

Topic 1

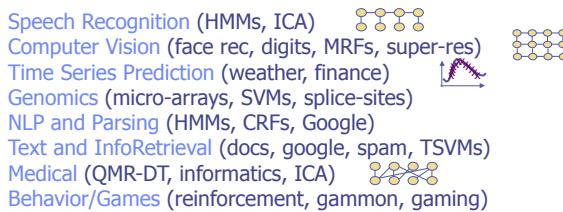
- Introduction
- Machine Learning: What, Why and Applications
- Syllabus, policies, texts, web page
- Historical Perspective
- Machine Learning Tasks and Tools
- Digit Recognition Example
- Machine Learning Approach
- Deterministic or Probabilistic Approach
- Why Probabilistic?

Machine Learning

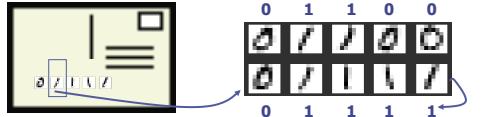
Slides use Tony Jebara's materials

Machine Learning Applications

- ML: Interdisciplinary (CS, Math, Stats, Physics, OR, Psych)
- Data-driven approach to AI
- Many domains are too hard to do manually



ML Example: Digit Recognition



- Want to automate zipcode reading in post office
- Look at an image and say if it is a '1' or '0'
- 8x8 pixels of gray-level (0.0=dark, 0.5=gray, 1.0=white)
- Learn from above labeled training images
- Predict labels on testing images
- Binary Classification [0,1]
- What to do?



Course Details & Requirements

- Probability/Stats, Linear Algebra, Calculus, AI
- Mathematical & Data Driven approach to AI
- Lots of Equations!
- Required Text: Introduction to Graphical Models by M. Jordan & C. Bishop (Online)
Pattern Recognition & Machine Learning by C. Bishop (Spring 2006 Edition)
- Reference Text: Pattern Classification (3rd Edition) by Duda, Hart and Stork
- Homework: Every 2-3 weeks
- Grading: homework, midterm, 2 quizzes & final examination
- Software Requirements: Matlab software & Acis account

Machine Learning: What/Why

Statistical Data-Driven Computational Models

Real domains (vision, speech, behavior):

no $E=MC^2$

noisy, complex, nonlinear

have many variables

non-deterministic

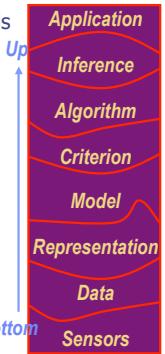
incomplete, approximate models

Need: statistical models driven by data & sensors, a.k.a Machine Learning

Bottom-Up: use data to form a model

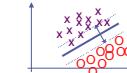
Why? Complex data everywhere,
audio, video, internet

Intelligence = Learning = Prediction



Machine Learning Tasks

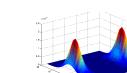
Classification $y=\text{sign}(f(x))$



Regression $y=f(x)$

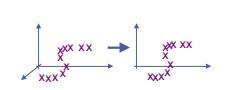


Modeling $p(x)$

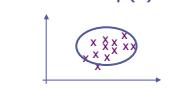


Clustering

Feature Selection



Detection $p(x) < t$



Supervised

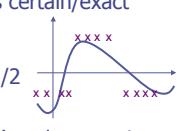
Unsupervised

Ex: Two Approaches

In ML, we will consider two complementary approaches:

1) Deterministic:

All variables/observables are treated as certain/exact
Find/fit a function $f(X)$ on an image X
Output 0 or 1 depending on input
Class label given by $y=\text{sign}(f(X))/2 + 1/2$

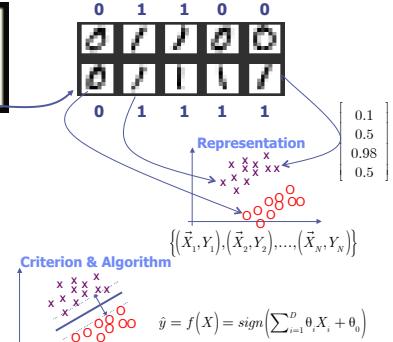
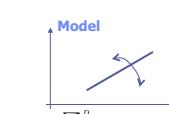
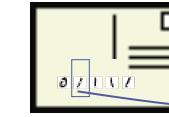


2) Probabilistic/Bayesian/Stochastic:

Variables/observables are random (R.V.) and uncertain
Probability image is a '0' digit: $p(y=0|X) = 0.43$
Probability image is a '1' digit: $p(y=1|X) = 0.57$
Output label with larger $p(y=0|\text{image})$ or $p(y=1|\text{image})$

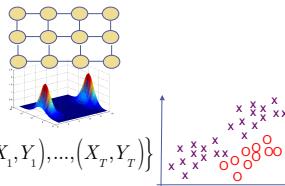
These are interconnected! Deterministic approaches can be generated from (more general) probabilistic approaches

Ex: 1) Deterministic Approach



Ex: 2) Probabilistic Approach

a) Provide Prior Model
Parameters & Structure
e.g. nearby pixels are co-dependent



b) Obtain Data and Labels $\{(X_1, Y_1), \dots, (X_t, Y_t)\}$

c) Learn a probability model with data
p(all system variables)

$$p(X, Y)$$

d) Use model for inference (classify/predict)

Probability image is '0': $p(y=0|X)$

Probability image is '1': $p(y=1|X)$

Output: $\arg \max_i p(y=i|X)$



Monty Hall Solution

Probabilistic Interpretation is Best

Bayesian Solution: Change your mind!

Assume we always start by picking A.

If prize behind A: Opens B/C \rightarrow Change A to C/B \rightarrow Lose

If prize behind B: Opens C \rightarrow Change A to B \rightarrow Win

If prize behind C: Opens B \rightarrow Change A to C \rightarrow Win

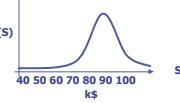
Probability of winning if change your mind = 66%

Probability of winning if stick to your guns = 33%

Why Probabilistic Approach?

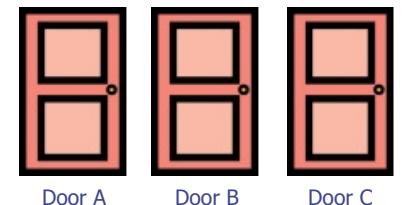
- Decision making often involves uncertainty
- Hidden variables, complexity, randomness in system
- Input data is noisy and uncertain
- Estimated model is noisy and uncertain
- Output data is uncertain (no single correct answer)

- Example: Predict your salary in the future
- Inputs: Field, Degree, University, City, IQ
- Output: \$Amount
- There is uncertainty and hidden variables
- No one answer (I.e. \$84K) is correct
- Answer = a distribution over salaries



Why Probabilistic? Monty Hall

- Behind one door is a prize (car? 1\$?)
- Pick a door

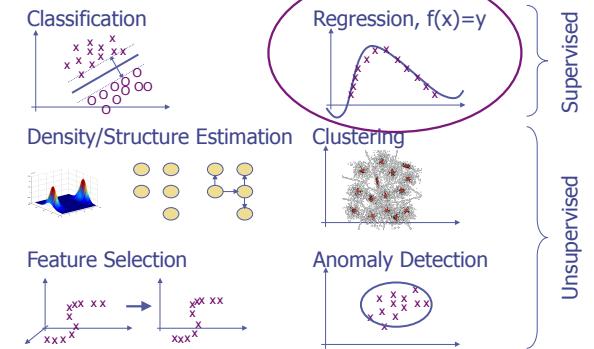


Machine Learning

Topic 2

- Regression
- Empirical Risk Minimization
- Least Squares
- Higher Order Polynomials
- Under-fitting / Over-fitting
- Cross-Validation

Regression



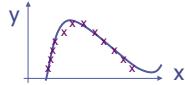
Function Approximation

- Start with training dataset

$$\mathcal{X} = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\} \quad x \in \mathbb{R}^D = \begin{bmatrix} x(1) \\ x(2) \\ \vdots \\ x(D) \end{bmatrix} \quad y \in \mathbb{R}^1$$

- Have N (input, output) pairs

- Find a function $f(x)$ to predict y from x
That fits the training data well



- Example: predict the price of house in

dollars y using $x = [\# \text{rooms}; \text{latitude}; \text{longitude}; \dots]$

- Need:
 - a) Way to evaluate how good a fit we have
 - b) Class of functions in which to search for $f(x)$

Empirical Risk Minimization

- Idea: minimize 'loss' on the training data set
- Empirical = use the training set to find the best fit

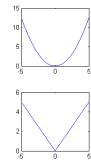
- Define a loss function of how good we fit a single point:
 $L(y, f(x))$

- Empirical Risk = the average loss over the dataset

$$R = \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i))$$

- Simplest loss: squared error from y value

$$L(y_i, f(x_i)) = \frac{1}{2} (y_i - f(x_i))^2$$



- Other possible loss: absolute error

$$L(y_i, f(x_i)) = |y_i - f(x_i)|$$



Min by Gradient=0

- Gradient=0 means the partial derivatives are all 0
 $\nabla_{\theta} R = \begin{bmatrix} \frac{\partial R}{\partial \theta_0} \\ \frac{\partial R}{\partial \theta_1} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

- Take partials of empirical risk:

$$R(\theta) = \frac{1}{2N} \sum_{i=1}^N (y_i - \theta_1 x_i - \theta_0)^2$$

$$\frac{\partial R}{\partial \theta_0} = \frac{1}{N} \sum_{i=1}^N (y_i - \theta_1 x_i - \theta_0)(-1) = 0$$

$$\frac{\partial R}{\partial \theta_1} = \frac{1}{N} \sum_{i=1}^N (y_i - \theta_1 x_i - \theta_0)(-x_i) = 0$$

$$\begin{aligned} \theta_0 &= \frac{1}{N} \sum y_i - \theta_1 \frac{1}{N} \sum x_i \\ \theta_1 &= \frac{1}{N} \sum y_i x_i - \theta_0 \frac{1}{N} \sum x_i \\ \theta_1 &= \frac{\sum y_i x_i - \frac{1}{N} \sum y_i \sum x_i}{\sum x_i^2 - \frac{1}{N} \sum x_i \sum x_i} \end{aligned}$$

Properties of the Solution

- Setting θ^* as before gives least squared error

- Define error on each data point as:
 $e_i = y_i - \theta^* x_i - \theta^*$

- Note property #1:

$$\frac{\partial R}{\partial \theta_0} = \frac{1}{N} \sum_{i=1}^N (y_i - \theta_1 x_i - \theta_0) = 0$$

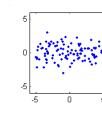
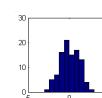
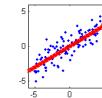
$$\dots \text{average error is zero} \quad \frac{1}{N} \sum e_i = 0$$

