## Bias/Variance Trade-off, Some Practical Issues, Semi-supervised and Active Learning

Piyush Rai

#### Introduction to Machine Learning (CS771A)

November 13, 2018

Intro to Machine Learning (CS771A)

Bias/Variance Trade-off, Some Practical Issues, Semi-supervised and Active Learning

- Bias/Variance Trade-off of Learning Algorithms
- Some basic guidelines for debugging ML algorithms
- What if training data and test data distributions are NOT the same
- How to utilize unlabeled data: Semi-supervised Learning
- How to selectively ask for the most informative data: Active Learning

イロト 人間 とくほとく ほど

## Bias/Variance Trade-off



< ロ > < 回 > < 回 > < 回 > < 回 >

• Assume  $\mathcal{F}$  to be a class of models (e.g., set of linear classifiers with some pre-defined features)



・ロト ・日 ト ・日 ト ・日 ト

- Assume  $\mathcal{F}$  to be a class of models (e.g., set of linear classifiers with some pre-defined features)
- Suppose we've learned a model  $f \in \mathcal{F}$  learned using some (finite amount of) training data



・ロト ・日 ト ・日 ト ・日 ト

- Assume  $\mathcal{F}$  to be a class of models (e.g., set of linear classifiers with some pre-defined features)
- Suppose we've learned a model  $f \in \mathcal{F}$  learned using some (finite amount of) training data
- Can decompose the test error  $\epsilon(f)$  of f as follows

$$\epsilon(f) = \underbrace{\left[\min_{f^* \in \mathcal{F}} \epsilon(f^*)\right]}_{\text{approximation error}} + \underbrace{\left[\epsilon(f) - \min_{f^* \in \mathcal{F}} \epsilon(f^*)\right]}_{\text{estimation error}}$$

- Assume  $\mathcal{F}$  to be a class of models (e.g., set of linear classifiers with some pre-defined features)
- Suppose we've learned a model  $f \in \mathcal{F}$  learned using some (finite amount of) training data
- Can decompose the test error  $\epsilon(f)$  of f as follows

$$\epsilon(f) = \underbrace{\left[\min_{f^* \in \mathcal{F}} \epsilon(f^*)\right]}_{\text{approximation error}} + \underbrace{\left[\epsilon(f) - \min_{f^* \in \mathcal{F}} \epsilon(f^*)\right]}_{\text{estimation error}}$$

• Here  $f^*$  is the best possible model in  $\mathcal{F}$  assuming infinite amount of training data

(日) (四) (三) (三) (三)

- Assume  $\mathcal{F}$  to be a class of models (e.g., set of linear classifiers with some pre-defined features)
- Suppose we've learned a model  $f \in \mathcal{F}$  learned using some (finite amount of) training data
- Can decompose the test error  $\epsilon(f)$  of f as follows

$$\epsilon(f) = \underbrace{\left[\min_{f^* \in \mathcal{F}} \epsilon(f^*)\right]}_{\text{approximation error}} + \underbrace{\left[\epsilon(f) - \min_{f^* \in \mathcal{F}} \epsilon(f^*)\right]}_{\text{estimation error}}$$

- Here  $f^*$  is the best possible model in  $\mathcal{F}$  assuming infinite amount of training data
- Approximation error: Error of best model  $f^*$  because of the model class  $\mathcal F$  being too simple

- Assume  $\mathcal{F}$  to be a class of models (e.g., set of linear classifiers with some pre-defined features)
- Suppose we've learned a model  $f \in \mathcal{F}$  learned using some (finite amount of) training data
- Can decompose the test error  $\epsilon(f)$  of f as follows

$$\epsilon(f) = \underbrace{\left[\min_{f^* \in \mathcal{F}} \epsilon(f^*)\right]}_{\text{approximation error}} + \underbrace{\left[\epsilon(f) - \min_{f^* \in \mathcal{F}} \epsilon(f^*)\right]}_{\text{estimation error}}$$

- Here  $f^*$  is the best possible model in  $\mathcal{F}$  assuming infinite amount of training data
- Approximation error: Error of best model  $f^*$  because of the model class  $\mathcal{F}$  being too simple
- Estimation error: Error of learned model f (relative to  $f^*$ ) because we only had finite training data

4日 + 4日 + 4日 + 4日 + (単)の(で)

• There is a trade-off between the two components of the error

$$\epsilon(f) = \underbrace{\left[\min_{f^* \in \mathcal{F}} \epsilon(f^*)\right]}_{\text{approximation error}} + \underbrace{\left[\epsilon(f) - \min_{f^* \in \mathcal{F}} \epsilon(f^*)\right]}_{\text{estimation error}}$$

・ロト ・日 ・ ・ ヨ ・ ・ ヨ ・

• There is a trade-off between the two components of the error



 $\bullet$  Can reduce the approximation error by making  ${\cal F}$  more complex/richer

イロト イポト イヨト イヨ

• There is a trade-off between the two components of the error

$$\epsilon(f) = \underbrace{\left[\min_{f^* \in \mathcal{F}} \epsilon(f^*)\right]}_{\text{approximation error}} + \underbrace{\left[\epsilon(f) - \min_{f^* \in \mathcal{F}} \epsilon(f^*)\right]}_{\text{estimation error}}$$

- $\bullet$  Can reduce the approximation error by making  ${\cal F}$  more complex/richer, e.g.,
  - Replace the linear classifiers by a higher order polynomials

• There is a trade-off between the two components of the error

$$\epsilon(f) = \underbrace{\left[\min_{f^* \in \mathcal{F}} \epsilon(f^*)\right]}_{\text{approximation error}} + \underbrace{\left[\epsilon(f) - \min_{f^* \in \mathcal{F}} \epsilon(f^*)\right]}_{\text{estimation error}}$$

