

CSE 446 Logistic Regression Winter 2012

Dan Weld

Some slides from Carlos Guestrin, Luke Zettlemoyer

Gaussian Naïve Bayes

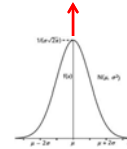
Sometimes Assume Variance

- is independent of Y (i.e., σ_i),
- or independent of X_i (i.e., σ_k)
- or both (i.e., σ)

$$P(Y | \mathbf{X}) \propto P(\mathbf{X} | Y) P(Y)$$

$$P(X_i = x | Y = y_k) = N(\mu_{ik}, \sigma_{ik})$$

$$N(\mu_{ik}, \sigma_{ik}) = \frac{1}{\sigma_{ik} \sqrt{2\pi}} e^{-\frac{(x - \mu_{ik})^2}{2\sigma_{ik}^2}}$$

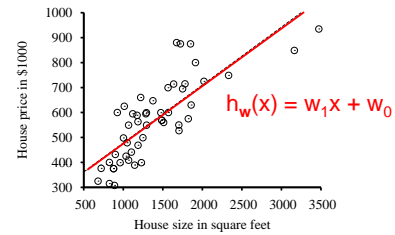


Generative vs. Discriminative Classifiers

- **Want to Learn:** $h: \mathbf{X} \mapsto Y$
 - \mathbf{X} - features
 - Y - target classes
- **Bayes optimal classifier** - $P(Y|\mathbf{X})$
- **Generative classifier**, e.g., Naïve Bayes: $P(Y | \mathbf{X}) \propto P(\mathbf{X} | Y) P(Y)$
 - Assume some **functional form** for $P(\mathbf{X}|Y)$, $P(Y)$
 - Estimate parameters of $P(\mathbf{X}|Y)$, $P(Y)$ directly from training data
 - Use Bayes rule to calculate $P(Y|\mathbf{X} = x)$
 - This is a **'generative' model**
 - Indirect computation of $P(Y|\mathbf{X})$ through Bayes rule
 - As a result, **can also generate a sample of the data**, $P(\mathbf{X}) = \sum_y P(y) P(\mathbf{X}|y)$
- **Discriminative classifiers**, e.g., Logistic Regression:
 - Assume some **functional form** for $P(Y|\mathbf{X})$
 - Estimate parameters of $P(Y|\mathbf{X})$ directly from training data
 - This is the **'discriminative' model**
 - Directly learn $P(Y|\mathbf{X})$
 - But **cannot obtain a sample of the data**, because $P(\mathbf{X})$ is not available

3

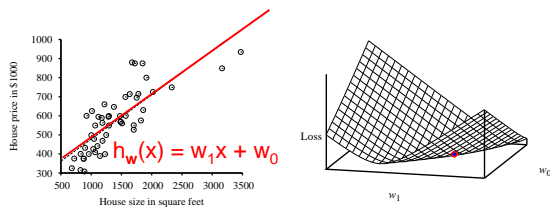
Univariate Linear Regression



$$\text{Loss}(h_w) = \sum_{j=1}^n L_2(y_j, h_w(x_j)) = \sum_{j=1}^n (y_j - h_w(x_j))^2 = \sum_{j=1}^n (y_j - (w_1x_j + w_0))^2$$