- Note property #2:

$$\frac{\partial R}{\partial \theta_1} = \frac{1}{N} \sum_{i=1}^N (y_i - \theta_1 x_i - \theta_0) x_i = 0$$

...error not correlated with data

$$\frac{1}{N} \sum e_i x_i = \frac{1}{N} e^T x = 0$$



Multi-Dimensional Regression

- More elegant/general to do $\nabla_{\theta} R = 0$ with linear algebra

- Rewrite empirical risk in vector-matrix notation:

$$\begin{aligned} R(\theta) &= \frac{1}{2N} \sum_{i=1}^N (y_i - \theta_1 x_i - \theta_0)^2 \\ &= \frac{1}{2N} \sum_{i=1}^N \left(y_i - \begin{bmatrix} 1 & x_i \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix} \right)^2 \\ &= \frac{1}{2N} \left\| \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} - \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix} \right\|^2 \\ &= \frac{1}{2N} \| \mathbf{y} - \mathbf{X}\theta \|^2 \end{aligned}$$

Can add more dimensions by adding columns to X matrix and rows to θ vector

Tony Jebara, Columbia University

Topic 3

Machine Learning

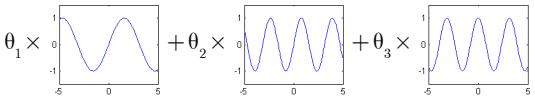
- Additive Models and Linear Regression
- Sinusoids and Radial Basis Functions
- Classification
- Logistic Regression
- Gradient Descent

Sinusoidal Basis Functions

- More generally, we don't just have to deal with polynomials, use any set of basis fn's:

$$f(x; \theta) = \sum_{p=1}^P \theta_p \phi_p(x) + \theta_0$$

- These are generally called Additive Models
- Regression adds linear combinations of the basis fn's
- For example: Fourier (sinusoidal) basis
 $\phi_{2k}(x_i) = \sin(kx_i)$ $\phi_{2k+1}(x_i) = \cos(kx_i)$
- Note, don't have to be a basis per se, usually subset



Evaluating Our Learned Function

- We minimized empirical risk to get θ^*
- How well does $f(x; \theta^*)$ perform on future data?
- It should *Generalize* and have low *True Risk*:

$$R_{\text{true}}(\theta) = \int P(x, y) L(y, f(x; \theta)) dx dy$$

- Can't compute true risk, instead use *Testing Empirical Risk*
- We randomly split data into training and testing portions

$$\{(x_1, y_1), \dots, (x_N, y_N)\} \quad \{(x_{N+1}, y_{N+1}), \dots, (x_{N+M}, y_{N+M})\}$$

- Find θ^* with *training data*: $R_{\text{train}}(\theta) = \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i; \theta))$

- Evaluate it with *testing data*: $R_{\text{test}}(\theta) = \frac{1}{M} \sum_{i=N+1}^{N+M} L(y_i, f(x_i; \theta))$

Radial Basis Functions

- Can act as prototypes of the data itself

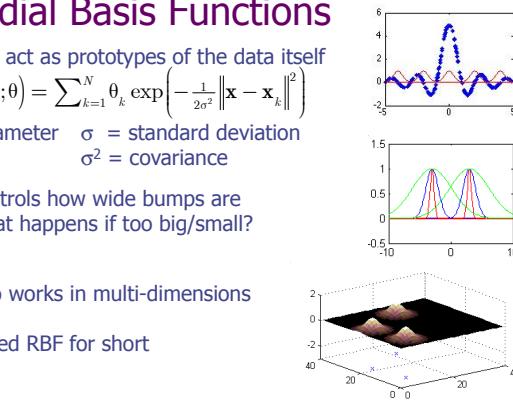
$$f(\mathbf{x}; \theta) = \sum_{k=1}^N \theta_k \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x} - \mathbf{x}_k\|^2\right)$$

- Parameter σ = standard deviation
 σ^2 = covariance

controls how wide bumps are
what happens if too big/small?

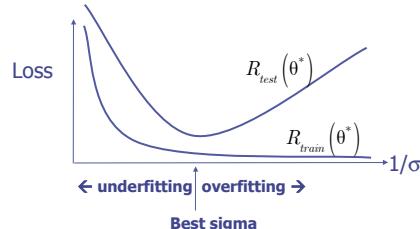
- Also works in multi-dimensions

- Called RBF for short



Crossvalidation

- Try fitting with different sigma radial basis function widths
- Select sigma which gives lowest $R_{\text{test}}(\theta^*)$

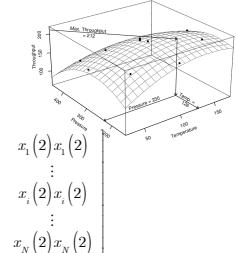


- Think of sigma as a measure of the simplicity of the model
Thinner RBFs are more flexible and complex

Polynomial Basis Functions

- To fit a P'th order polynomial function to multivariate data:
concatenate columns of all monomials up to power P
- E.g. 2 dimensional data and 2nd order polynomial (quadratic)

$$\begin{bmatrix} x_1(1) & x_1(2) \\ \vdots & \vdots \\ x_i(1) & x_i(2) \\ \vdots & \vdots \\ x_N(1) & x_N(2) \end{bmatrix} \xrightarrow{\text{Yellow Arrow}} \begin{bmatrix} 1 & x_1(1) & x_1(2) & x_1(1)x_1(1) & x_1(1)x_1(2) & x_1(2)x_1(2) \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_i(1) & x_i(2) & x_i(1)x_i(1) & x_i(1)x_i(2) & x_i(2)x_i(2) \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_N(1) & x_N(2) & x_N(1)x_N(1) & x_N(1)x_N(2) & x_N(2)x_N(2) \end{bmatrix}$$



Radial Basis Functions

- Each training point leads to a bump function

$$f(\mathbf{x}; \theta) = \sum_{k=1}^N \theta_k \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x} - \mathbf{x}_k\|^2\right)$$

- Reuse solution from linear regression: $\theta^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

- Can view the data instead as \mathbf{X} , a big matrix of size $N \times N$

$$\mathbf{X} = \begin{bmatrix} \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x}_1 - \mathbf{x}_1\|^2\right) & \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x}_1 - \mathbf{x}_2\|^2\right) & \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x}_1 - \mathbf{x}_3\|^2\right) \\ \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x}_2 - \mathbf{x}_1\|^2\right) & \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x}_2 - \mathbf{x}_2\|^2\right) & \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x}_2 - \mathbf{x}_3\|^2\right) \\ \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x}_3 - \mathbf{x}_1\|^2\right) & \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x}_3 - \mathbf{x}_2\|^2\right) & \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x}_3 - \mathbf{x}_3\|^2\right) \end{bmatrix}$$

- For RBFs, \mathbf{X} is square and symmetric, so solution is just

$$\nabla_{\theta} R = 0 \rightarrow \mathbf{X}^T \mathbf{X} \theta = \mathbf{X}^T \mathbf{y} \rightarrow \mathbf{X} \theta = \mathbf{y} \rightarrow \theta^* = \mathbf{X}^{-1} \mathbf{y}$$

Regularized Risk Minimization

- Empirical Risk Minimization gave overfitting & underfitting

- We want to add a penalty for using too many theta values

- This gives us the Regularized Risk

$$\begin{aligned} R_{\text{regularized}}(\theta) &= R_{\text{empirical}}(\theta) + \text{Penalty}(\theta) \\ &= \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i; \theta)) + \frac{\lambda}{2N} \|\theta\|^2 \end{aligned}$$

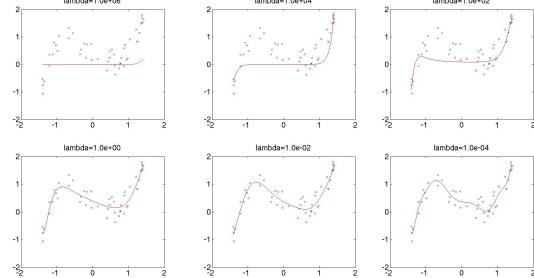
- Solution for Regularized Risk with Least Squares Loss:

$$\nabla_{\theta} R_{\text{regularized}} = 0 \Rightarrow \nabla_{\theta} \left(\frac{1}{2N} \|\mathbf{y} - \mathbf{X}\theta\|^2 + \frac{\lambda}{2N} \|\theta\|^2 \right) = 0$$

$$\theta^* = (\mathbf{X}^T \mathbf{X} + \lambda I)^{-1} \mathbf{X}^T \mathbf{y}$$

Regularized Risk Minimization

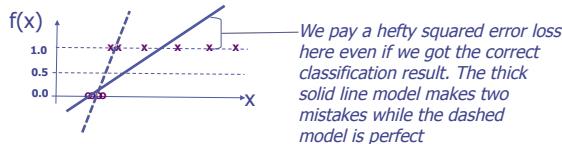
- Have D=16 features (or P=15 throughout)
- Try minimizing $R_{\text{regularized}}(\theta)$ to get θ^* with different λ
- Note that $\lambda=0$ give back Empirical Risk Minimization



Tony Jebara, Columbia University

Classification vs. Regression

- a) Classification needs binary answers like {0,1}
- b) Least squares is an unfair measure of risk here
 - e.g. Why penalize a correct but large positive y answer?
 - e.g. Why penalize a correct but large negative y answer?
- Example: not good to use regression output for a decision
 $f(x)>0.5 \rightarrow \text{Class 1}$ $f(x)<0.5 \rightarrow \text{Class 0}$
 if $f(x)=-3.8$ & correct class=0, squared error penalizes it...

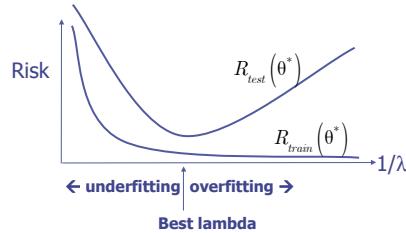


Logistic Regression

- Given a classification problem with binary outputs
 $\mathcal{X} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\} \quad \mathbf{x} \in \mathbb{R}^D \quad y \in \{0,1\}$
- Use this function and output 1 if $f(\mathbf{x})>0.5$ and 0 otherwise
 $f(\mathbf{x}; \theta) = (1 + \exp(-\theta^T \mathbf{x}))^{-1}$

Crossvalidation

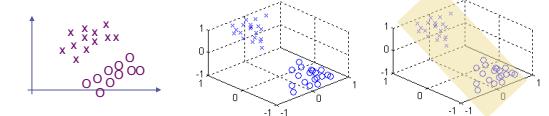
- Try fitting with different lambda regularization levels
- Select lambda which gives lowest $R_{\text{test}}(\theta^*)$



- Lambda measures simplicity of the model
- Models with low lambda are more flexible

From Regression To Classification

- Classification is another important learning problem
- Regression $\mathcal{X} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\} \quad \mathbf{x} \in \mathbb{R}^D \quad y \in \mathbb{R}^1$
- Classification $\mathcal{X} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\} \quad \mathbf{x} \in \mathbb{R}^D \quad y \in \{0,1\}$
- E.g. Given $\mathbf{x} = [\text{tumor size, tumor density}]$
 Predict y in {benign, malignant}
- Should we solve this as a least squares regression problem?

