

Module 4 - Linear Classifiers : Overfitting & Regularization

Training the classifier: ¶

- The dataset -> (labelled -> [(sentence1, +), (sentence2, -), ...]) is provided.
- The dataset is divided into **Training set** and **Validation set**.
- The training set is fed into the learn classifier and the algorithm is built, model is trained.
- The learnt model is run on the validation set and evaluate the classifier.

Classification Error

- 'Classification error' - It is the measure of the classifiers performance.
- If the 'true label' and 'predicted label' are the same then it is right else it is wrong.

Classification Error

- Error measures the fraction of mistakes.
- **error = #mistakes / total number of datapoints.**
- Best possible value = 0.0.

Accuracy

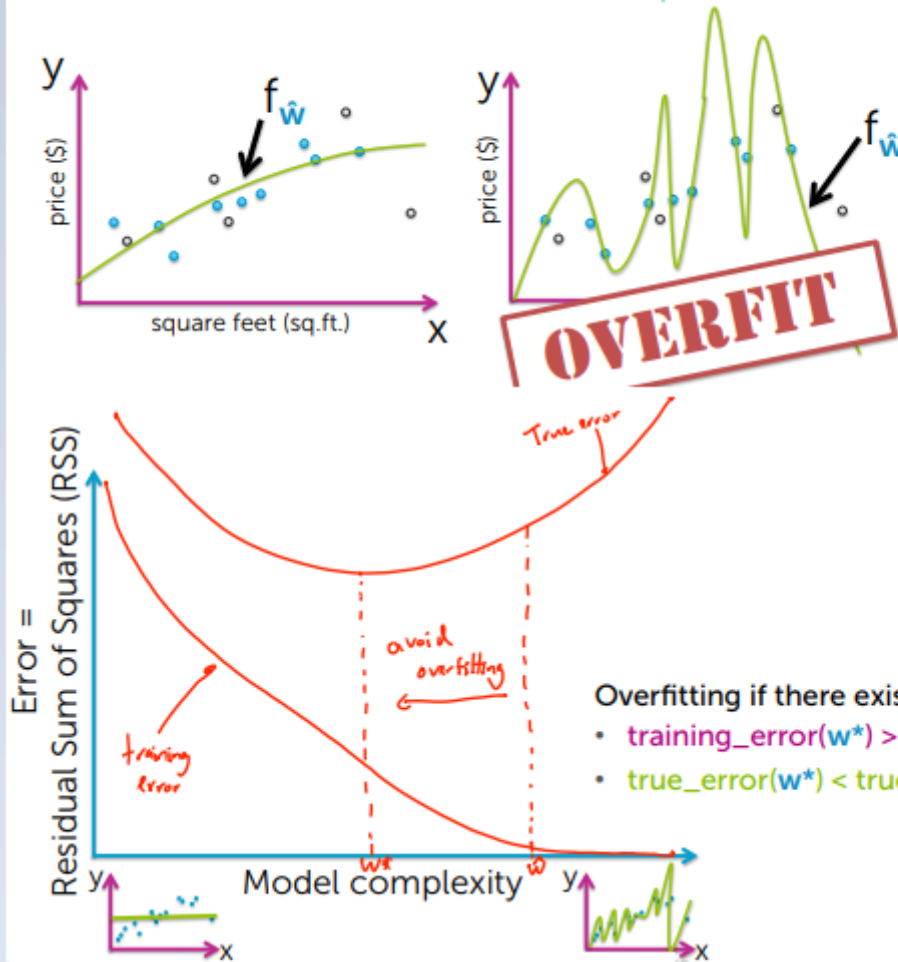
- Fraction of correct predictions.
- **accuracy = #correct / total number of datapoints.**
- Best possible value - 1.0.

