# Lecture 13 - Debugging ML Models and Error Analysis | Stanford CS229: Machine Learning (Autumn 2018)-Make Written Notes

## Debugging learning algorithms

Motivating example:

- Anti-spam. You carefully choose a small set of 100 words to use as features. (Instead of using all 50000+ words in English.)

- Logistic regression with regularization (Bayesian Logistic regression), implemented with gradient ascent, gets 20% test error, which is unacceptably high.

$$\max_{\theta} \sum_{i=1}^{m} \log p(y^{(i)}|x^{(i)}, \theta) - \lambda||\theta||^2$$

- What to do next?

## Fixing the learning algorithm

- Logistic regression (with regularization):

$$\max_{\theta} \sum_{i=1}^{m} \log p(y^{(i)}|x^{(i)}, \theta) - \lambda||\theta||^2$$

- Common approach: Try improving the algorithm in different ways.
  - Try getting more training examples.
  - Try a smaller set of features.
  - Try a larger set of features.
  - Try changing the features: Email header vs. email body features.
  - Run gradient descent for more iterations.
  - Try Newton's method.
  - Use a different value for $\lambda$.
  - Try using an SVM.

## Diagnostic for bias vs. variance

Better approach:
- Run diagnostics to figure out what the problem is.
- Fix whatever the problem is.

Logistic regression's test error is 20% (unacceptably high).

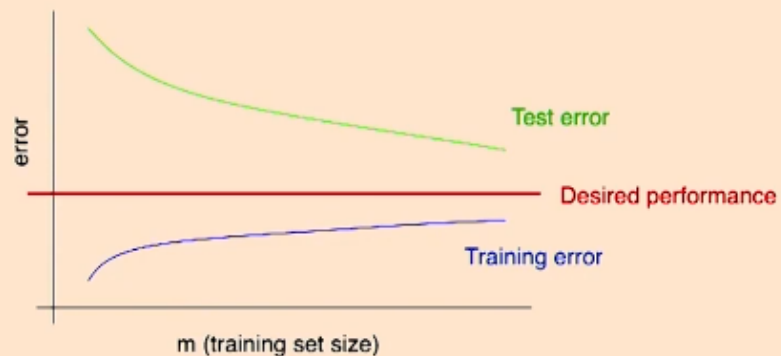Suppose you suspect the problem is either:
- Overfitting (high variance).
- Too few features to classify spam (high bias).

Diagnostic:
- Variance: Training error will be much lower than test error.
- Bias: Training error will also be high.
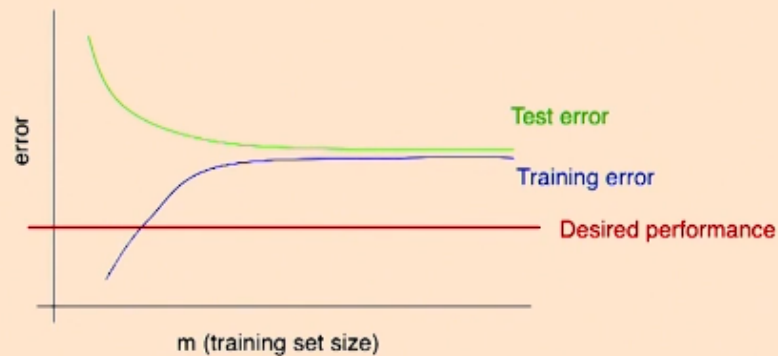
## More on bias vs. variance

Typical learning curve for high variance:



- Test error still decreasing as m increases. Suggests larger training set will help.
- Large gap between training and test error.

## More on bias vs. variance

Typical learning curve for high bias:



- error (y-axis)
- Test error
- Training error
- Desired performance
- m (training set size)

- Even training error is unacceptably high.
- Small gap between training and test error.
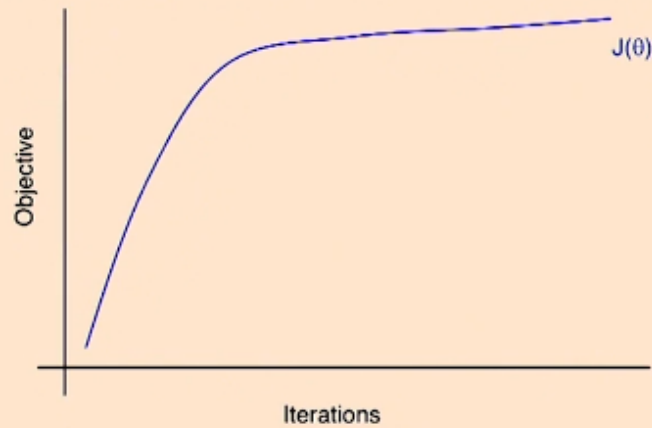
## Optimization algorithm diagnostics

- Bias vs. variance is one common diagnostic.