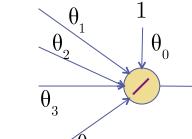
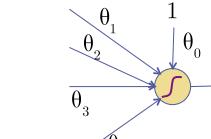


Classification vs. Regression

We will consider the following four steps to improve from naïve regression to get better classification learning:

- 1) Fix functions $f(\mathbf{x})$ to give binary output (logistic neuron)
- 2) Fix our definition of the Risk we will minimize so that we get good classification accuracy (logistic loss)
- ...and later on...
- 3) Make an even better fix on $f(\mathbf{x})$ to binarize (perceptron)
- 4) Make an even better risk (perceptron loss)

Logistic Neuron (McCullough-Pitts)

- To output binary, use squashing function $g()$.
- $f(\mathbf{x}; \theta) = \theta^T \mathbf{x}$
- $g(z) = (1 + \exp(-z))^{-1}$
- 
- 
- Linear neuron
- Logistic Neuron
- This squashing is called sigmoid or logistic function

Short hand for Linear Functions

- What happened to adding the intercept?

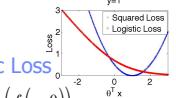
$$f(\mathbf{x}; \theta) = \theta^T \mathbf{x} + \theta_0$$

$$= \begin{bmatrix} \theta(0) \\ \theta(2) \\ \vdots \\ \theta(D) \end{bmatrix}^T \begin{bmatrix} \mathbf{x}(1) \\ \mathbf{x}(2) \\ \vdots \\ \mathbf{x}(D) \end{bmatrix} + \theta_0 = \begin{bmatrix} \theta_0 \\ \theta(1) \\ \theta(2) \\ \vdots \\ \theta(D) \end{bmatrix}^T \begin{bmatrix} 1 \\ \mathbf{x}(1) \\ \mathbf{x}(2) \\ \vdots \\ \mathbf{x}(D) \end{bmatrix} = \vec{\theta}^T \vec{\mathbf{x}}$$

Logistic Regression

- Given a classification problem with binary outputs
 $\mathcal{X} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\} \quad \mathbf{x} \in \mathbb{R}^D \quad y \in \{0,1\}$
- Fix#1: use $f(\mathbf{x})$ below, output 1 if $f(\mathbf{x})>0.5$ and 0 otherwise
 $f(\mathbf{x}; \theta) = (1 + \exp(-\theta^T \mathbf{x}))^{-1}$
- Fix#2: instead of squared loss, use Logistic Loss
 $L_{\text{log}}(y_i, f(\mathbf{x}; \theta)) = (y_i - 1) \log(1 - f(\mathbf{x}; \theta)) - y_i \log(f(\mathbf{x}; \theta))$
- This method is called Logistic Regression.
- But Empirical Risk Minimization has no closed-form sol'n:

$$R_{\text{emp}}(\theta) = \frac{1}{N} \sum_{i=1}^N (y_i - 1) \log(1 - f(\mathbf{x}_i; \theta)) - y_i \log(f(\mathbf{x}_i; \theta))$$



Logistic Regression

- With logistic squashing function, minimizing $R(\theta)$ is harder

$$R_{emp}(\theta) = \frac{1}{N} \sum_{i=1}^N (y_i - 1) \log(1 - f(\mathbf{x}_i; \theta)) - y_i \log(f(\mathbf{x}_i; \theta))$$

$$\nabla_{\theta} R = \frac{1}{N} \sum_{i=1}^N \left(\frac{1 - y_i}{1 - f(\mathbf{x}_i; \theta)} - \frac{y_i}{f(\mathbf{x}_i; \theta)} \right) f'(\mathbf{x}_i; \theta) = 0 \quad ???$$

- Can't minimize risk and find best theta analytically!
- Let's try finding best theta numerically.
- Use the following to compute gradient

$$f(\mathbf{x}; \theta) = (1 + \exp(-\theta^T \mathbf{x}))^{-1} = g(\theta^T \mathbf{x})$$

- Here, $g()$ is the logistic squashing function

$$g(z) = (1 + \exp(-z))^{-1} \quad g'(z) = g(z)(1 - g(z))$$

Gradient Descent

- Useful when we can't get minimum solution in closed form
- Gradient points in direction of fastest increase
- Take step in the opposite direction!

- Gradient Descent Algorithm

choose scalar step size η , & tolerance ε
initialize $\theta^0 = \text{small random vector}$

$$\begin{aligned} \theta^1 &= \theta^0 - \eta \nabla_{\theta} R_{emp}|_{\theta^0}, \quad t = 1 \\ \text{while } \|\theta^t - \theta^{t-1}\| &\geq \varepsilon \quad \{ \\ \theta^{t+1} &= \theta^t - \eta \nabla_{\theta} R_{emp}|_{\theta^t}, \quad t = t + 1 \end{aligned} \quad \}$$

- For appropriate η , this will converge to local minimum

Logistic Regression

- Logistic regression gives better classification performance

- Its empirical risk is

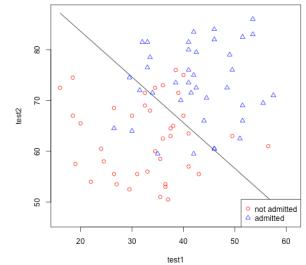
$$R_{emp}(\theta) = \frac{1}{N} \sum_{i=1}^N (y_i - 1) \log(1 - f(\mathbf{x}_i; \theta)) - y_i \log(f(\mathbf{x}_i; \theta))$$

- This $R(\theta)$ is convex so gradient descent always converges to the same solution

- Make predictions using

$$f(\mathbf{x}; \theta) = (1 + \exp(-\theta^T \mathbf{x}))^{-1}$$

- Output 1 if $f > 0.5$
- Output 0 otherwise



Topic 4

Machine Learning

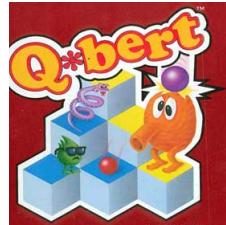
- Perceptron, Online & Stochastic Gradient Descent
- Convergence Guarantee
- Perceptron vs. Linear Regression
- Multi-Layer Neural Networks
- Back-Propagation
- Demo: LeNet
- Deep Learning

Perceptron & Classification Loss

- Classification loss for the Risk leads to hard minimization
- What does this $R(\theta)$ function look like?

$$R(\theta) = \frac{1}{N} \sum_{i=1}^N \text{step}(-y_i \theta^T x_i)$$

- Qbert-like, can't do gradient descent since the gradient is zero except at edges when a label flips



Perceptron & Perceptron Loss

- Instead of Classification Loss

$$R(\theta) = \frac{1}{N} \sum_{i=1}^N \text{step}(-y_i \theta^T x_i)$$

- Consider Perceptron Loss:

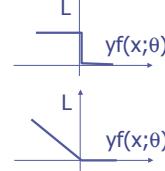
$$R^{per}(\theta) = -\frac{1}{N} \sum_{i \in \text{misclassified}} y_i (\theta^T x_i)$$

- Instead of staircase-shaped R get smooth piece-wise linear R

- Get reasonable gradients for gradient descent

$$\nabla_{\theta} R^{per}(\theta) = -\frac{1}{N} \sum_{i \in \text{misclassified}} y_i \mathbf{x}_i$$

$$\theta^{t+1} = \theta^t - \eta \nabla_{\theta} R^{per} \Big|_{\theta^t} = \theta^t + \eta \frac{1}{N} \sum_{i \in \text{misclassified}} y_i \mathbf{x}_i$$



Stochastic Gradient Descent

- Gradient Descent vs. Stochastic Gradient Descent

- Instead of computing the average gradient for all points and then taking a step

$$\nabla_{\theta} R^{per}(\theta) = -\frac{1}{N} \sum_{i \in \text{misclassified}} y_i \mathbf{x}_i$$

- Update the gradient for each mis-classified point by itself

$$\nabla_{\theta} R^{per}(\theta) = -y_i \mathbf{x}_i \quad \text{if } i \text{ mis-classified}$$

- Also, set η to 1 without loss of generality

$$\theta^{t+1} = \theta^t - \eta \nabla_{\theta} R^{per} \Big|_{\theta^t} = \theta^t + y_i \mathbf{x}_i \quad \text{if } i \text{ mis-classified}$$

Online Perceptron

- Apply stochastic gradient descent to a perceptron
- Get the "online perceptron" algorithm:

```

initialize t = 0 and θ⁰ = 0
while not converged {
    pick i ∈ {1,...,N}
    if (yᵢ xᵢ¹ θ⁰ ≤ 0)   {   θ¹⁺ = θ⁰ + yᵢ xᵢ
                                t = t + 1           }

```

- Either pick i randomly or use a "for $i=1$ to N " loop
- If the algorithm stops, we have a theta that separates data
- The total number of mistakes we made along the way is t

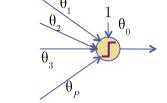
Perceptron (another Neuron)

- Classification scenario once again but consider +1, -1 labels

$$\mathcal{X} = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\} \quad x \in \mathbb{R}^D \quad y \in \{-1, 1\}$$

- A better choice for a classification squashing function is

$$g(z) = \begin{cases} -1 & \text{when } z < 0 \\ +1 & \text{when } z \geq 0 \end{cases}$$



- And a better choice is classification loss

$$L(y, f(\mathbf{x}; \theta)) = \text{step}(-y f(\mathbf{x}; \theta))$$

- Actually with above $g(z)$ any loss is like classification loss

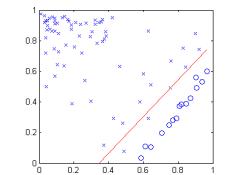
$$R(\theta) = \frac{1}{4N} \sum_{i=1}^N (y_i - g(\theta^T x_i))^2 \equiv \frac{1}{N} \sum_{i=1}^N \text{step}(-y_i \theta^T x_i)$$

- What does this $R(\theta)$ function look like?