Overfitting in regression

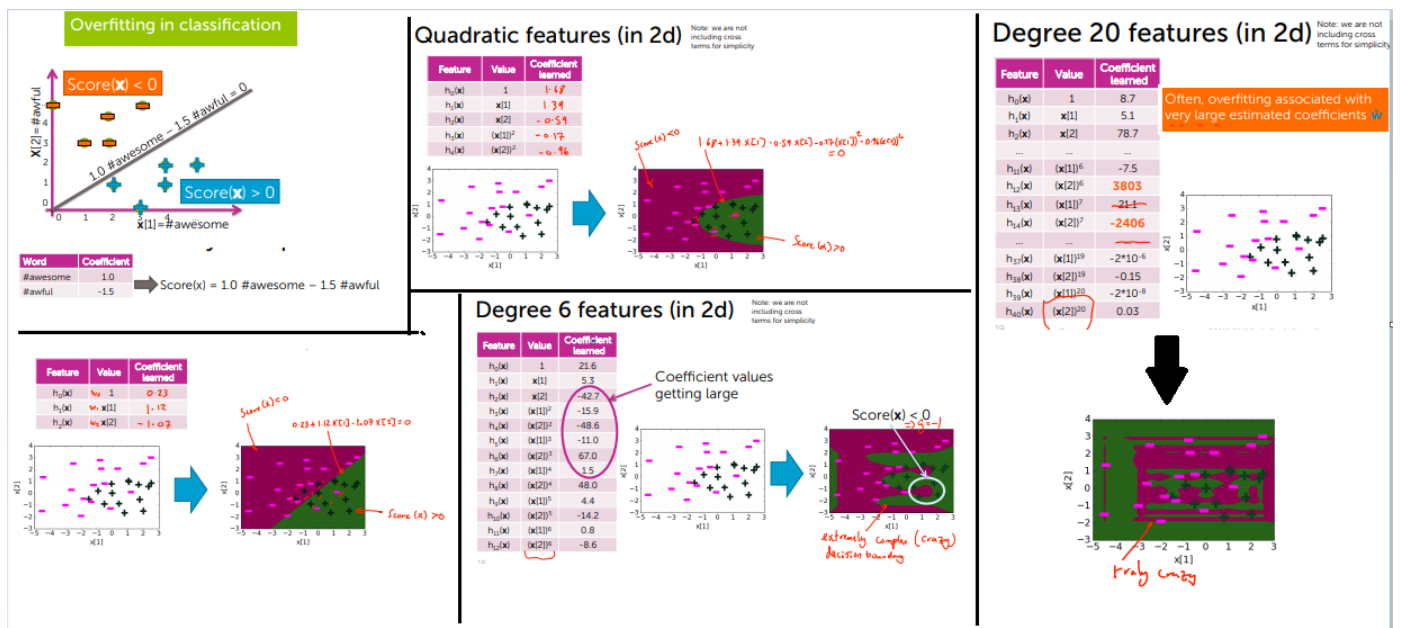
- Here as the model complexity increases, the model fits the training data very well, but might not be generalized leading to **overfitting**.

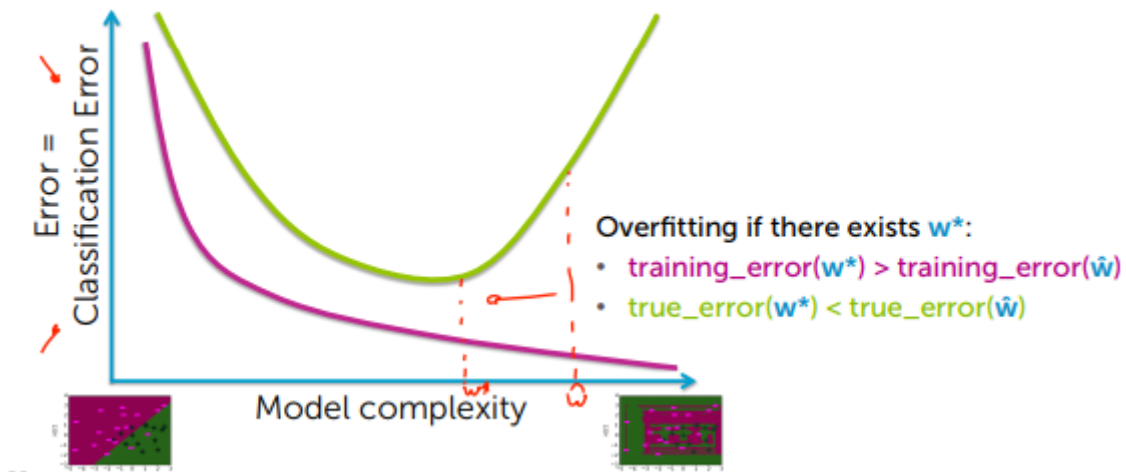
Overfitting in regression

$$y_i = w_0 + w_1 x_i + w_2 x_i^2 + \dots + w_p x_i^p + \epsilon_i$$



Overfitting in Classification





Overfitting in Classifiers -> Overconfident Predictions

- In the 'Logistic Regression Model' the score $\rightarrow w(\text{transpose}) * \text{features} \rightarrow \text{extends } [-\infty, \infty]$;
- This is reduce to their probability funtion of the output given input using the sigmoid/link function in the range $[0, 1] \rightarrow$ where the output ranges $y \in [-1, +1]$;

The subtle consequences of overfitting in logistic regression

- Overfitting leads to 'Large coefficient values'.

- $w^T * h(x_i) \rightarrow$ either very positive / very negative \rightarrow which makes the sigmoid($w^T * h(x_i)$) \rightarrow go to 0 or 1.

