# CS489/698: Intro to ML

Lecture 04: Logistic Regression

UNIVERSITY OF **WATERLOO**

# Outline

- Announcements

- Bernoulli model

- Logistic regression

- Computation

9/21/17
Yao-Liang Yu

UNIVERSITY OF
WATERLOO

# Outline

- **Announcements**


- Bernoulli model


- Logistic regression


- Computation

9/21/17 Yao-Liang Yu

# Announcements

- Assignment 1 <span style="color:red">due next Tuesday</span>

Yao-Liang Yu

**UNIVERSITY OF WATERLOO**

# Outline

- Announcements

- Bernoulli model

- Logistic regression

- Computation

9/21/17 Yao-Liang Yu

UNIVERSITY OF
**WATERLOO**

# Classification revisited

- $\hat{y} = \text{sign}(\ \mathbf{x}^{\top}\mathbf{w} + b\ )$

- How confident we are about $\hat{y}$?

- $|\mathbf{x}^{\top}\mathbf{w} + b|$ seems a good indicator
  - real-valued; hard to interpret
  - ways to transform into [0,1]

- Better(?) idea: learn confidence directly

9/21/17                                    Yao-Liang Yu

# Conditional probability

- P(Y=1 | X=$\mathbf{x}$): conditional on seeing $\mathbf{x}$, what is the chance of this instance being positive, i.e., Y=1?
  - obviously, value in [0,1]

- P(Y=0 | X=$\mathbf{x}$) = 1 − P(Y=1 | X=$\mathbf{x}$), if two classes
  - more generally, sum to 1

Notation (Simplex). $\Delta_{k-1} := \{\ \mathbf{p}\ \text{in}\ R^k : \mathbf{p} \geq 0,\ \Sigma_i\ p_i = 1\ \}$

Yao-Liang Yu

# Reduction to a harder problem

- $P(Y{=}1 \mid X{=}\mathbf{x}) = E(1_{Y=1} \mid X{=}\mathbf{x})$

$$1_A = \begin{cases} 1, & A \text{ is true} \\ 0, & A \text{ is false} \end{cases}$$

- Let $Z = 1_{Y=1}$, then regression function for $(X, Z)$
  - use linear regression for binary Z?

- Exploit structure!
  - conditional probabilities are in a simplex

- Never reduce to unnecessarily harder problem

9/21/17                                    Yao-Liang Yu

UNIVERSITY OF
WATERLOO

# Bernoulli model

- Let P(Y=1 | X=$\mathbf{x}$) = p($\mathbf{x}$; $\mathbf{w}$), parameterized by $\mathbf{w}$

- Conditional likelihood on {($\mathbf{x}_1$, $y_1$), … ($\mathbf{x}_n$, $y_n$)}:

$$\mathbf{P}(Y_1 = y_1, \ldots, Y_n = y_n | X_1 = \mathbf{x_1}, \ldots, X_n = \mathbf{x}_n)$$

  - simplifies if independence holds

$$\prod_{i=1}^{n} \mathbf{P}(Y_i = y_i | X_i = \mathbf{x}_i) = \prod_{i=1}^{n} p(\mathbf{x}_i; \mathbf{w})^{y_i} (1 - p(\mathbf{x}_i; \mathbf{w}))^{1-y_i}$$

  - Assuming $y_i$ is {0,1}-valued

# Naïve solution

$$\prod_{i=1}^{n} p(\mathbf{x}_i; \mathbf{w})^{y_i} (1 - p(\mathbf{x}_i; \mathbf{w}))^{1-y_i}$$

- Find **w** to maximize conditional likelihood

- What is the solution if p(**x**; **w**) does not depend on **x**?

- What is the solution if p(**x**; **w**) does not depend on ?

UNIVERSITY OF
**WATERLOO**

# Generalized linear models (GLM)

- y ~ Bernoulli(p);  p = p($\mathbf{x}$; $\mathbf{w}$) natural parameter
  - Logistic regression

- y ~ Normal($\boldsymbol{\mu}$, $\sigma^2$); $\boldsymbol{\mu}$ = $\boldsymbol{\mu}$($\mathbf{x}$; $\mathbf{w}$)
  - (weighted) least-squares regression

- GLM: y ~ exp( θ φ(y) – A(θ) )

log-partition function

sufficient statistics

UNIVERSITY OF
**WATERLOO**

# Outline

- Announcements

- Bernoulli model

- Logistic regression

- Computation

9/21/17    Yao-Liang Yu

UNIVERSITY OF
WATERLOO

# Logit transform

- p(**x**; **w**) = **w**$^\top$**x**?  p >=0 not guaranteed…

- log p(**x**; **w**) = **w**$^\top$**x**?  better!
  - LHS negative, RHS real-valued…

  odds ratio

- Logit transform  $\log \dfrac{p(\mathbf{x}; \mathbf{w})}{1 - p(\mathbf{x}; \mathbf{w})} = \mathbf{w}^\top \mathbf{x}$

