

Linear Models for Classification

CMPT 419/726

Mo Chen

SFU Computing Science

Jan. 22, 2020

Bishop PRML Ch. 4

Classification: Hand-written Digit Recognition

$$\mathbf{x}_i = \begin{array}{|c|} \hline \text{4} \\ \hline \end{array} \quad \mathbf{t}_i = (0, 0, 0, 1, 0, 0, 0, 0, 0, 0)$$

- Each input vector classified into one of K discrete classes
 - Denote classes by \mathcal{C}_k
- Represent input image as a vector $\mathbf{x}_i \in \mathbb{R}^{784}$.
- We have target vector $\mathbf{t}_i \in \{0, 1\}^{10}$
- Given a **training set** $\{(\mathbf{x}_1, \mathbf{t}_1), \dots, (\mathbf{x}_N, \mathbf{t}_N)\}$, learning problem is to construct a “good” function $\mathbf{y}(\mathbf{x})$ from these.
 - $\mathbf{y}: \mathbb{R}^{784} \rightarrow \mathbb{R}^{10}$

Generalized Linear Models

- Similar to previous chapter on linear models for regression, we will use a “linear” model for classification:

$$y(\mathbf{x}) = f(\mathbf{w}^T \mathbf{x} + w_0)$$

- This is called a **generalized linear model**
- $f(\cdot)$ is a fixed non-linear function
 - e.g.

$$f(u) = \begin{cases} 1, & \text{if } u \geq 0 \\ 0, & \text{otherwise} \end{cases}$$

- **Decision boundary** between classes will be linear function of \mathbf{x}
- Can also apply non-linearity to \mathbf{x} , as in $\phi_i(\mathbf{x})$ for regression

Outline

Discriminant Functions

Generative Models

Discriminative Models

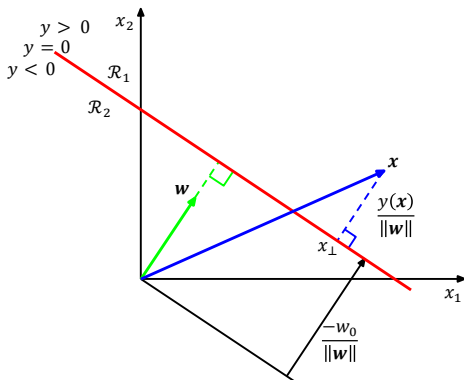
Outline

Discriminant Functions

Generative Models

Discriminative Models

Discriminant Functions with Two Classes

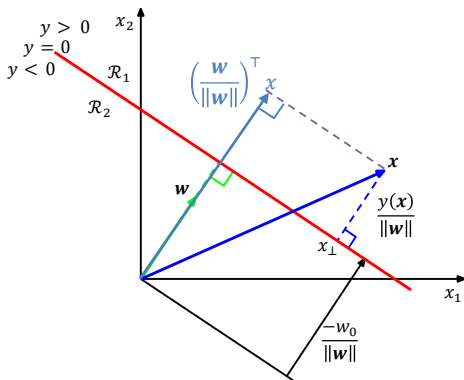


- Start with 2 class problem, $t_i \in \{0,1\}$
- Simple linear discriminant

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

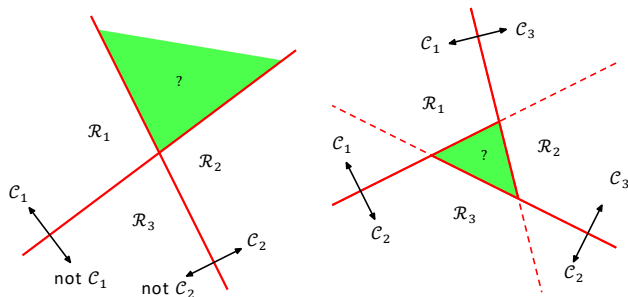
apply threshold function to get classification