- The model becomes extremely overconfident of predictions.
- Below 3 cases are listed for the same

* Input : #awesome = 2, #awful = 1
 * X-axis = #awesome - #awful = 2 - 1 = 1

- Coefficients are intercept, weight -awesome, weight -awful

* Y-axis = Score = $1 / (1 + e^{-(w^T * h(x_i))})$;

1. Coefficients $\rightarrow 0, +1, -1 \rightarrow 0.73$

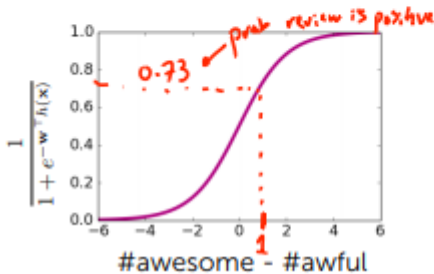
2. Coefficients $\rightarrow 0, +2, -2 \rightarrow 0.88$

3. Coefficients $\rightarrow 0, +6, -6 \rightarrow 0.997$

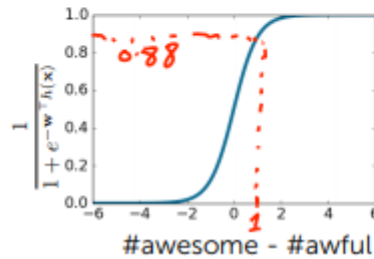
Effect of coefficients on logistic regression model

Input x : #awesome=2, #awful=1

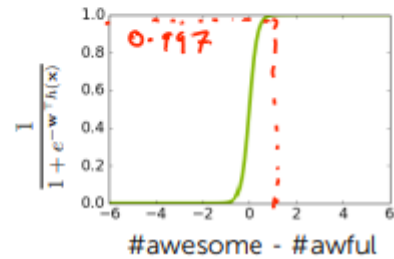
w_0	0
$w_{\text{#awesome}}$	+1
$w_{\text{#awful}}$	-1



w_0	0
$w_{\text{#awesome}}$	+2
$w_{\text{#awful}}$	-2

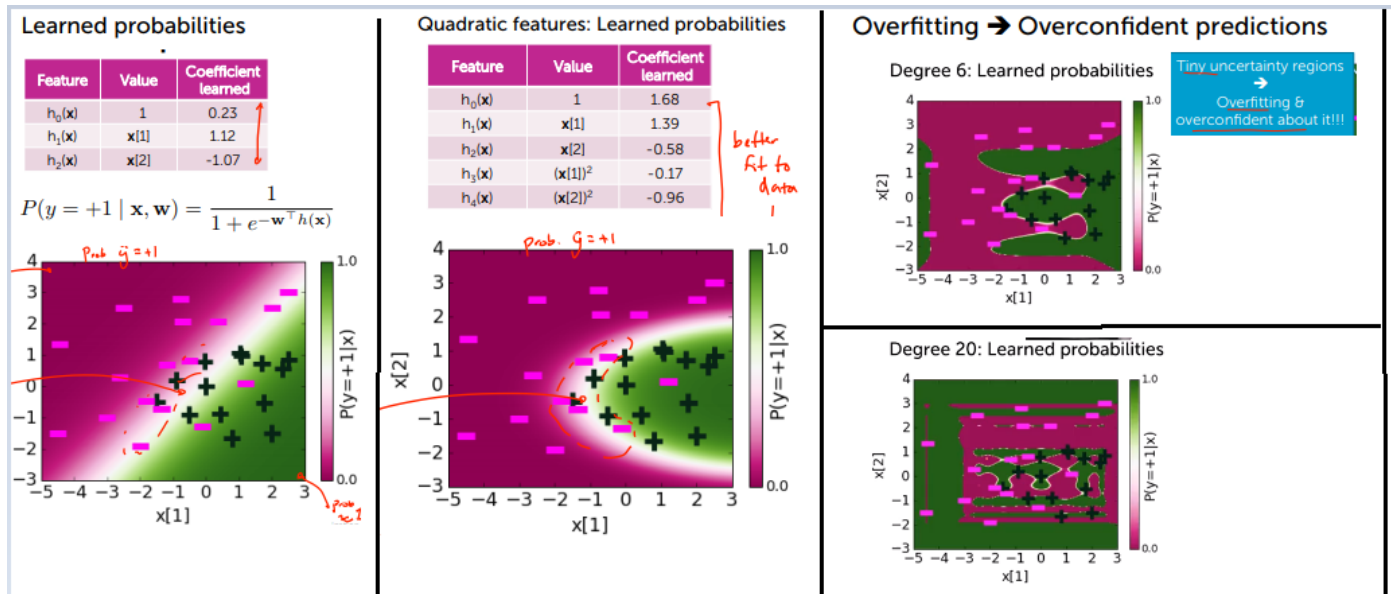


w_0	0
$w_{\text{#awesome}}$	+6
$w_{\text{#awful}}$	-6



Visualizing the dataset to identify Overfitting - Overconfident predictions

- As model complexity increases, the decision boundary between the classification data becomes very narrow but not very wide or extremely tiny. this indicates overfitting.
- The separation between the data must be narrow but not very wide or extremely tiny. **In case of tiny it is an indication of overfitting.**



Another perspective on overfitting logistic regression (ADVANCED)

Linearly-separable data: A dataset that can be classified into categories. A line can be drawn to segregate the data.