- Can reduce the approximation error by making  ${\cal F}$  more complex/richer, e.g.,
  - Replace the linear classifiers by a higher order polynomials
  - Add more features in the data

• There is a trade-off between the two components of the error

$$\epsilon(f) = \underbrace{\left[\min_{f^* \in \mathcal{F}} \epsilon(f^*)\right]}_{\text{approximation error}} + \underbrace{\left[\epsilon(f) - \min_{f^* \in \mathcal{F}} \epsilon(f^*)\right]}_{\text{estimation error}}$$

- Can reduce the approximation error by making  ${\cal F}$  more complex/richer, e.g.,
  - Replace the linear classifiers by a higher order polynomials
  - Add more features in the data
- However, making  $\mathcal{F}$  richer will usually cause estimation error to increase since f can now overfit

• There is a trade-off between the two components of the error

$$\epsilon(f) = \underbrace{\left[\min_{f^* \in \mathcal{F}} \epsilon(f^*)\right]}_{\text{approximation error}} + \underbrace{\left[\epsilon(f) - \min_{f^* \in \mathcal{F}} \epsilon(f^*)\right]}_{\text{estimation error}}$$

- Can reduce the approximation error by making  ${\cal F}$  more complex/richer, e.g.,
  - Replace the linear classifiers by a higher order polynomials
  - Add more features in the data
- However, making  $\mathcal{F}$  richer will usually cause estimation error to increase since f can now overfit
  - f may give small error on some datasets, large error on some datasets

• There is a trade-off between the two components of the error

$$\epsilon(f) = \underbrace{\left[\min_{f^* \in \mathcal{F}} \epsilon(f^*)\right]}_{\text{approximation error}} + \underbrace{\left[\epsilon(f) - \min_{f^* \in \mathcal{F}} \epsilon(f^*)\right]}_{\text{estimation error}}$$

- Can reduce the approximation error by making  ${\cal F}$  more complex/richer, e.g.,
  - Replace the linear classifiers by a higher order polynomials
  - Add more features in the data
- However, making  $\mathcal{F}$  richer will usually cause estimation error to increase since f can now overfit
  - f may give small error on some datasets, large error on some datasets
- Note: The above trade-off is popularly known as the bias/variance trade-off

イロト 人間 とくほとく ほど

• There is a trade-off between the two components of the error

$$\epsilon(f) = \underbrace{\left[\min_{f^* \in \mathcal{F}} \epsilon(f^*)\right]}_{\text{approximation error}} + \underbrace{\left[\epsilon(f) - \min_{f^* \in \mathcal{F}} \epsilon(f^*)\right]}_{\text{estimation error}}$$

- Can reduce the approximation error by making  ${\cal F}$  more complex/richer, e.g.,
  - Replace the linear classifiers by a higher order polynomials
  - Add more features in the data
- However, making  $\mathcal{F}$  richer will usually cause estimation error to increase since f can now overfit
  - f may give small error on some datasets, large error on some datasets
- Note: The above trade-off is popularly known as the bias/variance trade-off
  - Approximation error known as "bias"

• There is a trade-off between the two components of the error

$$\epsilon(f) = \underbrace{\left[\min_{f^* \in \mathcal{F}} \epsilon(f^*)\right]}_{\text{approximation error}} + \underbrace{\left[\epsilon(f) - \min_{f^* \in \mathcal{F}} \epsilon(f^*)\right]}_{\text{estimation error}}$$

- Can reduce the approximation error by making  ${\cal F}$  more complex/richer, e.g.,
  - Replace the linear classifiers by a higher order polynomials
  - Add more features in the data
- However, making  $\mathcal{F}$  richer will usually cause estimation error to increase since f can now overfit
  - f may give small error on some datasets, large error on some datasets
- Note: The above trade-off is popularly known as the bias/variance trade-off
  - Approximation error known as "bias" (high if the model is simple)

イロン スピン メヨン メヨン

• There is a trade-off between the two components of the error

$$\epsilon(f) = \underbrace{\left[\min_{f^* \in \mathcal{F}} \epsilon(f^*)\right]}_{\text{approximation error}} + \underbrace{\left[\epsilon(f) - \min_{f^* \in \mathcal{F}} \epsilon(f^*)\right]}_{\text{estimation error}}$$

- Can reduce the approximation error by making  ${\cal F}$  more complex/richer, e.g.,
  - Replace the linear classifiers by a higher order polynomials
  - Add more features in the data
- However, making  $\mathcal{F}$  richer will usually cause estimation error to increase since f can now overfit
  - f may give small error on some datasets, large error on some datasets
- Note: The above trade-off is popularly known as the bias/variance trade-off
  - Approximation error known as "bias" (high if the model is simple)
  - Estimation error known as "variance"

イロン スピン メヨン メヨン

• There is a trade-off between the two components of the error

$$\epsilon(f) = \underbrace{\left[\min_{f^* \in \mathcal{F}} \epsilon(f^*)\right]}_{\text{approximation error}} + \underbrace{\left[\epsilon(f) - \min_{f^* \in \mathcal{F}} \epsilon(f^*)\right]}_{\text{estimation error}}$$

- Can reduce the approximation error by making  ${\cal F}$  more complex/richer, e.g.,
  - Replace the linear classifiers by a higher order polynomials
  - Add more features in the data
- However, making  $\mathcal{F}$  richer will usually cause estimation error to increase since f can now overfit
  - f may give small error on some datasets, large error on some datasets
- Note: The above trade-off is popularly known as the bias/variance trade-off
  - Approximation error known as "bias" (high if the model is simple)
  - Estimation error known as "variance" (high if the model is complex)

イロト 人間 とくほとく ほど

#### A Visual Illustration of the Trade-off



# Debugging ML Algorithms



イロト 不得 とうせい うけん

- A notoriously hard problem in general
  - Note that code for ML algorithms is not procedural but data-driven

- A notoriously hard problem in general
  - Note that code for ML algorithms is not procedural but data-driven
- What to do when our model (say logistic regression) isn't doing well on test data



- A notoriously hard problem in general
  - Note that code for ML algorithms is not procedural but data-driven
- What to do when our model (say logistic regression) isn't doing well on test data
  - Use more training examples to train the model?