4

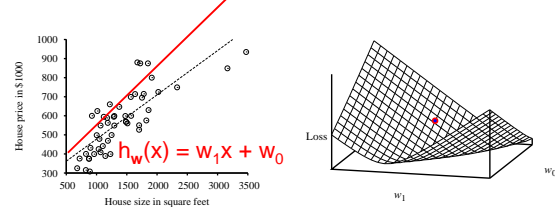
Understanding Weight Space



$$\text{Loss}(h_w) = \sum_{j=1}^n (y_j - (w_1x_j + w_0))^2$$

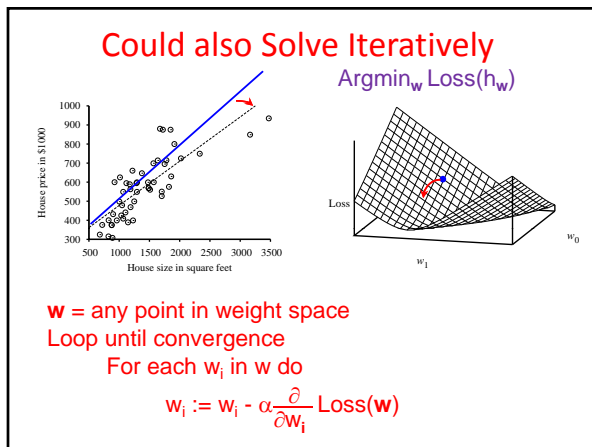
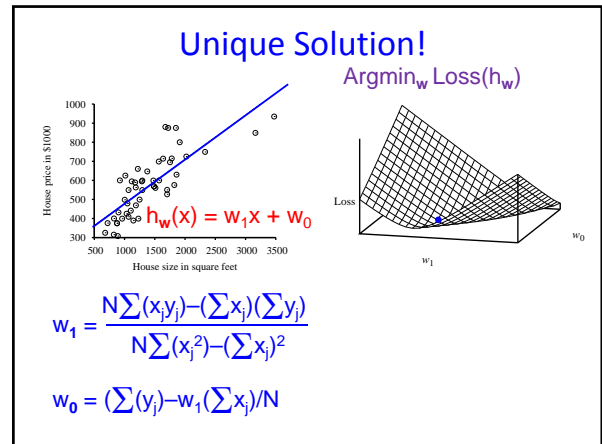
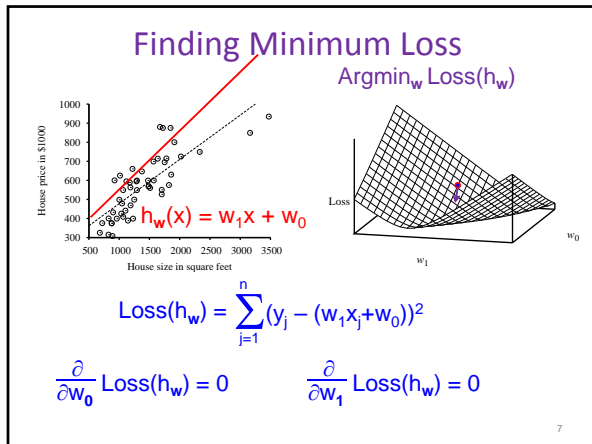
5

Understanding Weight Space



$$\text{Loss}(h_w) = \sum_{j=1}^n (y_j - (w_1x_j + w_0))^2$$

6



Multivariate Linear Regression

$h_w(x_j) = w_0 + \sum w_i x_{j,i} = \sum w_i x_{j,i} = w^T x_j$

Argmin_w Loss(h_w)

Unique Solution = $(x^T x)^{-1} x^T y$

Problem....

10

Overfitting

Regularize!!

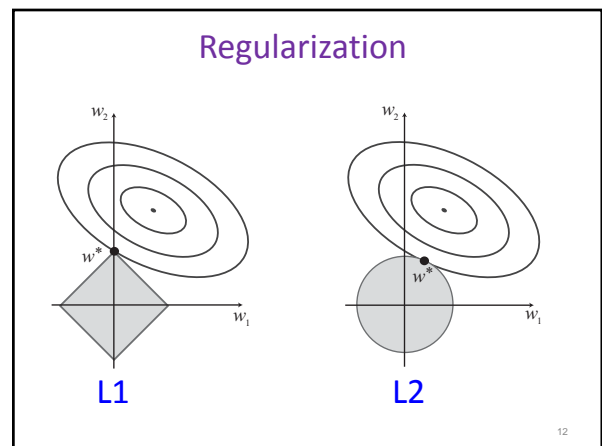
Penalize high weights

$$\text{Loss}(h_w) = \sum_{j=1}^n (y_j - (w_1x_j + w_0))^2 + \lambda \sum_{i=1}^k w_i^2$$

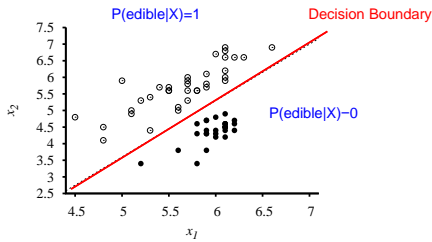
Alternatively....

