# LASSO: Big Picture

Machine Learning – CSE446

Carlos Guestrin

University of Washington
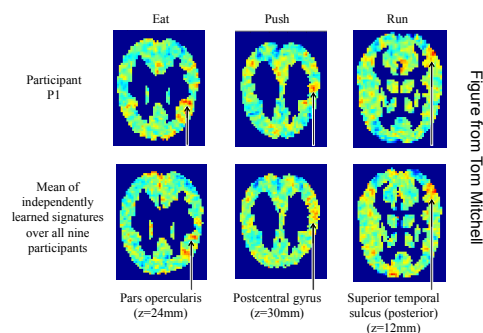
April 10, 2013

1

---

# Sparsity

- Vector **w** is sparse, if many entries are zero:

- Very useful for many tasks, e.g.,
  - □ **Efficiency**: If size(**w**) = 100B, each prediction is expensive:
    - If part of an online system, too slow
    - If **w** is sparse, prediction computation only depends on number of non-zeros
  - □ **Interpretability**: What are the relevant dimension to make a prediction?
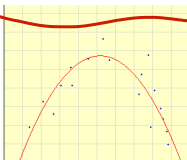    - E.g., what are the parts of the brain associated with particular words?



Eat    Push    Run

Participant P1

Mean of independently learned signatures over all nine participants

Pars opercularis (z=24mm)    Postcentral gyrus (z=30mm)    Superior temporal sulcus (posterior) (z=12mm)

Figure from Tom Mitchell
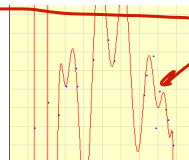
2

---

# Regularization in Linear Regression

- Overfitting usually leads to very large parameter choices, e.g.:

    $-2.2 + 3.1 X - 0.30 X^2$          $-1.1 + 4{,}700{,}910.7 X - 8{,}585{,}638.4 X^2 + \dots$

    *penalty for large weights*          *overfitting*

- ***Regularized*** or ***penalized*** regression aims to impose a "complexity" penalty by penalizing large weights
    - ☐ "Shrinkage" method

    $L_2$ *regularization* → *penalizes towards smoother functions*

©Carlos Guestrin 2005-2013                3

---

# LASSO Regression          $\lambda > 0$

- **LASSO:** least absolute shrinkage and selection operator

- New objective:

$$\min_{w} \; \sum_{j=1}^{N} \left( f(x_j) - \left( w_0 + \sum_{i} w_i h_i(x_j) \right) \right)^2 + \lambda \sum_{i=1}^{K} |w_i|$$

*don't regularize $w_0$*

©Carlos Guestrin 2005-2013                4

## Coordinate Descent for LASSO (aka Shooting Algorithm)

- **Repeat until convergence**
  - Pick a coordinate $j$ at (random or sequentially)    *round robin ?*
    - Set:
      $$\hat{w}_\ell = \begin{cases} (c_\ell + \lambda)/a_\ell & c_\ell < -\lambda \\ 0 & c_\ell \in [-\lambda, \lambda] \\ (c_\ell - \lambda)/a_\ell & c_\ell > \lambda \end{cases}$$
    - Where:
      $$a_\ell = 2\sum_{j=1}^{N}(h_\ell(\mathbf{x}_j))^2$$
      $$c_\ell = 2\sum_{j=1}^{N} h_\ell(\mathbf{x}_j)\left(t(\mathbf{x}_j) - (w_0 + \sum_{i\neq\ell} w_i h_i(\mathbf{x}_j))\right)$$

      *$w_0$ ??*
      *no regularization.*
      *$w_0 = c_0/a_0$*
      *$\Rightarrow$*
      *$w_0 = \frac{1}{N}\sum_{j=1}^{N}\left(t(x_j) - \sum_{i=1}^{K} w_i h_i(x_j)\right)$*
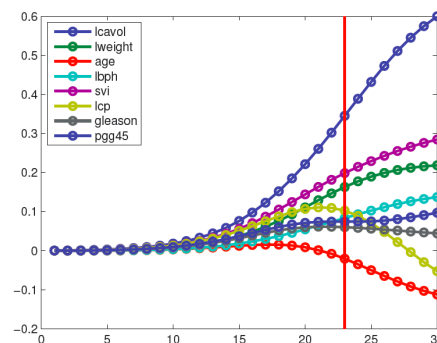
  - For convergence rates, see Shalev-Shwartz and Tewari 2009
- **Other common technique = LARS**
  - Least angle regression and shrinkage, Efron et al. 2004
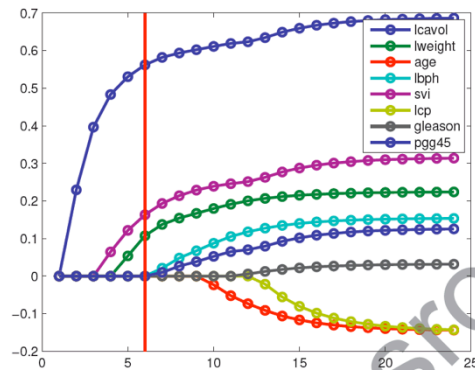
5

