

Lecture 4: Basic Principles Part II

*Instructor: Anant Sahai**Scribe: Austin Zane*

1 Regularization

1.1 Explicit Regularization

This is a recap of the previous lecture. For the ordinary least squares problem $\mathbf{X}\mathbf{w} \approx \mathbf{y}$, the solution is given by

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

However, this can give us very large and not very useful values for the parameter vector. In such cases, it is often beneficial to “regularize” the parameters when training so that they stay reasonable. A common choice is ridge regularization, which uses a modified version of the least-squares loss function:

$$\arg \min_{\mathbf{w}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_2^2.$$

A bit of vector calculus shows that the new solution for $\hat{\mathbf{w}}$ is given by

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

where the first formula is the classic ridge form and the second is kernel ridge form. The equivalence can be seen as follows,

$$\begin{aligned} (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y} &= (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top + \lambda \mathbf{I}) (\mathbf{X} \mathbf{X}^\top + \lambda \mathbf{I})^{-1} \mathbf{y} \\ &= (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} (\mathbf{X}^\top \mathbf{X} \mathbf{X}^\top + \lambda \mathbf{X}^\top) (\mathbf{X} \mathbf{X}^\top + \lambda \mathbf{I})^{-1} \mathbf{y} \\ &= (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}) \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top + \lambda \mathbf{I})^{-1} \mathbf{y} \\ &= \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top + \lambda \mathbf{I})^{-1} \mathbf{y}. \end{aligned}$$

There is a way of thinking of the first formulation as solving the primal and the second as solving the dual.

1.2 Data and Feature Augmentation

Instead of explicitly changing the loss function, we can add d “fake” data points to our data matrix to achieve the same regularizing effect (as proven in previous lecture):

$$\begin{bmatrix} \mathbf{X} \\ \sqrt{\lambda} \mathbf{I}_d \end{bmatrix} \mathbf{w} \approx \begin{bmatrix} \mathbf{y} \\ \mathbf{0}_d \end{bmatrix},$$

where the new data matrix is in $\mathbb{R}^{(n+d) \times d}$ and the new response is a vector in \mathbb{R}^{n+d} . This is an important concept and can be thought of as improving the conditioning number of the data matrix. Plugging the new data matrix and response vector into the classic OLS solution immediately gives us the solution for the ridge regularized problem.

Another equivalent option is adding n fake features:

$$\begin{bmatrix} \mathbf{X} & \sqrt{\lambda} \mathbf{I}_n \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ \mathbf{f} \end{bmatrix} = \mathbf{y}.$$

We use the Moore-Penrose pseudoinverse to solve for the new weight vector,

$$\begin{bmatrix} \mathbf{w} \\ \mathbf{f} \end{bmatrix} = \begin{bmatrix} \mathbf{X} & \sqrt{\lambda} \mathbf{I}_n \end{bmatrix}^\dagger \mathbf{y} = \begin{bmatrix} \mathbf{X}^\top \\ \sqrt{\lambda} \mathbf{I}_n \end{bmatrix} \left(\begin{bmatrix} \mathbf{X} & \sqrt{\lambda} \mathbf{I}_n \end{bmatrix} \begin{bmatrix} \mathbf{X}^\top \\ \sqrt{\lambda} \mathbf{I}_n \end{bmatrix} \right)^{-1} \mathbf{y}.$$

We are only interested in solving for \mathbf{w} , so we disregard the n rows corresponding to the false parameter \mathbf{f} . In the end, we are left with the minimum norm solution to the under-determined system specified above,

$$\mathbf{w} = \mathbf{X}^\top \left(\begin{bmatrix} \mathbf{X} & \sqrt{\lambda} \mathbf{I}_n \end{bmatrix} \begin{bmatrix} \mathbf{X}^\top \\ \sqrt{\lambda} \mathbf{I}_n \end{bmatrix} \right)^{-1} \mathbf{y} = \mathbf{X}^\top \left(\mathbf{X} \mathbf{X}^\top + \lambda \mathbf{I} \right)^{-1} \mathbf{y}.$$