Perceptron vs. Linear Regression

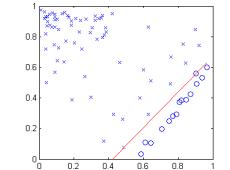
- Linear regression gets close but doesn't do perfectly

classification error = 2
squared error = 0.139



- Perceptron gets zero error

classification error = 0
perceptron err = 0



Online Perceptron Theorem

Theorem: the online perceptron algorithm converges to zero error in finite t if we assume

- 1) all data inside a sphere of radius r : $\|\mathbf{x}_i\| \leq r \quad \forall i$
- 2) data is separable with margin γ : $y_i (\theta^*)^T \mathbf{x}_i \geq \gamma \quad \forall i$

Proof:

- Part 1) Look at inner product of current θ^t with θ^*
assume we just updated a mistake on point i :

$$(\theta^*)^T \theta' = (\theta^*)^T \theta^{t-1} + y_i (\theta^*)^T \mathbf{x}_i \geq (\theta^*)^T \theta^{t-1} + \gamma$$

after applying t such updates, we must get:

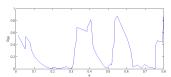
$$(\theta^*)^T \theta' = (\theta^*)^T \theta^t \geq t\gamma$$

Neural Networks Demo

- Again, take small step in direction opposite to gradient

Digits Demo: LeNet... <http://yann.lecun.com>

- Problems with back-prop
is that MLP over-fits...

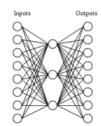


- Other problems: hard to interpret, black-box

- What are the hidden inner layers doing?

Auto-Encoders

A network:

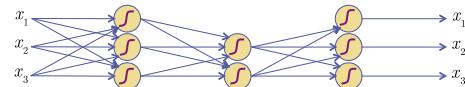


Learned hidden layer representation:

Input	Hidden Values	Output
10000000 → .89 .04 .08		→ 10000000
01000000 → .01 .11 .88		→ 01000000
00100000 → .01 .97 .27		→ 00100000
00010000 → .99 .97 .71		→ 00010000
00001000 → .03 .05 .02		→ 00001000
00000100 → .22 .99 .99		→ 00000100
00000010 → .80 .01 .98		→ 00000010
00000001 → .60 .94 .01		→ 00000001

Auto-Encoders

- Make the neural net reconstruct the input vector



- Set the target y vector to be the x vector again

- But, it gets narrow in the middle!

- So, there is some information "compression"

- This leads to better generalization

- This is unsupervised learning since we only use the x data

Auto-Encoders



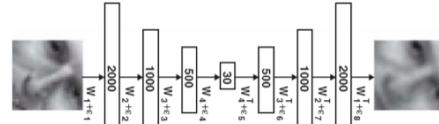
A target function:

Input	Output
10000000 → 10000000	
01000000 → 01000000	
00100000 → 00100000	
00010000 → 00010000	
00001000 → 00001000	
00000100 → 00000100	
00000010 → 00000010	
00000001 → 00000001	

Can this be learned??

Deep Learning

- We can stack several independently trained auto-encoders



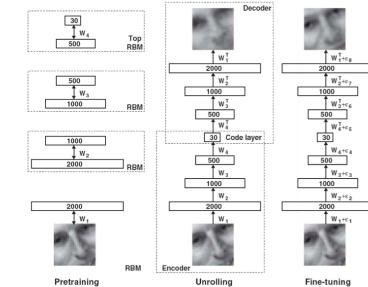
- Using back-propagation, we do *pre-training*

- Train Net 1 to go from 2000 inputs to 1000 to 2000 inputs
- Train Net 2 to go from 1000 hidden values to 500 to 1000
- Train Net 3 to go from 500 hidden to 30 to 500

- Then, do *unrolling* - link up the learned networks as above

Deep Learning

- Then do *fine-tuning* of the overall neural network by running back-propagation on the whole thing to reconstruct x from x



Deep Learning

- Does good reconstruction!
- Beats PCA on images of unaligned faces.
- PCA is better when face images are aligned...
- We will cover PCA in a few lectures...



Minimum Training Error?

- Is minimizing Empirical Risk the right thing?
- Are Perceptrons and Neural Networks giving the best classifier?

- We are getting: minimum training error
not minimum testing error

- Perceptrons are giving a bunch of solutions:



... a solution with *guarantees* → SVMs

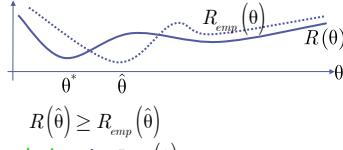
Machine Learning

Topic 5

- Generalization Guarantees
- VC-Dimension
- Nearest Neighbor Classification (infinite VC dimension)
- Structural Risk Minimization
- Support Vector Machines

Bounding the True Risk

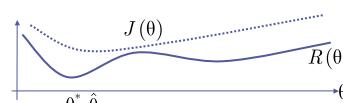
- ERM is inconsistent
not guaranteed
may do better
on training than
on test!



Idea: add a prior or regularizer to $R_{emp}(\theta)$

- Define capacity or confidence = $C(\theta)$ which favors simpler θ

$$J(\theta) = R_{emp}(\theta) + C(\theta)$$



- If $R(\theta) \leq J(\theta)$ we have bound $J(\theta)$ is a guaranteed risk
- After train, can guarantee future error rate is $\leq \min_{\theta} J(\theta)$

Bound the True Risk with VC

- But, how to find a guarantee? Difficult, but there is one...

- Theorem (Vapnik):** with probability $1-\eta$ where η is a number between $[0,1]$, the following bound holds:

$$R(\theta) \leq J(\theta) = R_{emp}(\theta) + \frac{2h \log(\frac{2eN}{h}) + 2 \log(\frac{1}{\eta})}{N} \left[1 + \sqrt{1 + \frac{NR_{emp}(\theta)}{h \log(\frac{2eN}{h}) + \log(\frac{1}{\eta})}} \right]$$

N = number of data points

h = Vapnik-Chervonenkis (VC) dimension (1970's)
= capacity of the classifier class $f(\cdot; \theta)$

- Note, above is independent of the true $P(x,y)$

- A worst-case scenario bound, guaranteed for all $P(x,y)$
- VC dimension not just the # of parameters a classifier has
- VC measures # of different datasets it can classify perfectly
- Structural Risk Minimization: minimize risk bound $J(\theta)$

VC Dimension & Shattering

- More generally for higher dimensional linear classifiers, a hyperplane in \mathbb{R}^d shatters any set of linearly independent points. Can choose $d+1$ linearly indep. points so $h=d+1$

- Note: VC is not necessarily proportional to # of parameters

- Example: sinusoidal 1d classifier $f(x; \theta) = \text{sign}(\sin(\theta x))$

number of parameters=1

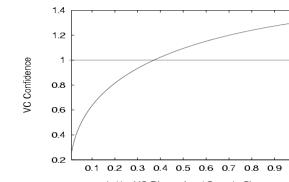
...but... $h=\infty$!

since I can choose: $x_i = 10^{-i}$ $i = 1, \dots, h$
no matter what labeling you challenge: $y_1, \dots, y_h \in \{\pm 1\}^h$
using $\theta = \pi \left(1 + \sum_{i=1}^h \frac{1}{2} (1 - y_i) 10^{-i} \right)$ shatters perfectly

But, as a side note, if I choose 4 equally spaced x's then cannot shatter

VC Dimension & Shattering

- Recall that VC dimension gives an upper bound
- We want to minimize h since that minimizes $C(\theta)$ & $J(\theta)$
- If can't compute h exactly but can compute h^+ can plug in h^+ in bound & still guarantee
- Also, sometimes bound is trivial
- Need $h/N = 0.3$ before $C(\theta) < 1$ (recall $R(\theta)$ in $[0,1]$)



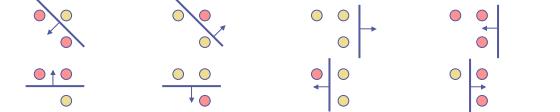
- Note: $h = \text{low} \Rightarrow \text{good performance}$ $h = \infty \times \text{poor performance}$ $h = \text{low} \Rightarrow \text{good performance}$

Empirical Risk Minimization

- Example: non-pdf linear classifiers $f(x; \theta) = \text{sign}(\theta^T x + \theta_0) \in \{-1, 1\}$
- Recall ERM: $R_{emp}(\theta) = \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i; \theta)) \in [0, 1]$
- Have loss function: quadratic: $L(y, x, \theta) = \frac{1}{2} |y - f(x; \theta)|^2$
linear: $L(y, x, \theta) = |y - f(x; \theta)|$
binary: $L(y, x, \theta) = \text{step}(-yf(x; \theta))$
- Empirical $R_{emp}(\theta)$ approximates the true risk (expected error)
 $R(\theta) = E_P \{L(x, y, \theta)\} = \int_{X \times Y} P(x, y) L(x, y, \theta) dx dy \in [0, 1]$
- But, we don't know the true $P(x,y)!$
- If infinite data, law of large numbers says:
 $\lim_{N \rightarrow \infty} \min_{\theta} R_{emp}(\theta) = \min_{\theta} R(\theta)$
- But, in general, can't make guarantees for ERM solution:
 $\arg \min_{\theta} R_{emp}(\theta) \neq \arg \min_{\theta} R(\theta)$

VC Dimension & Shattering

- How to compute h or VC for a family of functions $f(\cdot; \theta)$
 $h = \# \text{ of training points that can be shattered}$
- Recall, classifier maps input to output $f(x; \theta) \rightarrow y \in \{-1, 1\}$
- Shattering:** I pick h points & place them at x_1, \dots, x_h . You challenge me with 2^h possible labelings $y_1, \dots, y_h \in \{\pm 1\}^h$. VC dimension is maximum # of points I can place which a $f(x; \theta)$ can correctly classify for arbitrary labeling y_1, \dots, y_h
- Example: for 2d linear classifier $h=3$ $f(x; \theta) = x_1 \theta_1 + x_2 \theta_2 + \theta_0$

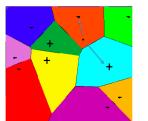


can't ever shatter 4 points! or 3 points on a straight line...