- For other problems, it's usually up to your own ingenuity to construct your own diagnostics to figure out what's wrong.

- Another example:
  - Logistic regression gets 2% error on spam, and 2% error on non-spam. (Unacceptably high error on non-spam.)
  - SVM using a linear kernel gets 10% error on spam, and 0.01% error on non-spam. (Acceptable performance.)
  - But you want to use logistic regression, because of computational efficiency, etc.

- What to do next?

## More diagnostics

- Other common questions:
  - Is the algorithm (gradient ascent for logistic regression) converging?

---

## More diagnostics

- Other common questions:
  - Is the algorithm (gradient ascent for logistic regression) converging?
  - Are you optimizing the right function?
  - I.e., what you care about:

$$a(\theta) = \sum_i w^{(i)} 1\{h_\theta(x^{(i)}) = y^{(i)}\}$$

  (weights $w^{(i)}$ higher for non-spam than for spam).
  - Logistic regression?  Correct value for $\lambda$?

$$\max_\theta J(\theta) = \sum_{i=1}^m \log p(y^{(i)}|x^{(i)}, \theta) - \lambda||\theta||^2$$

  - SVM?  Correct value for C?

$$\min_{w,b} \quad ||w||^2 + C\sum_{i=1}^m \xi_i$$
$$\text{s.t.} \quad y^{(i)}(w^T x^{(i)} - b) \geq 1 - \xi_i$$

## Diagnostic

An SVM outperforms logistic regression, but you really want to deploy logistic regression for your application.

Let $\theta_{SVM}$ be the parameters learned by an SVM.

Let $\theta_{BLR}$ be the parameters learned by logistic regression. (BLR = Bayesian logistic regression.)

You care about weighted accuracy:

$$a(\theta) = \max_{\theta} \sum_i w^{(i)} 1\{h_\theta(x^{(i)}) = y^{(i)}\}$$

$\theta_{SVM}$ outperforms $\theta_{BLR}$.  So:

$$a(\theta_{SVM}) > a(\theta_{BLR})$$

---

BLR tries to maximize:

$$J(\theta) = \sum_{i=1}^{m} \log p(y^{(i)}|x^{(i)}, \theta) - \lambda||\theta||^2$$

Diagnostic:

$$J(\theta_{SVM}) > J(\theta_{BLR})?$$

Stanfor

Andrew Y. Ng

- Very Important-New Insight-Please See it Future Mohtashim

## Two cases

Case 1:
$$a(\theta_{SVM}) > a(\theta_{BLR})$$
$$J(\theta_{SVM}) > J(\theta_{BLR})$$

But BLR was trying to maximize J(θ). This means that $\theta_{BLR}$ fails to maximize J, and the problem is with the convergence of the algorithm. Problem is with optimization algorithm.

Case 2:
$$a(\theta_{SVM}) > a(\theta_{BLR})$$
$$J(\theta_{SVM}) \leq J(\theta_{BLR})$$

This means that BLR succeeded at maximizing J(θ). But the SVM, which does worse on J(θ), actually does better on weighted accuracy a(θ).

This means that J(θ) is the wrong function to be maximizing, if you care about a(θ). Problem is with objective function of the maximization problem.

Stanfor

---
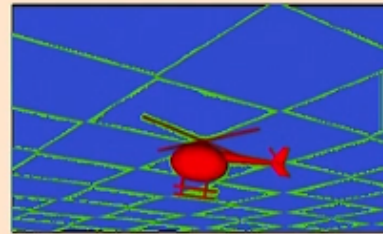
## The Stanford Autonomous Helicopter



**Payload: 14 pounds**
**Weight: 32 pounds**

**Machine learning algorithm**

1. Build a simulator of helicopter.

2. Choose a cost function. Say $J(\theta) = \|x - x_{desired}\|^2$      ($x$ = helicopter position)

3. Run reinforcement learning (RL) algorithm to fly helicopter in simulation, so as to try to minimize cost function:
$$\theta_{RL} = \arg\min_\theta J(\theta)$$

Suppose you do this, and the resulting controller parameters $\theta_{RL}$ gives much worse performance than your human pilot. What to do next?

Improve simulator?
Modify cost function J?
Modify RL algorithm?