- Data is linearly separable if: There are coefficients \hat{w} such that:
 - For all positive training data: $\text{Score}(x) = \hat{w}^T h(x) > 0$
 - For all negative training data: $\text{Score}(x) = \hat{w}^T h(x) < 0$
- training_error(\hat{w}) = 0.** For linearly-separable data the training_error = 0. This could be a situation of overfitting - especially w.r.t complex models.
- Note 1 :** If there are D features, linear separability happens in a D -dimensional space.

- **Note 2 : If you have enough features, data are (almost) always linearly separable. ***
- Polynomial to the degree 50, 100, 180, etc - data is gonna become linearly separated - lead to Problematic case.

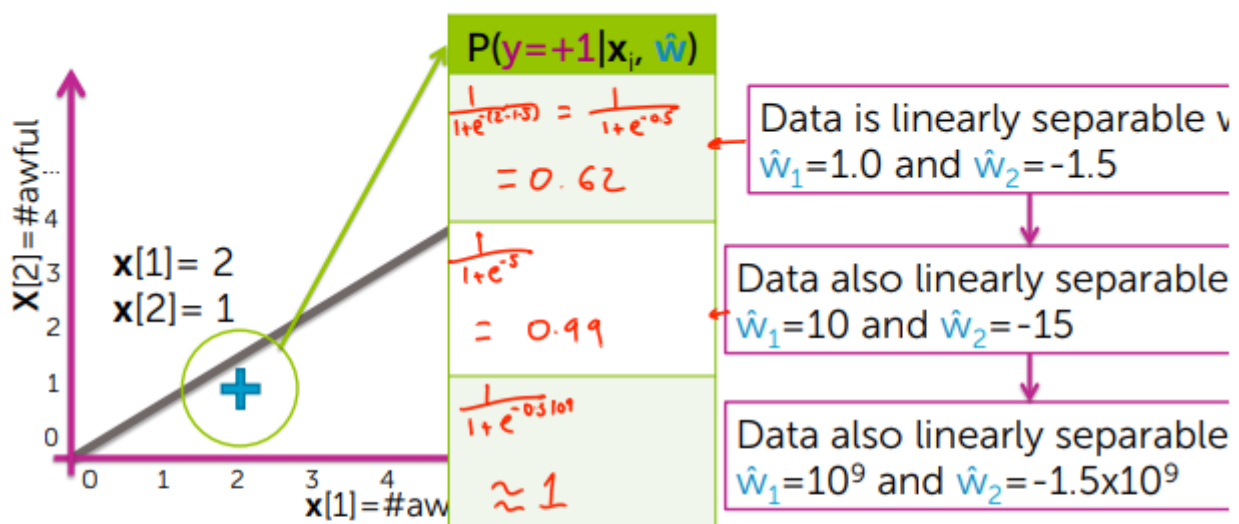
Effects of linear separability on coefficients

- Consider a plane that separates the data (positive and negative);
 - Plane $\rightarrow 1.0 \text{ #awesome} - 1.5 \text{ #awful} = 0$
 - Multiplying $\times 10 \rightarrow 10 \text{ #awesome} - 15 \text{ #awful} = 0$
 - Multiplying $\times 10^9 \rightarrow 10 \times 10^9 \text{ #awesome} - 15 \times 10^9 \text{ #awful} = 0$
 - In the above case \rightarrow although the values of the coefficients is increasing, the plane separating the positive and negative boundaries is still the same. Hence the prediction are not right if its a result of increase in magnitude of the coefficients.

Issue : MLE (Maximum likelihood estimation) - prefers most certain models, but here in case of overfit models / linearly-separable data the coefficients go to infinity, increasing the certainty of the prediction \rightarrow This is problem.

- The picture depicts the effects of high magnitude coefficients on probability.
- the point under consideration is near the boundary \rightarrow its probability is uncertain and closer to 0.5 and not 1, 0.
- But as the magnitude of the coefficients of the model increases its certainty raises. This leads to false results.

Effect of linear separability on coefficients



- The picture depicts the effects of high magnitude coefficients on probability.
- the point under consideration is near the boundary \rightarrow its probability is uncertain and closer to 0.5 and not 1, 0.
- But as the magnitude of the coefficients of the model increases its certainty raises. This leads to false results.

Maximum likelihood estimation (MLE)

prefers most certain model \rightarrow

Coefficients go to infinity for linearly-separable data!!!

Overfitting in logistic regression

- Learning tries to find decision boundary that separates data - 'Overly complex boundary'.
- If data is linearly separable -> coefficients go to infinity.

L2 regularized logistic regression

Penalizing large coefficients to mitigate overfitting

- Quality Metric is modified to handle - large coefficients and prevent overfitting.

