

CS 181 Spring 2019 Section 2 Notes (Model Selection)

1 Validation

1.1 Linear Regression

Suppose we have data $\{(x_i, y_i)\}_{i=1}^n$, with $x_i, y_i \in \mathbb{R}$, and we want to fit polynomial basis functions:

$$\phi(x)^\top = [\phi_1(x) = 1, \phi_2(x) = x, \dots, \phi_{d+1}(x) = x^d]$$

$$h(\mathbf{x}; \mathbf{w}) = \mathbf{w}^\top \phi(x)$$

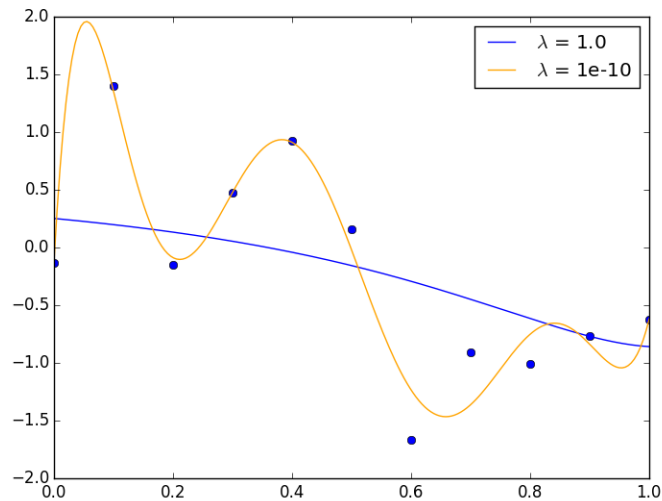
That is, we fit a degree d polynomial. With a small dataset and too high of a d , we get overfitting. Obviously, this will generalize poorly to new data points. How can we solve this problem?

1.2 Ridge Regression

One solution to overfitting linear regression is through ridge regression, which minimizes a modified least squares loss function:

$$\mathcal{L}(D) = \sum_{i=1}^n (y_i - h(x_i; \mathbf{w}))^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

Ridge regression is used to *regularize* a model, making it simpler and allowing it to generalize better to new data. Indeed, the extra term penalizes overly large weights in \mathbf{w} , leading to smaller coefficients for a “flatter” polynomial:



1.3 Validation Set: Model Selection

We can do model selection through a *validation set*, data that are separate from our training set used to fit the regression. By separating our full dataset into a training set and validation set (say in a 90/10 split), we can use our validation set to check our model's generalization ability on data it was not trained on. When tuning model parameters, we can train our models with different parameters on the training set and check their performance on the validation set in order to find the optimal value for the parameter.

1.4 Cross Validation

Cross validation is a more sophisticated technique for obtaining validation losses. Instead of splitting our data once into a 90/10 training/validation set, in k -fold cross validation, we split our data into k equal chunks. For each chunk, we set it to be the validation set and use the rest of the data to fit our model. Then, we obtain a validation loss on our current chunk, and averaging over the 10 chunks gives the final validation loss. Cross validation can also be used to find optimal parameter values as described in the previous section - we simply have an improved way of computing validation losses by averaging. This reduces the variance in the resulting validation loss, as each example is used in estimating the validation loss.

See this [notebook](#) for an interactive demo of how cross validation could be used.

2 Bias-Variance Decomposition

Bias-variance decomposition is a way of understanding how different sources of error (bias and variance) can affect the final performance of a model. A tradeoff between bias and variance is often made when selecting models to use, and can be informed by the results of the bias-variance decomposition.