Nearest Neighbors & VC

- Consider Nearest Neighbors classification algorithm:

Input a query example x
Find training example x_i in $\{x_1, \dots, x_N\}$ closest to x
Predict label for x as y_i of neighbor

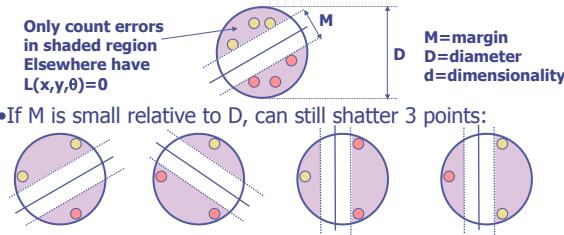


- Often use Euclidean distance $\|x - x_i\|$ to measure closeness
- Nearest neighbors shatters any set of points!
- So VC=infinity, $C(\theta)=\infty$, guaranteed risk=infinity
- But still works well in practice

$h = \infty \times \text{poor performance}$ $h = \text{low} \Rightarrow \text{good performance}$

VC Dimension & Large Margins

- Linear classifiers are too big a function class since $h=d+1$
- Can reduce VC dimension if we restrict them
- Constrain linear classifiers to data living inside a sphere
- Gap-Tolerant classifiers:** a linear classifier whose activity is constrained to a sphere & outside a margin



VC Dimension & Large Margins

- But, as M grows relative to D, can only shatter 2 points!
- Can't shatter 3
Can shatter 2
- For hyperplanes, as M grows vs. D, shatter fewer points!
- VC dimension h goes down if gap-tolerant classifier has larger margin, general formula is: $h \leq \min \left\{ \text{ceil} \left[\frac{D^2}{M^2} \right], d \right\} + 1$
- Before, just had $h=d+1$. Now we have a smaller h
- If data is anywhere, D is infinite and back to $h=d+1$
- Typically real data is bounded (by sphere), D is fixed
- Maximizing M reduces h, improving guaranteed risk $J(\theta)$
- Note: $R(\theta)$ doesn't count errors in margin or outside sphere

Support Vector Machines

- Support vector machines are (in the simplest case) linear classifiers that do structural risk minimization (SRM)
- Directly maximize margin to reduce guaranteed risk $J(\theta)$
- Assume first the 2-class data is linearly separable:
 - have $\{(x_1, y_1), \dots, (x_N, y_N)\}$ where $x_i \in \mathbb{R}^D$ and $y_i \in \{-1, 1\}$
 - $f(x; \theta) = \text{sign}(w^T x + b)$
- Decision boundary or hyperplane given by $w^T x + b = 0$
- Note: can scale w & b while keeping same boundary
- Many solutions exist which have empirical error $R_{\text{emp}}(\theta) = 0$
- Want widest or thickest one (max margin), also it's unique!



Support Vector Machines

Define: $w^T x + b = 0$

H_+ = positive margin hyperplane
 H_- = negative margin hyperplane
 q = distance from decision plane to origin

$$q = \min_x \|\vec{x} - \vec{0}\| \quad \text{subject to } w^T x + b = 0$$

$$\min_x \frac{1}{2} \|\vec{x} - \vec{0}\|^2 - \lambda(w^T x + b)$$

$$1) \text{ grad } \frac{\partial}{\partial x} \left(\frac{1}{2} \vec{x}^T \vec{x} - \lambda(w^T x + b) \right) = 0 \quad 2) \text{ plug into constraint } w^T x + b = 0$$

$$x - \lambda w = 0 \quad x = \lambda w \quad \lambda = -\frac{b}{w^T w}$$

$$3) \text{ Sol'n } \hat{x} = -\left(\frac{b}{w^T w}\right)w$$

$$4) \text{ distance } q = \|\hat{x} - \vec{0}\| = \left\| -\frac{b}{w^T w} w \right\| = \frac{|b|}{w^T w} \sqrt{w^T w} = \frac{|b|}{\|w\|}$$

$$5) \text{ Define without loss of generality } H \rightarrow w^T x + b = 0 \quad H_+ \rightarrow w^T x + b = +1 \quad H_- \rightarrow w^T x + b = -1$$

Support Vector Machines

- The constraints on the SVM for $R_{\text{emp}}(\theta)=0$ are thus:

$$w^T x_i + b \geq +1 \quad \forall y_i = +1$$

$$w^T x_i + b \leq -1 \quad \forall y_i = -1$$
- Or more simply: $y_i(w^T x_i + b) - 1 \geq 0$
- The margin of the SVM is:

$$m = d_+ + d_-$$
- Distance to origin: $H \rightarrow q = \frac{\|H\|}{\|w\|}$ $H_+ \rightarrow q_+ = \frac{\|H\|}{\|w\|}$ $H_- \rightarrow q_- = \frac{\|H\|}{\|w\|}$
- Therefore: $d_+ = d_- = \frac{1}{\|w\|}$ and margin $m = \frac{2}{\|w\|}$
- Want to max margin, or equivalently minimize: $\|w\|$ or $\|w\|^2$
- SVM Problem: $\min_w \frac{1}{2} \|w\|^2 \quad \text{subject to } y_i(w^T x_i + b) - 1 \geq 0$
- This is a quadratic program!
- Can plug this into a matlab function called "qp()", done!

Structural Risk Minimization

- Structural Risk Minimization:** minimize risk bound $J(\theta)$ reducing empirical error & reduce VC dimension h

$$R(\theta) \leq J(\theta) = R_{\text{emp}}(\theta) + \frac{2h \log(\frac{2eN}{h}) + 2 \log(\frac{4}{\eta})}{N} \left[1 + \sqrt{1 + \frac{NR_{\text{emp}}(\theta)}{h \log(\frac{2eN}{h}) + \log(\frac{4}{\eta})}} \right]$$

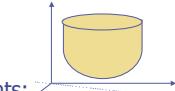
for each model i in list of hypothesis

- compute its $h = h_i$
 - $\theta^* = \arg \min_{\theta} R_{\text{emp}}(\theta)$
 - compute $J(\theta^*, h_i)$
- choose model with lowest $J(\theta^*, h_i)$

Or, directly optimize over both $(\theta^*, h) = \arg \min_{\theta, h} J(\theta, h)$

If possible, min empirical error while also minimizing VC

For gap-tolerant linear classifiers, minimize $R_{\text{emp}}(\theta)$ while maximizing margin, support vector machines do just that!

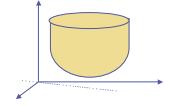


Side Note: Constraints

- Cost minimization with equality constraints:
 - Subtract each constraint times an extra variable (a Lagrange multiplier λ , like an adversary variable)

2) Take partials with respect to x and set to zero

3) Plug solution into constraint to find lambda



$$\begin{aligned} \min_{\vec{x}} \max_{\lambda} f(\vec{x}) - \lambda (\text{equality condition} = 0) \\ = \min_{\vec{x}} \max_{\lambda} b^T \vec{x} + \frac{1}{2} \vec{x}^T H \vec{x} - \lambda(x_1 - 2x_2) \\ \Rightarrow \frac{\partial f}{\partial \vec{x}} = \vec{b} + H \vec{x} - \lambda \begin{bmatrix} 1 \\ -2 \end{bmatrix} = 0 \Rightarrow \vec{x} = H^{-1} \lambda \begin{bmatrix} 1 \\ -2 \end{bmatrix} - H^{-1} \vec{b} \\ \Rightarrow \left(H^{-1} \lambda \begin{bmatrix} 1 \\ -2 \end{bmatrix} - H^{-1} \vec{b} \right)^T \begin{bmatrix} 1 \\ -2 \end{bmatrix} = 0 \Rightarrow \lambda = \frac{\vec{b}^T H^{-1} \begin{bmatrix} 1 \\ -2 \end{bmatrix}}{\left\| \begin{bmatrix} 1 \\ -2 \end{bmatrix} \right\|^2} \end{aligned}$$

Side Note: Optimization Tools

- A hierarchy of Matlab optimization packages to use:

Linear Programming $\min_{\vec{x}} \vec{b}^T \vec{x} \quad \text{s.t. } \vec{c}_i^T \vec{x} \geq \alpha_i \quad \forall i$

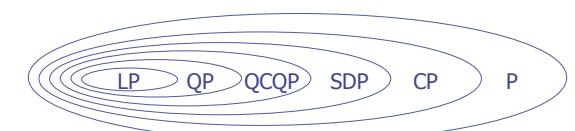
<Quadratic Programming $\min_{\vec{x}} \frac{1}{2} \vec{x}^T H \vec{x} + \vec{b}^T \vec{x} \quad \text{s.t. } \vec{c}_i^T \vec{x} \geq \alpha_i \quad \forall i$

<Quadratically Constrained Quadratic Programming

<Semidefinite Programming

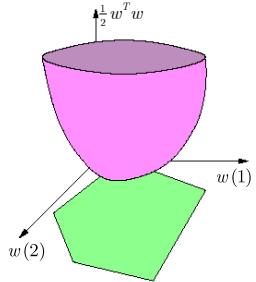
<Convex Programming

<Polynomial Time Algorithms



Side Note: Optimization Tools

- Each data point adds $y_i(w^T x_i + b) - 1 \geq 0$ linear inequality to QP
- Each point cuts a half plane of allowable SVMs and reduces green region
- The SVM is closest point to the origin that is still in the green region
- The perceptron algorithm just puts us randomly in green region
- QP runs in cubic polynomial time
- There are D values in the w vector
- Needs $O(D^3)$ run time... But, there is a DUAL SVM in $O(N^3)$!

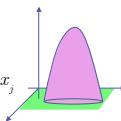
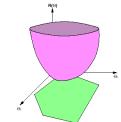


SVM in Dual Form

- We can also solve the problem via convex duality
- Primal SVM problem $L_P: \min_{w,b} \frac{1}{2} \|w\|^2$ subject to $y_i(w^T x_i + b) - 1 \geq 0$
- This is a quadratic program, quadratic cost function with multiple linear inequalities (these carve out a convex hull)
- Subtract from cost each inequality times an α : Lagrange multiplier, take derivatives of w & b:
$$L_P = \min_{w,b} \max_{\alpha \geq 0} \frac{1}{2} \|w\|^2 - \sum_i \alpha_i (y_i(w^T x_i + b) - 1)$$

$$\frac{\partial}{\partial w} L_P = w - \sum_i \alpha_i y_i x_i = 0 \rightarrow w = \sum_i \alpha_i y_i x_i$$

$$\frac{\partial}{\partial b} L_P = -\sum_i \alpha_i y_i = 0$$
- Plug back in, dual: $L_D = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j x_i^T x_j$
- Also have constraints: $\sum_i \alpha_i y_i = 0$ & $\alpha_i \geq 0$
- Above L_D must be maximized! convex duality... also qp()



Machine Learning

Topic 6

- Review: Support Vector Machines
- Primal & Dual Solution
- Non-separable SVMs
- Kernels
- SVM Demo

Support Vector Machines

• Define: $w^T x + b = 0$

H_+ = positive margin hyperplane
 H_- = negative margin hyperplane
 q = distance from decision plane to origin

$$q = \min_x \|\vec{x} - \vec{0}\|^2 \text{ subject to } w^T x + b = 0$$

$$\min_x \frac{1}{2} \|\vec{x} - \vec{0}\|^2 - \lambda(w^T x + b)$$

1) grad $\frac{\partial}{\partial x} \left(\frac{1}{2} x^T x - \lambda(w^T x + b) \right) = 0$ 2) plug into constraint $w^T x + b = 0$

$$x - \lambda w = 0$$

3) Sol'n $\hat{x} = -\left(\frac{b}{w^T w}\right)w$ $x = \lambda w$ $\lambda = -\frac{b}{w^T w}$

4) distance $q = \|\hat{x} - \vec{0}\| = \left\| -\frac{b}{w^T w}w \right\| = \frac{1}{w^T w} \sqrt{w^T w} = \frac{|b|}{\|w\|}$