Desired total cost format

- Want to balance:
 1. How well the function fits the data -> large.
 2. Magnitude of the coefficients -> small.
- Total quality = measure of fit - measure of magnitude of coefficients.
 - Measure of fit -> data likelihood -> large # = good fit for training data.
 - Measure of magnitude of coefficients -> large # = overfit.

Part 1 : Maximum likelihood estimate (MLE):

- Measure of fit = Data likelihood
 - Choose coefficients of w that maximize likelihood.

$$\ell(\mathbf{w}) = \prod_{i=1}^N P(y_i | \mathbf{x}_i, \mathbf{w})$$

- Typically, use log of likelihood function (simplifies math and has better gradient/convergence properties.)

$$\ell\ell(\mathbf{w}) = \ln \prod_{i=1}^N P(y_i | \mathbf{x}_i, \mathbf{w})$$

natural log

Data likelihood is supposed to be as big as possible.

Part 2 : Measure of magnitude of logistic regression coefficients:

- Sum of squares (L2 norm) -> penalize highly positive and highly negative numbers in the same way.
- Sum of absolute value (L1 norm) -> provides sparse solutions.
- Both the metrics penalize large coefficients.

L2 regularized logistic regression

- The mechanism of finding **lambda**/ tuning parameter that balances between the model fit and the coefficient magnitude is **L2 regularized logistic regression**. (In regression case -> term - Ridge Regression);

Picking lambda:

- Validation Set (for large dataset).
- Cross-validation (for smaller datasets).

Bias-Variance tradeoff

- Lambda controls the model complexity.
 - Large lambda : high bias, low variance. ($\hat{w} = 0$, $\lambda = \infty$);
 - Small lambda : low bias, high variance. (MLE for higher order polynomial, $\lambda = 0$).

Consider specific total cost

$$\text{Total quality} = \underbrace{\text{measure of fit}}_{\ell(w)} - \underbrace{\text{measure of magnitude of coefficients}}_{\|w\|_2^2}$$

$$\ell(w) - \lambda \|w\|_2^2$$

tuning parameter = balance of fit and magnitude

If $\lambda = 0$:

Reduces $\max_w \ell(w) \rightarrow$ Standard (unpenalized) MLE solution

If $\lambda = \infty$:

$\max_w \ell(w) - \infty \|w\|_2^2 \rightarrow$ only care about penalizing w , large coefficients $\rightarrow w = 0$

If λ in between:

Balance data fit against the magnitude of the coefficients

Bias-variance tradeoff

Large λ :

high bias, low variance

(e.g., $\hat{w} = 0$ for $\lambda = \infty$)

In essence, λ controls model complexity

Small λ :

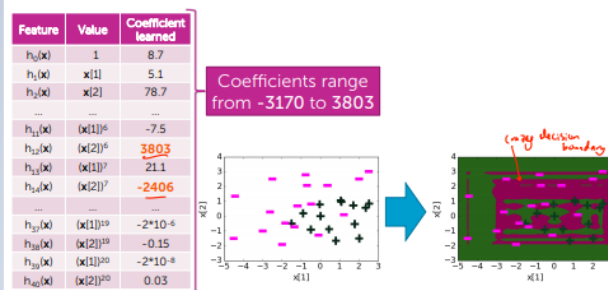
low bias, high variance

(e.g., maximum likelihood (MLE) fit of high-order polynomial for $\lambda = 0$)

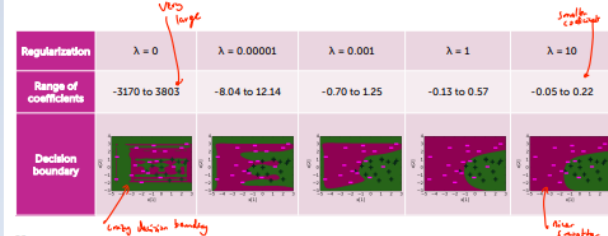
L2 regularization address overfitting issues.

- Choosing appropriate lambda value with higher order complex models
 - Provides a better decision boundary between the data.
 - The overconfidence predictions is reduced as a natural uncertainty region is obtained.

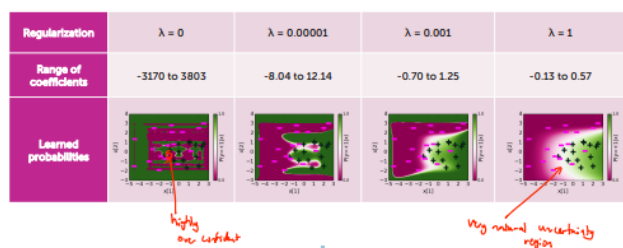
Degree 20 features, $\lambda = 0$



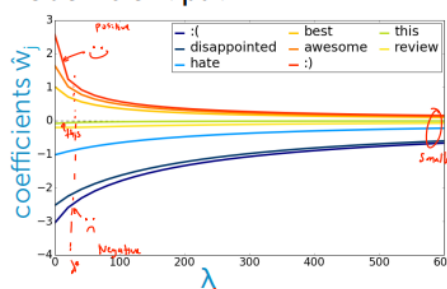
Degree 20 features, effect of regularization penalty λ