---

# Recall: *Ridge Coefficient Path*



From
Kevin Murphy
textbook

- Typical approach: select λ using cross validation

6

3

# Now: *LASSO Coefficient Path*



From
Kevin Murphy
textbook

# LASSO Example

| Term | Least Squares | Ridge | Lasso |
|---|---|---|---|
| Intercept | 2.465 | 2.452 | 2.468 |
| lcavol | 0.680 | 0.420 | 0.533 |
| lweight | 0.263 | 0.238 | 0.169 |
| age | $-0.141$ | $-0.046$ | |
| lbph | 0.210 | 0.162 | 0.002 |
| svi | 0.305 | 0.227 | 0.094 |
| lcp | $-0.288$ | 0.000 | |
| gleason | $-0.021$ | 0.040 | |
| pgg45 | 0.267 | 0.133 | |

From
Rob
Tibshirani
slides

# What you need to know

- Variable Selection: find a sparse solution to learning problem
- $L_1$ regularization is one way to do variable selection
  - Applies beyond regressions
  - Hundreds of other approaches out there
- LASSO objective non-differentiable, but convex ➔ Use subgradient
- No closed-form solution for minimization ➔ Use coordinate descent
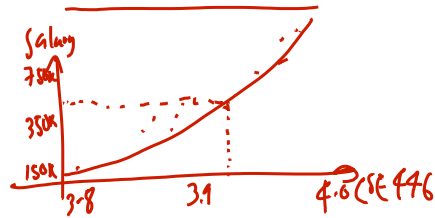- Shooting algorithm is very simple approach for solving LASSO

9

# Classification
# Logistic Regression

Machine Learning – CSE446

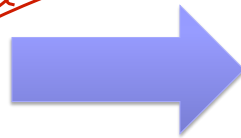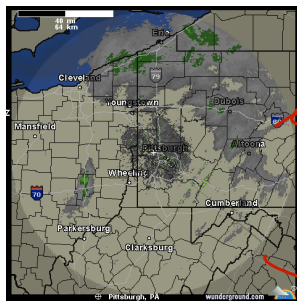Carlos Guestrin

University of Washington

April 15, 2013

10

# THUS FAR, REGRESSION: PREDICT A CONTINUOUS VALUE GIVEN SOME INPUTS

---

# Weather prediction revisted

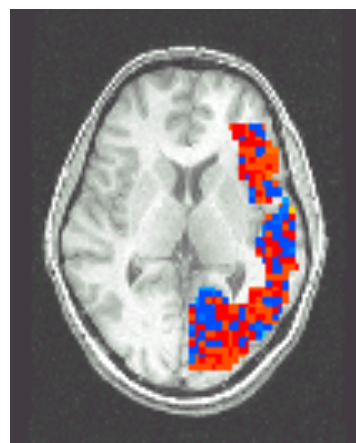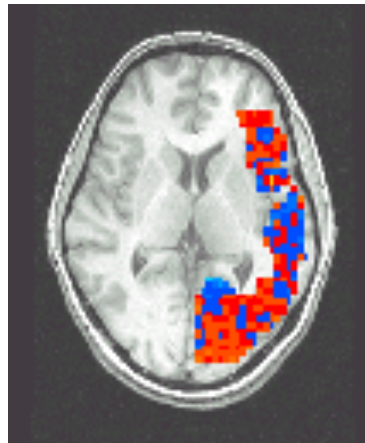## Reading Your Brain, Simple Example

[Mitchell et al.]

Pairwise classification accuracy: 85%

Person          Animal



©Carlos Guestrin 2005-2009          13

## Classification

- **Learn**: h:$\mathbf{X} \mapsto Y$
  - $\mathbf{X}$ – features
  - $Y$ – target classes

- Conditional probability: $P(Y|\mathbf{X})$

- Suppose you know $P(Y|\mathbf{X})$ exactly, how should you classify?
  - Bayes optimal classifier:

- **How do we estimate $P(Y|\mathbf{X})$?**

©Carlos Guestrin 2005-2013          14

7

# Link Functions

- Estimating P(Y|**X**): Why not use standard linear regression?



- Combing regression and probability?
  - Need a mapping from real values to [0,1]
  - A link function!