- A notoriously hard problem in general
  - Note that code for ML algorithms is not procedural but data-driven
- What to do when our model (say logistic regression) isn't doing well on test data
  - Use more training examples to train the model?
  - Use a smaller number of features?

イロン スピン メヨン メヨン

- A notoriously hard problem in general
  - Note that code for ML algorithms is not procedural but data-driven
- What to do when our model (say logistic regression) isn't doing well on test data
  - Use more training examples to train the model?
  - Use a smaller number of features?

Intro to Machine Learning (CS771A)

• Introduce new features (can be combinations of existing features)?

- A notoriously hard problem in general
  - Note that code for ML algorithms is not procedural but data-driven
- What to do when our model (say logistic regression) isn't doing well on test data
  - Use more training examples to train the model?
  - Use a smaller number of features?
  - Introduce new features (can be combinations of existing features)?
  - Try tuning the regularization parameter?

イロン スピン メヨン メヨン

- A notoriously hard problem in general
  - Note that code for ML algorithms is not procedural but data-driven
- What to do when our model (say logistic regression) isn't doing well on test data
  - Use more training examples to train the model?
  - Use a smaller number of features?
  - Introduce new features (can be combinations of existing features)?
  - Try tuning the regularization parameter?
  - Run (the iterative) optimizer longer, i.e., for more iterations?

- A notoriously hard problem in general
  - Note that code for ML algorithms is not procedural but data-driven
- What to do when our model (say logistic regression) isn't doing well on test data
  - Use more training examples to train the model?
  - Use a smaller number of features?
  - Introduce new features (can be combinations of existing features)?
  - Try tuning the regularization parameter?
  - Run (the iterative) optimizer longer, i.e., for more iterations?
  - Change the optimization algorithm (e.g., GD to SGD or Newton..)?

イロン スピン メヨン メヨン

- A notoriously hard problem in general
  - Note that code for ML algorithms is not procedural but data-driven
- What to do when our model (say logistic regression) isn't doing well on test data
  - Use more training examples to train the model?
  - Use a smaller number of features?
  - Introduce new features (can be combinations of existing features)?
  - Try tuning the regularization parameter?
  - Run (the iterative) optimizer longer, i.e., for more iterations?
  - Change the optimization algorithm (e.g., GD to SGD or Newton..)?
  - Give up and switch to a different model (e.g., SVM)?

イロン スピン メヨン メヨン

- A notoriously hard problem in general
  - Note that code for ML algorithms is not procedural but data-driven
- What to do when our model (say logistic regression) isn't doing well on test data
  - Use more training examples to train the model?
  - Use a smaller number of features?
  - Introduce new features (can be combinations of existing features)?
  - Try tuning the regularization parameter?
  - Run (the iterative) optimizer longer, i.e., for more iterations?
  - Change the optimization algorithm (e.g., GD to SGD or Newton..)?
  - Give up and switch to a different model (e.g., SVM)?
- How to know what might be going wrong and how to debug?

#### **High Bias or High Variance?**

- The bad performance (low accuracy on test data) could be due either
  - High Bias (Underfitting)
  - High Variance (Overfitting)
- Looking at the training and test error can tell which of the two is the case



### **High Bias or High Variance?**

- The bad performance (low accuracy on test data) could be due either
  - High Bias (Underfitting)
  - High Variance (Overfitting)
- Looking at the training and test error can tell which of the two is the case



• High Bias: Both training and test errors are large

イロト イポト イヨト イヨ

### **High Bias or High Variance?**

- The bad performance (low accuracy on test data) could be due either
  - High Bias (Underfitting)
  - High Variance (Overfitting)
- Looking at the training and test error can tell which of the two is the case



- High Bias: Both training and test errors are large
- High Variance: Small training error, large test error (and huge gap)

Intro to Machine Learning (CS771A)

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・

### Some Guidelines for Debugging Learning Algorithms



• If the model has high bias (underfitting) then the model class is weak




- If the model has high bias (underfitting) then the model class is weak
  - Adding more training data won't help



- If the model has high bias (underfitting) then the model class is weak
  - Adding more training data won't help
  - Instead, try making the model class richer, or add more features (makes the model richer)



- If the model has high bias (underfitting) then the model class is weak
  - Adding more training data won't help
  - Instead, try making the model class richer, or add more features (makes the model richer)
- If the model has high variance (overfitting) then the model class is sufficiently rich

A B > A B > A B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A



- If the model has high bias (underfitting) then the model class is weak
  - Adding more training data won't help
  - Instead, try making the model class richer, or add more features (makes the model richer)
- If the model has high variance (overfitting) then the model class is sufficiently rich
  - Adding more training data can help (overfitting might be due to limited training data)

