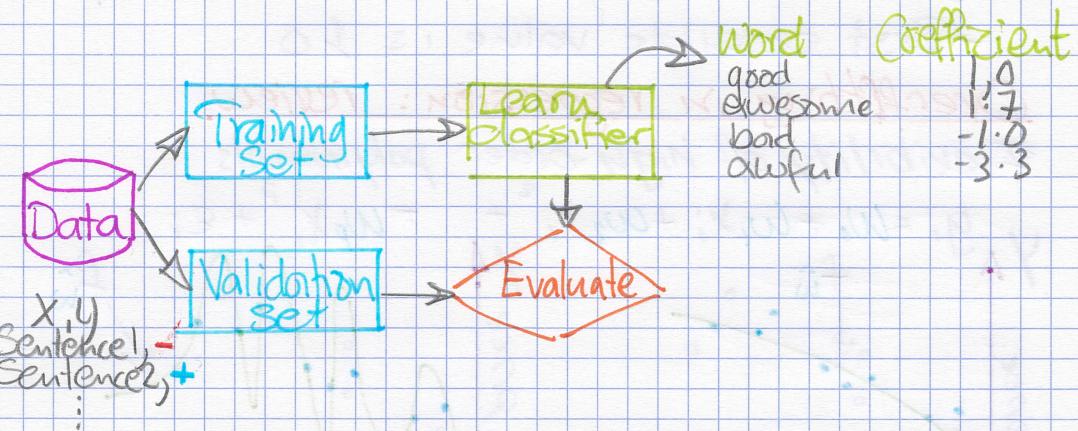


Overfitting & Regularization in Logistic Regression

Training and evaluating a classifier

Training a classifier = Learning the coefficients



Classification error

"Sushi was great" Learned classifier
"Food was OK" $\hat{y} = +, \hat{y} = +$

Test example

Sushi was great, +
Food was OK, -

Correct!
Mistake!
+ -

Hide label

Correct
Mistakes

1
1

Classification error & accuracy

- Error measures fraction of mistakes

$$\text{error} = \frac{\# \text{Mistakes}}{\text{Total number of data points}}$$

- Best possible value is 0.0

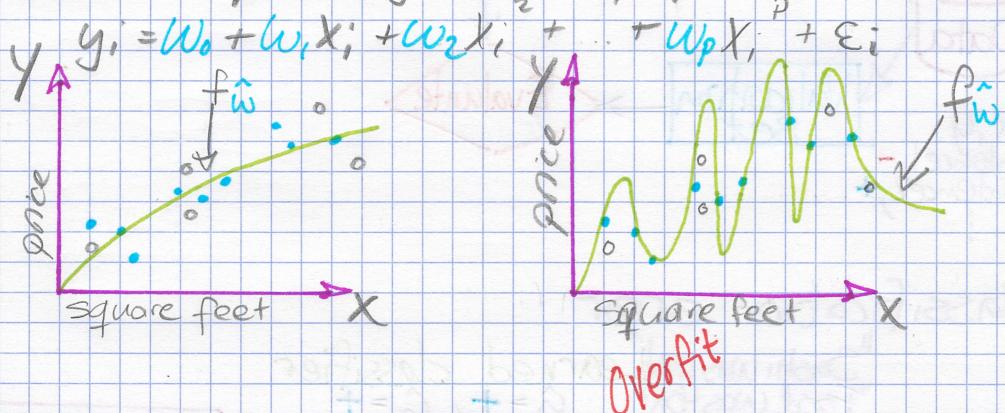
- Often, measure accuracy
 - Fraction of correct predictions

$$\text{accuracy} = \frac{\# \text{correct}}{\text{Total number of data points}}$$