15

---
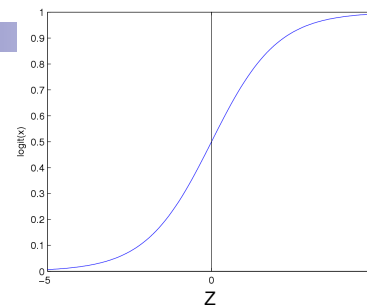
# Logistic Regression

**Logistic function (or Sigmoid):** $\dfrac{1}{1+exp(-z)}$



- Learn P(Y|**X**) directly
  - Assume a particular functional form for link function
  - Sigmoid applied to a linear function of the input features:

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

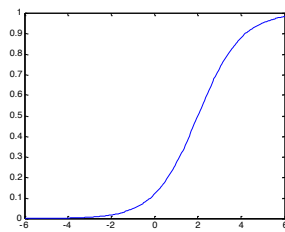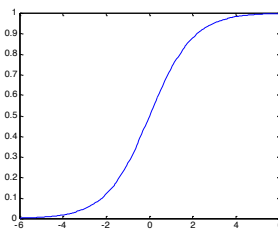**Features can be discrete or continuous!**

16

8

# Understanding the sigmoid

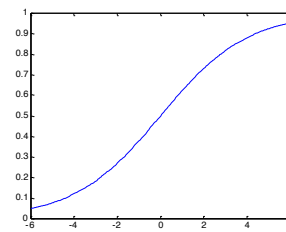$$g(w_0 + \sum_i w_i x_i) \;=\; \frac{1}{1 + e^{w_0 + \sum_i w_i x_i}}$$

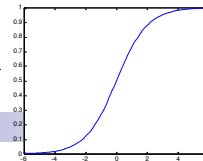$w_0$=-2, $w_1$=-1          $w_0$=0, $w_1$=-1          $w_0$=0, $w_1$=-0.5

**17**

# Logistic Regression – a Linear classifier

$$\frac{1}{1 + exp(-z)}$$

$$g(w_0 + \sum_i w_i x_i) \;=\; \frac{1}{1 + e^{w_0 + \sum_i w_i x_i}}$$

**18**

# Very convenient!

$$P(Y = 0 \mid X = < X_1, ... X_n >) = \frac{1}{1 + exp(w_0 + \sum_i w_i X_i)}$$

implies

$$P(Y = 1 \mid X = < X_1, ... X_n >) = \frac{exp(w_0 + \sum_i w_i X_i)}{1 + exp(w_0 + \sum_i w_i X_i)}$$

implies

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

linear classification rule!

implies

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

19

# Loss function: Conditional Likelihood

- Have a bunch of iid data of the form:

- Discriminative (logistic regression) loss function:
  **Conditional Data Likelihood**

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

20

# Expressing Conditional Log Likelihood

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

$$l(\mathbf{w}) \equiv \sum_j \ln P(y^j | \mathbf{x}^j, \mathbf{w})$$

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

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

---

# Maximizing Conditional Log Likelihood

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

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

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

Good news: $l(\mathbf{w})$ is concave function of $\mathbf{w}$, no local optima problems

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

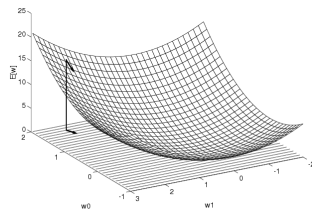Good news: concave functions easy to optimize

# Optimizing concave function – Gradient ascent

- Conditional likelihood for Logistic Regression is concave. Find optimum with gradient ascent

**Gradient:** $\nabla_{\mathbf{w}} l(\mathbf{w}) = [\frac{\partial l(\mathbf{w})}{\partial w_0}, \ldots, \frac{\partial l(\mathbf{w})}{\partial w_n}]'$

Step size, $\eta > 0$

**Update rule:** $\triangle \mathbf{w} = \eta \nabla_{\mathbf{w}} l(\mathbf{w})$

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

- Gradient ascent is simplest of optimization approaches
  - e.g., Conjugate gradient ascent can be much better

---

# Maximize Conditional Log Likelihood: Gradient ascent

$$l(\mathbf{w}) = \sum_j y^j (w_0 + \sum_i^n w_i x_i^j) - \ln(1 + exp(w_0 + \sum_i^n w_i x_i^j))$$

# Gradient Ascent for LR

Gradient ascent algorithm: iterate until change < ε

$$w_0^{(t+1)} \leftarrow w_0^{(t)} + \eta \sum_j [y^j - \hat{P}(Y^j = 1 \mid \mathbf{x}^j, \mathbf{w}^{(t)})]$$

For i=1,…,n,

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \sum_j x_i^j [y^j - \hat{P}(Y^j = 1 \mid \mathbf{x}^j, \mathbf{w}^{(t)})]$$
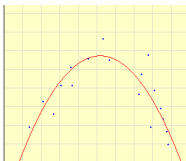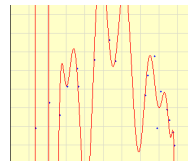
repeat