・ロト ・ 同ト ・ ヨト ・ ヨト



- If the model has high bias (underfitting) then the model class is weak
  - Adding more training data won't help
  - Instead, try making the model class richer, or add more features (makes the model richer)
- If the model has high variance (overfitting) then the model class is sufficiently rich
  - Adding more training data can help (overfitting might be due to limited training data)
  - Using a simpler model class or fewer features or regularization can help

# When Train and Test Distributions Differ



### When Train and Test Conditions are Different..

• Training and test inputs typically assumed to be drawn from same distribution, i.e.,  $p_{tr}(\mathbf{x}) = p_{te}(\mathbf{x})$ 

・ロト ・四ト ・日ト ・日ト

# When Train and Test Conditions are Different..

- Training and test inputs typically assumed to be drawn from same distribution, i.e.,  $p_{tr}(\mathbf{x}) = p_{te}(\mathbf{x})$
- However, this is rarely true in real-world applications of ML where  $p_{tr}(\mathbf{x}) \neq p_{te}(\mathbf{x})$



• So the training and test data are from different "domains". How can we "adapt"?

# When Train and Test Conditions are Different..

- Training and test inputs typically assumed to be drawn from same distribution, i.e.,  $p_{tr}(\mathbf{x}) = p_{te}(\mathbf{x})$
- However, this is rarely true in real-world applications of ML where  $p_{tr}(\mathbf{x}) \neq p_{te}(\mathbf{x})$



- So the training and test data are from different "domains". How can we "adapt"?
- Known as "domain adaptation" (also "covariate-shift", since features x are also called covariates)

・ロト ・日ト ・ヨト ・ヨト (道) りんの

## Handling Domain Differences: Some Approaches

• Learning a transformation that makes the training and test data distributions "close"





## Handling Domain Differences: Some Approaches

• Learning a transformation that makes the training and test data distributions "close"



• Use an "importance weighted" correction of the loss function defined on the training data

$$\sum_{n=1}^{N} \frac{p_{te}(\boldsymbol{x}_n)}{p_{tr}(\boldsymbol{x}_n)} \ell(y_n, f(\boldsymbol{x}_n))$$

Intro to Machine Learning (CS771A)

# Handling Domain Differences: Some Approaches

• Learning a transformation that makes the training and test data distributions "close"



• Use an "importance weighted" correction of the loss function defined on the training data

$$\sum_{n=1}^{N} \frac{p_{te}(\boldsymbol{x}_n)}{p_{tr}(\boldsymbol{x}_n)} \ell(\boldsymbol{y}_n, f(\boldsymbol{x}_n))$$

.. basically, assign a weight to each term of the loss function, depending on how "close" the training example  $(x_n, y_n)$  is from the test distribution

イロト イヨト イヨト

• Note that, since we care about test data, ideally, we would like to learn f by minimizing

$$\mathbb{E}_{p_{te}(\boldsymbol{x})}[\ell(y,f(\boldsymbol{x}))] = \int p_{te}(\boldsymbol{x})\ell(y,f(\boldsymbol{x}))d\boldsymbol{x}$$



• Note that, since we care about test data, ideally, we would like to learn f by minimizing

$$\mathbb{E}_{p_{te}(\boldsymbol{x})}[\ell(y,f(\boldsymbol{x}))] = \int p_{te}(\boldsymbol{x})\ell(y,f(\boldsymbol{x}))d\boldsymbol{x}$$

• However, we don't have labeled data from the test domain. :-(

• Note that, since we care about test data, ideally, we would like to learn f by minimizing

$$\mathbb{E}_{p_{te}(\boldsymbol{x})}[\ell(y,f(\boldsymbol{x}))] = \int p_{te}(\boldsymbol{x})\ell(y,f(\boldsymbol{x}))d\boldsymbol{x}$$

• However, we don't have labeled data from the test domain. :-( What can we do now?



• Note that, since we care about test data, ideally, we would like to learn f by minimizing

$$\mathbb{E}_{p_{te}(\boldsymbol{x})}[\ell(y,f(\boldsymbol{x}))] = \int p_{te}(\boldsymbol{x})\ell(y,f(\boldsymbol{x}))d\boldsymbol{x}$$

- However, we don't have labeled data from the test domain. :-( What can we do now?
- Let's re-express the above in terms of the loss on training data

$$\int p_{te}(x)\ell(y,f(x))dx$$



A B > A B > A B >
 A
 B >
 A
 B > A B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B

• Note that, since we care about test data, ideally, we would like to learn f by minimizing

$$\mathbb{E}_{p_{te}(\boldsymbol{x})}[\ell(y,f(\boldsymbol{x}))] = \int p_{te}(\boldsymbol{x})\ell(y,f(\boldsymbol{x}))d\boldsymbol{x}$$

- However, we don't have labeled data from the test domain. :-( What can we do now?
- Let's re-express the above in terms of the loss on training data

$$\int p_{te}(\mathbf{x})\ell(y,f(\mathbf{x}))d\mathbf{x} = \int p_{tr}(\mathbf{x})\frac{p_{te}(\mathbf{x})}{p_{tr}(\mathbf{x})}\ell(y,f(\mathbf{x}))d\mathbf{x}$$

• Note that, since we care about test data, ideally, we would like to learn f by minimizing

$$\mathbb{E}_{p_{te}(\mathbf{x})}[\ell(y, f(\mathbf{x}))] = \int p_{te}(\mathbf{x})\ell(y, f(\mathbf{x}))d\mathbf{x}$$