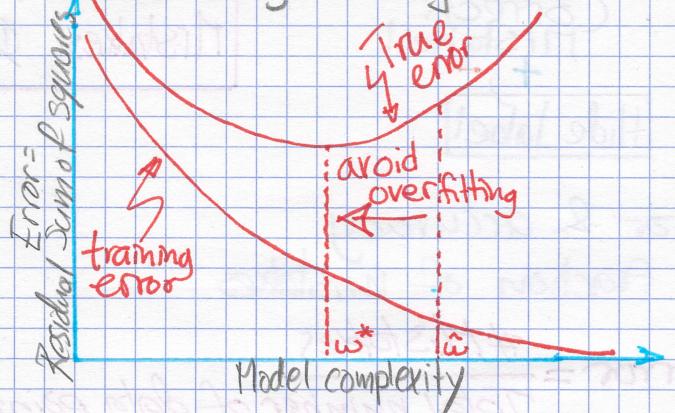
- Best possible value is 1.0

Overfitting in regression: review

Flexibility of high-order polynomials



Overfitting in regression



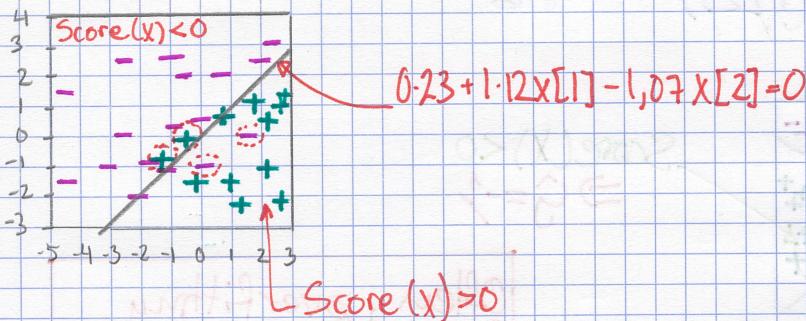
Overfitting if there exists \hat{w} :

- $\text{training_error}(w) > \text{training_error}(\hat{w})$
- $\text{true_error}(w^*) < \text{true_error}(\hat{w})$

Overfitting in classification

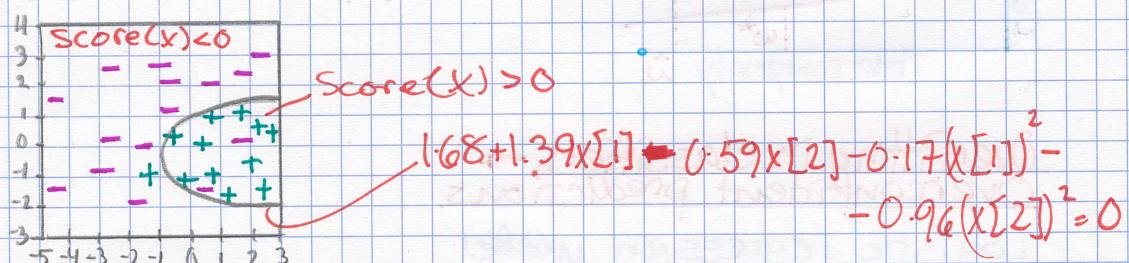
Learned decision boundary

Feature	Value	Coefficient Learned
$h_0(x)$	w_0	1
$h_1(x)$	w_1	$x[1]$
$h_2(x)$	w_2	$x[2]$



Quadratic features (in 2d)

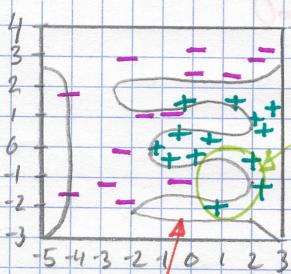
Feature	Value	Coefficient Learned
$h_0(x)$	1	1.68
$h_1(x)$	$x[1]$	1.39
$h_2(x)$	$x[2]$	-0.59
$h_3(x)$	$(x[1])^2$	-0.17
$h_4(x)$	$(x[2])^2$	-0.96



Degree 6 features (in 2d)

Feature	Value	Coefficients learned
$h_0(x)$	1	91.6
$h_1(x)$	$x[1]$	5.3
$h_2(x)$	$x[2]$	-42.7
$h_3(x)$	$(x[1])^2$	-15.9
:		!
$h_{11}(x)$	$(x[1])^6$	5.8
$h_{12}(x)$	$(x[2])^6$	-8.6

Coefficient values getting large

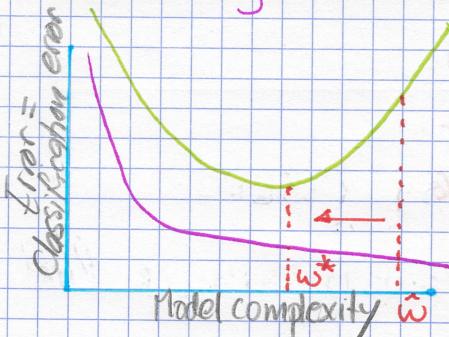


$$\text{Score}(x) < 0 \Rightarrow \hat{y} = -1$$

extremely complex decision boundary

Often, overfitting associated with very large estimated coefficients \hat{w}

Overfitting in classification



Overfitting if there exists w^*

- $\text{training_error}(w^*) > \text{training_error}(\hat{w})$
- $\text{true_error}(w^*) < \text{true_error}(\hat{w})$