Discriminant Functions with Two Classes



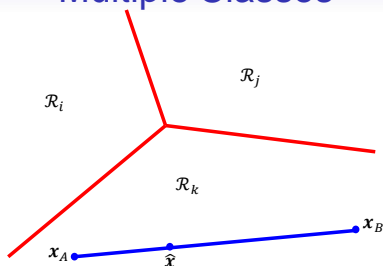
- $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$
 - Gradient of y is \mathbf{w}
 - Constant y values \Rightarrow parallel lines
- If $y = 0$ (decision boundary),
 $\mathbf{w}^T \mathbf{x} = -w_0 \Rightarrow \left(\frac{\mathbf{w}}{\|\mathbf{w}\|} \right)^T \mathbf{x} = -\frac{w_0}{\|\mathbf{w}\|}$
- In general, $\frac{y}{\|\mathbf{w}\|} = \frac{\mathbf{w}^T \mathbf{x}}{\|\mathbf{w}\|} + \frac{w_0}{\|\mathbf{w}\|}$, or
 $\left(\frac{\mathbf{w}}{\|\mathbf{w}\|} \right)^T \mathbf{x} = \frac{y(\mathbf{x})}{\|\mathbf{w}\|} - \frac{w_0}{\|\mathbf{w}\|}$

Multiple Classes



- A linear discriminant between two classes separates with a hyperplane
- How to use this for multiple classes?
- **One-versus-the-rest** method: build $K - 1$ classifiers, between \mathcal{C}_k and all others
- **One-versus-one** method: build $K(K - 1)/2$ classifiers, between all pairs

Multiple Classes



- A solution is to build K linear functions:

$$y_k(\mathbf{x}) = \mathbf{w}_k^\top \mathbf{x} + w_{k0}$$

assign \mathbf{x} to class $\arg \max_k y_k(\mathbf{x})$

- Gives connected, convex **decision regions**

$$\hat{\mathbf{x}} = \lambda \mathbf{x}_A + (1 - \lambda) \mathbf{x}_B$$

$$y_k(\hat{\mathbf{x}}) = \lambda y_k(\mathbf{x}_A) + (1 - \lambda) y_k(\mathbf{x}_B)$$

$$\Rightarrow y_k(\hat{\mathbf{x}}) > y_j(\hat{\mathbf{x}}), \quad \forall j \neq k$$

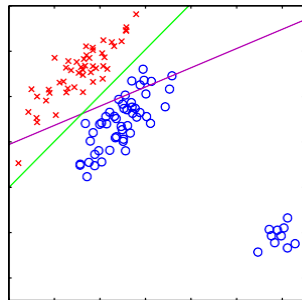
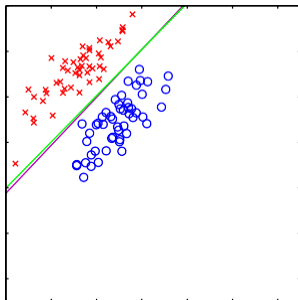
Least Squares for Classification

- How do we learn the decision boundaries (\mathbf{w}_k, w_{k0})?
- One approach is to use least squares, similar to regression
- Find \mathbf{W} to minimize squared error over all examples and all components of the label vector:

$$E(\mathbf{W}) = \frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K (y_k(\mathbf{x}_n) - t_{nk})^2$$

- Some algebra, we get a solution using the pseudo-inverse as in regression

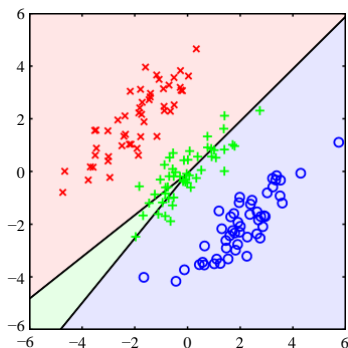
Problems with Least Squares



- Looks okay... **least squares decision boundary**
 - Similar to **logistic regression decision boundary** (more later)