1.3 Using Singular Value Decomposition to Simplify Regularization

Using the singular value decomposition (SVD) in place of \mathbf{X} in the above equations allows us to simplify things and make the algorithms far easier to handle. Namely, we are able to update each weight individually instead of solving systems of equations. This can be seen by taking the SVD of \mathbf{X} and solving the unregularized problem

$$\begin{aligned} \mathbf{X} \mathbf{w} &= \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top \mathbf{w} \approx \mathbf{y} \\ \Rightarrow \mathbf{\Sigma} \tilde{\mathbf{w}} &\approx \tilde{\mathbf{y}}, \end{aligned}$$

where $\tilde{\mathbf{w}} := \mathbf{V}^\top \mathbf{w}$, $\tilde{\mathbf{y}} := \mathbf{U}^\top \mathbf{y}$, and $\mathbf{\Sigma} \in \mathbb{R}^{n \times d}$. Note that \mathbf{U} , \mathbf{V} are orthogonal matrices so they merely rotate our vectors while preserving their norms. Let σ_i denote the i^{th} singular value. Because $\mathbf{\Sigma}$ is diagonal and we are assuming that $n > d$, only the first d equations will be meaningful here:

$$\begin{aligned} \text{for } d \text{ equations: } \sigma_i \tilde{\mathbf{w}}[i] &\approx \tilde{\mathbf{y}}[i] \Rightarrow \tilde{\mathbf{w}}[i] \approx \frac{1}{\sigma_i} \tilde{\mathbf{y}}[i], \\ \text{for } n - d \text{ equations: } 0 &\approx \tilde{\mathbf{y}}[i]. \end{aligned}$$

As previously stated, using these coordinates removes the system of equations and allows us to simply solve for the individual weight components. Observe that we are dividing by the singular values so problems may arise if they become too small (i.e. if the matrix is ill-conditioned). This is the underlying cause of OLS sometimes giving us wild values.

Next, we consider using ridge regression. In this setting, the solution is obtained by plugging the SVD of \mathbf{X}

into the previously mentioned classic solution for ridge regression.

$$\begin{aligned}
\hat{\mathbf{w}} &= \left(\mathbf{V} \mathbf{\Sigma}^\top \mathbf{\Sigma} \mathbf{V}^\top + \lambda \mathbf{I} \right)^{-1} \mathbf{V} \mathbf{\Sigma}^\top \mathbf{U}^\top \mathbf{y} \\
&= \mathbf{V} \left(\mathbf{\Sigma}^\top \mathbf{\Sigma} + \lambda \mathbf{I} \right)^{-1} \mathbf{V}^\top \mathbf{V} \mathbf{\Sigma}^\top \tilde{\mathbf{y}} \\
&= \mathbf{V} \left(\mathbf{\Sigma}^\top \mathbf{\Sigma} + \lambda \mathbf{I} \right)^{-1} \mathbf{\Sigma}^\top \tilde{\mathbf{y}} \\
\Rightarrow \hat{\mathbf{w}} &= \left(\mathbf{\Sigma}^\top \mathbf{\Sigma} + \lambda \mathbf{I} \right)^{-1} \mathbf{\Sigma}^\top \tilde{\mathbf{y}}.
\end{aligned}$$

In the end, we are left with solutions of the form

$$\tilde{\mathbf{w}}[i] = \frac{\sigma_i}{\sigma_i^2 + \lambda} \tilde{\mathbf{y}}[i].$$

If $\lambda \ll \sigma_i^2$, then we are in the same situation as the unregularized case. If $\lambda \gg \sigma_i^2$, then the weights are forced to stay small instead of behaving wildly.

1.4 Implicit Regularization

Implicit regularization is the regularization that occurs when we aren't consciously doing any regularization. When performing explicit regularization as above, we must specify a specific regularization hyperparameter and modify the training data, model architecture, or loss function. In contrast, implicit regularization is an unexpected benefit stemming from our choice of optimizer. We will see that choosing gradient descent as our optimization algorithm, combined with the large size of DNNs, provides enough regularization for DNNs to generalize well without us intentionally restricting the parameters.

To gain intuition, let's look at gradient descent (GD) updates for OLS in SVD coordinates:

$$\tilde{\mathbf{w}}_{t+1} = \tilde{\mathbf{w}}_t + 2\eta \mathbf{\Sigma}^\top (\tilde{\mathbf{y}} - \mathbf{\Sigma} \tilde{\mathbf{w}}_t),$$

where $\tilde{\mathbf{w}}_t$ represents the parameter vector at the current step, $\tilde{\mathbf{w}}_{t+1}$ represents the updated parameter vector, $\mathbf{\Sigma}$ represents the diagonal matrix of singular values from the above SVD, $\tilde{\mathbf{y}} := \mathbf{U}^\top \mathbf{y}$ as defined above, and η is our learning rate hyperparameter.

Note that because we are dealing with a diagonal matrix, this works out to updating each component of the weight vector individually. They don't interact with each other at all during GD, so each is being modified as follows:

$$\tilde{w}_{t+1}[i] = \tilde{w}_t[i] + 2\eta \sigma_i (\tilde{y}[i] - \sigma_i \tilde{w}_t[i]).$$

This is potentially unstable because we are not reducing $\tilde{w}_t[i]$ at each step as we did in the last lecture. This means that, subject to a bounded input, we might get an unbounded output if we allow the algorithm to run forever. Observe that the stationary point is the solution we discussed earlier, $\tilde{\mathbf{w}}[i] = \frac{1}{\sigma_i} \tilde{\mathbf{y}}[i]$. If σ_i is tiny, then we find ourselves in a bad situation.