Exercise: Decompose the generalization error into the sum of bias squared (systematic error), variance (sensitivity of prediction), and noise (irreducible error) by following the steps below (**try not to peek at your notes!**). You will find the following notation useful:

- h_D : The trained model, $h_D : \mathcal{X} \mapsto \mathbb{R}$.
 - D : The data, a random variable sampled $D \sim F^n$.
 - \mathbf{x} : A new input.
 - y : The true result of input \mathbf{x} . Conditioned on \mathbf{x} , y is a r.v. (may be noise.)
 - \bar{y} : The true conditional mean, $\bar{y} = \mathbb{E}_{y|\mathbf{x}}[y]$.
 - $\bar{h}(\mathbf{x})$: The prediction mean, $\bar{h}(\mathbf{x}) = \mathbb{E}_D[h_D(\mathbf{x})]$.
1. Start with the equation for the generalization error - the expected error, in least squares terms, on an unseen sample:

$$\mathbb{E}_{D, y|\mathbf{x}}[(y - h_D(\mathbf{x}))^2]$$

and use the linearity of expectation to derive an equation of the form:

$$\underbrace{\mathbb{E}_{y|\mathbf{x}}[(y - \bar{y})^2]}_{\text{noise}} + \underbrace{\mathbb{E}_D[(\bar{y} - h_D(\mathbf{x}))^2]}_{\text{bias+var}} + \text{*****} \quad (1)$$

where the *s denote a third term. What is this third term? (**Hint:** add and subtract \bar{y}).

2. Show that this third term is equal to 0 (**Hint:** take advantage of the fact that \bar{y} and $h_D(\mathbf{x})$ do not depend on $y|x$).
3. The first term in (1) is the noise. We therefore want to decompose the second term into the bias and variance. Again, using the linearity of expectation, re-write the second term in equation (1) in the form:

$$\underbrace{(\bar{y} - \bar{h}(\mathbf{x}))^2}_{\text{bias squared}} + \underbrace{\mathbb{E}_D[(\bar{h}(\mathbf{x}) - h_D(\mathbf{x}))^2]}_{\text{variance}} + 2\mathbb{E}_D[(\bar{y} - \bar{h}(\mathbf{x}))(\bar{h}(\mathbf{x}) - h_D(\mathbf{x}))] \quad (2)$$

show that the third term is equal to 0.

4. Plug the results of part 3 back into (1) to show that we have decomposed the error into noise, bias, and variance.

Solution:

1. Follow the hint:

$$\begin{aligned}
 & \mathbb{E}_{D, y|\mathbf{x}}[(y - h_D(\mathbf{x}))^2] \\
 &= \mathbb{E}_{D, y|\mathbf{x}}[(y - \bar{y} + \bar{y} - h_D(\mathbf{x}))^2] \\
 &= \underbrace{\mathbb{E}_{y|\mathbf{x}}[(y - \bar{y})^2]}_{\text{noise}} + \underbrace{\mathbb{E}_D[(\bar{y} - h_D(\mathbf{x}))^2]}_{\text{bias+var}} + \underbrace{2\mathbb{E}_{D, y|\mathbf{x}}[(y - \bar{y})(\bar{y} - h_D(\mathbf{x}))]}_0
 \end{aligned}$$

2. Using the hint:

$$2\mathbb{E}_D[\bar{y} - h_D(\mathbf{x}) \cdot \mathbb{E}_{y|\mathbf{x}}[y - \bar{y}]] = 2\mathbb{E}_D[\bar{y} - h_D(\mathbf{x}) \cdot 0] = 0.$$

3. Following a similar procedure as in part 1:

$$\begin{aligned}
 & \mathbb{E}_D[(\bar{y} - h_D(\mathbf{x}))^2] \\
 &= \mathbb{E}_D[(\bar{y} - \bar{h}(\mathbf{x}) + \bar{h}(\mathbf{x}) - h_D(\mathbf{x}))^2] \\
 &= \underbrace{(\bar{y} - \bar{h}(\mathbf{x}))^2}_{\text{bias squared}} + \underbrace{\mathbb{E}_D[(\bar{h}(\mathbf{x}) - h_D(\mathbf{x}))^2]}_{\text{variance}} + \underbrace{2\mathbb{E}_D[(\bar{y} - \bar{h}(\mathbf{x}))(\bar{h}(\mathbf{x}) - h_D(\mathbf{x}))]}_0
 \end{aligned}$$

where the third term is 0 by:

$$2(\bar{y} - \bar{h}(\mathbf{x}))\mathbb{E}_D[\bar{h}(\mathbf{x}) - h_D(\mathbf{x})] = 2(\bar{y} - \bar{h}(\mathbf{x}))(0) = 0.$$