Overfitting in classifiers \rightarrow Overconfident predictions

Logistic regression model

$$P(y=+1 | X_i, w) = \text{Sigmoid}(w^\top h(x_i))$$

The subtle (negative) consequence of overfitting in logistic regression

Overfitting \rightarrow Large coefficient value

$\hat{w}^T h(x_i)$ is very positive (or very negative)

$\Rightarrow \text{Sigmoid}(\hat{w}^T h(x_i))$ goes to 1 (or to 0)

Model becomes extremely overconfident of predictions

Penalizing large coefficients to mitigate overfitting.

Desired total cost format

Want to balance:

- i How well function fits data
- ii Magnitude of coefficients

Total quality = $\underbrace{\text{measure of fit}}_{\substack{\uparrow \\ (\text{data likelihood})}} - \underbrace{\text{measure of magnitude of coefficients}}_{\substack{\downarrow \\ \text{large # = overfit}}} \quad \text{want to balance}$

\uparrow \downarrow

large # = good fit to training data

large # = overfit

Maximum likelihood estimation (MLE):

Measure of fit = Data likelihood

Choose coefficients so that maximize likelihood:

$$\prod_{i=1}^N P(y_i | x_i; w)$$

• Typically we use the log of likelihood function
(Simplifies math and has better convergence properties)

$$Q(w) = \ln \prod_{i=1}^N P(y_i | x_i; w)$$

Measure of magnitude of logistic regression coefficients

What summary # is indicative of size of logistic regression coefficients?

- Sum of squares (L_2 norm)

$$\|w\|_2^2 = w_0^2 + w_1^2 + w_2^2 + \dots + w_d^2$$

large weights

- Sum of absolute value (L_1 norm)

$$\|w\|_1 = |w_0| + |w_1| + |w_2| + \dots + |w_d|$$

penalize large weights

Consider specific total cost

\max_w

Total quality =

measure of fit - measure of magnitude of coefficients

$\ell(w)$

↑
log data likelihood

$\|w\|_2^2$

↑
 L_2 penalty

L_2 regularized logistic regression

Consider resulting objective

What if w selected to minimize:

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

↑ tuning parameter = balance of fit and magnitude

If $\lambda = 0$: Reduces $\max_w \ell(w) \rightarrow$ Standard MLE Solution

If $\lambda = \infty$: $\max_w \ell(w) = \infty \|w\|_2^2 \rightarrow$ Only care about penalizing w , large coefficient $\rightarrow w = 0$

If λ in between

Balance data fit against the magnitude of the coefficients

L_2 regularized logistic regression

Pick λ using:

- Validation set (for large datasets)
- Cross-validation (for smaller datasets)

Bias-variance tradeoff

Large λ :

high bias low variance
(e.g., $w = 0$ for $\lambda = \infty$)