- Gets worse by adding easy points?!
- Why?
 - If target value is 1, points far from boundary will have high value, say 10; this is a large error so the boundary is moved

More Least Squares Problems



- Easily separated by hyperplanes, but not found using least squares!
- We'll address these problems later with better models
- First, a look at a different criterion for linear discriminant

Fisher's Linear Discriminant

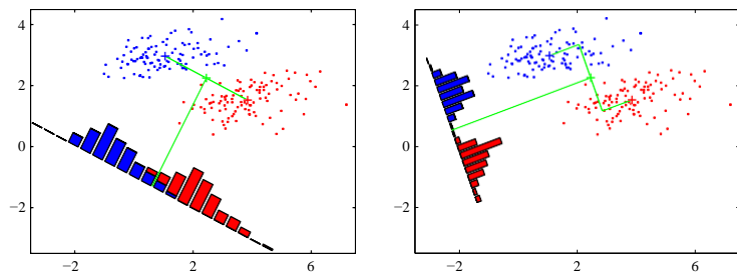
- The two-class linear discriminant acts as a projection:

$$y = \mathbf{w}^T \mathbf{x} \geq -w_0$$

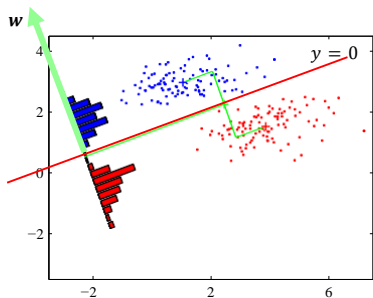
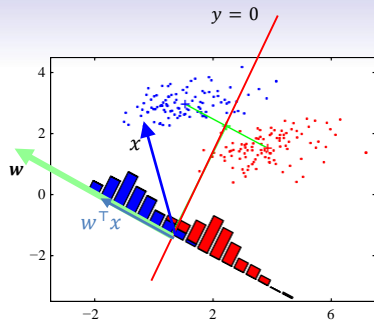
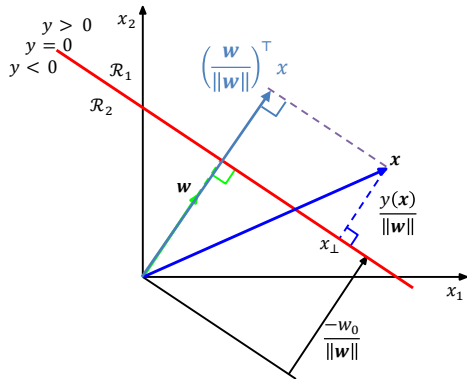
followed by a threshold

- In which direction \mathbf{w} should we project?
- One which separates classes “well”

Fisher's Linear Discriminant



- A natural idea would be to project in the direction of the line connecting class means
- However, problematic if classes have variance in this direction
- Fisher criterion: maximize ratio of inter-class separation (between) to intra-class variance (inside)



Math time - FLD

- Projection $\mathbf{y}_n = \mathbf{w}^\top \mathbf{x}_n$
- Inter-class separation is distance between class means (good):

$$m_k = \frac{1}{N_k} \sum_{n \in \mathcal{C}_k} \mathbf{w}^\top \mathbf{x}_n$$

- Intra-class variance (bad):

$$s_k^2 = \sum_{n \in \mathcal{C}_k} (y_n - m_k)^2$$

- Fisher criterion:

$$J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2}$$

maximize wrt \mathbf{w}

Math time - FLD

$$J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2} = \frac{\mathbf{w}^\top S_B \mathbf{w}}{\mathbf{w}^\top S_W \mathbf{w}}$$

Between-class covariance:

$$S_B = (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^\top$$

Within-class covariance:

$$S_W = \sum_{n \in \mathcal{C}_1} (\mathbf{x}_n - \mathbf{m}_1)(\mathbf{x}_n - \mathbf{m}_1)^\top + \sum_{n \in \mathcal{C}_2} (\mathbf{x}_n - \mathbf{m}_2)(\mathbf{x}_n - \mathbf{m}_2)^\top$$