- Or equivalently  $p(\mathbf{x}; \mathbf{w}) = \dfrac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})}$

UNIVERSITY OF
WATERLOO

# Prediction with confidence

$$p(\mathbf{x}; \mathbf{w}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})}$$

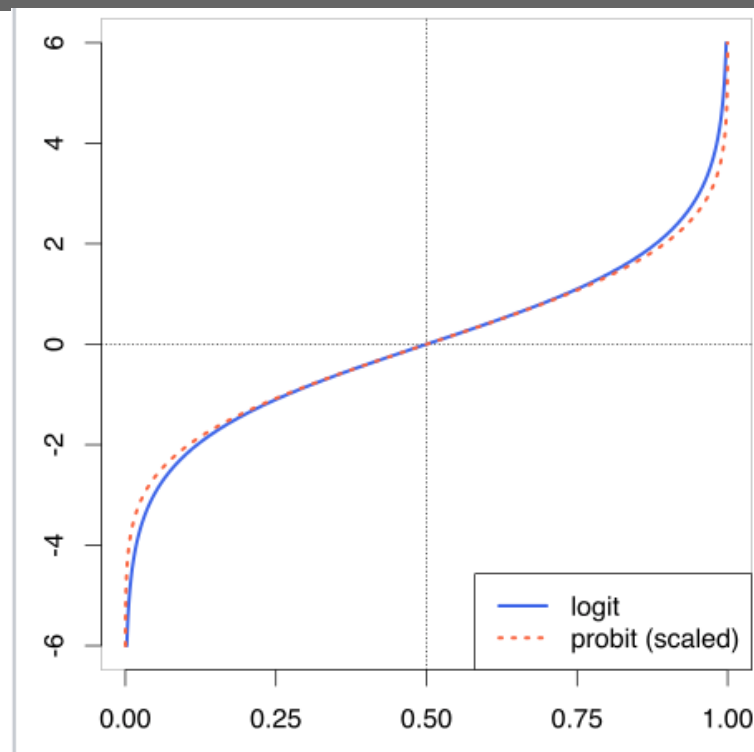- ŷ = 1 if p = P(Y=1 | X=x) > ½ iff $\mathbf{w}^\top\mathbf{x}$ > 0

- Decision boundary $\mathbf{w}^\top\mathbf{x}$ = 0

- ŷ = sign($\mathbf{w}^\top\mathbf{x}$) as before, but with confidence p(x; w)

Yao-Liang Yu

UNIVERSITY OF
WATERLOO

# Not just a classification algorithm

- Logistic regression does more than classification
  - it estimates conditional probabilities
  - under the logit transform assumption

- Having confidence in prediction is nice
  - the price is an assumption that may or may not hold

- If classification is the sole goal, then doing extra work
  - as shall see, SVM only estimates decision boundary

# More than logistic regression

- F(p) transforms p from [0,1] to R

- Then, equating F(p) to a linear function $\mathbf{w}^\top\mathbf{x}$

- But, there are many other choices for F!
  - precisely the inverse of any distribution function!
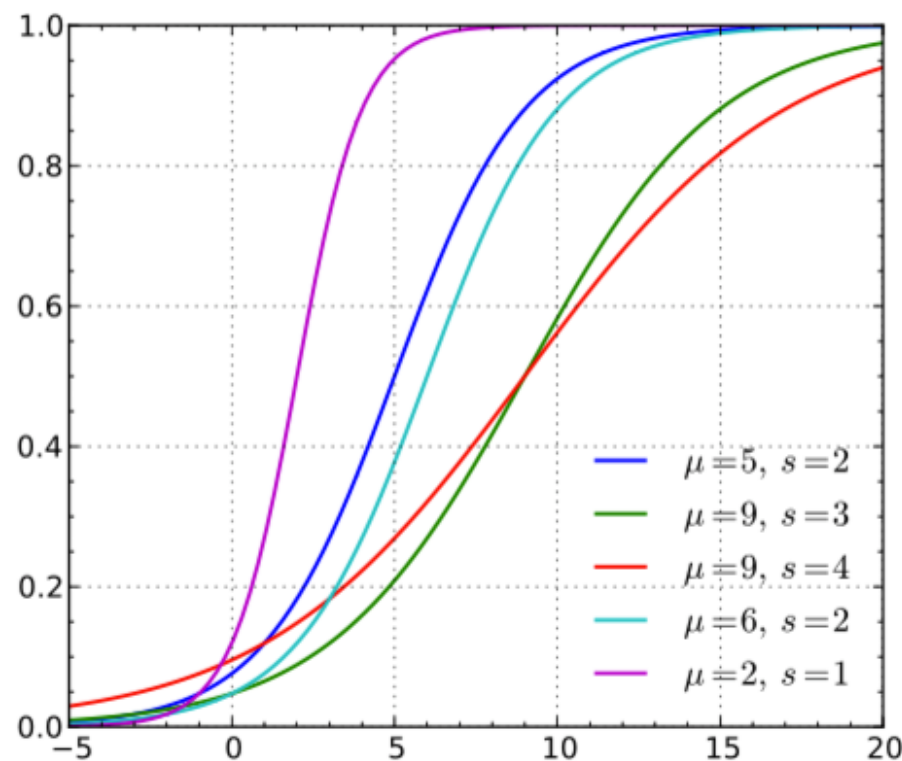