5) Define without loss of generality since can scale b & w $H \rightarrow w^T x + b = 0$ $H_+ \rightarrow w^T x + b = +1$ $H_- \rightarrow w^T x + b = -1$

Support Vector Machines

- The constraints on the SVM for $R_{\text{emp}}(\theta)=0$ are thus:
- $$w^T x_i + b \geq +1 \quad \forall y_i = +1$$
- $$w^T x_i + b \leq -1 \quad \forall y_i = -1$$
- Or more simply: $y_i(w^T x_i + b) - 1 \geq 0$
- The margin of the SVM is:
- $$m = d_+ + d_-$$
- Distance to origin: $H \rightarrow q = \frac{\|w\|}{\|H\|}$ $H_+ \rightarrow q_+ = \frac{\|w\|}{\|H\|}$ $H_- \rightarrow q_- = \frac{\|w\|}{\|H\|}$
- Therefore: $d_+ = d_- = \frac{1}{\|w\|}$ and margin $m = \frac{2}{\|w\|}$
- Want to max margin, or equivalently minimize: $\|w\|$ or $\|w\|^2$
- SVM Problem: $\min \frac{1}{2} \|w\|^2 \text{ subject to } y_i(w^T x_i + b) - 1 \geq 0$
- This is a quadratic program!
- Can plug this into a matlab function called "qp()", done!
-

SVM Dual Solution Properties

- We have dual convex program:
- $$\sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i^T x_j \text{ subject to } \sum_i \alpha_i y_i = 0 \quad \& \quad \alpha_i \geq 0$$
- Solve for N alphas (one per data point) instead of D w's
- Still convex (qp) so unique solution, gives alphas
- Alphas can be used to get w: $w = \sum_i \alpha_i y_i x_i$
- Support Vectors: have non-zero alphas shown with thicker circles, all live on the margin: $w^T x_i + b = \pm 1$
- Solution is sparse, most alphas=0 these are non-support vectors SVM ignores them if they move (without crossing margin) or if they are deleted, SVM doesn't change (stays robust)
-

SVM Dual Solution Properties

- Primal & Dual Illustration:
-
- Recall we could get w from alphas: $w = \sum_i \alpha_i y_i x_i$
- Or could use as is: $f(x) = \text{sign}(w^T x + b) = \text{sign}(\sum_i \alpha_i y_i x_i^T x + b)$
- Karush-Kuhn-Tucker Conditions (KKT): solve value of b on margin (for nonzero alphas) have: $w^T x_i + b = y_i$ using known w, compute b for each support vector $\bar{b}_i = y_i - w^T x_i \quad \forall i : \alpha_i > 0$ then... $b = \text{average}(\bar{b}_i)$
- Sparsity (few nonzero alphas) is useful for several reasons
- Means SVM only uses some of training data to learn
- Should help improve its ability to generalize to test data
- Computationally faster when using final learned classifier

Review: SVM

- Support vector machines are (in the simplest case) linear classifiers that do structural risk minimization (SRM)
 - Directly maximize margin to reduce guaranteed risk $J(\theta)$
 - Assume first the 2-class data is linearly separable: have $\{(x_1, y_1), \dots, (x_N, y_N)\}$ where $x_i \in \mathbb{R}^D$ and $y_i \in \{-1, 1\}$
 $f(x; \theta) = \text{sign}(w^T x + b)$
 - Decision boundary or hyperplane given by $w^T x + b = 0$
 - Note: can scale w & b while keeping same boundary
 - Many solutions exist which have empirical error $R_{\text{emp}}(\theta) = 0$
 - Want widest or thickest one (max margin), also it's unique!
-

SVM in Dual Form

- We can also solve the problem via convex duality
- Primal SVM problem $L_p: \min \frac{1}{2} \|w\|^2 \text{ subject to } y_i(w^T x_i + b) - 1 \geq 0$
- This is a quadratic program, quadratic cost function with multiple linear inequalities (these carve out a convex hull)
- Subtract from cost each inequality times an α Lagrange multiplier, take derivatives of w & b:
- $$L_p = \min_{w, b} \max_{\alpha \geq 0} \frac{1}{2} \|w\|^2 - \sum_i \alpha_i (y_i(w^T x_i + b) - 1)$$
- $$\frac{\partial}{\partial w} L_p = w - \sum_i \alpha_i y_i x_i = 0 \rightarrow w = \sum_i \alpha_i y_i x_i$$
- $$\frac{\partial}{\partial b} L_p = -\sum_i \alpha_i y_i = 0$$
- Plug back in, dual: $L_d = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j x_i^T x_j$
- Also have constraints: $\sum_i \alpha_i y_i = 0$ & $\alpha_i \geq 0$
- Above L_d must be maximized! convex duality... also qp()
-

Non-Separable SVMs

- What happens when non-separable?
- There is no solution and convex hull shrinks to nothing
-
- Not all constraints can be resolved, their alphas go to ∞
- Instead of perfectly classifying each point: $y_i(w^T x_i + b) \geq 1$ we "Relax" the problem with (positive) slack variables ξ_i 's allow data to (sometimes) fall on wrong side, for example:
- $$w^T x_i + b \geq -0.03 \quad \text{if } y_i = +1$$
- New constraints: $w^T x_i + b \geq 1 - \xi_i \quad \text{if } y_i = +1 \quad \text{where } \xi_i \geq 0$

$$w^T x_i + b \leq -1 + \xi_i \quad \text{if } y_i = -1 \quad \text{where } \xi_i \geq 0$$
- But too much ξ_i 's means too much slack, so penalize them
- $$L_p : \min \frac{1}{2} \|w\|^2 + C \sum_i \xi_i \text{ subject to } y_i(w^T x_i + b) - 1 + \xi_i \geq 0$$

Non-Separable SVMs

- This new problem is still convex, still $\text{qp}()$!
- User chooses scalar C (or cross-validates) which controls how much slack ξ_i to use (how non-separable) and how robust to outliers or bad points on the wrong side

Large margin ↗ Low slack ↗ On right side ↗ For x_i positivity ↗

$$L_p : \min_{\frac{1}{2}\|w\|^2 + C \sum_i \xi_i - \sum_i \alpha_i (y_i(w^T x_i + b) - 1 + \xi_i) - \sum_i \beta_i \xi_i}$$

$$\frac{\partial}{\partial b} L_p \text{ and } \frac{\partial}{\partial w} L_p \text{ as before...}$$

$$\frac{\partial}{\partial \xi_i} L_p = C - \alpha_i - \beta_i = 0$$

$$\alpha_i = C - \beta_i \text{ but... } \alpha_i \text{ and } \beta_i \geq 0$$

$$\therefore 0 \leq \alpha_i \leq C$$

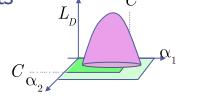
- Can now write dual problem (to maximize):

$$L_D : \max \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i^T x_j \text{ subject to } \sum_i \alpha_i y_i = 0 \text{ and } \alpha_i \in [0, C]$$

- Same dual as before but alphas can't grow beyond C

Non-Separable SVMs

- As we try to enforce a classification for a data point its Lagrange multiplier alpha keeps growing endlessly
- Clamping alpha to stop growing at C makes SVM "give up" on those non-separable points
- The dual program is now:
- Solve as before with extra constraints that alphas positive AND less than C ... gives alphas... from alphas get $w = \sum_i \alpha_i y_i x_i$
- Karush-Kuhn-Tucker Conditions (KKT):** solve value of b on margin for not=zero alphas AND not= C alphas for all others have support vectors, assume $\xi_i = 0$ and use formula $y_i(w^T x_i + b) - 1 + \xi_i = 0$ to get b_i and $b = \text{average}(b_i)$
- Mechanical analogy: support vector forces & torques



Kernels (see <http://www.youtube.com/watch?v=3liCbRZPrZA>)

- One important aspect of SVMs: all math involves only the *inner products* between the phi features!

$$f(x) = \text{sign} \left(\sum_i \alpha_i y_i \phi(x)^T \phi(x) + b \right)$$

- Replace all inner products with a general kernel function

- Mercer kernel: accepts 2 inputs and outputs a scalar via:

$$k(x, \tilde{x}) = \langle \phi(x), \phi(\tilde{x}) \rangle = \begin{cases} \phi(x)^T \phi(\tilde{x}) & \text{for finite } \phi \\ \int_0^1 \phi(x, t) \phi(\tilde{x}, t) dt & \text{otherwise} \end{cases}$$

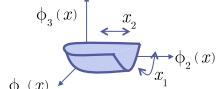
- Example: quadratic polynomial

$$\phi(x) = \begin{bmatrix} x_1^2 & \sqrt{2}x_1x_2 & x_2^2 \end{bmatrix}^T$$

$$k(x, \tilde{x}) = \phi(x)^T \phi(\tilde{x})$$

$$= x_1^2 \tilde{x}_1^2 + 2x_1 \tilde{x}_1 x_2 \tilde{x}_2 + x_2^2 \tilde{x}_2^2$$

$$= (x_1 \tilde{x}_1 + x_2 \tilde{x}_2)^2$$



Kernels

- Sometimes, many $\Phi(x)$ will produce the same $k(x, x')$
- Sometimes $k(x, x')$ computable but features huge or infinite!

- Example: polynomials

If explicit polynomial mapping, feature space $\Phi(x)$ is huge d-dimensional data, p-th order polynomial, $\dim(H) = \binom{d+p-1}{p}$

images of size 16x16 with p=4 have $\dim(H)=183$ million but can equivalently just use kernel: $k(x, y) = (x^T y)^p$

$$k(x, \tilde{x}) = (x^T \tilde{x})^p = (\sum_i x_i \tilde{x}_i)^p$$

Multinomial Theorem

$$\propto \sum_r \frac{p!}{r_1! r_2! \dots (p - r_1 - r_2 - \dots)!} x_1^{r_1} \tilde{x}_1^{r_1} \dots x_d^{r_d} \tilde{x}_d^{r_d}$$

$$\propto \sum_r (\sqrt{w_r} x_1^{r_1} \tilde{x}_1^{r_1} \dots x_d^{r_d})^p$$

$$\propto \phi(x) \phi(\tilde{x})$$

w=weight on term

Equivalent!