- However, we don't have labeled data from the test domain. :-( What can we do now?
- Let's re-express the above in terms of the loss on training data

$$\int p_{te}(\mathbf{x})\ell(y,f(\mathbf{x}))d\mathbf{x} = \int p_{tr}(\mathbf{x})\frac{p_{te}(\mathbf{x})}{p_{tr}(\mathbf{x})}\ell(y,f(\mathbf{x}))d\mathbf{x}$$
$$= \mathbb{E}_{p_{tr}(\mathbf{x})}\left[\frac{p_{te}(\mathbf{x})}{p_{tr}(\mathbf{x})}\ell(y,f(\mathbf{x}))\right]$$

• Note that, since we care about test data, ideally, we would like to learn f by minimizing

$$\mathbb{E}_{p_{te}(\mathbf{x})}[\ell(y, f(\mathbf{x}))] = \int p_{te}(\mathbf{x})\ell(y, f(\mathbf{x}))d\mathbf{x}$$

- However, we don't have labeled data from the test domain. :-( What can we do now?
- Let's re-express the above in terms of the loss on training data

$$\begin{split} \int p_{te}(x)\ell(y,f(x))dx &= \int p_{tr}(x)\frac{p_{te}(x)}{p_{tr}(x)}\ell(y,f(x))dx \\ &= \mathbb{E}_{p_{tr}(x)}\left[\frac{p_{te}(x)}{p_{tr}(x)}\ell(y,f(x))\right] \\ &\approx \sum_{n=1}^{N}\frac{p_{te}(x_{n})}{p_{tr}(x_{n})}\ell(y_{n},f(x_{n})) \quad \text{(uses only training data from } p_{tr}) \end{split}$$

• Note that, since we care about test data, ideally, we would like to learn f by minimizing

$$\mathbb{E}_{p_{te}(\boldsymbol{x})}[\ell(y,f(\boldsymbol{x}))] = \int p_{te}(\boldsymbol{x})\ell(y,f(\boldsymbol{x}))d\boldsymbol{x}$$

- However, we don't have labeled data from the test domain. :-( What can we do now?
- Let's re-express the above in terms of the loss on training data

$$\begin{split} \int p_{te}(x)\ell(y,f(x))dx &= \int p_{tr}(x)\frac{p_{te}(x)}{p_{tr}(x)}\ell(y,f(x))dx \\ &= \mathbb{E}_{p_{tr}(x)}\left[\frac{p_{te}(x)}{p_{tr}(x)}\ell(y,f(x))\right] \\ &\approx \sum_{n=1}^{N}\frac{p_{te}(x_n)}{p_{tr}(x_n)}\ell(y_n,f(x_n)) \quad \text{(uses only training data from } p_{tr}) \end{split}$$

• In order to do this, we need to estimate the two densities  $p_{tr}(x)$  and  $p_{te}(x)$  (can assume some form for these and estimate them using unlabeled data from both domains)

# Semi-supervised Learning



# Labeled vs Unlabeled Data

- Supervised Learning models require labeled data
- Learning a reliable model usually requires plenty of labeled data
- Labeled Data: Expensive and Scarce (someone has to do the labeling)
  - Often labeling is very difficult too (e.g., in speech analysis or NLP problems)
- Unlabeled Data: Abundant and Free/Cheap
  - E.g., can easily crawl the web and download webpages/images



• Usually such problems come in one of the following two flavors



Intro to Machine Learning (CS771A)

- Usually such problems come in one of the following two flavors
- Semi-supervised Learning

• Semi-Unsupervised Learning



- Usually such problems come in one of the following two flavors
- Semi-supervised Learning
  - Training set containslabeled data  $\mathcal{L} = \{x_i, y_i\}_{i=1}^L$  and unlabeled data  $\mathcal{U} = \{x_j\}_{j=L+1}^{L+U}$  (usually  $U \gg L$ )

• Semi-Unsupervised Learning



- Usually such problems come in one of the following two flavors
- Semi-supervised Learning
  - Training set containslabeled data  $\mathcal{L} = \{x_i, y_i\}_{i=1}^L$  and unlabeled data  $\mathcal{U} = \{x_j\}_{j=L+1}^{L+U}$  (usually  $U \gg L$ )
  - Note: In some cases,  $\mathcal U$  is also our test data (this is called transductive setting)
- Semi-Unsupervised Learning

- Usually such problems come in one of the following two flavors
- Semi-supervised Learning
  - Training set containslabeled data  $\mathcal{L} = \{x_i, y_i\}_{i=1}^L$  and unlabeled data  $\mathcal{U} = \{x_j\}_{j=L+1}^{L+U}$  (usually  $U \gg L$ )
  - Note: In some cases,  $\mathcal{U}$  is also our test data (this is called transductive setting)
  - Inductive setting is more standard: Want to learn a classifier for future test data
- Semi-Unsupervised Learning

Intro to Machine Learning (CS771A)

・ロト ・日下・ ・日下・ ・日下

- Usually such problems come in one of the following two flavors
- Semi-supervised Learning
  - Training set containslabeled data  $\mathcal{L} = \{x_i, y_i\}_{i=1}^L$  and unlabeled data  $\mathcal{U} = \{x_j\}_{j=L+1}^{L+U}$  (usually  $U \gg L$ )
  - Note: In some cases,  $\mathcal{U}$  is also our test data (this is called transductive setting)
  - Inductive setting is more standard: Want to learn a classifier for future test data
- Semi-Unsupervised Learning
  - We are given unlabeled data  $\mathcal{D} = \{ \pmb{x}_i \}_{i=1}^N$

- Usually such problems come in one of the following two flavors
- Semi-supervised Learning
  - Training set containslabeled data  $\mathcal{L} = \{x_i, y_i\}_{i=1}^L$  and unlabeled data  $\mathcal{U} = \{x_j\}_{j=L+1}^{L+U}$  (usually  $U \gg L$ )
  - Note: In some cases,  $\mathcal{U}$  is also our test data (this is called transductive setting)
  - Inductive setting is more standard: Want to learn a classifier for future test data
- Semi-Unsupervised Learning
  - We are given unlabeled data  $\mathcal{D} = \{ \pmb{x}_i \}_{i=1}^N$
  - Additionally, we're given supervision in form of some constraints on data (e.g., points  $x_n$  and  $x_m$  belong to the same cluster) or "labels" of some points.