Comparison of the logit function with a scaled probit (i.e. the inverse CDF of the normal distribution), comparing $\mathrm{logit}(x)$ vs. $\Phi^{-1}(x)/\sqrt{\dfrac{\pi}{8}}$, which makes the slopes the same at the origin.

Yao-Liang Yu

# Logistic distribution

- Cumulative Distribution Function

$$F(x; \mu, s) = \frac{1}{1 + \exp(-\frac{x-\mu}{s})}$$

- Mean mu, variance $s^2\pi^2/3$



Legend:
- $\mu = 5,\ s = 2$
- $\mu = 9,\ s = 3$
- $\mu = 9,\ s = 4$
- $\mu = 6,\ s = 2$
- $\mu = 2,\ s = 1$

# Outline

- Announcements

- Bernoulli model

- Logistic regression

- Computation

                                  Yao-Liang Yu

UNIVERSITY OF
**WATERLOO**

# Maximum likelihood

$$\prod_{i=1}^{n} p(\mathbf{x}_i; \mathbf{w})^{y_i} (1 - p(\mathbf{x}_i; \mathbf{w}))^{1-y_i}$$

$$p(\mathbf{x}; \mathbf{w}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})}$$

- Minimize negative log-likelihood

$$\sum_i \log(e^{(1-y_i)\mathbf{w}^\top \mathbf{x}_i} + e^{-y_i \mathbf{w}^\top \mathbf{x}_i}) \equiv \sum_i \log(1 + e^{-\tilde{y}_i \mathbf{w}^\top \mathbf{x}_i})$$

Yao-Liang Yu

# Newton's algorithm

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \eta_t [\nabla^2 f(\mathbf{w}_t)]^{-1} \cdot \nabla f(\mathbf{w}_t)$$

$$\nabla f(\mathbf{w}_t) = X^\top (\mathbf{p} - \mathbf{y})$$

$$\nabla^2 f(\mathbf{w}_t) = \sum_i p_i (1 - p_i) \mathbf{x}_i \mathbf{x}_i^\top$$ PSD

$$p_i = \frac{1}{1 + e^{-\mathbf{w}_t^\top \mathbf{x}_i}}$$

Uncertain predictions get bigger weight

- η = 1: iterative weighted least-squares

# A word about implementation

- Numerically computing exponential can be tricky
  - easily underflows or overflows

- The usual trick
  - estimate the range of the exponents
  - shift the mean of the exponents to 0

UNIVERSITY OF
**WATERLOO**

# Robustness

$$\ell(t) = \log(1 + e^t) \quad L(\hat{y}, y) = \ell(-\hat{y}y) \quad \hat{y} = \mathbf{w}^\top \mathbf{x}$$

- Bounded derivative

$$\ell'(t) = \frac{e^t}{1 + e^t} = \frac{1}{1 + e^{-t}}$$

- Variational exponential

Larger exp loss gets smaller weights

$$\log(1 + e^t) = \min_{0 \le \eta \le 1} \boxed{\eta e^t} - \log(\eta) + \eta - 1$$

9/21/17                                    Yao-Liang Yu

UNIVERSITY OF
WATERLOO

# More than 2 classes

- Softmax

$$\mathbf{P}(Y = c | \mathbf{x}, W) = \frac{\exp(\mathbf{w}_c^\top \mathbf{x})}{\sum_{q=1}^{k} \exp(\mathbf{w}_q^\top \mathbf{x})}$$

- Again, nonnegative and sum to 1

- Negative log-likelihood (y is one-hot)

$$-\log \prod_{i=1}^{n} \prod_{c=1}^{k} p_{ic}^{y_{ic}} = -\sum_{i=1}^{n} \sum_{c=1}^{k} y_{ic} \log p_{ic}$$

9/21/17 Yao-Liang Yu

UNIVERSITY OF
**WATERLOO**

# Questions?

         Yao-Liang Yu

UNIVERSITY OF
**WATERLOO**