$$\text{Loss}(h_w) = \sum_{j=1}^n (y_j - (w_1x_j + w_0))^2 + \lambda \sum_{i=1}^k |w_i|$$

11



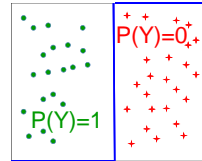
Back to Classification



13

Logistic Regression

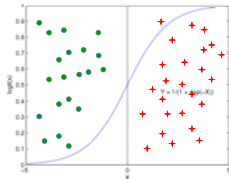
- Learn $P(Y|X)$ directly!
 - Assume a particular functional form
 - Not differentiable...**



14

Logistic Regression

- Learn $P(Y|X)$ directly!
 - Assume a particular functional form
 - Logistic Function
 - Aka Sigmoid
- $$\frac{1}{1 + \exp(-z)}$$

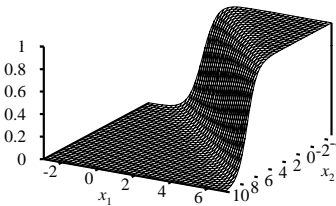


15

Logistic Function in n Dimensions

$$P(Y = 1|X) = \frac{1}{1 + \exp(w_0 + \sum_{i=1}^n w_i X_i)}$$

Sigmoid applied to a linear function of the data:



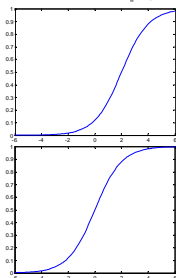
Features can be discrete or continuous!

16

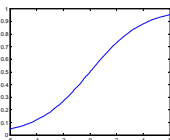
Understanding Sigmoids

$$g(w_0 + \sum_i w_i x_i) = \frac{1}{1 + \exp(w_0 + \sum_i w_i x_i)}$$

$w_0 = -2, w_1 = -1$



$w_0 = 0, w_1 = -1$



$w_0 = 0, w_1 = -0.5$

Very convenient!

$$P(Y = 1|X = \langle X_1, \dots, X_n \rangle) = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

implies

$$P(Y = 0|X = \langle X_1, \dots, X_n \rangle) = \frac{\exp(w_0 + \sum_i w_i X_i)}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

implies

$$\frac{P(Y = 0|X)}{P(Y = 1|X)} = \exp(w_0 + \sum_i w_i X_i)$$

implies

$$\ln \frac{P(Y = 0|X)}{P(Y = 1|X)} = w_0 + \sum_i w_i X_i$$

linear classification rule!

©Carlos Guestrin 2005-2009

Logistic regression more generally

Logistic regression in more general case, where $Y \in \{y_1, \dots, y_R\}$

for $k < R$

$$P(Y = y_k | X) = \frac{\exp(w_{k0} + \sum_{i=1}^n w_{ki} X_i)}{1 + \sum_{j=1}^{R-1} \exp(w_{j0} + \sum_{i=1}^n w_{ji} X_i)}$$

for $k=R$ (normalization, so no weights for this class)

$$P(Y = y_R | X) = \frac{1}{1 + \sum_{j=1}^{R-1} \exp(w_{j0} + \sum_{i=1}^n w_{ji} X_i)}$$

Features can be discrete or continuous!

©Carlos Guestrin 2005-2009

19

Loss Functions: Likelihood vs. Conditional Likelihood

Generative (Naïve Bayes) Loss function: Data likelihood

$$\begin{aligned} \ln P(\mathcal{D} | \mathbf{w}) &= \sum_{j=1}^N \ln P(\mathbf{x}^j, y^j | \mathbf{w}) \\ &= \sum_{j=1}^N \ln P(y^j | \mathbf{x}^j, \mathbf{w}) + \sum_{j=1}^N \ln P(\mathbf{x}^j | \mathbf{w}) \end{aligned}$$

Discriminative (Logistic Regr.) Loss function: Conditional Data Likelihood

$$\ln P(\mathcal{D}_Y | \mathcal{D}_X, \mathbf{w}) = \sum_{j=1}^N \ln P(y^j | \mathbf{x}^j, \mathbf{w})$$