Let's carefully calculate the first few steps of GD to see what's going on:

$$\begin{aligned}
\tilde{w}_0[i] &= 0 \\
\tilde{w}_1[i] &= 2\eta \sigma_i \tilde{y}[i] \\
\tilde{w}_2[i] &= 2\eta \sigma_i \tilde{y}[i] + 2\eta \sigma_i (\tilde{y}[i] - \sigma_i 2\eta \sigma_i \tilde{y}[i]) \approx 4\eta \sigma_i \tilde{y}[i].
\end{aligned}$$

Observe that this is roughly a linear function with an extremely small slope if σ_i is tiny. In this situation, GD barely moves in the early stages even though it will eventually converge to a very large value, as previously discussed. Together with early stopping, this means that GD is trying to do something like ridge regularization for us because it will resist enlarging the directions corresponding to small singular values. Early stopping is when we stop the training process because validation performance has gotten worse or has not improved for a long time. It is important to note that GD, when initialized at zero, will converge to the minimum-norm solution. This is a good exercise for the reader to verify.

To summarize, there are three kinds of regularization: **explicit regularization**, **data augmentation (adding fake observations or features)**, and **implicit regularization (optimizer has implicit regularizing effect)**. Regarding DNNs, the combination of the min-norm seeking behavior of gradient descent and the feature augmentation that is implicit when using large networks gives a lot of regularization even if we weren't thinking about it.

2 Trade-offs Between Qualitatively Different Sources of Error

Suppose we have learned a model $\hat{\theta} \rightarrow f_{\hat{\theta}}$. At “test time”, we look at the error, $(Y_{\text{observed}} - f_{\hat{\theta}}(x))$. There are three main sources of error:

1. **Irreducible error**: This is due to **noise or randomness from $Y|X$ itself**. In short, there is some level of noise that is impossible for our model to account for. One possible situation is that the underlying response is a deterministic function of X (i.e. not random), but there is randomness in the measurement. It is possible that our model perfectly predicts the underlying signal, but the model will still disagree with Y_{observed} and contribute to the error. In the classic setup of $y = f_{\text{true}}(x) + \epsilon_{\text{noise}}$, the ϵ_{noise} term contributes to irreducible error.
2. **Approximation error**: This error **comes from limited expressive power of f_{θ} as a finitely-parameterized model**. In other words, our model isn't “flexible” enough to capture the true signal. For example, trying to fit the function $y = \cos(x)$ using a single 6th degree polynomial model, $f_{\theta}(x) = \sum_{i=0}^6 a_i x^i$.
3. **Estimation error**: There are two components to this type of error, both of which are well-covered in prerequisite machine learning courses:
 - (a) **Bias**: Bias captures the systematic error of our learning algorithm and training data in terms of making predictions. We write this mathematically by $\mathbb{E}_{\epsilon, \mathcal{D}}[f_{\hat{\theta}(\mathcal{D})}(X) - Y | X]$. Note that ϵ represents the randomness in $Y|X$ and \mathcal{D} represents the randomness in the training process used to get $\hat{\theta}$. For example, this could include the training data set we used or the splits made in random forest trees. Traditionally, X is separate from \mathcal{D} and is not random. It is common to look at the squared bias to prevent positive and negative biases for different observations from canceling.
 - (b) **Variance**: **This is the variable part of error**. We write it as $\mathbb{E}_{\mathcal{D}}[(f_{\hat{\theta}(\mathcal{D})}(X) - \mathbb{E}_{\mathcal{D}}[f_{\hat{\theta}(\mathcal{D})}(X)])^2 | X]$, where \mathcal{D} again represents the randomness of the training process. It describes how much our prediction “shakes” as a function of the randomness in the training.

This perspective can be useful and many papers utilize the bias-variance decomposition in their derivations. However, it doesn't always match what our intuition might be, especially in deep learning settings. The following figure is intended to help us understand this point.