Lots of math:

$$\mathbf{w} \propto S_W^{-1}(\mathbf{m}_2 - \mathbf{m}_1)$$

If covariance S_W is isotropic, reduces to class mean difference vector

FLD Summary

- FLD is a **dimensionality reduction** technique (more later in the course)
- Criterion for choosing projection based on class labels
 - Still suffers from outliers (e.g. earlier least squares example)

Perceptrons

- **Perceptrons** is used to refer to many neural network structures (more coming up)
- The classic type is a fixed non-linear transformation of input, one layer of adaptive weights, and a threshold:

$$y(\mathbf{x}) = f(\mathbf{w}^T \phi(\mathbf{x}))$$

- Developed by Rosenblatt in the 50s
- The main difference compared to the methods we've seen so far is the learning algorithm

Perceptron Learning

- Two class problem
- For ease of notation, we will use $t = 1$ for class \mathcal{C}_1 and $t = -1$ for class \mathcal{C}_2 ; we choose f such that $f(a) = 1$ if $a \geq 0$ and $f(a) = -1$ otherwise
- We saw that squared error was problematic
- Instead, we'd like to minimize the number of misclassified examples
 - An example is mis-classified if $\mathbf{w}^\top \phi(\mathbf{x}_n)t_n < 0$
 - **Perceptron criterion:**

$$E_P(\mathbf{w}) = - \sum_{n \in \mathcal{M}} \mathbf{w}^\top \phi(\mathbf{x}_n)t_n$$

sum over mis-classified examples only

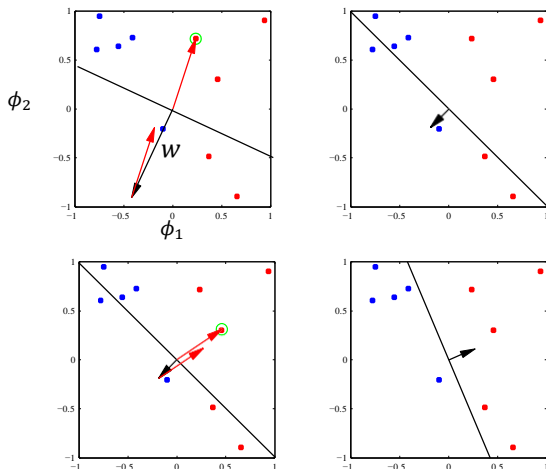
Perceptron Learning Algorithm

- Minimize the error function using stochastic gradient descent (gradient descent per example):

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_P(\mathbf{w}) = \mathbf{w}^{(\tau)} + \underbrace{\eta \phi(\mathbf{x}_n) t_n}_{\text{if incorrect}}$$

- Iterate over all training examples, only change \mathbf{w} if the example is mis-classified
- Guaranteed to converge if data are **linearly separable**
- Will not converge if not
- May take many iterations
- Sensitive to initialization

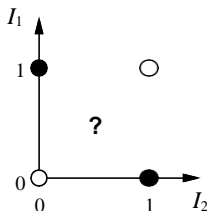
Perceptron Learning Illustration



- Note there are many hyperplanes with 0 error
 - Support vector machines have a nice way of choosing one

Limitations of Perceptrons

- Perceptrons can only solve linearly separable problems in feature space
 - Same as the other models in this chapter
- Canonical example of non-separable problem is X-OR
 - Real datasets can look like this too



Outline

Discriminant Functions

Generative Models

Discriminative Models

Probabilistic Generative Models

- Up to now we've looked at learning classification by choosing parameters to minimize an error function
- We'll now develop a probabilistic approach
- With 2 classes, \mathcal{C}_1 and \mathcal{C}_2 :

$$p(\mathcal{C}_1|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x})} \quad \text{Bayes' Rule}$$

$$p(\mathcal{C}_1|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}, \mathcal{C}_1) + p(\mathbf{x}, \mathcal{C}_2)} \quad \text{Sum rule}$$

$$p(\mathcal{C}_1|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1) + p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)} \quad \text{Product rule}$$

- In **generative models** we specify the distribution $p(\mathbf{x}|\mathcal{C}_k)$ which generates the data for each class

Probabilistic Generative Models - Example

- Let's say we observe x which is the current temperature
- Determine if we are in Vancouver (\mathcal{C}_1) or Honolulu (\mathcal{C}_2)
- Generative model:

$$p(\mathcal{C}_1|x) = \frac{p(x|\mathcal{C}_1)p(\mathcal{C}_1)}{p(x|\mathcal{C}_1)p(\mathcal{C}_1) + p(x|\mathcal{C}_2)p(\mathcal{C}_2)}$$

- $p(x|\mathcal{C}_1)$ is distribution over typical temperatures in Vancouver
 - e.g. $p(x|\mathcal{C}_1) = \mathcal{N}(x; 10, 5)$
- $p(x|\mathcal{C}_2)$ is distribution over typical temperatures in Honolulu
 - e.g. $p(x|\mathcal{C}_1) = \mathcal{N}(x; 25, 5)$
 - Class priors $p(\mathcal{C}_1) = 0.1, p(\mathcal{C}_2) = 0.9$
- $p(\mathcal{C}_1|x = 15) = \frac{0.0484 \times 0.1}{0.0484 \times 0.1 + 0.0108 \times 0.9} \approx 0.33$

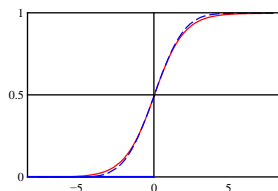
Generalized Linear Models

- We can write the classifier in another form

$$\begin{aligned} p(\mathcal{C}_1|\mathbf{x}) &= \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1) + p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)} \\ &= \frac{1}{1 + \frac{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)}{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}} \\ &= \frac{1}{1 + \exp(-a)} \equiv \sigma(a) \end{aligned}$$

$$\text{where } a = \ln \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)}$$

Logistic Sigmoid



- The function $\sigma(a) = \frac{1}{1 + \exp(-a)}$ is known as the logistic sigmoid
- It squashes the real axis down to $[0, 1]$
- It is continuous and differentiable
- It avoids the problems encountered with the *too correct* least-squares error fitting

Multi-class Extension

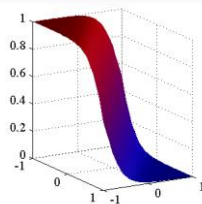
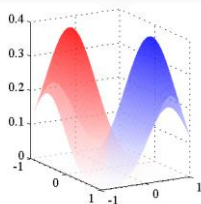
- There is a generalization of the logistic sigmoid to $K > 2$ classes:

$$\begin{aligned} p(\mathcal{C}_k|\mathbf{x}) &= \frac{p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)}{\sum_j p(\mathbf{x}|\mathcal{C}_j)p(\mathcal{C}_j)} \\ &= \frac{\exp(a_k)}{\sum_j \exp(a_j)} \end{aligned}$$

where $a_k = \ln p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)$

- a.k.a. **softmax function**
 - If some $a_k \gg a_j$, $p(\mathcal{C}_k|\mathbf{x})$ goes to 1

Gaussian Class-Conditional Densities



- Back to that a in the logistic sigmoid for 2 classes
- Let's assume the class-conditional densities $p(x|\mathcal{C}_k)$ are Gaussians, and have the same covariance matrix Σ :

$$p(x|\mathcal{C}_k) = \frac{1}{(2\pi)^{D/2}|\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^\top \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}_k)\right\}$$