Discriminative models *can't* compute $P(\mathbf{x}^j | \mathbf{w})!$

Or, ... "They don't waste effort learning $P(\mathbf{X})$ "

Focus only on $P(Y|X)$ - all that matters for classification

©Carlos Guestrin 2005-2009

Expressing Conditional Log Likelihood

$$\begin{aligned} P(Y=0|X, \mathbf{w}) &= \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)} \\ \ln P(Y=0|X, \mathbf{w}) &= -\ln(1 + \exp(w_0 + \sum_i w_i X_i)) \\ P(Y=1|X, \mathbf{w}) &= \frac{\exp(w_0 + \sum_i w_i X_i)}{1 + \exp(w_0 + \sum_i w_i X_i)} \\ \ln P(Y=1|X, \mathbf{w}) &= w_0 + \sum_i w_i X_i - \ln(1 + \exp(w_0 + \sum_i w_i X_i)) \end{aligned}$$

$$l(\mathbf{w}) = \sum_j y^j \ln P(y^j = 1 | \mathbf{x}^j, \mathbf{w}) + (1 - y^j) \ln P(y^j = 0 | \mathbf{x}^j, \mathbf{w})$$

1 when correct answer is 1

Probability of predicting 1

1 when correct answer is 0

Probability of predicting 0

©Carlos Guestrin 2005-2009

21

Expressing Conditional Log Likelihood

$$l(\mathbf{w}) = \sum_j \ln P(y^j | \mathbf{x}^j, \mathbf{w}) \quad \ln P(Y=0|X, \mathbf{w}) = -\ln(1 + \exp(w_0 + \sum_i w_i X_i))$$

$$\ln P(Y=1|X, \mathbf{w}) = w_0 + \sum_i w_i X_i - \ln(1 + \exp(w_0 + \sum_i w_i X_i))$$

$$l(\mathbf{w}) = \sum_j y^j \ln P(y^j = 1 | \mathbf{x}^j, \mathbf{w}) + (1 - y^j) \ln P(y^j = 0 | \mathbf{x}^j, \mathbf{w})$$

$$= \sum_j y^j (w_0 + \sum_i w_i x_i^j) - \ln(1 + \exp(w_0 + \sum_i w_i x_i^j))$$

Maximizing Conditional Log Likelihood

$$\begin{aligned} P(Y=0|X, \mathbf{w}) &= \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)} \\ P(Y=1|X, \mathbf{w}) &= \frac{\exp(w_0 + \sum_i w_i X_i)}{1 + \exp(w_0 + \sum_i w_i X_i)} \\ l(\mathbf{w}) &\equiv \ln \prod_j P(y^j | \mathbf{x}^j, \mathbf{w}) \\ &= \sum_j y^j (w_0 + \sum_i w_i x_i^j) - \ln(1 + \exp(w_0 + \sum_i w_i x_i^j)) \end{aligned}$$

Bad news: no closed-form solution to maximize $l(\mathbf{w})$

Good news: $l(\mathbf{w})$ is concave function of $\mathbf{w}!$

No local minima

Concave functions easy to optimize

©Carlos Guestrin 2005-2009

23

Optimizing Concave Functions Gradient Ascent

Conditional likelihood for Logistic Regression is concave!

Find optimum with gradient ascent

$$\text{Gradient: } \nabla_{\mathbf{w}} l(\mathbf{w}) = \left[\frac{\partial l(\mathbf{w})}{\partial w_0}, \dots, \frac{\partial l(\mathbf{w})}{\partial w_n} \right]'$$

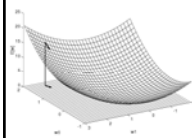
Learning rate, $\eta > 0$

$$\text{Update rule: } \Delta \mathbf{w} = \eta \nabla_{\mathbf{w}} l(\mathbf{w})$$

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \frac{\partial l(\mathbf{w})}{\partial w_i}$$

Gradient ascent is simplest of optimization approaches

e.g., Conjugate gradient ascent much better (see reading)



©Carlos Guestrin 2005-2009