4. Substituting (2) back into (1), we have:

$$\begin{aligned}
 & \mathbb{E}_{D, y|\mathbf{x}}[(y - h_D(\mathbf{x}))^2] \\
 &= \mathbb{E}_{y|\mathbf{x}}[(y - \bar{y})^2] + (\bar{y} - \bar{h}(\mathbf{x}))^2 + \mathbb{E}_D[(\bar{h}(\mathbf{x}) - h_D(\mathbf{x}))^2] \\
 &= \text{noise}(\mathbf{x}) + \text{bias}^2(h(\mathbf{x})) + \text{Var}_D(h_D(\mathbf{x})).
 \end{aligned}$$

Considering the expectation over \mathbf{x} (you are not asked to do this in the exercises), the generalization error is:

$$\mathbb{E}_{\mathbf{x}} [\text{noise}(\mathbf{x}) + \text{bias}^2(h(\mathbf{x})) + \text{Var}_D(h_D(\mathbf{x}))]$$

2.1 Limitations

Although the bias-variance decomposition provides some interesting insights into model selection from a complexity perspective, it has limited practical value, as it is based on

averages of independent data sets drawn from some distribution. In practice, however, we only have a single observed data set. The bias and variance can be estimated through “bootstrap” style approaches where we sample with replacement to form additional data sets, but still— why not more directly compute validation loss and use this to find the best model? The main interest in the bias-variance decomposition is to gain conceptual insight.

3 Ensemble Methods

Ensemble methods take advantage of multiple models to obtain better predictive accuracy than with a single model alone. The two most common types of ensemble methods are bagging and boosting.

3.1 Bootstrap aggregating (Bagging)

In bagging, we fit each individual model on a random sample of the training set. To predict data in the test set, we either use an average of the predictions from the individual models (for regression) or take the majority vote (for classification). As an average of models, bagging tends to decrease the variance of a learning algorithm without changing the bias. An example is a random forest, which trains multiple decision trees and takes the average prediction from the ensemble of learned models.

3.2 Boosting

In boosting, we train the individual models sequentially. Thus, after training the i^{th} model on a sample of the training set, we train the $(i + 1)^{th}$ model on a new sample based on the performance of the i^{th} model. Examples classified incorrectly in the previous step receive higher weights in the new sample, encouraging the new model to focus on those examples. During testing, we take a weighted average or weighted majority vote of the models’ predictions based on their respective training accuracies on their reweighted training data (i.e. higher models have larger weights). A common example is the Adaboost algorithm.

4 Practice Questions

1. Ridge Regression

Suppose we have some data matrix $\mathbf{X} \in \mathbb{R}^{n \times m}$ and targets $\mathbf{y} \in \mathbb{R}^n$. Suppose the data are orthogonal*, i.e. satisfies $\mathbf{X}^\top \mathbf{X} = \mathbf{I}$. Show that if $\hat{\mathbf{w}}$ is the solution to linear regression, and $\hat{\mathbf{w}}_{ridge}$ is the solution to ridge regression, then

$$\hat{\mathbf{w}}_{ridge} = \frac{1}{1 + \lambda} \hat{\mathbf{w}}$$

This explicitly illustrates the phenomenon of weight shrinkage.