Stanfor

Andrew Y. Ng

---

**Debugging an RL algorithm**

The controller given by $\theta_{RL}$ performs poorly.
Suppose that:
1. The helicopter simulator is accurate.
2. The RL algorithm correctly controls the helicopter (in simulation) so as to minimize $J(\theta)$.
3. Minimizing $J(\theta)$ corresponds to correct autonomous flight.

Then: The learned parameters $\theta_{RL}$ should fly well on the actual helicopter.

Diagnostics:≥
1. If $\theta_{RL}$ flies well in simulation, but not in real life, then the problem is in the simulator. Otherwise:
2. Let $\theta_{human}$ be the human control policy. If $J(\theta_{human}) < J(\theta_{RL})$, then the problem is in the reinforcement learning algorithm. (Failing to minimize the cost function J.)
3. If $J(\theta_{human}) \geq J(\theta_{RL})$, then the problem is in the cost function. (Maximizing it doesn't correspond to good autonomous flight.)
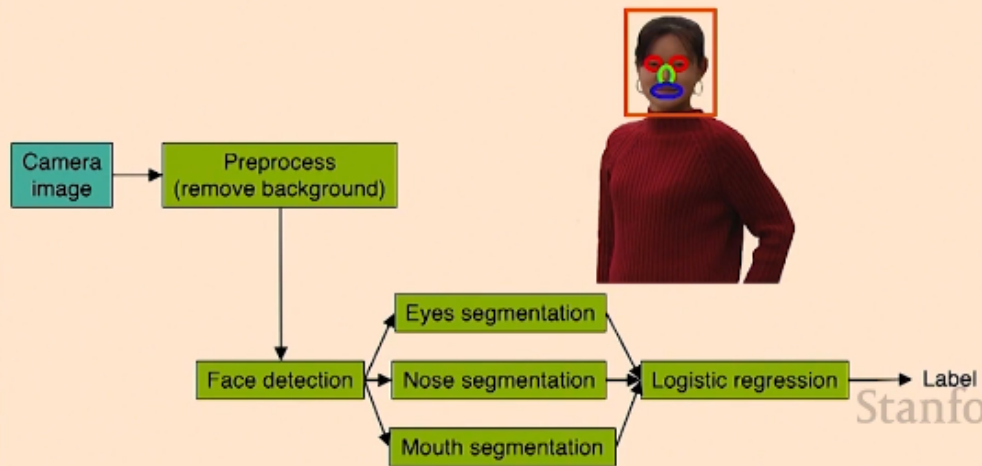
Stanfor

Andrew Y. Ng

- Remember two things:
  - Bias/Variance Trade-off
  - Analyse whether the problem is in algorithm or objective that you are trying to achieve
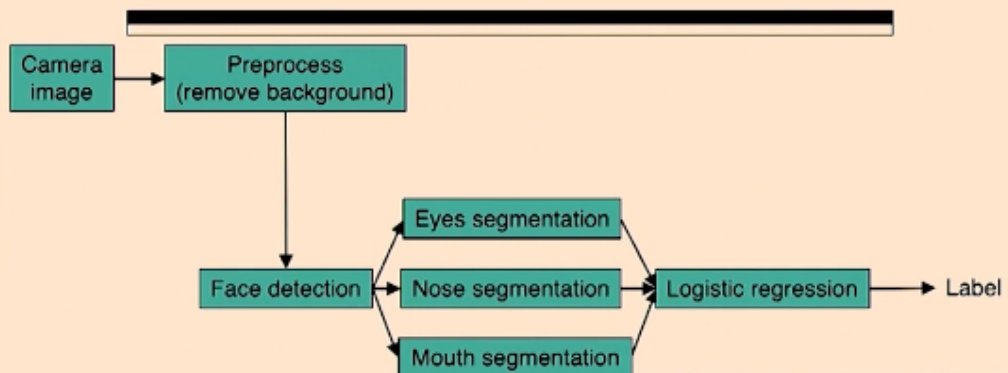
# Ablative analysis

Simple logistic regression without any clever features get 94% performance.

Just what accounts for your improvement from 94 to 99.9%?

Ablative analysis: Remove components from your system one at a time, to see how it breaks.

| Component | Accuracy | |
|---|---|---|
| Overall system | 99.9% | |
| Spelling correction | 99.0 | |
| Sender host features | 98.9% | |
| Email header features | 98.9% | |
| Email text parser features | 95% | |
| Javascript parser | 94.5% | |
| Features from images | 94.0% | [baseline] |

Conclusion: The email text parser features account for most of the improvement.

Stanfor