

# CS229 Lecture notes

Andrew Ng

## Part VI

# Regularization and model selection

Suppose we are trying select among several different models for a learning problem. For instance, we might be using a polynomial regression model  $h_{\theta}(x) = g(\theta_0 + \theta_1x + \theta_2x^2 + \dots + \theta_kx^k)$ , and wish to decide if  $k$  should be 0, 1, ..., or 10. How can we automatically select a model that represents a good tradeoff between the twin evils of bias and variance<sup>1</sup>? Alternatively, suppose we want to automatically choose the bandwidth parameter  $\tau$  for locally weighted regression, or the parameter  $C$  for our  $\ell_1$ -regularized SVM. How can we do that?

For the sake of concreteness, in these notes we assume we have some finite set of models  $\mathcal{M} = \{M_1, \dots, M_d\}$  that we're trying to select among. For instance, in our first example above, the model  $M_i$  would be an  $i$ -th order polynomial regression model. (The generalization to infinite  $\mathcal{M}$  is not hard.<sup>2</sup>) Alternatively, if we are trying to decide between using an SVM, a neural network or logistic regression, then  $\mathcal{M}$  may contain these models.

---

<sup>1</sup>Given that we said in the previous set of notes that bias and variance are two very different beasts, some readers may be wondering if we should be calling them “twin” evils here. Perhaps it'd be better to think of them as non-identical twins. The phrase “the fraternal twin evils of bias and variance” doesn't have the same ring to it, though.

<sup>2</sup>If we are trying to choose from an infinite set of models, say corresponding to the possible values of the bandwidth  $\tau \in \mathbb{R}^+$ , we may discretize  $\tau$  and consider only a finite number of possible values for it. More generally, most of the algorithms described here can all be viewed as performing optimization search in the space of models, and we can perform this search over infinite model classes as well.

# 1 Cross validation

Lets suppose we are, as usual, given a training set  $S$ . Given what we know about empirical risk minimization, here's what might initially seem like a algorithm, resulting from using empirical risk minimization for model selection:

1. Train each model  $M_i$  on  $S$ , to get some hypothesis  $h_i$ .
2. Pick the hypotheses with the smallest training error.

This algorithm does *not* work. Consider choosing the order of a polynomial. The higher the order of the polynomial, the better it will fit the training set  $S$ , and thus the lower the training error. Hence, this method will always select a high-variance, high-degree polynomial model, which we saw previously is often poor choice.

Here's an algorithm that works better. In **hold-out cross validation** (also called **simple cross validation**), we do the following:

1. Randomly split  $S$  into  $S_{\text{train}}$  (say, 70% of the data) and  $S_{\text{cv}}$  (the remaining 30%). Here,  $S_{\text{cv}}$  is called the hold-out cross validation set.
2. Train each model  $M_i$  on  $S_{\text{train}}$  only, to get some hypothesis  $h_i$ .
3. Select and output the hypothesis  $h_i$  that had the smallest error  $\hat{\epsilon}_{S_{\text{cv}}}(h_i)$  on the hold out cross validation set. (Recall,  $\hat{\epsilon}_{S_{\text{cv}}}(h)$  denotes the empirical error of  $h$  on the set of examples in  $S_{\text{cv}}$ .)

By testing on a set of examples  $S_{\text{cv}}$  that the models were not trained on, we obtain a better estimate of each hypothesis  $h_i$ 's true generalization error, and can then pick the one with the smallest estimated generalization error. Usually, somewhere between 1/4 – 1/3 of the data is used in the hold out cross validation set, and 30% is a typical choice.

Optionally, step 3 in the algorithm may also be replaced with selecting the model  $M_i$  according to  $\arg \min_i \hat{\epsilon}_{S_{\text{cv}}}(h_i)$ , and then retraining  $M_i$  on the entire training set  $S$ . (This is often a good idea, with one exception being learning algorithms that are be very sensitive to perturbations of the initial conditions and/or data. For these methods,  $M_i$  doing well on  $S_{\text{train}}$  does not necessarily mean it will also do well on  $S_{\text{cv}}$ , and it might be better to forgo this retraining step.)