25

---

# Regularization in linear regression

- Overfitting usually leads to very large parameter choices, e.g.:

  -2.2 + 3.1 X – 0.30 X$^2$         -1.1 + 4,700,910.7 X – 8,585,638.4 X$^2$ + …

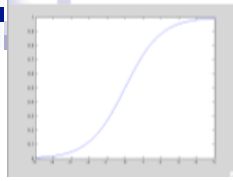- Regularized least-squares (a.k.a. ridge regression), for λ>0:

$$\mathbf{w}^* \;=\; \arg\min_{\mathbf{w}} \sum_j \left( t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2 + \lambda \sum_{i=1}^k w_i^2$$
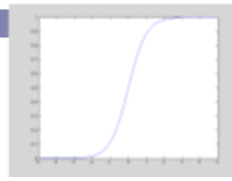
26

13

# Large parameters → Overfitting



$$\frac{1}{1 + e^{-x}}$$

$$\frac{1}{1 + e^{-2x}}$$

$$\frac{1}{1 + e^{-100x}}$$

- If data is linearly separable, weights go to infinity
- Leads to overfitting:

- Penalizing high weights can prevent overfitting…

27

---

# Regularized Conditional Log Likelihood

- Add regularization penalty, e.g., $L_2$:

$$\ell(\mathbf{w}) = \ln \prod_j P(y^j | \mathbf{x}^j, \mathbf{w})) - \lambda ||\mathbf{w}||_2^2$$

- Practical note about $w_0$:

- Gradient of regularized likelihood:

28

14

# Standard v. Regularized Updates

- Maximum conditional likelihood estimate

$$\mathbf{w}^* = \arg\max_{\mathbf{w}} \ln \left[ \prod_{j=1}^{N} P(y^j \mid \mathbf{x}^j, \mathbf{w}) \right]$$

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \sum_j x_i^j [y^j - \hat{P}(Y^j = 1 \mid \mathbf{x}^j, \mathbf{w}^{(t)})]$$

- Regularized maximum conditional likelihood estimate

$$\mathbf{w}^* = \arg\max_{\mathbf{w}} \ln \left[ \prod_j P(y^j \mid \mathbf{x}^j, \mathbf{w})) \right] - \lambda \sum_{i>0} w_i^2$$

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \left\{ -\lambda w_i^{(t)} + \sum_j x_i^j [y^j - \hat{P}(Y^j = 1 \mid \mathbf{x}^j, \mathbf{w}^{(t)})] \right\}$$

29

# Please Stop!! Stopping criterion

$$\ell(\mathbf{w}) = \ln \prod_j P(y^j \mid \mathbf{x}^j, \mathbf{w})) - \lambda ||\mathbf{w}||_2^2$$

- When do we stop doing gradient descent?

- Because $l(\mathbf{w})$ is strongly concave:
  - □ i.e., because of some technical condition

$$\ell(\mathbf{w}^*) - \ell(\mathbf{w}) \leq \frac{1}{2\lambda} ||\nabla \ell(\mathbf{w})||_2^2$$

- Thus, stop when:

30

15

# Stopping criterion

$$\ell(\mathbf{w}) = \ln \prod_j P(y^j | \mathbf{x}^j, \mathbf{w})) - \lambda ||\mathbf{w}||_2^2$$

- Regularized logistic regression is strongly concave
    - Negative second derivative bounded away from zero:

- Strong concavity (convexity) is super helpful!!

- For example, for strongly concave *l*(**w**):

$$\ell(\mathbf{w}^*) - \ell(\mathbf{w}) \le \frac{1}{2\lambda} ||\nabla \ell(\mathbf{w})||_2^2$$

---

# Convergence rates for gradient descent/ascent

- Number of Iterations to get to accuracy

$$\ell(\mathbf{w}^*) - \ell(\mathbf{w}) \le \epsilon$$

- If func Lipschitz: $O(1/\epsilon^2)$

- If gradient of func Lipschitz: $O(1/\epsilon)$

- If func is strongly convex: $O(\ln(1/\epsilon))$

# Digression: Logistic regression for more than 2 classes

- Logistic regression in more general case (k+1 classes), where *Y in* $\{y_1,\ldots,y_R\}$

# Digression: Logistic regression more generally

- Logistic regression in more general case, where $Y$ *in* $\{y_1,\ldots,y_R\}$

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

**Learning procedure is basically the same as what we derived!**

# What you should know…

- Classification: predict discrete classes rather than real values
- Logistic regression model: Linear model
  - Logistic function maps real values to [0,1]
- Optimize conditional likelihood
- Gradient computation
- Overfitting
- Regularization
- Regularized optimization

35