Degree 20 features: regularization reduces "overconfidence"



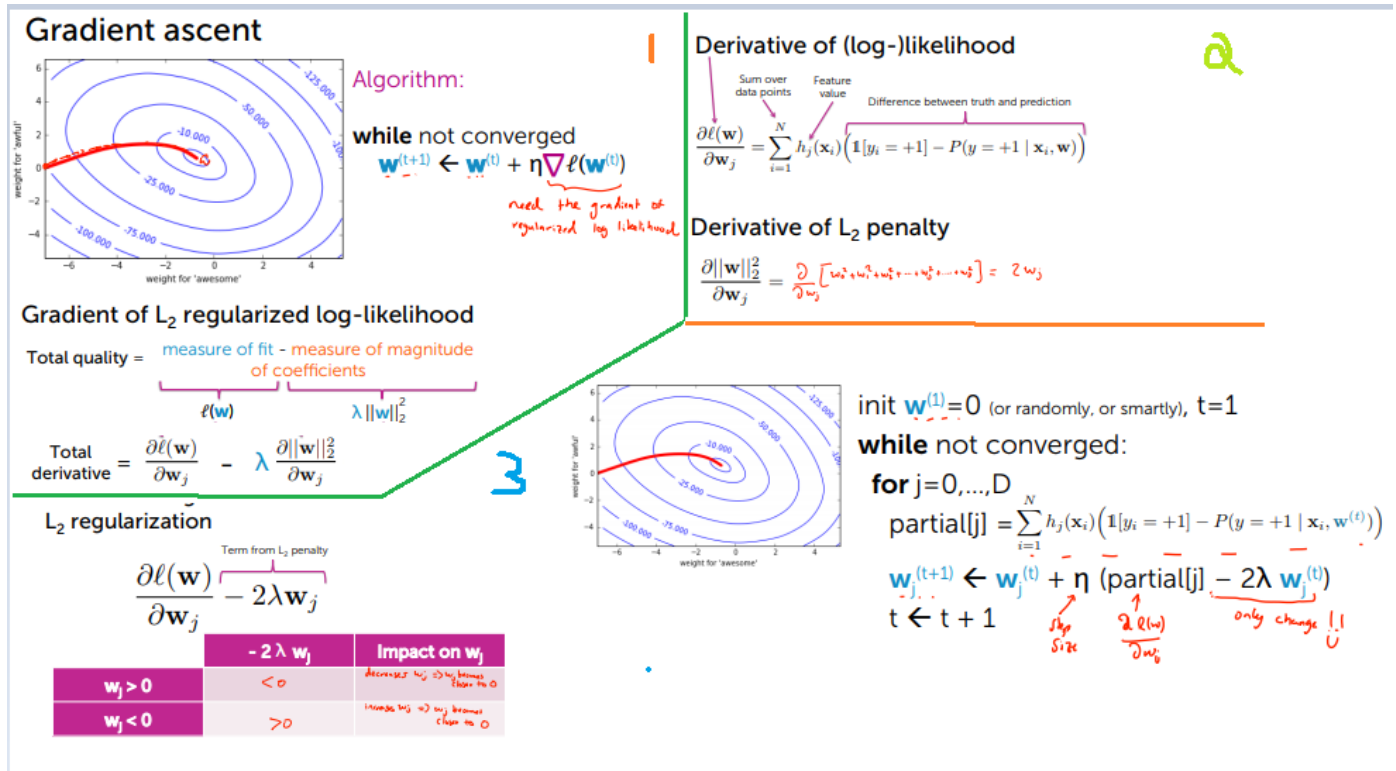
Coefficient path



Coefficient path
As lambda increases coefficient decreases. Choose lambda* that fits model best.

Learning L2 regularized logistic regression with gradient ascent

- Algorithm to optimize to get \hat{w} .
- Employ gradient ascent algorithm to find the \hat{w} .



Sparse logistic regression with L1 regularization

- This provides Efficiency and Interpretability.

Efficiency:

- If $\text{size}(\hat{\mathbf{w}}) = 100\text{B}$, each prediction is expensive.
- If $\hat{\mathbf{w}}$ is sparse, computation only depends on the # non-zeros.

$$\hat{y}(i) = \text{sign}(\text{Sum}(\mathbf{w}_j \neq 0) \hat{\mathbf{w}}_j * h_j(\mathbf{x}_i))$$

Interpretability:

- Can decipher the features truly relevant for the prediction.

Sparse Logistic Regression

- Total quality = measure of fit - measure of magnitude of coefficients
- Total quality = $\ell(\mathbf{w}) - \|\mathbf{w}\|_1$
- L1 norm leads to sparse solutions.

L1 regularized logistic regression

- Lambda is a tuning parameter that provides a balance between the model fit and sparsity.