The disadvantage of using hold out cross validation is that it “wastes” about 30% of the data. Even if we were to take the optional step of retraining

the model on the entire training set, it's still as if we're trying to find a good model for a learning problem in which we had  $0.7m$  training examples, rather than  $n$  training examples, since we're testing models that were trained on only  $0.7m$  examples each time. While this is fine if data is abundant and/or cheap, in learning problems in which data is scarce (consider a problem with  $m = 20$ , say), we'd like to do something better.

Here is a method, called  **$k$ -fold cross validation**, that holds out less data each time:

1. Randomly split  $S$  into  $k$  disjoint subsets of  $m/k$  training examples each. Lets call these subsets  $S_1, \dots, S_k$ .
2. For each model  $M_i$ , we evaluate it as follows:

For  $j = 1, \dots, k$

Train the model  $M_i$  on  $S_1 \cup \dots \cup S_{j-1} \cup S_{j+1} \cup \dots \cup S_k$  (i.e., train on all the data except  $S_j$ ) to get some hypothesis  $h_{ij}$ .

Test the hypothesis  $h_{ij}$  on  $S_j$ , to get  $\hat{\epsilon}_{S_j}(h_{ij})$ .

The estimated generalization error of model  $M_i$  is then calculated as the average of the  $\hat{\epsilon}_{S_j}(h_{ij})$ 's (averaged over  $j$ ).

3. Pick the model  $M_i$  with the lowest estimated generalization error, and retrain that model on the entire training set  $S$ . The resulting hypothesis is then output as our final answer.

A typical choice for the number of folds to use here would be  $k = 10$ . While the fraction of data held out each time is now  $1/k$ —much smaller than before—this procedure may also be more computationally expensive than hold-out cross validation, since we now need train to each model  $k$  times.

While  $k = 10$  is a commonly used choice, in problems in which data is really scarce, sometimes we will use the extreme choice of  $k = m$  in order to leave out as little data as possible each time. In this setting, we would repeatedly train on all but one of the training examples in  $S$ , and test on that held-out example. The resulting  $m = k$  errors are then averaged together to obtain our estimate of the generalization error of a model. This method has its own name; since we're holding out one training example at a time, this method is called **leave-one-out cross validation**.

Finally, even though we have described the different versions of cross validation as methods for selecting a model, they can also be used more simply to evaluate a *single* model or algorithm. For example, if you have implemented

some learning algorithm and want to estimate how well it performs for your application (or if you have invented a novel learning algorithm and want to report in a technical paper how well it performs on various test sets), cross validation would give a reasonable way of doing so.

## 2 Feature Selection

One special and important case of model selection is called feature selection. To motivate this, imagine that you have a supervised learning problem where the number of features  $d$  is very large (perhaps  $n \gg n$ ), but you suspect that there is only a small number of features that are “relevant” to the learning task. Even if you use the a simple linear classifier (such as the perceptron) over the  $d$  input features, the VC dimension of your hypothesis class would still be  $O(n)$ , and thus overfitting would be a potential problem unless the training set is fairly large.

In such a setting, you can apply a feature selection algorithm to reduce the number of features. Given  $d$  features, there are  $2^d$  possible feature subsets (since each of the  $d$  features can either be included or excluded from the subset), and thus feature selection can be posed as a model selection problem over  $2^d$  possible models. For large values of  $d$ , it’s usually too expensive to explicitly enumerate over and compare all  $2^d$  models, and so typically some heuristic search procedure is used to find a good feature subset. The following search procedure is called **forward search**:

1. Initialize  $\mathcal{F} = \emptyset$ .
2. Repeat {
  - (a) For  $i = 1, \dots, d$  if  $i \notin \mathcal{F}$ , let  $\mathcal{F}_i = \mathcal{F} \cup \{i\}$ , and use some version of cross validation to evaluate features  $\mathcal{F}_i$ . (I.e., train your learning algorithm using only the features in  $\mathcal{F}_i$ , and estimate its generalization error.)
  - (b) Set  $\mathcal{F}$  to be the best feature subset found on step (a).
3. Select and output the best feature subset that was evaluated during the entire search procedure.

The outer loop of the algorithm can be terminated either when  $\mathcal{F} = \{1, \dots, d\}$  is the set of all features, or when  $|\mathcal{F}|$  exceeds some pre-set threshold (corresponding to the maximum number of features that you want the algorithm to consider using).

This algorithm described above one instantiation of **wrapper model feature selection**, since it is a procedure that “wraps” around your learning algorithm, and repeatedly makes calls to the learning algorithm to evaluate how well it does using different feature subsets. Aside from forward search, other search procedures can also be used. For example, **backward search** starts off with  $\mathcal{F} = \{1, \dots, d\}$  as the set of all features, and repeatedly deletes features one at a time (evaluating single-feature deletions in a similar manner to how forward search evaluates single-feature additions) until  $\mathcal{F} = \emptyset$ .

Wrapper feature selection algorithms often work quite well, but can be computationally expensive given how that they need to make many calls to the learning algorithm. Indeed, complete forward search (terminating when  $\mathcal{F} = \{1, \dots, d\}$ ) would take about  $O(n^2)$  calls to the learning algorithm.

**Filter feature selection** methods give heuristic, but computationally much cheaper, ways of choosing a feature subset. The idea here is to compute some simple score  $S(i)$  that measures how informative each feature  $x_i$  is about the class labels  $y$ . Then, we simply pick the  $k$  features with the largest scores  $S(i)$ .

One possible choice of the score would be define  $S(i)$  to be (the absolute value of) the correlation between  $x_i$  and  $y$ , as measured on the training data. This would result in our choosing the features that are the most strongly correlated with the class labels. In practice, it is more common (particularly for discrete-valued features  $x_i$ ) to choose  $S(i)$  to be the **mutual information**  $\text{MI}(x_i, y)$  between  $x_i$  and  $y$ :

$$\text{MI}(x_i, y) = \sum_{x_i \in \{0,1\}} \sum_{y \in \{0,1\}} p(x_i, y) \log \frac{p(x_i, y)}{p(x_i)p(y)}.$$

(The equation above assumes that  $x_i$  and  $y$  are binary-valued; more generally the summations would be over the domains of the variables.) The probabilities above  $p(x_i, y)$ ,  $p(x_i)$  and  $p(y)$  can all be estimated according to their empirical distributions on the training set.

To gain intuition about what this score does, note that the mutual information can also be expressed as a Kullback-Leibler (KL) divergence:

$$\text{MI}(x_i, y) = \text{KL}(p(x_i, y) || p(x_i)p(y))$$

You’ll get to play more with KL-divergence in Problem set #3, but informally, this gives a measure of how different the probability distributions

$p(x_i, y)$  and  $p(x_i)p(y)$  are. If  $x_i$  and  $y$  are independent random variables, then we would have  $p(x_i, y) = p(x_i)p(y)$ , and the KL-divergence between the two distributions will be zero. This is consistent with the idea if  $x_i$  and  $y$  are independent, then  $x_i$  is clearly very “non-informative” about  $y$ , and thus the score  $S(i)$  should be small. Conversely, if  $x_i$  is very “informative” about  $y$ , then their mutual information  $\text{MI}(x_i, y)$  would be large.

One final detail: Now that you’ve ranked the features according to their scores  $S(i)$ , how do you decide how many features  $k$  to choose? Well, one standard way to do so is to use cross validation to select among the possible values of  $k$ . For example, when applying naive Bayes to text classification—a problem where  $d$ , the vocabulary size, is usually very large—using this method to select a feature subset often results in increased classifier accuracy.

### 3 Bayesian statistics and regularization

In this section, we will talk about one more tool in our arsenal for our battle against overfitting.

At the beginning of the quarter, we talked about parameter fitting using maximum likelihood estimation (MLE), and chose our parameters according to

$$\theta_{\text{MLE}} = \arg \max_{\theta} \prod_{i=1}^n p(y^{(i)} | x^{(i)}; \theta).$$

Throughout our subsequent discussions, we viewed  $\theta$  as an unknown parameter of the world. This view of the  $\theta$  as being *constant-valued but unknown* is taken in **frequentist** statistics. In the frequentist this view of the world,  $\theta$  is not random—it just happens to be unknown—and it’s our job to come up with statistical procedures (such as maximum likelihood) to try to estimate this parameter.

An alternative way to approach our parameter estimation problems is to take the **Bayesian** view of the world, and think of  $\theta$  as being a *random variable* whose value is unknown. In this approach, we would specify a **prior distribution**  $p(\theta)$  on  $\theta$  that expresses our “prior beliefs” about the parameters. Given a training set  $S = \{(x^{(i)}, y^{(i)})\}_{i=1}^n$ , when we are asked to make a prediction on a new value of  $x$ , we can then compute the posterior

distribution on the parameters

$$\begin{aligned} p(\theta|S) &= \frac{p(S|\theta)p(\theta)}{p(S)} \\ &= \frac{\left(\prod_{i=1}^n p(y^{(i)}|x^{(i)}, \theta)\right) p(\theta)}{\int_{\theta} \left(\prod_{i=1}^n p(y^{(i)}|x^{(i)}, \theta)\right) p(\theta) d\theta} \end{aligned} \quad (1)$$

In the equation above,  $p(y^{(i)}|x^{(i)}, \theta)$  comes from whatever model you're using for your learning problem. For example, if you are using Bayesian logistic regression, then you might choose  $p(y^{(i)}|x^{(i)}, \theta) = h_{\theta}(x^{(i)})^{y^{(i)}}(1-h_{\theta}(x^{(i)}))^{(1-y^{(i)})}$ , where  $h_{\theta}(x^{(i)}) = 1/(1 + \exp(-\theta^T x^{(i)}))$ .<sup>3</sup>

When we are given a new test example  $x$  and asked to make it prediction on it, we can compute our posterior distribution on the class label using the posterior distribution on  $\theta$ :

$$p(y|x, S) = \int_{\theta} p(y|x, \theta)p(\theta|S)d\theta \quad (2)$$

In the equation above,  $p(\theta|S)$  comes from Equation (1). Thus, for example, if the goal is to the predict the expected value of  $y$  given  $x$ , then we would output<sup>4</sup>

$$E[y|x, S] = \int_y yp(y|x, S)dy$$

The procedure that we've outlined here can be thought of as doing "fully Bayesian" prediction, where our prediction is computed by taking an average with respect to the posterior  $p(\theta|S)$  over  $\theta$ . Unfortunately, in general it is computationally very difficult to compute this posterior distribution. This is because it requires taking integrals over the (usually high-dimensional)  $\theta$  as in Equation (1), and this typically cannot be done in closed-form.

Thus, in practice we will instead approximate the posterior distribution for  $\theta$ . One common approximation is to replace our posterior distribution for  $\theta$  (as in Equation 2) with a single point estimate. The **MAP (maximum a posteriori)** estimate for  $\theta$  is given by

$$\theta_{\text{MAP}} = \arg \max_{\theta} \prod_{i=1}^n p(y^{(i)}|x^{(i)}, \theta)p(\theta). \quad (3)$$

---

<sup>3</sup>Since we are now viewing  $\theta$  as a random variable, it is okay to condition on it value, and write " $p(y|x, \theta)$ " instead of " $p(y|x; \theta)$ ."

<sup>4</sup>The integral below would be replaced by a summation if  $y$  is discrete-valued.

Note that this is the same formulas as for the MLE (maximum likelihood) estimate for  $\theta$ , except for the prior  $p(\theta)$  term at the end.

In practical applications, a common choice for the prior  $p(\theta)$  is to assume that  $\theta \sim \mathcal{N}(0, \tau^2 I)$ . Using this choice of prior, the fitted parameters  $\theta_{\text{MAP}}$  will have smaller norm than that selected by maximum likelihood. (See Problem Set #3.) In practice, this causes the Bayesian MAP estimate to be less susceptible to overfitting than the ML estimate of the parameters. For example, Bayesian logistic regression turns out to be an effective algorithm for text classification, even though in text classification we usually have  $d \gg n$ .