Kernelized SVMs

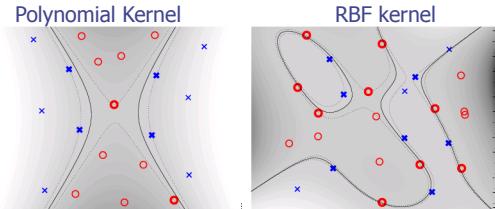
- Polynomial kernel:

$$k(x_i, x_j) = (x_i^T x_j + 1)^p$$

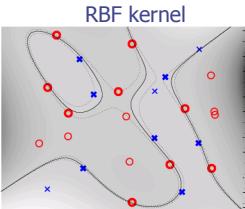
- Radial basis function kernel:

$$k(x_i, x_j) = \exp\left(-\frac{1}{2\sigma^2} \|x_i - x_j\|^2\right)$$

Polynomial Kernel



RBF kernel



- Least-squares, logistic-regression, perceptron are also kernelizable

SVM Demo

- SVM Demo by Steve Gunn:

<http://www.isis.ecs.soton.ac.uk/resources/svminfo/>

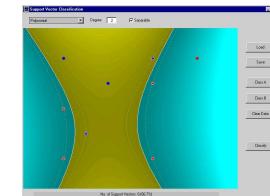
- In svc.m replace

[alpha lambda how] = qp(...);

with

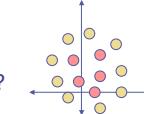
[alpha lambda how] = quadprog(H,c,[],[],A,b,vlb,vub,x0);

This replaces the old Matlab command qp (quadratic programming) with the new one for more recent versions



Nonlinear SVMs

- What if the problem is not linear?



- We can use our old trick...

- Map d-dimensional x data from L-space to high dimensional H (Hilbert) feature-space via basis functions $\Phi(x)$

- For example, quadratic classifier:

$$x_i \rightarrow \Phi(x_i) \text{ via } \Phi(\tilde{x}) = \begin{bmatrix} \vec{x} \\ \text{vec}(\tilde{x} \tilde{x}^T) \end{bmatrix}$$

- Call phi's **feature vectors** computed from original x inputs

- Replace all x 's in the SVM equations with phi's

- Now solve the following learning problem:

$$L_D : \max \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \Phi(x_i)^T \Phi(x_j) \text{ s.t. } \alpha_i \in [0, C], \sum_i \alpha_i y_i = 0$$

- Which gives a **nonlinear classifier** in original space:

$$f(x) = \text{sign} \left[\sum_i \alpha_i y_i \Phi(x_i)^T \Phi(x) + b \right]$$

Kernels

- Replace each $x_i^T x_j \rightarrow k(x_i, x_j)$, for example:

P-th Order Polynomial Kernel: $k(x, \tilde{x}) = (x^T \tilde{x} + 1)^p$

RBF Kernel (infinite!): $k(x, \tilde{x}) = \exp\left(-\frac{1}{2\sigma^2} \|x - \tilde{x}\|^2\right)$

Sigmoid (hyperbolic tan) Kernel: $k(x, \tilde{x}) = \tanh(\kappa x^T \tilde{x} - \delta)$

- Using kernels we get generalized inner product SVM:

$$L_D : \max \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j k(x_i, x_j) \text{ s.t. } \alpha_i \in [0, C], \sum_i \alpha_i y_i = 0$$

$$f(x) = \text{sign} \left(\sum_i \alpha_i y_i k(x, x_i) + b \right)$$

- Still qp solver, just use **Gram matrix K** (positive definite)

$$K = \begin{bmatrix} k(x_1, x_1) & k(x_1, x_2) & k(x_1, x_3) \\ k(x_2, x_1) & k(x_2, x_2) & k(x_2, x_3) \\ k(x_3, x_1) & k(x_3, x_2) & k(x_3, x_3) \end{bmatrix} \quad K_{ij} = k(x_i, x_j)$$

Machine Learning

Topic 7

- Unsupervised Learning
- Statistical Perspective
- Probability Models
- Discrete & Continuous: Gaussian, Bernoulli, Multinomial
- Maximum Likelihood \rightarrow Logistic Regression
- Conditioning, Marginalizing, Bayes Rule, Expectations
- Classification, Regression, Detection
- Dependence/Independence
- Maximum Likelihood \rightarrow Naïve Bayes

Statistical Perspective

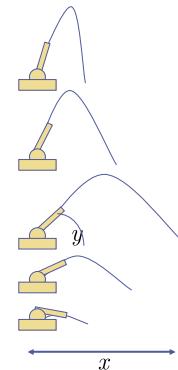
- Several problems with framework so far:
 - Only have input-output approaches (SVM, Neural Net)
 - Pulled non-linear squashing functions out of a hat
 - Pulled loss functions (squared error, etc.) out of a hat
- Better approach for classification?
- What if we have multi-class classification?
- What if other problems, i.e. unobserved values of x, y , etc...
- Also, what if we don't have a true function?
- Example of Projectile Cannon (c.f. Distal Learning)



- Would like to train a regression function to control a cannon's angle of fire (y) given target distance (x)

Statistical Perspective

- Example of Projectile Cannon (45 degree problem)
 - x = input target distance
 - y = output cannon angle
$$x = \frac{v(0)^2}{g} \sin(2y) + \text{noise}$$
- What does least squares do?
- Conditional statistical models address this problem...



Probability Models

- Now: our output is a probability density function (pdf) $p(y)$
- Probability Model: a family of pdf's with adjustable parameters which lets us select one of many $p(y) \rightarrow p(y | \Theta)$
- E.g.: 1-dim Gaussian distribution 'given' 'mean' parameter μ :

$$p(y | \mu) = N(y | \mu) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(y-\mu)^2}$$
- Want mean centered on $f(x)$'s value $p(y) = N(y | f(x))$
- Now, linear regression is:

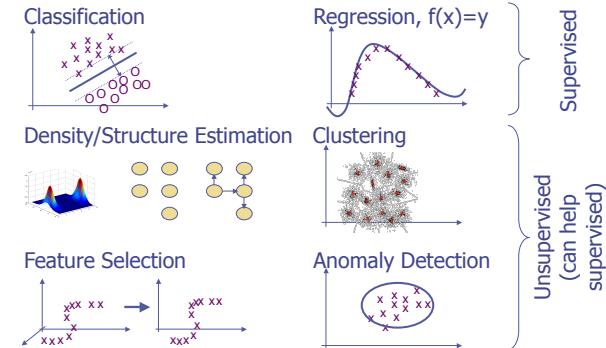
$$\begin{aligned} N(y | f(x)) &= \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(y-f(x))^2} \\ &= \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(y-b^T x - b)^2} \end{aligned}$$

Probability Models

- To fit to data, we typically "maximize likelihood" of the probability model
- Log-likelihood = objective function (i.e. negative of cost) for probability models which we want to maximize
- Define (conditional) likelihood as $L(\Theta) = \prod_{i=1}^N p(y_i | x_i)$
- or log-Likelihood as $l(\Theta) = \log(L(\Theta)) = \sum_{i=1}^N \log p(y_i | x_i)$
- For Gaussian $p(y|x)$, maximum likelihood is least squares!

$$\begin{aligned} \sum_{i=1}^N \log p(y_i | x_i) &= \sum_{i=1}^N \log N(y_i | f(x_i)) = \sum_{i=1}^N \log \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(y_i - f(x_i))^2} \\ &= -N \log(\sqrt{2\pi}) - \sum_{i=1}^N \frac{1}{2} (y_i - f(x_i))^2 \end{aligned}$$

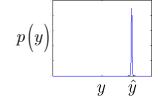
Unsupervised Learning



Probability Models

- Instead of deterministic functions, output is a probability
- Previously: our output was a scalar $\hat{y} = f(x) = \theta^T x + b$
- Now: our output is a probability $p(y)$
 - e.g. a probability bump:
- $p(y)$ subsumes or is a superset of \hat{y}
- Why is this representation for our answer more general?
 - A deterministic answer \hat{y} with complete confidence is like putting a probability $p(y)$ where all the mass is at \hat{y} !

$$\hat{y} \Leftrightarrow p(y) = \delta(y - \hat{y})$$



Probability Models

- Can extend probability model to 2 bumps:

$$p(y | \Theta) = \frac{1}{2} N(y | \mu_1) + \frac{1}{2} N(y | \mu_2)$$
- Each mean can be a linear regression fn.

$$\begin{aligned} p(y | x, \Theta) &= \frac{1}{2} N(y | f_1(x)) + \frac{1}{2} N(y | f_2(x)) \\ &= \frac{1}{2} N(y | \theta_1^T x + b_1) + \frac{1}{2} N(y | \theta_2^T x + b_2) \end{aligned}$$
- Therefore the (conditional) log-likelihood to maximize is:

$$l(\Theta) = \sum_{i=1}^N \log \left(\frac{1}{2} N(y_i | \theta_1^T x_i + b_1) + \frac{1}{2} N(y_i | \theta_2^T x_i + b_2) \right)$$
- Maximize $l(\Theta)$ using gradient ascent
 - Nicely handles the "cannon firing" data

Probability Models

- Now classification: can also go beyond deterministic!

- Previously: wanted output to be binary $\hat{y} = \{0,1\}$

- Now: our output is a probability $p(y)$

e.g. a probability table:

$y=0$	$y=1$
0.73	0.27

- This subsumes or is a superset again...

- Consider probability over binary events (coin flips!):

e.g. Bernoulli distribution (i.e 1x2 probability table) with parameter α

$$p(y|\alpha) = \alpha^y (1-\alpha)^{1-y} \quad \alpha \in [0,1]$$

- Linear classification can be done by setting α equal to $f(x)$:

$$p(y|x) = f(x)^y (1-f(x))^{1-y} \quad f(x) \in [0,1]$$

Probability Models

- Now linear classification is:

$$p(y|x) = f(x)^y (1-f(x))^{1-y} \quad f(x) \equiv \alpha \in [0,1]$$

- Log-likelihood is (negative of cost function):

$$\begin{aligned} \sum_{i=1}^N \log p(y_i|x_i) &= \sum_{i=1}^N \log f(x_i)^{y_i} (1-f(x_i))^{1-y_i} \\ &= \sum_{i=1}^N y_i \log f(x_i) + (1-y_i) \log (1-f(x_i)) \\ &= \sum_{i \in \text{class1}} \log f(x_i) + \sum_{i \in \text{class0}} \log (1-f(x_i)) \end{aligned}$$

- But, need a squashing function since $f(x)$ in $[0,1]$

- Use sigmoid or logistic again...

$$f(x) = \text{sigmoid}(\theta^T x + b) \in [0,1]$$

- Called logistic regression \rightarrow new loss function

- Do gradient descent, similar to logistic output neural net!

- Can also handle multi-layer $f(x)$ and do backprop again!

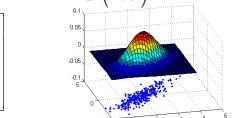
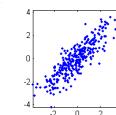
Generative Probability Models

- Idea: Extend probability to describe both X and Y

- Find probability density function over both: $p(x,y)$

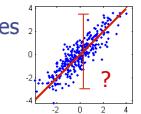
E.g. *describe* data

with Multi-Dim.
Gaussian (later...)



