

Emily Fox & Carlos Guestrin Machine Learning Specialization University of Washington

©2015-2016 Emily Fox & Carlos Guestrir

Training and evaluating a classifier

©2015-2016 Emily Fox & Carlos Guestrin

Training a classifier = Learning the coefficients



Classification error

Learned classifier

ŷ = 🕂



Miastradat!





©2015-2016 Emily Fox & Carlos Guestrin

Classification error & accuracy

• Error measures fraction of mistakes



- Best possible value is 0.0
- Often, measure accuracy
 - Fraction of correct predictions

– Best possible value is 1.0

6

©2015-2016 Emily Fox & Carlos Guestrin

Overfitting in regression: review

©2015-2016 Emily Fox & Carlos Guestrin

Flexibility of high-order polynomials



Machine Learning Specialization



Overfitting in classification

©2015-2016 Emily Fox & Carlos Guestrin

Decision boundary example





©2015-2016 Emily Fox & Carlos Guestrin

Learned decision boundary



©2015-2016 Emily Fox & Carlos Guestrin



Quadratic features (in 2d)

Note: we are not including cross terms for simplicity

17

©2015-2016 Emily Fox & Carlos Guestrin

Degree 6 features (in 2d)

Note: we are not including cross terms for simplicity



Degree 20 features (in 2d) Note: we are not including cross terms for simplicity



Often, overfitting associated with very large estimated coefficients ŵ





©2015-2016 Emily Fox & Carlos Guestrin



Overfitting in classifiers → Overconfident predictions

©2015-2016 Emily Fox & Carlos Guestrin



The subtle (negative) consequence of overfitting in logistic regression

Overfitting \rightarrow Large coefficient values

→ sigmoid($\mathbf{M}^{T}h(\mathbf{x}_{i})$) is very positive (or very negative) → sigmoid($\mathbf{M}^{T}h(\mathbf{x}_{i})$) goes to 1 (or to 0)

Model becomes extremely overconfident of predictions

©2015-2016 Emily Fox & Carlos Guestrin

Effect of coefficients on logistic regression model

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







©2015-2016 Emily Fox & Carlos Guestrin

Learned probabilities



Quadratic features: Learned probabilities



Overfitting Overconfident predictions



©2015-2016 Emily Fox & Carlos Guestrin

Machine Learning Specialization

29

Overfitting in logistic regression: Another perspective



@2015-2016 Emily Fox & Carlos Guestrin

Linearly-separable data



Note 1: If you are using D features, linear separability happens in a D-dimensional space

Note 2: If you have enough features, data are (almost) always linearly separable

Data are linearly separable if:

- There exist coefficients ŵ such that:
 - For all positive training data

 $S_{core}(x) = \hat{w}^T h(x) > 0$

- For all negative training data $\int core(x) = \hat{w}^T h(x) < 0$



Machine Learning Specialization

33



©2015-2016 Emily Fox & Carlos Guestrin

Maximum likelihood estimation (MLE) prefers most certain model → Coefficients go to infinity for linearly-separable data!!!



©2015-2016 Emily Fox & Carlos Guestrin

Machine Learning Specialization

35

Overfitting in logistic regression is "twice as bad"



Penalizing large coefficients to mitigate overfitting

©2015 Emily Fox & Carlos Guestrin



Desired total cost format

Want to balance:

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



Maximum likelihood estimation (MLE): Measure of fit = Data likelihood

• Choose coefficients **w** that maximize likelihood:

NΤ

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

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

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

©2015-2016 Emily Fox & Carlos Guestrin

Machine Learning Specialization

41

Measure of magnitude of logistic regression coefficients

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



©2015 Emily Fox & Carlos Guestrin

Consider specific total cost



Machine Learning Specialization

©2015 Emily Fox & Carlos Guestrin

Consider resulting objective

 $\ell(\mathbf{w}) - \lambda ||\mathbf{w}||_{2}^{2}$ $lf \lambda = 0:$ $lf \lambda = 0:$ $lf \lambda = 0:$ $lf \lambda = \infty:$ $lf \lambda = \infty:$ - 9 max R(w) - ob ||w||2 -> only care about penalizing w, large coefficients -> W=o If λ in between: Balance Anta hit against the magnitude of the coefficients 45 ©2015 Emily Fox & Carlos Guestrin Machine Learning Specialization

Consider resulting objective

What if $\hat{\mathbf{w}}$ selected to minimize

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

tuning parameter = balance of fit and magnitude

L₂ regularized logistic regression

Pick λ using:

- Validation set (for large datasets)
- Cross-validation (for smaller datasets) (see regression course)

©2015 Emily Fox & Carlos Guestrin

Bias-variance tradeoff

Large λ :

high bias, low variance

(e.g., $\hat{\mathbf{w}} = 0$ 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 λ =0)

Visualizing effect of regularization on logistic regression

©2015-2016 Emily Fox & Carlos Guestrin

Degree 20 features, $\lambda = 0$



©2015-2016 Emily Fox & Carlos Guestrin

Machine Learning Specialization

51

Degree 20 features, effect of regularization penalty λ



Coefficient path



Machine Learning Specialization

53

Degree 20 features: regularization reduces "overconfidence"



Finding best L₂ regularized linear classifier with gradient ascent



Gradient ascent



Algorithm:

while not converged $\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} + \eta \nabla \ell(\mathbf{w}^{(t)})$ need the gradient of regularized log likelihood

©2015 Emily Fox & Carlos Guestrin

Gradient of L₂ regularized log-likelihood



Machine Learning Specialization

61

©2015 Emily Fox & Carlos Guestrin

Derivative of (log-)likelihood



Derivative of L₂ penalty

$$\frac{\partial ||\mathbf{w}||_2^2}{\partial \mathbf{w}_j} = \frac{\partial}{\partial \mathbf{w}_j} \left[\frac{\omega_0^2 + \omega_1^2 + \omega_2^2 + \dots + \omega_j^2}{\partial \mathbf{w}_j} \right] = 2 \omega_j$$

Understanding contribution of L₂ regularization



Summary of gradient ascent for logistic regression with L₂ Regularization



init $\mathbf{w}^{(1)} = 0$ (or randomly, or smartly), t=1while not converged: for j=0,...,Dpartial[j] $= \sum_{i=1}^{N} h_j(\mathbf{x}_i) \left(\mathbb{1}[y_i = +1] - P(y = +1 \mid \mathbf{x}_i, \mathbf{w}^{(t)}) \right)$ $\mathbf{w}^{(t+1)}_j \leftarrow \mathbf{w}^{(t)}_j + \eta$ (partial[j] $- 2\lambda \mathbf{w}^{(t)}_j$) $t \leftarrow t+1$

©2015 Emily Fox & Carlos Guestrin

Sparse logistic regression with L₁ regularization

©2015-2016 Emily Fox & Carlos Guestrin

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

Efficiency:

- If size(w) = 100B, each prediction is expensive
- If $\hat{\mathbf{w}}$ sparse, computation only depends on # of non-zeros $\hat{y}_i = sign\left(\sum_{\hat{\mathbf{w}}_j \neq 0} \hat{\mathbf{w}}_j h_j(\mathbf{x}_i)\right)$

Interpretability:

- Which features are relevant for prediction?

Sparse logistic regression



©2015 Emily Fox & Carlos Guestrin

Machine Learning Specialization

69

L₁ regularized logistic regression

Just like L2 regularization, solution is governed by a continuous parameter λ

 $\ell(\mathbf{w}) - \lambda \| \mathbf{w} \|_{1}$ tuning parameter = balance of fit and sparsity $ko cycle risebon \rightarrow shadad KLE solution$ $lf \lambda = \infty:$ $s all weight is on regularization \rightarrow \hat{\omega} = 0$ $lf \lambda \text{ in between:}$ $Space solutions: Some \tilde{w}_{j} \neq 0, may other \tilde{w}_{j} = 0$ 21



©2015 Emily Fox & Carlos Guestrin



73

©2015 Emily Fox & Carlos Guestrin

Summary of overfitting in logistic regression

©2015-2016 Emily Fox & Carlos Guestrin

What you can do now...

- Identify when overfitting is happening
- Relate large learned coefficients to overfitting
- Describe the impact of overfitting on decision boundaries and predicted probabilities of linear classifiers
- Motivate the form of L₂ regularized logistic regression quality metric
- Describe what happens to estimated coefficients as tuning parameter λ is varied
- Interpret coefficient path plot
- Estimate L₂ regularized logistic regression coefficients using gradient ascent
- Describe the use of L₁ regularization to obtain sparse logistic regression solutions