↓ □ ▶ ↓ @ ▶ ↓ E ▶ ↓ E ▶ ↓

- Usually such problems come in one of the following two flavors
- Semi-supervised Learning
  - Training set containslabeled data  $\mathcal{L} = \{x_i, y_i\}_{i=1}^L$  and unlabeled data  $\mathcal{U} = \{x_j\}_{j=L+1}^{L+U}$  (usually  $U \gg L$ )
  - Note: In some cases,  $\mathcal{U}$  is also our test data (this is called transductive setting)
  - Inductive setting is more standard: Want to learn a classifier for future test data
- Semi-Unsupervised Learning
  - We are given unlabeled data  $\mathcal{D} = \{ \pmb{x}_i \}_{i=1}^N$
  - Additionally, we're given supervision in form of some constraints on data (e.g., points  $x_n$  and  $x_m$  belong to the same cluster) or "labels" of some points.
  - Want to learn an unsupervised learning model combining both data sources

- Usually such problems come in one of the following two flavors
- Semi-supervised Learning
  - Training set containslabeled data  $\mathcal{L} = \{x_i, y_i\}_{i=1}^L$  and unlabeled data  $\mathcal{U} = \{x_j\}_{j=L+1}^{L+U}$  (usually  $U \gg L$ )
  - Note: In some cases,  $\mathcal{U}$  is also our test data (this is called transductive setting)
  - Inductive setting is more standard: Want to learn a classifier for future test data
- Semi-Unsupervised Learning
  - We are given unlabeled data  $\mathcal{D} = \{ \pmb{x}_i \}_{i=1}^N$
  - Additionally, we're given supervision in form of some constraints on data (e.g., points  $x_n$  and  $x_m$  belong to the same cluster) or "labels" of some points.
  - Want to learn an unsupervised learning model combining both data sources
- Here, we will focus on Semi-supervised Learning (SSL)

• Red: + 1, Dark Blue: -1





• Red: + 1, Dark Blue: -1





• Red: + 1, Dark Blue: -1



• Let's include some additional unlabeled points (Light Blue points)



• Red: + 1, Dark Blue: -1



• Let's include some additional unlabeled points (Light Blue points)



• Red: + 1, Dark Blue: -1



• Let's include some additional unlabeled points (Light Blue points)

• Assumption: Examples from the same class are clustered together
• Red: + 1, Dark Blue: -1



• Let's include some additional unlabeled points (Light Blue points)

- Assumption: Examples from the same class are clustered together
- Assumption: Decision boundary lies in the region where data has low density

Intro to Machine Learning (CS771A)

・ロト ・戸ト ・ヨト ・ヨト







Ĭ.



乏



Intro to Machine Learning (CS771A)

Ξi



Knowledge of the input distribution helps even if the data is unlabeled

Intro to Machine Learning (CS771A)

Bias/Variance Trade-off, Some Practical Issues, Semi-supervised and Active Learning

イロト イヨト イヨト

- Nearby points (may) have the same label
  - Smoothness assumption



- Nearby points (may) have the same label
  - Smoothness assumption
- Points in the same cluster (may) have same label
  - Cluster assumption

- Nearby points (may) have the same label
  - Smoothness assumption
- Points in the same cluster (may) have same label
  - Cluster assumption
- Decision boundary lies in areas where data has low density ۲
  - Low-density separation

- Nearby points (may) have the same label
  - Smoothness assumption
- Points in the same cluster (may) have same label
  - Cluster assumption
- Decision boundary lies in areas where data has low density
  - Low-density separation
- Note: Generative models already take some of these assumptions into account (recall HW3 problem on semi-supervised learning using EM)

- Based on smoothness assumption
- Requires constructing a graph between all the examples (labeled/unlabeled)



- Based on smoothness assumption
- Requires constructing a graph between all the examples (labeled/unlabeled)



• The graph can be constructed using several ways

- Based on smoothness assumption
- Requires constructing a graph between all the examples (labeled/unlabeled)



- The graph can be constructed using several ways
  - Connecting every example with its top k neighbors (labeled/unlabeled)

- Based on smoothness assumption
- Requires constructing a graph between all the examples (labeled/unlabeled)



- The graph can be constructed using several ways
  - Connecting every example with its top k neighbors (labeled/unlabeled)
  - Constructing an all-connected weighted graph

- Based on smoothness assumption
- Requires constructing a graph between all the examples (labeled/unlabeled)



- The graph can be constructed using several ways
  - Connecting every example with its top k neighbors (labeled/unlabeled)
  - Constructing an all-connected weighted graph
  - Weighted case: weight of edge connecting examples x<sub>i</sub> and x<sub>j</sub>

$$a_{ij} = \exp(-||\boldsymbol{x}_i - \boldsymbol{x}_j||^2 / \sigma^2)$$

.. where **A** is the  $(L + U) \times (L + U)$  matrix of pairwise similarities

• Suppose we want to learn a function f using labeled+unlabeld data



- Suppose we want to learn a function f using labeled+unlabeld data  $\mathcal{L} \cup \mathcal{U}$
- Suppose  $f_i$  denotes the prediction on example  $x_i$