Coefficient path - L1 penalty

- The coefficients that contribute to the model become 0 eventually with increasing lambda. While those hardly contributing to the model become 0 initially.

Sparse logistic regression

$$\hat{y}_i = \text{sign} \left(\sum_{\hat{w}_j \neq 0} \hat{w}_j h_j(\mathbf{x}_i) \right)$$

Total quality =

measure of fit - measure of magnitude of coefficients

$\ell(\mathbf{w})$

$\|\mathbf{w}\|_1 = |w_0| + \dots + |w_D|$

L_1 regularized logistic regression

Leads to sparse solutions!

L_1 regularized logistic regression

$$\ell(\mathbf{w}) - \lambda \|\mathbf{w}\|_1$$

tuning parameter = balance of fit and sparsity

If $\lambda=0$:

No regularization \rightarrow standard MLE solution

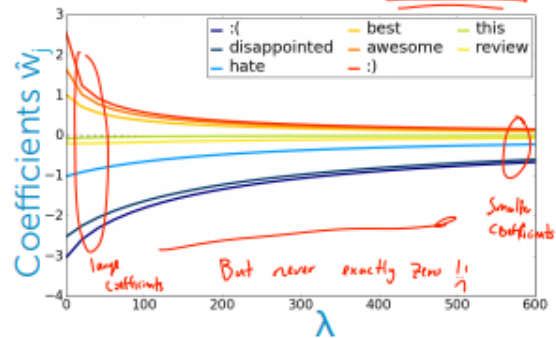
If $\lambda=\infty$:

all weight is on regularization $\rightarrow \hat{\mathbf{w}} = 0$

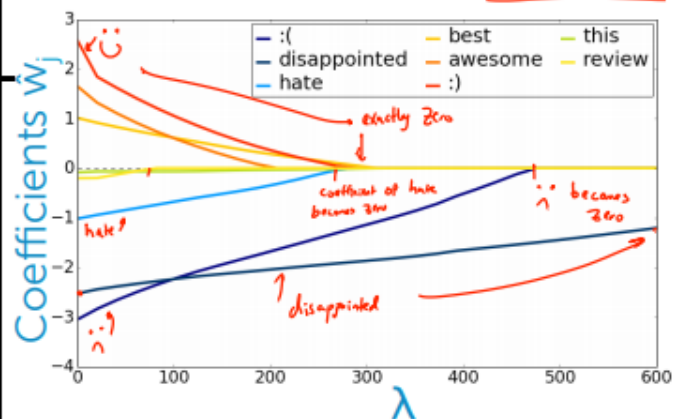
If λ in between:

Sparse solutions: some $\hat{w}_j \neq 0$, many other $\hat{w}_j = 0$

Regularization path – L_2 penalty



Regularization path – L_1 penalty



Quiz

1. Consider four classifiers, whose classification performance is given by the following table:

	Classification error on training set	Classification error on validation set
Classifier 1	0.2	0.6
Classifier 2	0.8	0.6
Classifier 3	0.2	0.2
Classifier 4	0.5	0.4

Which of the four classifiers is most likely overfit?

- ☒ Classifier 1
☐ Classifier 2
☐ Classifier 3
☐ Classifier 4

2. Suppose a classifier classifies 23100 examples correctly and 1900 examples incorrectly. Compute error by hand. Round your answer to 3 decimal places.

0.076

3. (True/False) Accuracy and error measured on the same dataset always sum to 1.

- ☒ True
☐ False

4. Which of the following is NOT a correct description of complex models?

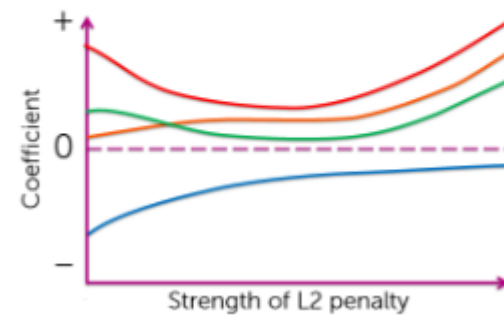
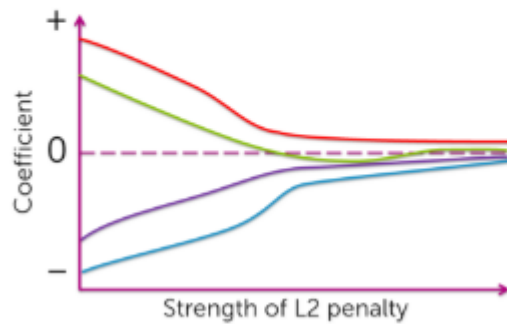
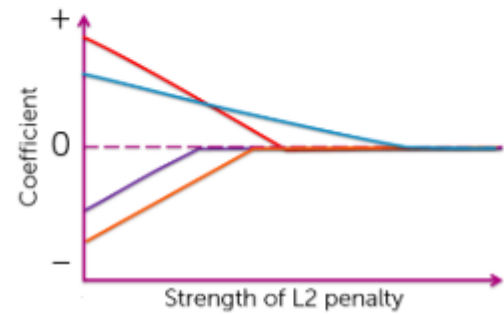
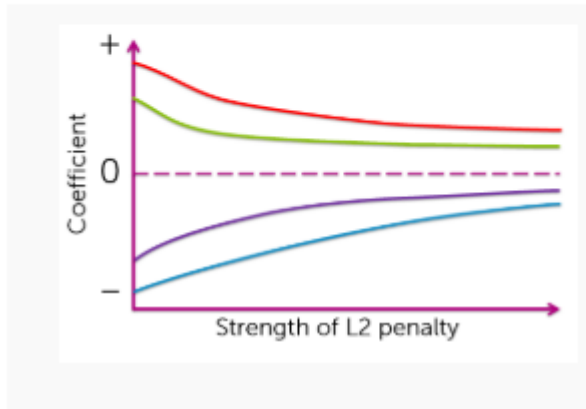
- ☐ Complex models accommodate many features.
☐ Complex models tend to produce lower training error than simple models.
☒ Complex models tend to generalize better than simple models.
☐ Complex models tend to exhibit high variance in response to perturbation in the training data.
☐ Complex models tend to exhibit low bias, capturing many patterns in the training data that simple models may have missed.