- a takes a simple form:

$$a = \ln \frac{p(x|\mathcal{C}_1)p(\mathcal{C}_1)}{p(x|\mathcal{C}_2)p(\mathcal{C}_2)} = \mathbf{w}^\top \mathbf{x} + w_0$$

- Note that quadratic terms $\mathbf{x}^\top \Sigma^{-1} \mathbf{x}$ cancel

Maximum Likelihood Learning

- We can fit the parameters to this model using **maximum likelihood**
 - Parameters are $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \Sigma^{-1}, p(\mathcal{C}_1) \equiv \pi, p(\mathcal{C}_2) \equiv 1 - \pi$
 - Refer to as θ
- For a datapoint x_n from class \mathcal{C}_1 ($t_n = 1$):

$$p(\mathbf{x}_n, \mathcal{C}_1) = p(\mathcal{C}_1)p(\mathbf{x}_n|\mathcal{C}_1) = \pi\mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_1, \Sigma)$$

- For a datapoint x_n from class \mathcal{C}_2 ($t_n = 0$):

$$p(\mathbf{x}_n, \mathcal{C}_2) = p(\mathcal{C}_2)p(\mathbf{x}_n|\mathcal{C}_2) = (1 - \pi)\mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_2, \Sigma)$$

Maximum Likelihood Learning

- The likelihood of the training data is:

$$p(\mathbf{t}|\pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}) = \prod_{n=1}^N [\pi \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_1, \boldsymbol{\Sigma})]^{t_n} [(1 - \pi) \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_2, \boldsymbol{\Sigma})]^{1-t_n}$$

- As usual, $\ln(\cdot)$ is our friend:

$$l(\mathbf{t}|\theta) = \sum_{n=1}^N \underbrace{(t_n \ln \pi + (1 - t_n) \ln(1 - \pi))}_{\pi} + \underbrace{t_n \ln \mathcal{N}_1 + (1 - t_n) \ln \mathcal{N}_2}_{\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}}$$

- Maximize for each separately

Maximum Likelihood Learning - Class Priors

- Maximization with respect to the class priors parameter π is straightforward:

$$\frac{\partial}{\partial \pi} l(t|\theta) = \sum_{n=1}^N \left(\frac{t_n}{\pi} - \frac{1-t_n}{1-\pi} \right)$$
$$\Rightarrow \pi = \frac{N_1}{N_1 + N_2}$$

- N_1 and N_2 are the number of training points in each class
- Prior is simply the fraction of points in each class

Maximum Likelihood Learning - Gaussian Parameters

- The other parameters can also be found in the same fashion
- Class means:

$$\mu_1 = \frac{1}{N_1} \sum_{n=1}^N t_n \mathbf{x}_n$$

$$\mu_2 = \frac{1}{N_2} \sum_{n=1}^N (1 - t_n) \mathbf{x}_n$$

- Means of training examples from each class
- Shared covariance matrix:

$$\Sigma = \frac{N_1}{N} \frac{1}{N_1} \sum_{n \in \mathcal{C}_1} (\mathbf{x}_n - \boldsymbol{\mu}_1)(\mathbf{x}_n - \boldsymbol{\mu}_1)^\top + \frac{N_2}{N} \frac{1}{N_2} \sum_{n \in \mathcal{C}_2} (\mathbf{x}_n - \boldsymbol{\mu}_2)(\mathbf{x}_n - \boldsymbol{\mu}_2)^\top$$

- Weighted average of class covariances

Probabilistic Generative Models Summary

- Fitting Gaussian using ML criterion is sensitive to outliers
- Simple linear form for a in logistic sigmoid occurs for more than just Gaussian distributions
 - Arises for any distribution in the [exponential family](#), a large class of distributions