- Suppose we want to learn a function f using labeled+unlabeld data  $\mathcal{L} \cup \mathcal{U}$
- Suppose  $f_i$  denotes the prediction on example  $x_i$ 
  - Similar examples  $x_i$  and  $x_i$  (thus high  $a_{ii}$ ) should have similar  $f_i$  and  $f_i$

- Suppose we want to learn a function f using labeled+unlabeld data  $\mathcal{L} \cup \mathcal{U}$
- Suppose  $f_i$  denotes the prediction on example  $x_i$ 
  - Similar examples  $x_i$  and  $x_j$  (thus high  $a_{ij}$ ) should have similar  $f_i$  and  $f_j$
  - Graph-based regularization optimizes the following objective:



・ロト ・ 日 ・ ・ 日 ・ ・ 日

• Suppose data  $\mathcal{D}$  is labeled  $\mathcal{L} = \{ \mathbf{x}_i, y_i \}_{i=1}^{L}$  and unlabeled  $\mathcal{U} = \{ \mathbf{x}_j \}_{i=L+1}^{L+U}$ 



- Suppose data  $\mathcal{D}$  is labeled  $\mathcal{L} = \{ \mathbf{x}_i, y_i \}_{i=1}^{L}$  and unlabeled  $\mathcal{U} = \{ \mathbf{x}_j \}_{i=L+1}^{L+U}$
- Assume the following model

$$p(\mathcal{D}|\theta) = \underbrace{\prod_{i=1}^{L} p(\mathbf{x}_i, y_i|\theta)}_{\text{labeled}} \underbrace{\prod_{j=L+1}^{L+U} p(\mathbf{x}_j|\theta)}_{\text{unlabeled}} = \prod_{i=1}^{L} p(\mathbf{x}_i, y_i|\theta) \prod_{j=L+1}^{L+U} \sum_{\mathbf{y}_j} p(\mathbf{x}_j, y_j|\theta)$$

< ロ > < 回 > < 回 > < 回 > < 回 >

- Suppose data  $\mathcal{D}$  is labeled  $\mathcal{L} = \{ \mathbf{x}_i, y_i \}_{i=1}^{L}$  and unlabeled  $\mathcal{U} = \{ \mathbf{x}_j \}_{i=L+1}^{L+U}$
- Assume the following model

$$p(\mathcal{D}|\theta) = \prod_{i=1}^{L} p(\mathbf{x}_i, y_i|\theta) \prod_{\substack{j=L+1\\ \text{unlabeled}}}^{L+U} p(\mathbf{x}_j|\theta) = \prod_{i=1}^{L} p(\mathbf{x}_i, y_i|\theta) \prod_{j=L+1}^{L+U} \sum_{\mathbf{y}_j} p(\mathbf{x}_j, \mathbf{y}_j|\theta)$$

• The unknowns are  $\{y_j\}_{j=L+1}^{L+U}$  (latent variables) and  $\theta$  (parameters)

- Suppose data  $\mathcal{D}$  is labeled  $\mathcal{L} = \{ \mathbf{x}_i, y_i \}_{i=1}^L$  and unlabeled  $\mathcal{U} = \{ \mathbf{x}_j \}_{j=L+1}^{L+U}$
- Assume the following model

$$p(\mathcal{D}|\theta) = \prod_{i=1}^{L} p(\mathbf{x}_i, y_i|\theta) \prod_{\substack{j=L+1\\ \text{unlabeled}}}^{L+U} p(\mathbf{x}_j|\theta) = \prod_{i=1}^{L} p(\mathbf{x}_i, y_i|\theta) \prod_{j=L+1}^{L+U} \sum_{\mathbf{y}_j} p(\mathbf{x}_j, \mathbf{y}_j|\theta)$$

- The unknowns are  $\{y_j\}_{j=L+1}^{L+U}$  (latent variables) and  $\theta$  (parameters)
- We can use EM to estimate the unknowns

・ロト ・ (日) ・ (1) \cdot (1) \cdot

- Suppose data  $\mathcal{D}$  is labeled  $\mathcal{L} = \{ \mathbf{x}_i, y_i \}_{i=1}^L$  and unlabeled  $\mathcal{U} = \{ \mathbf{x}_j \}_{j=L+1}^{L+U}$
- Assume the following model

$$p(\mathcal{D}|\theta) = \prod_{i=1}^{L} p(\mathbf{x}_i, y_i|\theta) \prod_{\substack{j=L+1\\ \text{unlabeled}}}^{L+U} p(\mathbf{x}_j|\theta) = \prod_{i=1}^{L} p(\mathbf{x}_i, y_i|\theta) \prod_{j=L+1}^{L+U} \sum_{\mathbf{y}_j} p(\mathbf{x}_j, y_j|\theta)$$

- The unknowns are  $\{y_j\}_{j=L+1}^{L+U}$  (latent variables) and  $\theta$  (parameters)
- We can use EM to estimate the unknowns
  - Given  $\theta = \hat{\theta}$ , E step will compute expected labels for unlabeled examples

$$\mathbb{E}[y_j] = +1 imes P(y_j = +1 | \hat{ heta}, \mathbf{x}_j) + (-1) imes P(y_j = -1 | \hat{ heta}, \mathbf{x}_j)$$

A B > A B > A B >

- Suppose data  $\mathcal{D}$  is labeled  $\mathcal{L} = \{ \mathbf{x}_i, y_i \}_{i=1}^L$  and unlabeled  $\mathcal{U} = \{ \mathbf{x}_j \}_{j=L+1}^{L+U}$
- Assume the following model

$$p(\mathcal{D}|\theta) = \prod_{i=1}^{L} p(\mathbf{x}_i, y_i|\theta) \prod_{\substack{j=L+1\\ \text{unlabeled}}}^{L+U} p(\mathbf{x}_j|\theta) = \prod_{i=1}^{L} p(\mathbf{x}_i, y_i|\theta) \prod_{j=L+1}^{L+U} \sum_{\mathbf{y}_j} p(\mathbf{x}_j, y_j|\theta)$$

- The unknowns are  $\{y_j\}_{j=L+1}^{L+U}$  (latent variables) and  $\theta$  (parameters)
- We can use EM to estimate the unknowns
  - $\bullet~{\rm Given}~\theta=\hat{\theta},~{\rm E}$  step will compute expected labels for unlabeled examples

$$\mathbb{E}[y_j] = +1 imes P(y_j = +1|\hat{ heta}, oldsymbol{x}_j) + (-1) imes P(y_j = -1|\hat{ heta}, oldsymbol{x}_j)$$

• M step can then perform standard MLE for re-estimating the parameters  $\theta$ 

A D > A D > A D > A D > A D >

- Suppose data  $\mathcal{D}$  is labeled  $\mathcal{L} = \{ \mathbf{x}_i, y_i \}_{i=1}^L$  and unlabeled  $\mathcal{U} = \{ \mathbf{x}_j \}_{j=L+1}^{L+U}$
- Assume the following model

$$p(\mathcal{D}|\theta) = \prod_{i=1}^{L} p(\mathbf{x}_i, y_i|\theta) \prod_{\substack{j=L+1\\ \text{unlabeled}}}^{L+U} p(\mathbf{x}_j|\theta) = \prod_{i=1}^{L} p(\mathbf{x}_i, y_i|\theta) \prod_{j=L+1}^{L+U} \sum_{\mathbf{y}_j} p(\mathbf{x}_j, y_j|\theta)$$

- The unknowns are  $\{y_j\}_{j=L+1}^{L+U}$  (latent variables) and  $\theta$  (parameters)
- We can use EM to estimate the unknowns
  - Given  $\theta = \hat{\theta}$ , E step will compute expected labels for unlabeled examples

$$\mathbb{E}[y_j] = +1 imes P(y_j = +1 | \hat{ heta}, \mathbf{x}_j) + (-1) imes P(y_j = -1 | \hat{ heta}, \mathbf{x}_j)$$

- $\bullet\,$  M step can then perform standard MLE for re-estimating the parameters  $\theta\,$
- A fairly general framework for semi-supervised learning. Can be used for different types of data (by choosing the appropriate  $p(\mathbf{x}|y)$  distribution)

## Things can go wrong..

If assumptions are not appropriate for the data (e.g., incorrectly specified class conditional distributions)



Thus need to be careful/flexible about the choice of class conditional distributions



イロト イロト イモト イモト

#### **Passive vs Active Learning**

Standard supervised learning assumes a "passive" learner



#### **Passive vs Active Learning**

Standard supervised learning assumes a "passive" learner



Called "passive" because the learner doesn't have any control over the labeled data it gets to learn

・ロト ・日 ・ モ ・ ・ モ ・ (話) りんの

An "active" learner can specifically request labels of "hard" examples that are most useful for learning



raw unlabeled data  $x_1, x_2, x_3, \ldots$ 



Assumes some small amount of initial labeled training data

active learner induces a classifier



イロト イポト イヨト イヨト

expert / oracle analyzes experiments to determine labels

An "active" learner can specifically request labels of "hard" examples that are most useful for learning



An "active" learner can specifically request labels of "hard" examples that are most useful for learning



An "active" learner can specifically request labels of "hard" examples that are most useful for learning



An "active" learner can specifically request labels of "hard" examples that are most useful for learning



An "active" learner can specifically request labels of "hard" examples that are most useful for learning


An "active" learner can specifically request labels of "hard" examples that are most useful for learning



How to quantify the "hardness" or "usefulness" of an unlabeled example?

Intro to Machine Learning (CS771A)

- Various ways to define the usefulness/hardness of an unlabeled example  $\boldsymbol{x}$
- A probabilistic approach can use distribution of y given x to quantify usefulness of x



イロト イロト イヨト イヨト

Intro to Machine Learning (CS771A)

- Various ways to define the usefulness/hardness of an unlabeled example  ${m x}$
- A probabilistic approach can use distribution of y given x to quantify usefulness of x, e.g.,
  - In classification, can look at the entropy of the distribution p(y|x, w) or posterior predictive p(y|x)



.. and choose example(s) for which the distribution of y given x has the largest entropy

- Various ways to define the usefulness/hardness of an unlabeled example  ${m x}$
- A probabilistic approach can use distribution of y given x to quantify usefulness of x, e.g.,
  - In classification, can look at the entropy of the distribution p(y|x, w) or posterior predictive p(y|x)



.. and choose example(s) for which the distribution of y given x has the largest entropy

• Look at the variance of the posterior predictive p(y|x). E.g. for regression,



 $\ldots$  and choose choose example(s) for which variance of posterior predictive is the largest

イロト イポト イヨト イヨト

- Various ways to define the usefulness/hardness of an unlabeled example  ${m x}$
- A probabilistic approach can use distribution of y given x to quantify usefulness of x, e.g.,
  - In classification, can look at the entropy of the distribution p(y|x, w) or posterior predictive p(y|x)



.. and choose example(s) for which the distribution of y given x has the largest entropy

• Look at the variance of the posterior predictive p(y|x). E.g. for regression,



 $\ldots$  and choose choose example(s) for which variance of posterior predictive is the largest

• ... many other ways to utilize the information in p(y|x)