of folds $F_{-j} = \bigcup_{j' \neq j} F_{j'}$ and validate on fold F_j
 - For each value of the tuning parameter $\theta \in \{\theta_1, \dots, \theta_m\}$, train on F_{-j} to obtain predictor \hat{f}_θ^{-j} , and record the loss on the validation set $\hat{\mathcal{R}}_j(\hat{f}_\theta^{-j})$.
2. For each parameter θ , compute the average loss over all folds

$$\hat{\mathcal{R}}_{\text{CV}}(\theta) = \frac{1}{k} \sum_{j=1}^k \hat{\mathcal{R}}_j(\hat{f}_\theta^{-j})$$

Then we will choose the parameter $\hat{\theta}$ that minimize $\hat{\mathcal{R}}_{\text{CV}}(\theta)$.