Outline

Discriminant Functions

Generative Models

Discriminative Models

Probabilistic Discriminative Models

- Generative model made assumptions about form of class-conditional distributions (e.g. Gaussian)
 - Resulted in logistic sigmoid of linear function of \mathbf{x}
- Discriminative model - explicitly use functional form

$$p(\mathcal{C}_1|\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x} + w_0)}$$

and find \mathbf{w} directly

- For the generative model we had $\underbrace{2M + M(M + 1)/2 + 1}_{\text{Means, variance, prior}}$ parameters
 - M is dimensionality of \mathbf{x}
- Discriminative model will have $M + 1$ parameters

Generative vs. Discriminative

- Generative models
 - Can generate synthetic example data
 - Perhaps accurate classification is equivalent to accurate synthesis
 - e.g. vision and graphics
 - Tend to have more parameters
 - Require good model of class distributions
- Discriminative models
 - Only usable for classification
 - Don't solve a harder problem than you need to
 - Tend to have fewer parameters
 - Require good model of decision boundary

Maximum Likelihood Learning - Discriminative Model

- As usual we can use the maximum likelihood criterion for learning

$$p(\mathbf{t}|\mathbf{w}) = \prod_{n=1}^N y_n^{t_n} \{1 - y_n\}^{1-t_n}, \text{ where } y_n = p(\mathcal{C}_1|\mathbf{x}_n)$$

- Taking ln and derivative gives:

$$\nabla l(\mathbf{w}) = \sum_{n=1}^N (t_n - y_n) \mathbf{x}_n$$

- This time no closed-form solution since $y_n = \sigma(\mathbf{w}^T \mathbf{x})$
- Could use (stochastic) gradient descent
 - But there's a better iterative technique

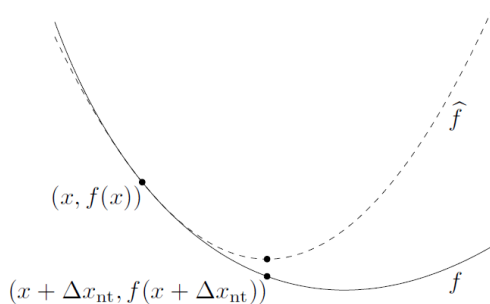
Iterative Reweighted Least Squares

- **Iterative reweighted least squares (IRLS)** is a descent method
 - As in **gradient descent**, start with an initial guess, improve it
 - Gradient descent - take a step (how large?) in the gradient direction
- IRLS is a special case of a **Newton-Raphson** method
 - Approximate function using second-order Taylor expansion:

$$\hat{f}(\mathbf{w} + \mathbf{v}) = f(\mathbf{w}) + \nabla f(\mathbf{w})^T (\mathbf{v} - \mathbf{w}) + \frac{1}{2} (\mathbf{v} - \mathbf{w})^T H f(\mathbf{w}) (\mathbf{v} - \mathbf{w})$$

- Closed-form solution to minimize this is straight-forward: quadratic, derivatives linear
- In IRLS this second-order Taylor expansion ends up being a weighted least-squares problem, as in the regression case from last week
 - Hence the name IRLS

Newton-Raphson



- Figure from Boyd and Vandenberghe, *Convex Optimization*
 - Excellent reference, free for download online
<http://www.stanford.edu/~boyd/cvxbook/>

Conclusion

- Readings: Ch. 4.1.1-4.1.4, 4.1.7, 4.2.1-4.2.2, 4.3.1-4.3.3
- Generalized linear models $y(\mathbf{x}) = f(\mathbf{w}^T \mathbf{x} + w_0)$
- Threshold/max function for $f(\cdot)$
 - Minimize with least squares
 - Fisher criterion - class separation
 - Perceptron criterion - mis-classified examples
- Probabilistic models: logistic sigmoid / softmax for $f(\cdot)$
 - Generative model - assume class conditional densities in exponential family; obtain sigmoid
 - Discriminative model - directly model posterior using sigmoid (a. k. a. **logistic regression**, though classification)
 - Can learn either using maximum likelihood
- **All of these models are limited to linear decision boundaries in feature space**