

## Lecture 7 → Regularization and Model Selection

• MLE → Maximum Likelihood Estimation

Objective is to maximize  $\log P(D|\omega)$

Example is Linear Regression log-likelihood

• MAP → Maximum a Posteriori

Objective is to maximize  $\log P(D|\omega) + \log P(\omega)$  log prior

Example is Regularized Linear Regression

• MLE discriminative model →  $\log P(D|\omega) = \log \prod_{n=1}^N P(y^{(n)} | \phi(x^{(n)}), \omega)$

$$= \sum_{n=1}^N \log P(y^{(n)} | \phi(x^{(n)}), \omega)$$

• MLE generative model →  $\log P(D|\omega) = \sum_{n=1}^N \log P(y^{(n)}, x^{(n)} | \omega)$

• Issue with MLE is that we run the risk of overfitting

• MAP assumes prior dist<sup>n</sup>  $P(\omega)$

• Point estimate using Bayes rule →  $\arg \max_{\omega} P(\omega | D) = \arg \max_{\omega} P(D | \omega) P(\omega)$

$$\log P(D|\omega) P(\omega) = \log \prod_{n=1}^N P(y^{(n)} | \phi(x^{(n)}), \omega) + \log P(\omega)$$

$$= \sum_{n=1}^N \log P(y^{(n)} | \phi(x^{(n)}), \omega) + \log P(\omega)$$

• Typically L2-norm is used as prior →  $P(\omega) = \mathcal{N}(0, \lambda^{-1} I)$

$$\log P(\omega) = -\frac{\lambda}{2} \|\omega\|^2 + \text{constant}$$

$$P(\omega) = \mathcal{N}(0, \lambda^{-1} I)$$

$$= \text{const} * \exp\left(-\frac{1}{2} \omega^T (\lambda^{-1} I)^{-1} \omega\right)$$

$$= \text{const} * \exp\left(-\frac{\lambda}{2} \omega^T \omega\right)$$

$$\log P(\omega) = \text{const} - \frac{\lambda}{2} \omega^T \omega$$

$$= \text{const} - \frac{\lambda}{2} \|\omega\|^2$$

• MAP derivation →  $\log P(D|\omega) P(\omega) = \sum_{n=1}^N \log P(y^{(n)} | \phi(x^{(n)}), \omega) + \log P(\omega)$

$$= \sum_{n=1}^N \log \left( \sqrt{\frac{\beta}{2\pi}} \exp\left(-\frac{\beta}{2} \|y^{(n)} - \omega^T \phi(x^{(n)})\|^2\right) \right) + \text{const} - \frac{\lambda}{2} \|\omega\|^2$$

$$= \frac{N}{2} \log \beta - \frac{N}{2} \log 2\pi - \sum_{n=1}^N \frac{\beta}{2} \|y^{(n)} - \omega^T \phi(x^{(n)})\|^2 + \text{const} - \frac{\lambda}{2} \|\omega\|^2$$

• Minimize →  $\tilde{E}[\omega] = \sum_{n=1}^N \frac{\beta}{2} \|y^{(n)} - \omega^T \phi(x^{(n)})\|^2 + \frac{\lambda}{2} \|\omega\|^2 + \text{const}$  w.r.t  $\omega$  which is minimized by,

$$\omega_{ML} = (\lambda I + \Phi^T \Phi)^{-1} \Phi^T y$$

• Overfitting → when the test accuracy is significantly higher than the training accuracy due to the model fitting too strongly to the training set and failing to generalize

• Underfitting → when the model fails to capture the relationship between the input and output, leading to high error on both the training set and the unseen data



$$E_{rms} = \sqrt{2E[w^2]/N}$$

- L2 regularization controls the tradeoff between fitting error and complexity
- Small L2 regularization results in complex models, but with the risk of overfitting
- Large L2 regularization results in simple models, but with the risk of underfitting

To avoid overfitting, ① more training data and ② regularization

## Bias Variance Tradeoff

We want to learn a model with,

- ① Small bias  $\rightarrow$  how well a model fits the data on average?
- ② Small variance  $\rightarrow$  how stable a model is w.r.t data samples?

If we have multiple datasets, each of size  $N$ , then any particular dataset,  $D$ , will give a particular function  $h(x; D)$

We have,

$$\begin{aligned} E_D \int \{h(x; D) - E[y|x]\}^2 p(x) dx \\ = \underbrace{\int \{E_D[h(x; D)] - E[y|x]\}^2 p(x) dx}_{\text{bias}^2} + \underbrace{\int E_D[\{h(x; D) - E_D[h(x; D)]\}^2] p(x) dx}_{\text{variance}} \end{aligned}$$

Expected loss  $\rightarrow E[L] = \iint \{h(x) - y\}^2 p(x, y) dx dy$   
 $= (\text{bias})^2 + \text{variance} + \text{noise}$

$$\begin{aligned} \text{where } (\text{bias})^2 &= \int \{E_D[h(x; D)] - E[y|x]\}^2 p(x) dx \\ \text{variance} &= \int E_D[\{h(x; D) - E_D[h(x; D)]\}^2] p(x) dx \\ \text{noise} &= \iint \{E[y|x] - y\}^2 p(x, y) dx dy \end{aligned}$$

- Over-regularized model (large  $\lambda$ ) will have a high bias and low variance
- Under-regularized model (small  $\lambda$ ) will have high variance and low bias

## Model Selection

- Hold out cross validation  $\rightarrow$  ① Randomly split  $D$  into  $D_{\text{train}}$  and  $D_{\text{val}}$   
70%:30% split is typical
  - ② train each model  $M_i$  on  $D_{\text{train}}$  to get some hypothesis  $h_i$
  - ③ select and output hypothesis  $h_i$  that has the smallest error on the holdout validation set
- Disadvantages  $\rightarrow$  ① lose data in  $D_{\text{val}}$  as we cannot train on this  
② some data is only used for training, some is only used for validation

- K-fold cross validation  $\rightarrow$  create a K-fold partition of the dataset  
for each of the K experiments, use K-1 folds for training and the remaining one for validating  
true error is estimated as the average error rate
- Special case,  $K=N \rightarrow$  leave one out cross validation (LOOCV)  
- expensive, but wastes least amount of training data for cross validation  
- might be useful when data size is not big
- Popular values of  $K=3, 5, 10$



• Training set → set of examples used for learning

to fit the parameters to the classifier

• Validation set → set of examples used to tune hyperparameters of a classifier

use to find the "optimal" hyperparameters

• Test set → set of examples used to assess performance of a fully trained classifier

• Test set is used for evaluation, not for tuning your parameters