This drawing is for high-level intuition, so we mustn't allow ourselves to become confused over the exact dimensions and projections. We first turn our attention to the figure on the left. The line passing through

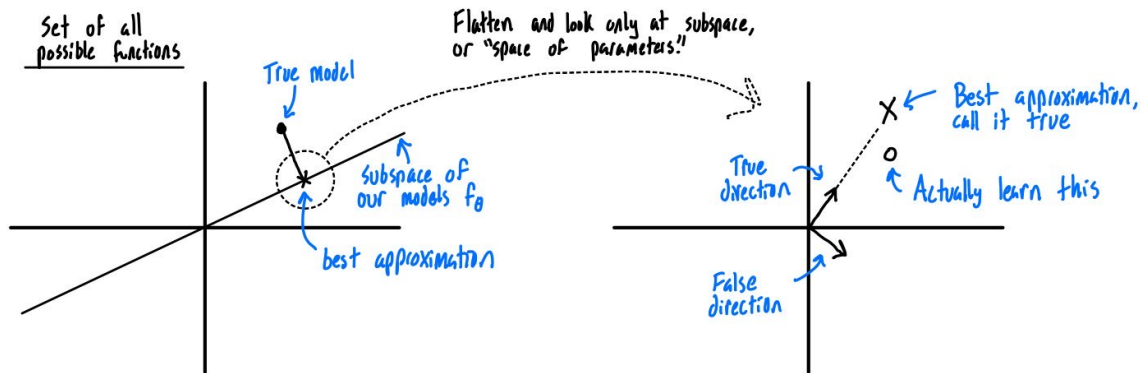


Figure 1: Approximation and Estimation Error

the origin represents the subspace models that we are considering. Note that we are only capable of providing estimates that fall along this line. As we can see, the true model is not in this subspace so we will suffer a certain amount of approximation error.

Turning our attention to the figure on the right, we project the true model onto the subspace of models that we are considering. The result is the best possible approximation and is the goal of our learning algorithm. We see in the figure that our predicted model is not quite equal to the best approximation. In this perspective, instead of thinking of things in terms of bias and variance, we consider how much of the “true” and “false” directions we are incorporating into our prediction.

We make this a bit more formal by discussing *survival* and *contamination*. Survival reflects, in expectation, how much of the true pattern survives the estimation/learning process. Contamination is how much useless information get “learned”, e.g. spurious features that our model is capable of picking up but don’t help with prediction. Survival and contamination can be thought of as the intuition behind bias and variance, respectively.

3 What are “Features”?

The following is a simplified sketch of a neural network with ℓ hidden layers.

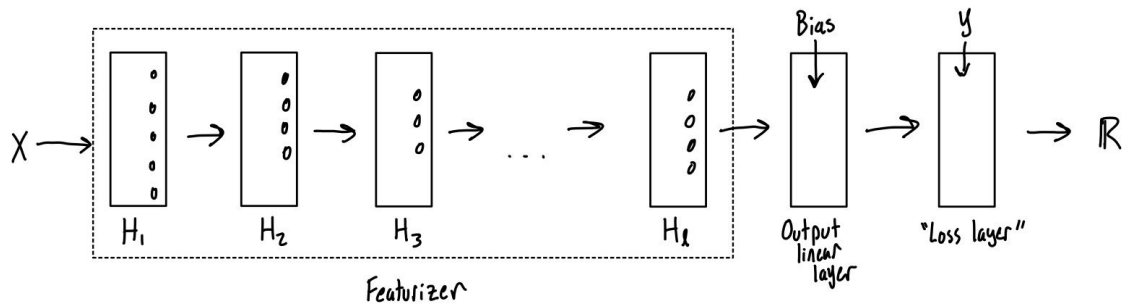


Figure 2: Simplified DNN

If we treat everything before the output of H_ℓ as a black box, we can think of things from the perspective of a generalized linear model. We have a featurizer, some linear function of those new features, and a loss that we are optimizing. The featurizer lifts or distills the input \mathbf{X} into a nicer feature space. In this perspective, the “learned” features are the outputs of the penultimate layer in the featurizer, H_ℓ . We want the featurization to be data-driven instead of hand-picked, so the layers essentially find a representation of the data that allows the generalized linear model (GLM) to work well.

However, there is another important point of view. Suppose a generalized linear model is given by $\hat{y} = \sum_{i=1}^{\lambda} w_i \phi_i(\mathbf{x})$, where $\phi(\mathbf{x})$ is the output of the “featurizer” in the figure above. When we are actually using the model on new data, these are simply the features. However, from the training perspective, they also determine the gradient. The derivative of the linear model with respect to the i^{th} parameter is $\frac{\partial \hat{y}}{\partial w_i} = \phi_i(\mathbf{x})$.

In short, we don’t understand nonlinear systems well, so our standard approach to understanding a nonlinear system is local linearization. We zoom in until things are roughly linear around a certain point. In terms of DNNs, we say that the deep network is Taylor expanded around the features in such a way that it is some constant term plus a local GLM in which small increments of the features change the predictions in a small way. This will be covered in greater detail during the next lecture.