5. Which of the following is a symptom of overfitting in the context of logistic regression? Select all that apply.

- ☒ Large estimated coefficients
☐ Good generalization to previously unseen data
☐ Simple decision boundary
☒ Complex decision boundary
☒ Overconfident predictions of class probabilities

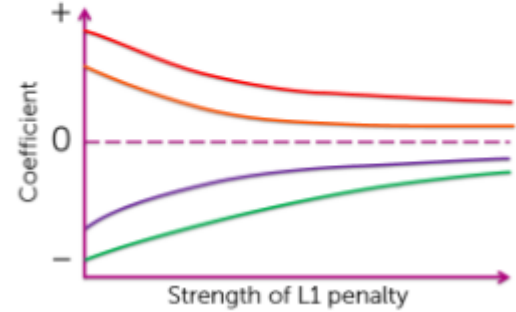
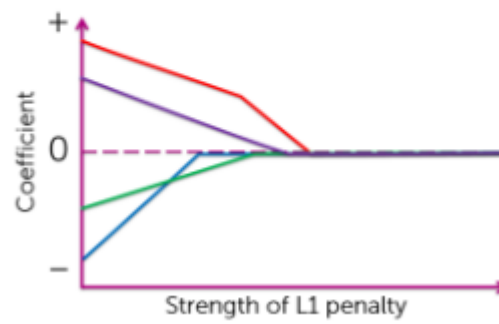
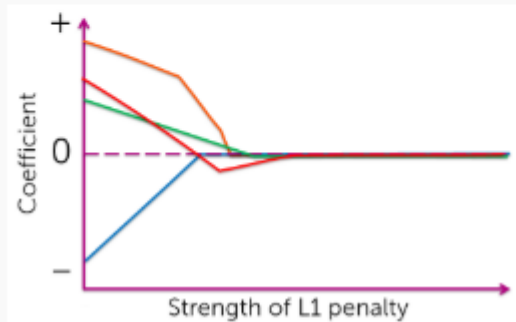
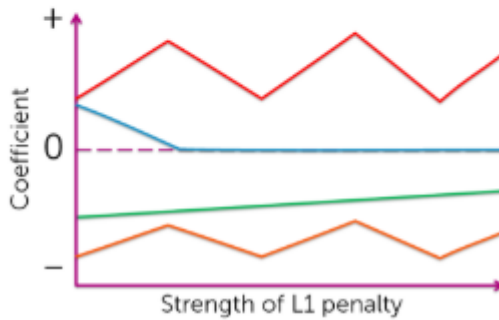
6. Suppose we perform L2 regularized logistic regression to fit a sentiment classifier. Which of the following plots does NOT describe a possible coefficient path? Choose all that apply.

Note. Assume that the algorithm runs for a wide range of L2 penalty values and each coefficient plot is zoomed out enough to capture all long-term trends.



7. Suppose we perform L1 regularized logistic regression to fit a sentiment classifier. Which of the following plots does NOT describe a possible coefficient path? Choose all that apply.

Note. Assume that the algorithm runs for a wide range of L1 penalty values and each coefficient plot is zoomed out enough to capture all long-term trends.



8. In the context of L2 regularized logistic regression, which of the following occurs as we increase the L2 penalty λ ? Choose all that apply.

- ☒ The L2 norm of the set of coefficients gets smaller
- ☐ Region of uncertainty becomes narrower, i.e., the classifier makes predictions with higher confidence.
- ☒ Decision boundary becomes less complex
- ☐ Training error decreases
- ☒ The classifier has lower variance
- ☐ Some features are excluded from the classifier

In []: