

Supervised Learning

Regression and Classification

Shan-Hung Wu
shwu@cs.nthu.edu.tw

Department of Computer Science,
National Tsing Hua University, Taiwan

NetDB-ML, Spring 2015

1 Regression

- Linear Regression
- Interpolation vs. Regression
- Probability Interpretation

2 Two-Class Classification

- Logistic Regression
- Perceptron

3 Multiclass Classification

- Wrapper Methods
- Direct Models

4 Non-Parametric Methods

1 Regression

- Linear Regression
- Interpolation vs. Regression
- Probability Interpretation

2 Two-Class Classification

- Logistic Regression
- Perceptron

3 Multiclass Classification

- Wrapper Methods
- Direct Models

4 Non-Parametric Methods

The Regression Problem

- Given

- a **training dataset** $\mathcal{X} = \{(\mathbf{x}^{(t)}, r^{(t)})\}_{t=1}^N$, where $\mathbf{x}^{(t)} \in \mathbb{R}^d$'s are **examples** (or **instances** or **observations**) consisting of **attributes** (or **inputs** or **features**) and $r^{(t)} \in \mathbb{R}$'s are **labels**, and
- a **testing instance** \mathbf{x}' ,

predict the label y' of \mathbf{x}'

- Example: stock price forecasting

ML Process Revisited

- 1 Data collection and preprocessing (e.g., integration, cleaning, etc.)
- 2 Model development
 - 1 Assume a **model** that represents the posteriori knowledge we want to discover. The model has parameters
 - 2 Define an **objective** that measures “how good the model with a particular combination of parameters can explain the data”
- 3 **Training**: employ an algorithm that optimizes the objective by finding the best (or good enough) parameters
- 4 **Testing**: evaluate the model performance on hold-out data
- 5 Using the model

Modeling a Regressor

- Model: Let the model be a collection of functions, called **hypothesis class** and denoted as $\mathcal{H} = \{h: \mathcal{J} \times \Theta \rightarrow \mathbb{R}\}$, where \mathcal{J} is the **input space** (or **feature space**) and Θ is the set of all possible parameters
 - A particular $\theta \in \Theta$ instantiates a **hypothesis** h that makes the **prediction** (or **output**) $y' = h(\mathbf{x}'; \theta) > 0$
- Objective: $\arg_{\theta} \min \sum_{t=1}^N l(h(\mathbf{x}^{(t)}; \theta), r^{(t)})$, where l is some **loss function** which penalizes the error of predictions made on the training dataset
 - We want the hypothesis to have the minimal **empirical error**:
$$emp(\theta; \mathcal{X}) = \sum_{t=1}^N l(h(\mathbf{x}^{(t)}; \theta), r^{(t)})$$

The Objective

- Common choice: $\arg_{\theta} \min \sum_{t=1}^N [r^{(t)} - h(\mathbf{x}^{(t)}; \theta)]^2$
 - $emp(\theta; \mathcal{X}) = \sum_{t=1}^N [r^{(t)} - h(\mathbf{x}^{(t)}; \theta)]^2$ has a specific name called the **Sum of Square Errors (SSE)**
- Alternatively, the objective can be formed using the absolute error: $\arg_{\theta} \min \sum_{t=1}^N |r^{(t)} - h(\mathbf{x}^{(t)}; \theta)|$
 - What is the difference? [Homework]

1 Regression

- Linear Regression
- Interpolation vs. Regression
- Probability Interpretation

2 Two-Class Classification

- Logistic Regression
- Perceptron

3 Multiclass Classification

- Wrapper Methods
- Direct Models

4 Non-Parametric Methods

- Suppose x is a scalar and h is a line, i.e., $h(x; \theta) = w_1x + w_0$, we have the objective:

- To find w_0 and w_1 that minimizes

$$emp(\theta; \mathcal{X}) = \sum_{t=1}^N \left(r^{(t)} - \mathbf{w}^\top \begin{bmatrix} 1 \\ x^{(t)} \end{bmatrix} \right)^2,$$

- where $\mathbf{w} = [w_0, w_1]^\top \in \mathbb{R}^2$

Training: Analytic Solution (1)

- We take the partial derivatives of *emp* with respect to w_0 and w_1 and set them to 0
 - We have a system of linear equations

$$\begin{cases} \sum_{t=1}^N r^{(t)} = Nw_0 + w_1 \sum_{t=1}^N x^{(t)} \\ \sum_{t=1}^N x^{(t)} r^{(t)} = w_0 \sum_{t=1}^N x^{(t)} + w_1 \sum_{t=1}^N (x^{(t)})^2 \end{cases}$$

- Let $\mathbf{A} = \begin{bmatrix} N & \sum_{t=1}^N x^{(t)} \\ \sum_{t=1}^N x^{(t)} & \sum_{t=1}^N (x^{(t)})^2 \end{bmatrix}$, $\mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$, and $\mathbf{y} = \begin{bmatrix} \sum_{t=1}^N r^{(t)} \\ \sum_{t=1}^N r^{(t)} x^{(t)} \end{bmatrix}$, we can solve \mathbf{w} by $\mathbf{w} = \mathbf{A}^{-1} \mathbf{y}$

Training: Analytic Solution (2)

- A bit of arithmetic leads to

$$\begin{cases} w_0 = \bar{r} - w_1 \bar{x} \\ w_1 = \left(\sum_{t=1}^N x^{(t)} r^{(t)} - \bar{x} \bar{r} N \right) / \left(\sum_{t=1}^N (x^{(t)})^2 - N \bar{x}^2 \right) \end{cases} ,$$

where $\bar{x} = \frac{1}{N} \sum_{t=1}^N x^{(t)}$ and $\bar{r} = \frac{1}{N} \sum_{t=1}^N r^{(t)}$ [Proof]

Multivariate Linear Regression

- Given $\mathbf{x} \in \mathbb{R}^d$, suppose h is linear: $h(\mathbf{x}; \theta) = \mathbf{w}^\top \begin{bmatrix} 1 \\ \mathbf{x}^{(t)} \end{bmatrix}$, where $\mathbf{w} = [w_0, w_1, \dots, w_d]^\top \in \mathbb{R}^{d+1}$

- We can solve \mathbf{w} by $\mathbf{w} = \mathbf{A}^{-1}\mathbf{y}$, where $\mathbf{y} = \begin{bmatrix} \sum_{t=1}^N r^{(t)} \\ \sum_{t=1}^N r^{(t)} x_1^{(t)} \\ \vdots \\ \sum_{t=1}^N r^{(t)} x_d^{(t)} \end{bmatrix}$ and

$$\mathbf{A} = \begin{bmatrix} N & \sum_{t=1}^N x_1^{(t)} & \cdots & \sum_{t=1}^N x_d^{(t)} \\ \sum_{t=1}^N x_1^{(t)} & \sum_{t=1}^N x_1^{(t)2} & \cdots & \sum_{t=1}^N x_1^{(t)} x_d^{(t)} \\ \vdots & \vdots & \ddots & \vdots \\ \sum_{t=1}^N x_d^{(t)} & \sum_{t=1}^N x_d^{(t)} x_1^{(t)} & \cdots & \sum_{t=1}^N x_d^{(t)2} \end{bmatrix} \quad [\text{Proof}]$$

From Least Squares to Linear Regression

- Let $\mathbf{X} = \begin{bmatrix} 1 & x_1^{(1)} & \cdots & x_d^{(1)} \\ 1 & x_1^{(2)} & \cdots & x_d^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(N)} & \cdots & x_d^{(N)} \end{bmatrix}$, $\mathbf{w} = [w_0, w_1, \dots, w_d]^\top$, and $\mathbf{r} = [r^{(1)}, r^{(2)}, \dots, r^{(N)}]^\top$.

- Ideally, we want to solve \mathbf{w} such that $\mathbf{X}\mathbf{w} = \mathbf{r}$, but impossible if $N > d$
- We can instead solve the “closest approximation:” $\arg \min_{\mathbf{w}} \|\mathbf{r} - \mathbf{X}\mathbf{w}\|^2$
 - $\|\mathbf{r} - \mathbf{X}\mathbf{w}\|^2$ is exactly the SSE!
- The **least square problem**: find \mathbf{w} such that $\|\mathbf{r} - \mathbf{X}\mathbf{w}\|^2$ is minimized. Solution?

From Least Squares to Linear Regression

- Let $\mathbf{X} = \begin{bmatrix} 1 & x_1^{(1)} & \cdots & x_d^{(1)} \\ 1 & x_1^{(2)} & \cdots & x_d^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(N)} & \cdots & x_d^{(N)} \end{bmatrix}$, $\mathbf{w} = [w_0, w_1, \dots, w_d]^\top$, and $\mathbf{r} = [r^{(1)}, r^{(2)}, \dots, r^{(N)}]^\top$.

- Ideally, we want to solve \mathbf{w} such that $\mathbf{X}\mathbf{w} = \mathbf{r}$, but impossible if $N > d$
- We can instead solve the “closest approximation:” $\arg \min_{\mathbf{w}} \|\mathbf{r} - \mathbf{X}\mathbf{w}\|^2$
 - $\|\mathbf{r} - \mathbf{X}\mathbf{w}\|^2$ is exactly the SSE!
- The **least square problem**: find \mathbf{w} such that $\|\mathbf{r} - \mathbf{X}\mathbf{w}\|^2$ is minimized. Solution?
 - $\mathbf{w}^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{r}$ if \mathbf{X} is full column rank (remember the normal equations?)
 - $(\mathbf{X}^\top \mathbf{X})^{-1}$ and $\mathbf{X}^\top \mathbf{r}$ are exactly \mathbf{A}^{-1} and \mathbf{y} seen previously

Analytic Solution Revisited

- What if X is not full column rank?

Analytic Solution Revisited

- What if \mathbf{X} is not full column rank?
- ① Anyone in the set $\mathbf{X}^\dagger \mathbf{r} + \mathcal{N}(\mathbf{X})$ is the solution (remember the SVD solution to least squares?)
- ② Make \mathbf{X} full column rank by changing the objective (to be explained later)

Training: Numeric Methods

- Machine learning solutions need not be accurate
 - Close-to-optimal solutions enough for making good predictions
- Numeric methods suffice
 - E.g., gradient descent:

Repeat until convergence {

$$\mathbf{w} := \mathbf{w} - \eta \nabla emp(\mathbf{w}; \mathcal{X}) = \mathbf{w} + 2\eta \sum_{t=1}^N (r^{(t)} - \mathbf{w}^\top \begin{bmatrix} 1 \\ \mathbf{x}^{(t)} \end{bmatrix}) \begin{bmatrix} 1 \\ \mathbf{x}^{(t)} \end{bmatrix};$$

}

- The step size η is called the *learning rate*

- 1 Regression**
 - Linear Regression
 - Interpolation vs. Regression
 - Probability Interpretation
- 2 Two-Class Classification**
 - Logistic Regression
 - Perceptron
- 3 Multiclass Classification**
 - Wrapper Methods
 - Direct Models
- 4 Non-Parametric Methods**

Interpolation vs. Regression (1)

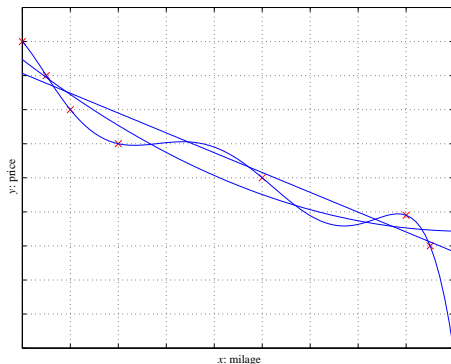
- Instead of regression, we can perform the **interpolation** that fits a hypothesis $h: \mathbb{R} \times \Theta \rightarrow \mathbb{R}$ to examples, i.e., $h(x^{(t)}; \theta) = r^{(t)}$
 - In polynomial interpolation, we can always fit a polynomial of degree $(N-1)$ to N 1-D points
 - Let $\theta = (w_0, \dots, w_{N-1})$ and $h(x; \theta) = w_0 + w_1x + \dots + w_{N-1}x^{N-1}$, $x \in \mathbb{R}$
 - Obtain θ by solving

$$\begin{bmatrix} (x^{(1)})^0 & \dots & (x^{(1)})^{N-1} \\ \vdots & \ddots & \vdots \\ (x^{(N)})^0 & \dots & (x^{(N)})^{N-1} \end{bmatrix} \begin{bmatrix} w_0 \\ \vdots \\ w_{N-1} \end{bmatrix} = \begin{bmatrix} r^{(1)} \\ \vdots \\ r^{(N)} \end{bmatrix}$$

- The label of a new instance x' can be predicted by $y' = h(x'; \theta)$

Interpolation vs. Regression (2)

- Given 7 examples, the right shows the regression results using polynomials of degrees 1, 2, and 6
 - $x^{(t)}$ is the mileage of a used car and $r^{(t)}$ is the price
- It is unlikely that the real curve shapes like the 6th-order polynomial



Interpolation vs. Regression (3)

- In the presence of noise, we don't need an exact fitting
- The target of regression is to catch the trend
 - Differs from interpolation in finding a “simple” hypothesis (e.g., low degree polynomial) that is “close enough” to the examples

How About Nonlinear Trend/Regression? (1)

- In the case of univariate regression (where x 's are scalars), we can assume a polynomial hypothesis with an arbitrary degree k :

$$h(x; \theta) = w_0 + w_1x + \dots + w_kx^k,$$

- We can solve $\mathbf{w} = \begin{bmatrix} w_0 \\ \vdots \\ w_k \end{bmatrix}$ by $\mathbf{w} = \mathbf{A}^{-1}\mathbf{y}$, where

$$\mathbf{A} = \begin{bmatrix} \sum_{t=1}^N x^{(t)} & \sum_{t=1}^N x^{(t)2} & \dots & \sum_{t=1}^N x^{(t)k} \\ \vdots & \vdots & \ddots & \vdots \\ \sum_{t=1}^N x^{(t)k} & \sum_{t=1}^N x^{(t)(k+1)} & \dots & \sum_{t=1}^N x^{(t)2k} \end{bmatrix},$$
$$\mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_k \end{bmatrix}, \text{ and } \mathbf{y} = \begin{bmatrix} \sum_{t=1}^N r^{(t)} \\ \sum_{t=1}^N r^{(t)}x^{(t)} \\ \vdots \\ \sum_{t=1}^N r^{(t)}x^{(t)k} \end{bmatrix} \quad [\text{Proof}]$$

How About Nonlinear Trend/Regression? (2)

- In multivariate regression, we seldom assume h to be a polynomial with degree higher than 1
 - Why?

How About Nonlinear Trend/Regression? (2)

- In multivariate regression, we seldom assume h to be a polynomial with degree higher than 1
 - Why?
- ① Analytical simplicity
- ② More descriptive model:
 - The sign of w_j tells whether x_j has positive or negative effect on the prediction
 - The absolute value of w_j indicates how important the feature is (provided that features are in the same range); if w_j is close to 0, the feature can even be removed
- ③ We can instead augment the inputs to achieve the effect of nonlinear regression (to be explained later)

1 Regression

- Linear Regression
- Interpolation vs. Regression
- Probability Interpretation

2 Two-Class Classification

- Logistic Regression
- Perceptron

3 Multiclass Classification

- Wrapper Methods
- Direct Models

4 Non-Parametric Methods

Probability Interpretation (1)

- Given $\mathcal{X} = \{\mathbf{x}^{(t)}, r^{(t)}\}_{t=1}^N$, where $r^{(t)} \in \mathbb{R}$. Assume
 - $(\mathbf{x}^{(t)}, r^{(t)})$ are i.i.d samples drawn from some joint distribution of \mathbf{x} and r (otherwise can never learn r from \mathbf{x})
 - In particular, $r^{(t)} = f(\mathbf{x}^{(t)}; \theta) + \epsilon$, $\epsilon \sim \mathcal{N}(0, \beta^{-1})$ for some **hyperparameter** (i.e., constant fixed during the objective solving) β
 - The marginal distribution $p(r|\mathbf{x})$ follows: $p(r|\mathbf{x}) = p_{N_{h(\mathbf{x};\theta), \beta^{-1}}}(r)$
- We want to estimate f using \mathcal{X}
 - Hypothesis: $h(\mathbf{x}; w_0, w_1, \dots, w_d) = w_0 + w_1x_1 + \dots + w_dx_d$, a line
 - Once getting w_0, w_1, \dots, w_d , we can predict the unknown r' of a new instance \mathbf{x}' by $y' = \arg_y \max p(y|\mathbf{x}') = \arg_y \max p_{N_{h(\mathbf{x}';\theta), \beta^{-1}}}(y) = h(\mathbf{x}'; \theta)$
 - Note that we don't need to know β to make prediction

Probability Interpretation (2)

- How to obtain the estimate h of f ? How to obtain θ ?
- Now let θ be a random variable, we can pick θ maximizing $p(\theta|\mathcal{X})$, the **posterior** probability
- Or, by Baye's theorem, θ maximizing the **likelihood** $p(\mathcal{X}|\theta)$ (if we assume $p(\theta)$ remains the same for all θ)
- Or, θ maximizing the **log likelihood** $\log p(\mathcal{X}|\theta) = \log \left(\prod_{t=1}^N p(\mathbf{x}^{(t)}, r^{(t)}|\theta) \right) = \log \left(\prod_{t=1}^N p(r^{(t)}|\mathbf{x}^{(t)}, \theta) p(\mathbf{x}^{(t)}|\theta) \right) = \log \left(\prod_{t=1}^N p(h(\mathbf{x}^{(t)}; \theta) + \epsilon|\mathbf{x}^{(t)}, \theta) p(\mathbf{x}^{(t)}|\theta) \right)$
- Ignoring $p(\mathbf{x}^{(t)}|\theta) = p(\mathbf{x}^{(t)})$ (since it is irrelevant to θ) and constants we have $\log p(\mathcal{X}|\theta) \propto -N \log \left(\sqrt{\frac{2\pi}{\beta}} \right) - \frac{\beta}{2} \sum_{t=1}^N (r^{(t)} - h(\mathbf{x}^{(t)}; \theta))^2$
- Dropping the first term and constants we have $\log p(\mathcal{X}|\theta) \propto - \sum_{t=1}^N (r^{(t)} - h(\mathbf{x}^{(t)}; \theta))^2$; that is, we seek for θ minimizing the SSE (sum of square errors)!

1 Regression

- Linear Regression
- Interpolation vs. Regression
- Probability Interpretation

2 Two-Class Classification

- Logistic Regression
- Perceptron

3 Multiclass Classification

- Wrapper Methods
- Direct Models

4 Non-Parametric Methods

Two-Class Classification Problem

- Given a **training dataset** $\mathcal{X} = \{(\mathbf{x}^{(t)}, r^{(t)})\}_{t=1}^N$, where $r^{(t)} \in \{1, -1\}$, and a testing instance \mathbf{x}' , predict the label of \mathbf{x}'
- Model (or **hypothesis class**): $\mathcal{H} = \{h : \mathcal{J} \times \Theta \rightarrow \{1, -1\}\}$
 - Or $\mathcal{H} = \{h : \mathcal{J} \times \Theta \rightarrow \mathbb{R}\}$ with prediction $\text{sgn}(h(\mathbf{x}'; \theta))$
- Objective: $\arg_{\theta} \min \sum_{t=1}^N l(h(\mathbf{x}^{(t)}; \theta), r^{(t)})$ with some loss function l
 - Example: the **0-1 loss function**: $l(a, b) = 1$ if $a \neq b$; 0 otherwise

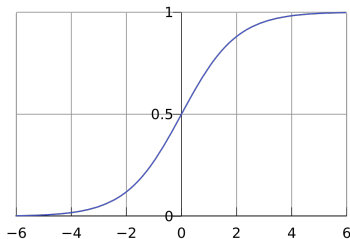
- 1 **Regression**
 - Linear Regression
 - Interpolation vs. Regression
 - Probability Interpretation
- 2 **Two-Class Classification**
 - Logistic Regression
 - Perceptron
- 3 **Multiclass Classification**
 - Wrapper Methods
 - Direct Models
- 4 **Non-Parametric Methods**

Logistic Function

- The **logistic function** (a special case of **sigmoid functions**) is defined as

$$g(z) = \frac{e^z}{e^z + 1} = \frac{1}{1 + e^{-z}}$$

- Always gives values between $(0, 1)$
- The larger the z , the higher the $g(z)$
- The smaller the z , the higher the $1 - g(z)$



Logistic Regression

- In regression, we learn $p(r|\mathbf{x}; \theta)$ from \mathcal{X} and make predictions by $y' = \arg \max_y p(y|\mathbf{x}'; \theta)$
- In **logistic regression** everything is the same except that $P(r|\mathbf{x}; \theta)$ is modeled by a Bernoulli distribution parametrized by ϕ :

$$P(r|\mathbf{x}; \theta) = \begin{cases} \phi, & \text{if } r = 1, \\ 1 - \phi, & \text{otherwise,} \end{cases}$$

- We can simply write $P(r|\mathbf{x}; \theta) = \phi^q(1 - \phi)^{(1-q)}$, where $q = \frac{r+1}{2}$
- Furthermore, $\phi = \pi(\mathbf{x}; \boldsymbol{\beta}) = \frac{e^{\boldsymbol{\beta}^\top \tilde{\mathbf{x}}}}{e^{\boldsymbol{\beta}^\top \tilde{\mathbf{x}}} + 1} = \frac{1}{1 + e^{-\boldsymbol{\beta}^\top \tilde{\mathbf{x}}}}$ is a deterministic function, where $\tilde{\mathbf{x}} = [1, \mathbf{x}]^\top$
 - So the larger the projection of $\tilde{\mathbf{x}}$ onto a line, the higher the ϕ
- Prediction: $y' = \arg \max_y p(y|\mathbf{x}'; \theta) = \arg \max_y \{\phi, 1 - \phi\} = \text{sgn}(\boldsymbol{\beta}^\top \tilde{\mathbf{x}}') = \text{sgn}(\mathbf{w}^\top \mathbf{x}' + b)$

Fitting Logistic Regression Models (1)

- How to obtain β ?

Fitting Logistic Regression Models (1)

- How to obtain β ?
 - By β maximizing $p(\beta|\mathcal{X})$
 - Or, by Bayes' Rule and assuming uniform $p(\beta)$, β maximizing $p(\mathcal{X}|\beta)$
- Log-likelihood:

$$\begin{aligned}l(\beta) &= \log \prod_{t=1}^N p(\mathbf{x}^{(t)}, r^{(t)}|\beta) \\ &= \log \prod_{t=1}^N P(r^{(t)}|\mathbf{x}^{(t)}, \beta) p(\mathbf{x}^{(t)}|\beta) \\ &\propto \log \prod_{t=1}^N \pi(\mathbf{x}^{(t)}; \beta)^{q^{(t)}} (1 - \pi(\mathbf{x}^{(t)}; \beta))^{(1-q^{(t)})}\end{aligned}$$

- $p(\mathbf{x}^{(t)}|\beta) = p(\mathbf{x}^{(t)})$ can be dropped
- $l(\beta) = \sum_{t=1}^N \{q^{(t)} \log \pi(\mathbf{x}^{(t)}; \beta) + (1 - q^{(t)}) \log (1 - \pi(\mathbf{x}^{(t)}; \beta))\} = \sum_{t=1}^N \left\{ q^{(t)} \beta^\top \tilde{\mathbf{x}}^{(t)} - \log \left(1 + e^{\beta^\top \tilde{\mathbf{x}}^{(t)}} \right) \right\}$ [Homework]

Fitting Logistic Regression Models (2)

- To maximize the log-likelihood, we set its derivative to zero:

$$\frac{\partial l(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = \sum_{t=1}^N \tilde{\mathbf{x}}^{(t)\top} \left(q^{(t)} - \pi(\mathbf{x}^{(t)}; \boldsymbol{\beta}) \right) = \mathbf{0}^\top$$

- Unlike the linear regression, we cannot solve $\boldsymbol{\beta}$ analytically in a closed-form
- How to obtain $\boldsymbol{\beta}$?

Fitting Logistic Regression Models (2)

- To maximize the log-likelihood, we set its derivative to zero:

$$\frac{\partial l(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = \sum_{t=1}^N \tilde{\mathbf{x}}^{(t)\top} \left(q^{(t)} - \pi(\mathbf{x}^{(t)}; \boldsymbol{\beta}) \right) = \mathbf{0}^\top$$

- Unlike the linear regression, we cannot solve $\boldsymbol{\beta}$ analytically in a closed-form
- How to obtain $\boldsymbol{\beta}$? Iterative algorithms
- Gradient descent:

Repeat until convergence {
 $\boldsymbol{\beta} := \boldsymbol{\beta} + \eta \nabla l(\boldsymbol{\beta}) = \boldsymbol{\beta} + \eta \sum_{t=1}^N \tilde{\mathbf{x}}^{(t)\top} (q^{(t)} - \pi(\mathbf{x}^{(t)}; \boldsymbol{\beta}))$;
}

- Observe that $l(\boldsymbol{\beta})$ is concave [Homework]
 - So iterative algorithms approach to global optimal

Newton's Method for Logistic Regression*

Update rule: $\boldsymbol{\beta} := \boldsymbol{\beta} - (\nabla^2 l(\boldsymbol{\beta}))^{-1} \nabla l(\boldsymbol{\beta}) = \boldsymbol{\beta} - (\nabla^2 l(\boldsymbol{\beta}))^{-1} \nabla l(\boldsymbol{\beta})$

- Given $\mathbf{q} \in \mathbb{R}^N$ the vector of $q^{(t)}$'s, $\mathbf{X} \in \mathbb{R}^{N \times (d+1)}$ the row matrix of $\tilde{\mathbf{x}}^{(t)}$'s, $\boldsymbol{\pi} \in \mathbb{R}^N$ with the t th element $\pi(\mathbf{x}^{(t)}; \boldsymbol{\beta})$, and $\mathbf{W} \in \mathbb{R}^{N \times N}$ a diagonal matrix with the t th diagonal element $\pi(\mathbf{x}^{(t)}; \boldsymbol{\beta})(1 - \pi(\mathbf{x}^{(t)}; \boldsymbol{\beta}))$, then

$$\nabla l(\boldsymbol{\beta}) = \left(\frac{\partial l(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \right)^\top = \mathbf{X}^\top (\mathbf{q} - \boldsymbol{\pi}),$$

$$\begin{aligned} \nabla^2 l(\boldsymbol{\beta}) &= \left(\frac{\partial \nabla l(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \right)^\top = \left(\frac{\partial \sum_{t=1}^N \tilde{\mathbf{x}}^{(t)} (q^{(t)} - \pi(\mathbf{x}^{(t)}; \boldsymbol{\beta}))}{\partial \boldsymbol{\beta}} \right)^\top \\ &= \left(- \sum_{t=1}^N \tilde{\mathbf{x}}^{(t)} \pi(\mathbf{x}^{(t)}; \boldsymbol{\beta})(1 - \pi(\mathbf{x}^{(t)}; \boldsymbol{\beta})) \mathbf{x}^{(t)\top} \right)^\top = -\mathbf{X}^\top \mathbf{W} \mathbf{X} \end{aligned}$$

- Note that $g'(z) = \frac{1}{(1+e^{-z})^2} e^{-z} = \frac{1}{1+e^{-z}} \left(1 - \frac{1}{1+e^{-z}} \right) = g(z)(1 - g(z))$.

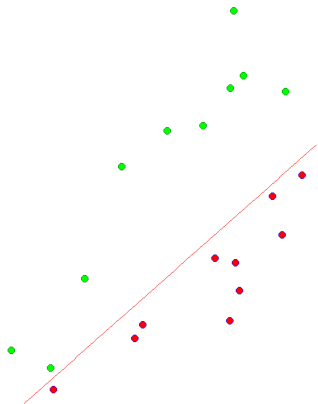
- 1 **Regression**
 - Linear Regression
 - Interpolation vs. Regression
 - Probability Interpretation
- 2 **Two-Class Classification**
 - Logistic Regression
 - Perceptron
- 3 **Multiclass Classification**
 - Wrapper Methods
 - Direct Models
- 4 **Non-Parametric Methods**

Perceptron (1)

- Recall that in logistic regression, we make prediction by $y' = \arg \max_y p(y|\mathbf{x}'; \theta) = \arg \max_y \{\phi, 1 - \phi\} = \text{sgn}(\boldsymbol{\beta}^\top \tilde{\mathbf{x}}') = \text{sgn}(\mathbf{w}^\top \mathbf{x}' + b)$
- Why not just making prediction based on $\text{sgn}(\mathbf{w}^\top \mathbf{x}' + b)$ directly?

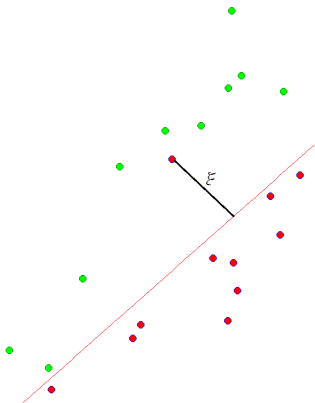
Perceptron (2)

- Model: $\mathcal{H} = \{f: f: \mathbb{R}^d \rightarrow \mathbb{R}, f(\mathbf{x}; \theta) = \mathbf{w}^\top \mathbf{x} + b\}$
 $\Theta = \{\mathbf{w}, b: \mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}\}$
 - A collection of hyperplanes
- Prediction: $y' = \text{sgn}(f(\mathbf{x}'))$
- Objective 1: any $f \in \mathcal{H}$ such that
 - $\mathbf{w}^\top \mathbf{x}^{(t)} + b > 0$ if $r^{(t)} = 1$;
 $\mathbf{w}^\top \mathbf{x}^{(t)} + b < 0$ otherwise
 - or simply $r^{(t)}(\mathbf{w}^\top \mathbf{x}^{(t)} + b) > 0, \forall t$



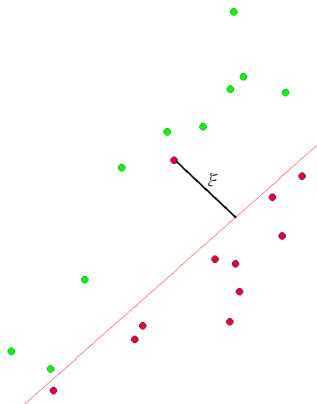
Non Separable Datasets

- What if the examples are *not* separable by a hyperplane?



Non Separable Datasets

- What if the examples are *not* separable by a hyperplane?



- Don't insist perfect separation as in Objective 1

- Objective 2:

$$\arg \min_{\mathbf{w}, b, \xi} \sum_{t=1}^N \xi_t$$

subject to $r^{(t)}(\mathbf{w}^\top \mathbf{x}^{(t)} + b) > -\xi_t$ and $\xi_t \geq 0, \forall t = 1, \dots, N$

- ξ_t 's are called the **slacks**
- We minimize $\sum_{t=1}^N \xi_t$ instead of $\sum_{t=1}^N \xi_t^2$ to make the hypothesis robust to outliers
- Alternative form: $\arg \min_{\mathbf{w}, b} \sum_{t=1}^N \max(0, -r^{(t)}(\mathbf{w}^\top \mathbf{x}^{(t)} + b))$
 - No slack to solve, no constraint, **convex**
- $l(h(\mathbf{x}^{(t)}; \theta), r^{(t)}) := \max(0, -r^{(t)}(\mathbf{w}^\top \mathbf{x}^{(t)} + b))$ is called the **hinge loss** function (why?)
 - $emp(\theta; \mathcal{X}) = \sum_{t=1}^N l(h(\mathbf{x}^{(t)}; \theta), r^{(t)})$

Training the Perceptron Classifier

- Let $\tilde{\mathbf{x}}^{(t)} = \begin{bmatrix} \mathbf{x}^{(t)} \\ 1 \end{bmatrix}$ and $\tilde{\mathbf{w}}^{(t)} = \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix}$, we need to solve

$$\arg \min_{\tilde{\mathbf{w}}} emp(\tilde{\mathbf{w}}) := \sum_{t=1}^N \max(0, -r^{(t)} \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}^{(t)})$$

- Let's consider the Gradient descent method due to its simplicity

- $\nabla emp(\tilde{\mathbf{w}}) = \sum_{t=1}^N \nabla l^{(t)}(\tilde{\mathbf{w}})$, where

$$\nabla l^{(t)}(\tilde{\mathbf{w}}) = \begin{cases} 0, & \text{if } r^{(t)} \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}^{(t)} > 0 \\ -r^{(t)} \tilde{\mathbf{x}}^{(t)}, & \text{otherwise} \end{cases}$$

- Can be also written as $\nabla emp(\tilde{\mathbf{w}}) = -\frac{1}{2} \sum_{t=1}^N \tilde{\mathbf{x}}^{(t)T} \left(r^{(t)} - \text{sgn}(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}) \right)$

Repeat until convergence {

$$\tilde{\mathbf{w}} := \tilde{\mathbf{w}} - \eta \nabla emp(\tilde{\mathbf{w}})$$

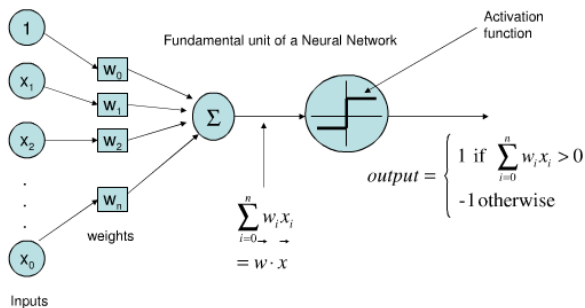
$$= \tilde{\mathbf{w}} + \eta' \sum_{t=1}^N \tilde{\mathbf{x}}^{(t)T} \left(r^{(t)} - \text{sgn}(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}) \right),$$

$$\text{where } \eta' = \eta/2$$

}

Remarks (1)

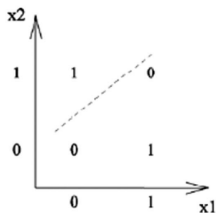
- A model for how individual neurons in human brain work



- Not good at recognizing non-linear classes/patterns
 - E.g., identifying an object in an image
- Improvements:
 - Chained to form a **neural network**
 - Make instances linearly separable (to be discussed later)

Remarks (2)

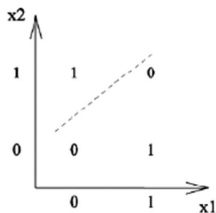
- If positive and negative examples are not linearly separable by $\text{sgn}(\beta^\top \tilde{\mathbf{x}}')$, the training algorithm will **not** converge



- Solution?

Remarks (2)

- If positive and negative examples are not linearly separable by $\text{sgn}(\boldsymbol{\beta}^\top \tilde{\mathbf{x}}')$, the training algorithm will **not** converge



- Solution?
 - Limit the maximum number of iterations, or
 - Stop if $|\text{emp}(\tilde{\mathbf{w}}^{(new)}) - \text{emp}(\tilde{\mathbf{w}}^{(old)})| < \epsilon$

Update Rules: Perceptron vs. Logistic Regression

- If we “harden” the logistic function $\pi(\mathbf{x}; \boldsymbol{\beta}) = \frac{1}{1+e^{-\boldsymbol{\beta}^\top \tilde{\mathbf{x}}}}$ to

$$\pi(\mathbf{x}; \boldsymbol{\beta}) = \begin{cases} 1, & \text{if } \boldsymbol{\beta}^\top \tilde{\mathbf{x}} \geq 0, \\ 0, & \text{otherwise,} \end{cases}$$

so that the gradient descent update rule becomes:

Repeat until convergence {

$$\begin{aligned} \boldsymbol{\beta} &:= \boldsymbol{\beta} + \eta \nabla l(\boldsymbol{\beta}) = \boldsymbol{\beta} + \eta \sum_{t=1}^N \tilde{\mathbf{x}}^{(t)\top} (q^{(t)} - \pi(\mathbf{x}^{(t)}; \boldsymbol{\beta})) \\ &= \boldsymbol{\beta} + \eta' \sum_{t=1}^N \tilde{\mathbf{x}}^{(t)\top} \left(r^{(t)} - \text{sgn}(\boldsymbol{\beta}^\top \tilde{\mathbf{x}}') \right), \end{aligned}$$

$$\text{where } q^{(t)} = \frac{r^{(t)}+1}{2} \text{ and } \eta' = \eta/2$$

}

- Despite its cosmetic similarity with logistic regression, perceptron learning has no simple probabilistic interpretation

Outline

- 1 Regression
 - Linear Regression
 - Interpolation vs. Regression
 - Probability Interpretation
- 2 Two-Class Classification
 - Logistic Regression
 - Perceptron
- 3 Multiclass Classification
 - Wrapper Methods
 - Direct Models
- 4 Non-Parametric Methods

Learning Multiple Classes

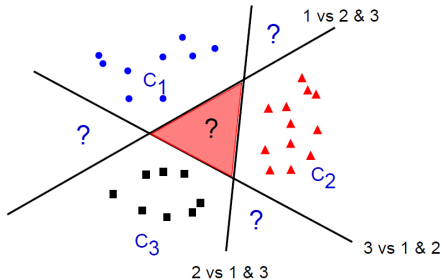
- What if we have K classes instead of 2?
- Applications:
 - OCR (Optical Character Recognition)
 - Medical diagnosis
 - Surveillance, etc.
- Training set: $\mathcal{X} = \{\mathbf{x}^{(t)}, \mathbf{r}^{(t)}\}_{t=1}^N$, where $\mathbf{r}^{(t)} \in \mathbb{R}^K$ and
$$r_i^{(t)} = \begin{cases} 1, & \mathbf{x}^{(t)} \in C_i \\ -1, & \textit{otherwise} \end{cases}$$

Outline

- 1 Regression
 - Linear Regression
 - Interpolation vs. Regression
 - Probability Interpretation
- 2 Two-Class Classification
 - Logistic Regression
 - Perceptron
- 3 Multiclass Classification
 - Wrapper Methods
 - Direct Models
- 4 Non-Parametric Methods

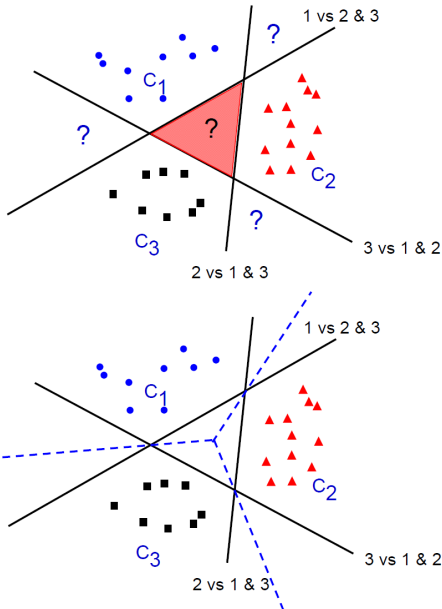
1 vs. All (1)

- Assume a model consisting of K hypotheses h_i
 - There is no need for these K hypotheses to belong to the same hypotheses class
- Perform the two-class classification K times
 - Each time treat the examples of a certain class as positive and the rest as negative
- How to handle the *cases of doubt?*



1 vs. All (1)

- Assume a model consisting of K hypotheses h_i
 - There is no need for these K hypotheses to belong to the same hypotheses class
- Perform the two-class classification K times
 - Each time treat the examples of a certain class as positive and the rest as negative
- How to handle the *cases of doubt*?
 - Define decision boundaries, e.g., $y' := \arg \max_i h_i(\mathbf{x}'; \theta_i)$



1 vs. All (2)

- Pros:
 - Easy to implement
 - # classifiers grows with K
- Cons:
 - Time consuming (each of the K classifiers takes the whole dataset as input)
 - Each classifier deals with imbalance dataset

- Perform 1 vs. 1 classification $\binom{K}{2}$ times, and predict by voting
- Pros:
 - Avoid creating imbalanced dataset for each classifier
 - Faster and memory economic (each classifier takes only two classes in the dataset as input)
- Cons:
 - # classifiers grows with K^2 , not suitable for datasets with massive classes

Wrappers based on Error-Correcting Codes

	h_1		\dots		h_L
C_1	-1	-1	-1	-1	1
\vdots	1	-1	1	1	-1
C_K	1	1	-1	-1	1

- Rows: predefined codewords of length L
- Columns: a particular grouping of examples for training a two-class classifier

- To make prediction:
 - 1 Obtain a codeword for x' based on the predictions of L classifiers
 - 2 Assign x' to the label with the most similar codeword
- If codewords are designed such that each pair has Hamming distance at least s , then $\lfloor \frac{s-1}{2} \rfloor$ wrong predictions can be tolerated

Outline

- 1 Regression
 - Linear Regression
 - Interpolation vs. Regression
 - Probability Interpretation
- 2 Two-Class Classification
 - Logistic Regression
 - Perceptron
- 3 **Multiclass Classification**
 - Wrapper Methods
 - Direct Models
- 4 Non-Parametric Methods

Generalized Linear Models

TBA

Multi-Hyperplane Classifier

- Learn K separating hyperplanes simultaneously:

$$\begin{aligned} & \arg \min_{\{\mathbf{w}_i, b_i\}_{i=1}^K, \xi} \sum_{t=1}^N \xi_{t,r} \\ \text{subject to } & (\mathbf{w}_{idx(\mathbf{r}^{(t)})}^\top \mathbf{x}^{(t)} - b_{idx(\mathbf{r}^{(t)})}) - (\mathbf{w}_r^\top \mathbf{x}^{(t)} - b_r) > -\xi_{t,r} \\ & \text{and } \xi_{t,r} \geq 0, \forall t, r \neq idx(\mathbf{r}^{(t)}) \end{aligned}$$

- For an example of class r , the corresponding hyperplane should give value higher than those given by other hyperplanes
- Prediction: $y' := \arg \max_i \mathbf{w}_i^\top \mathbf{x}' - b_i$
- Hyperplanes are correlated
 - No one will give values significantly higher than the others
- In practice,
 - There is little or no performance improvement over the wrappers
 - Very slow and memory hungry

- 1 **Regression**
 - Linear Regression
 - Interpolation vs. Regression
 - Probability Interpretation
- 2 **Two-Class Classification**
 - Logistic Regression
 - Perceptron
- 3 **Multiclass Classification**
 - Wrapper Methods
 - Direct Models
- 4 **Non-Parametric Methods**

- There are another simple ways, call k -NN methods, to make predictions
- Given a test instance \mathbf{x}' , predict its label by the (weighted) average of labels of k examples in \mathcal{X} most similar to \mathbf{x}'
 - Applies to both continuous and discrete labels
- Needs a similarity metric $k(\mathbf{x}, \mathbf{y})$ between any two instances
 - E.g., *cosin* similarity: $k(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x}^\top \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|} \in [-1, 1]$
- Training: simply remember \mathcal{X}

Non-Parametric Methods

- k -NN methods are special cases of *non-parametric* (or *memory-based*) methods
 - Non-parametric in the sense that f cannot be described by parameters
 - Memory-based in that all data (rather than just parameters) need to be memorized during the training process
- *Lazy* since the hypothesis is obtained only before the prediction
- This allows the development of *local models*

Local Weighted Linear Regression

- Recall in (eager) linear regression, we fit $\mathbf{w} \in \mathbb{R}^{d+1}$ to minimize the SSE: $\sum_i (r^{(i)} - \mathbf{w}^\top \begin{bmatrix} 1 \\ \mathbf{x}^{(i)} \end{bmatrix})^2$
- Local model: fit \mathbf{w} to minimize SSE **local to the instance \mathbf{x}' we want to predict**:

$$\sum_i l(\mathbf{x}^{(i)}; \mathbf{x}') (r^{(i)} - \mathbf{w}^\top \begin{bmatrix} 1 \\ \mathbf{x}^{(i)} \end{bmatrix})^2$$

where $l: \mathbb{R}^d \rightarrow \mathbb{R}$ is a weighting function

- Idea: only examples nearby (or local to) \mathbf{x}' should be taken into account in $emp(\theta; \mathcal{X})$
- Possible choice for l : $l(\mathbf{x}^{(i)}; \mathbf{x}') := \exp\left(-\frac{(\mathbf{x}^{(i)} - \mathbf{x}')^2}{2\tau^2}\right)$ for some τ (mimics k -NN)

Summary of Supervised Learning Models

- Three main categories (either parametric or non-parametric):
 - 1 Those learning the **discriminant functions** f 's (no probability interpretation)
 - E.g., perceptron, k NN, etc.
 - 2 Those based on probability and learn $p(r|\mathbf{x})$ directly
 - E.g., linear regression, logistic regression, etc.
 - $p(r|\mathbf{x};\theta)$ with θ (constant) estimated from \mathcal{X}
 - Methods in 1 and 2 are called **discriminative methods**
 - 3 Those learn $p(r|\mathbf{x})$ indirectly from $p(\mathbf{x}|r)p(r)$
 - To be discussed later
 - These are called **generative methods**, as $p(\mathbf{x}|r)p(r)$ explains how \mathcal{X} is generated