Recall that the linear regression solution is

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

and recall from homework that the ridge regression solution is

$$\hat{\mathbf{w}}_{ridge} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$$

If $\mathbf{X}^\top \mathbf{X} = \mathbf{I}$, we see that

$$\begin{aligned} \hat{\mathbf{w}} &= \mathbf{X}^\top \mathbf{y} \\ \hat{\mathbf{w}}_{ridge} &= \frac{1}{1 + \lambda} \mathbf{X}^\top \mathbf{y} \end{aligned}$$

hence giving us the result.

* Orthogonal data is a very special case in which the inner product between any two distinct features is zero. Normally we expect features to be correlated. But it is used to gain this clean illustration of the effect of ridge regression. Technically, we have $\mathbf{X} = [\mathbf{v}_1, \dots, \mathbf{v}_m]$ where $\mathbf{v}_1, \dots, \mathbf{v}_m$ are n dimensional, orthogonal column vectors.

2. Bias and Variance

We consider a very simple example where the data is a univariate Gaussian, with $x_i \sim \mathcal{N}(\mu, 1)$ with known variance but unknown mean. In this case, there are no features, and the hypothesis doesn't depend on x . A very simple hypothesis, for example, is the sample mean

$$h_D = \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

for data $(x_1, \dots, x_n) \in \mathbb{R}^n$. Calculate the bias and variance for the following two hypotheses:

- (a) Estimate 1: Use the same mean of data D .
- (b) Estimate 2: Use the constant hypothesis, 0.

- (a) The bias is $\mu - E_D[\bar{x}] = \mu - E_D[(1/n) \sum x_i] = 0$. The variance is $E_D[(\bar{x} - \mu)^2] = \sigma^2/n$, the variance of the sample mean on n examples (as is standard).
- (b) The bias equals μ , the expected difference between the estimator and the true value. This prediction is constant, so the variance of our prediction is 0.

3. Deriving Lasso Regularization with Lagrange Multipliers

Show that minimization of the unregularized sum-of-squares error function given by

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (t_n - \mathbf{w}^\top \phi(\mathbf{x}_n))^2,$$

subject to the constraint

$$\sum_{j=1}^M |w_j| \leq \eta,$$

is equivalent to minimizing the regularized error function

$$\frac{1}{2} \sum_{n=1}^N (t_n - \mathbf{w}^\top \phi(\mathbf{x}_n))^2 + \frac{\lambda}{2} \sum_{j=1}^M |w_j|$$

Rewrite the constraint as

$$\sum_{j=1}^M |w_j| - \eta \leq 0$$

We get the Lagrangian function

$$L(\mathbf{w}, \lambda) = \frac{1}{2} \sum_{n=1}^N (t_n - \mathbf{w}^\top \phi(\mathbf{x}_n))^2 + \frac{\lambda}{2} \sum_{j=1}^M (|w_j| - \eta)$$

where we introduce the factor of $1/2$ in front of the second term for convenience. We see immediately that the above function is equal to the regularized error function plus the terms of η which do not depend on \mathbf{w} . Therefore, minimizing the Lagrangian with respect to \mathbf{w} will give the same \mathbf{w}^* as minimizing the regularized error function.

4. Priors for Model Selection

Suppose you had three models, M_1, M_2, M_3 , each increasing in complexity. For example, you could imagine that the models represented unregularized polynomial regression, with M_1 linear regression, M_2 quadratic regression, and M_3 cubic regression. Within the context of Bayesian model selection, come up with a way to penalize the complexity of a model so you do not always choose M_3 . Additionally, explain why, in many cases, Bayesian model selection will recover the simplest model to explain the data without explicit penalization.

This is the same thing that we did to encourage less complex parameters in Bayesian linear regression; we put a prior over the weight vector assigning the highest probability to weight vectors that lie close to zero. This is a phenomenon known as Bayesian Occam's Razor and is discussed on in Bishop (pg. 164, fig. 3.13). Briefly, since simpler models can explain a smaller subset of data than more complex models, and each model has to integrate to 1, the simpler models can put higher probabilities on the data that they do describe than the complex models, which must put some of their mass on the data that the simpler model cannot describe.