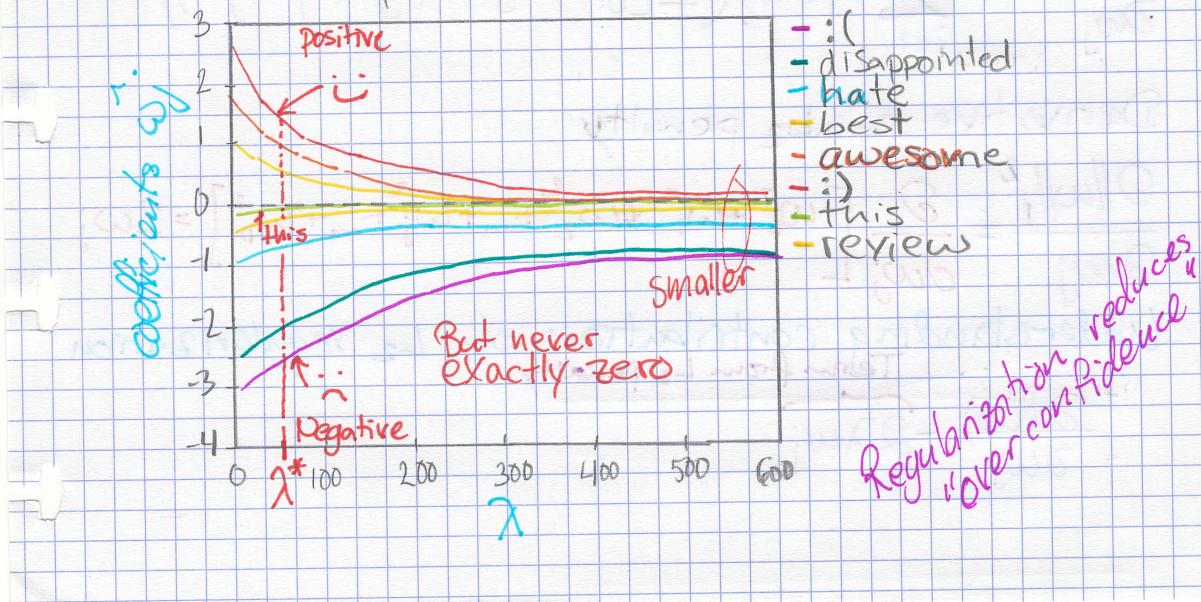
In essence, λ controls model complexity

Small λ :

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

Visualizing effect of L_2 regularization in logistic regression

Coefficient path



Finding best L₂ regularized linear classifier with gradient ascent

Algorithm:

while not converged:

$$w^{(t+1)} \leftarrow w^{(t)} + \eta \nabla \ell(w^{(t)})$$

Gradient of L₂ regularized log-likelihood

Total quality =

measure of fit - measure of magnitude of coefficients

$$\text{Total derivative} = \frac{\partial \ell(w)}{\partial w_j} - \lambda \frac{\partial \|w\|_2^2}{\partial w_j}$$

$$\frac{\partial \ell(w)}{\partial w_j} = \sum_{i=1}^N h_j(x_i) \left(\mathbb{1}[y_i=+1] - P(y=+1|x_i, w) \right)$$

Derivative of L₂ penalty

$$\frac{\partial \|w\|_2^2}{\partial w_j} = \frac{\partial}{\partial w_j} [w_0^2 + w_1^2 + w_2^2 + \dots + w_j^2 + \dots + w_d^2] = 2w_j$$

Understanding contribution of L₂ regularization

$$\frac{\partial \ell(w)}{\partial w_j} - 2\lambda w_j$$

$w_j > 0$ $-2\lambda w_j < 0$ Impact on w_j
 decreases $w_j \Rightarrow w_j$ becomes closer to 0

$w_j < 0$ > 0 Increasing $w_j \Rightarrow w_j$ becomes closer to 0

Summary of gradient ascent for logistic regression with L_2 regularization

init $\underline{w}^{(0)} = 0$ (or randomly, or smartly), $t=1$

while not converged:

for $j = 0, \dots, D$:
 $\text{partial}[j] = \sum_{i=1}^n h_j(x_i) (\mathbb{I}[y_i = +1] - P(y=+1|X_i, \underline{w}^{(t)}))$

$$w_j^{(t+1)} \leftarrow w_j^{(t)} + \eta (\text{partial}[j] - 2\lambda w_j^{(t)})$$

↓ stepsize ↓ $\frac{\partial Q(w)}{\partial w_j}$ ↓ only change!

Sparse logistic regression with L_1 regularization

Recall sparsity (many $\hat{w}_j = 0$) gives efficiency and interpretability

Efficiency

- If $\text{size}(w) = 100B$, each prediction is expensive

- If \hat{w} sparse, computation only depends on # of non-zeros
 many zeros

$$\hat{y}_i = \text{Sign} \left(\sum_{\substack{j \\ \hat{w}_j \neq 0}} \hat{w}_j h_j(x_i) \right)$$

Interpretability:

- which features are relevant for predictions?

Sparse logistic regression

Total quality =

measure of fit - measure of magnitude
of coefficients

$$\underbrace{\varphi(w)}$$

$$\|w\|_1 = |w_0| + \dots + |w_d|$$

L_1 regularized
logistic regression

↳ Leads to
sparse
solutions!

L_1 regularized logistic regression

Just like L_2 regularization; solution is governed by a 'continuous' parameter ↗

$$\varphi(w) - \lambda \|w\|_1$$

↑ tuning parameter =
balance of fit and sparsity

If $\lambda = 0$: No regularization → standard MLE solution

If $\lambda = \infty$: all weight is on regularization → $\hat{w} = 0$

If λ in between: Sparse solutions:

Some $\hat{w}_j \neq 0$, many other $\hat{w}_j = 0$

Regularization path - L_1 penalty