- Called a 'Generative Model' because we can use it to synthesize or re-generate data similar to the training data we learned from

- Regression models & classification boundaries are not as flexible
don't keep info about X
don't model noise/uncertainty



Properties of PDFs

- Let's review some basics of probability theory

- First, pdf is a function, multiple inputs, one output:

$$p(x_1, \dots, x_n) \quad p(x_1 = 0.3, \dots, x_n = 1) = 0.2$$

- Function's output is always non-negative:

$$p(x_1, \dots, x_n) \geq 0$$

- Can have discrete or continuous or both inputs:

$$p(x_1 = 1, x_2 = 0, x_3 = 0, x_4 = 3.1415)$$

- Summing over the domain of all inputs gives unity:

$$\int_{y=-\infty}^{\infty} \int_{x=-\infty}^{\infty} p(x,y) dx dy = 1 \quad \sum_y \sum_x p(x,y) = 1$$

Continuous \rightarrow integral, Discrete \rightarrow sum

Properties of PDFs

- Marginalizing: integrate/sum out a variable leaves a marginal distribution over the remaining ones...

$$\sum_y p(x,y) = p(x)$$

- Conditioning: if a variable 'y' is 'given' we get a conditional distribution over the remaining ones...

$$p(x|y) = \frac{p(x,y)}{p(y)}$$

- Bayes Rule: mathematically just redo conditioning but has a deeper meaning (1764)... if we have X being data and θ being a model

$$\text{posterior} \xrightarrow{\text{likelihood}} p(\theta | \mathcal{X}) = \frac{p(\mathcal{X} | \theta) p(\theta)}{p(\mathcal{X})} \xrightarrow{\text{evidence}} \text{prior}$$



Properties of PDFs

- Covariance: how strongly x and y vary together

$$\text{Cov}\{x,y\} = E\{(x - E\{x\})(y - E\{y\})\} = E\{xy\} - E\{x\}E\{y\}$$

- Conditional Expectation: $E\{y|x\} = \int_y p(y|x) y dy$

$$E\{E\{y|x\}\} = \int_x p(x) \int_y p(y|x) y dy dx = E\{y\}$$

- Sample Expectation: If we don't have pdf $p(x,y)$ can approximate expectations using samples of data

$$E_{p(x)}\{f(x)\} \simeq \frac{1}{N} \sum_{i=1}^N f(x_i)$$

- Sample Mean: $E\{x\} \simeq \bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$

- Sample Var: $E\{(x - E\{x\})^2\} \simeq \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2$

- Sample Cov: $E\{(x - E\{x\})(y - E\{y\})\} \simeq \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})$



More Properties of PDFs

- Independence: probabilities of independent variables multiply. Denote with the following notation:

$$x \perp\!\!\!\perp y \rightarrow p(x,y) = p(x)p(y)$$

$$x \perp\!\!\!\perp y \rightarrow p(x|y) = p(x)$$

also note in this case:

$$\begin{aligned} E_{p(x,y)}\{xy\} &= \int_x \int_y p(x,y) xy dx dy \\ &= \int_x p(x) dx \int_y p(y) y dy = E_{p(x)}\{x\} E_{p(y)}\{y\} \end{aligned}$$

- Conditional independence: when two variables become independent only if another is observed

$$x \perp\!\!\!\perp y | z \rightarrow p(x|y,z) = p(x|z)$$

$$x \perp\!\!\!\perp y | z \rightarrow p(x|y) \neq p(x)$$

The IID Assumption

- Most of the time, we will assume that a dataset is independent and identically distributed (IID)

- In many real situations, data is generated by some black box phenomenon in an arbitrary order.

- Assume we are given a dataset:

$$\mathcal{X} = \{x_1, \dots, x_N\}$$

"Independent" means that (given the model Θ) the probability of our data multiplies:

$$p(x_1, \dots, x_N | \Theta) = \prod_{i=1}^N p_i(x_i | \Theta)$$

"Identically distributed" means that each marginal probability is the same for each data point

$$p(x_1, \dots, x_N | \Theta) = \prod_{i=1}^N p_i(x_i | \Theta) = \prod_{i=1}^N p(x_i | \Theta)$$

The IID Assumption

- Bayes rule says likelihood is probability of data given model

$$\text{posterior} \rightarrow p(\theta | \mathcal{X}) = \frac{p(\mathcal{X} | \theta) p(\theta)}{p(\mathcal{X})} \quad \begin{matrix} \text{likelihood} \\ \text{prior} \\ \text{evidence} \end{matrix}$$

- The likelihood of $\mathcal{X} = \{x_1, \dots, x_N\}$ under IID assumptions is:

$$p(\mathcal{X} | \Theta) = p(x_1, \dots, x_N | \Theta) = \prod_{i=1}^N p_i(x_i | \Theta) = \prod_{i=1}^N p(x_i | \Theta)$$

- Learn joint distribution $p(x | \Theta)$ by maximum likelihood:

$$\Theta^* = \arg \max_{\Theta} \prod_{i=1}^N p(x_i | \Theta) = \arg \max_{\Theta} \sum_{i=1}^N \log p(x_i | \Theta)$$

- Learn conditional $p(y | x, \Theta)$ by max conditional likelihood:

$$\Theta^* = \arg \max_{\Theta} \prod_{i=1}^N p(y_i | x_i, \Theta) = \arg \max_{\Theta} \sum_{i=1}^N \log p(y_i | x_i, \Theta)$$

Uses of PDFs

- Classification: have $p(x, y)$ and given x . Asked for discrete y output, give most probable one

$$p(x, y) \rightarrow p(y | x) \rightarrow \hat{y} = \arg \max_m p(y = m | x)$$

- Regression: have $p(x, y)$ and given x . Asked for a scalar y output, give most probable or expected one

$$\hat{y} = \begin{cases} \arg \max_y p(y | x) \\ E_{p(y|x)} \{y\} \end{cases}$$

$\{(x_i, y_i)\} \rightarrow p(x, y) \rightarrow p(y | x)$

- Anomaly Detection: if have $p(x, y)$ and given both x, y . Asked if it is similar \rightarrow threshold

$$p(x, y) \geq \text{threshold} \rightarrow \{\text{normal, anomaly}\}$$





Machine Learning

Topic 8

- Discrete Probability Models
- Independence
- Bernoulli Distribution
- Text: Naïve Bayes
- Categorical / Multinomial Distribution
- Text: Bag of Words

Bernoulli Probability Models

- **Bernoulli:** recall binary (coin flip) probability, just 1×2 table

$$p(x) = \alpha^x (1 - \alpha)^{1-x} \quad \alpha \in [0,1] \quad x \in \{0,1\}$$

x=0	x=1
0.73	0.27

- **Multidimensional Bernoulli:** multiple binary events

$$p(x_1, x_2)$$

0	x ₂ =0	x ₂ =1
x ₁	0.4	0.1
1	0.3	0.2
x ₁		

$$p(x_1, x_2, x_3)$$

x=T	x=H
0.4	0.6

• Why do we write these as an equations instead of tables?

• To do things like... maximum likelihood...

• Fill in the table so that it matches real data...

• Example: coin flips H,H,T,TT,T,H,T,H,H,H ???

x=T	x=H

Bernoulli Probability Models



- **Bernoulli:** recall binary (coin flip) probability, just 1×2 table

$$p(x) = \alpha^x (1 - \alpha)^{1-x} \quad \alpha \in [0,1] \quad x \in \{0,1\}$$

x=0	x=1
0.73	0.27

- **Multidimensional Probability Table:** multiple binary events

$$p(x_1, x_2)$$

0	x ₂ =0	x ₂ =1
x ₁	0.4	0.1
1	0.3	0.2
x ₁		

$$p(x_1, x_2, x_3)$$

x=T	x=H
0.4	0.6

• Why do we write these as an equations instead of tables?

• To do things like... maximum likelihood...

• Fill in the table so that it matches real data...

• Example: coin flips H,H,T,TT,T,H,T,H,H,H

• Why is this correct?

Bernoulli Maximum Likelihood

- **Bernoulli:**

$$p(x) = \alpha^x (1 - \alpha)^{1-x} \quad \alpha \in [0,1] \quad x \in \{0,1\}$$

- **Log-Likelihood (IID):** $\sum_{i=1}^N \log p(x_i | \alpha) = \sum_{i=1}^N \log \alpha^{x_i} (1 - \alpha)^{1-x_i}$

- Gradient=0:

$$\begin{aligned} \frac{\partial}{\partial \alpha} \sum_{i=1}^N \log \alpha^{x_i} (1 - \alpha)^{1-x_i} &= 0 \\ \frac{\partial}{\partial \alpha} \sum_{i=1}^N x_i \log \alpha + (1 - x_i) \log (1 - \alpha) &= 0 \\ \frac{\partial}{\partial \alpha} \sum_{i \in \text{class1}} \log \alpha + \sum_{i \in \text{class0}} \log (1 - \alpha) &= 0 \\ \sum_{i \in \text{class1}} \frac{1}{\alpha} - \sum_{i \in \text{class0}} \frac{1}{1-\alpha} &= 0 \\ N_1 \frac{1}{\alpha} - N_0 \frac{1}{1-\alpha} &= 0 \\ N_1 (1 - \alpha) - N_0 \alpha &= 0 \\ N_1 - (N_1 + N_0) \alpha &= 0 \\ \alpha &= \frac{N_1}{N_1 + N_0} \end{aligned}$$

$$\begin{array}{c|c} x=0 & x=1 \\ \hline \frac{N_0}{N_0 + N_1} & \frac{N_1}{N_0 + N_1} \end{array}$$

Text Modeling via Naïve Bayes

- Maximum likelihood: assume we have several IID vectors
- Have N documents, each a 50,000 dimension binary vector
- Each dimension is a word, set to 1 if word in the document

$$\begin{array}{ccccc} \vec{x}_1 & \vec{x}_2 & \vec{x}_3 & \vec{x}_4 & \vec{x}_5 \\ \text{Dim1: "the"} & = & 1 & 0 & 1 & 1 \\ \text{Dim2: "hello"} & = & 0 & 1 & 0 & 1 \\ \text{Dim3: "and"} & = & 1 & 1 & 0 & 1 \\ \text{Dim4: "happy"} & = & 1 & 0 & 0 & 1 \end{array}$$

$$\text{Likelihood} = \prod_{i=1}^N p(\vec{x}_i | \vec{\alpha}) = \prod_{i=1}^N \prod_{d=1}^{50000} \vec{\alpha}(d)^{\vec{x}_i(d)} (1 - \vec{\alpha}(d))^{(1 - \vec{x}_i(d))}$$

• Max likelihood solution: for each word d count number of documents it appears in divided by total N documents

• To classify a new document x , build two models α_{+1}, α_{-1} & compare $\text{prediction} = \arg \max_{y \in \{\pm 1\}} p(\vec{x} | \vec{\alpha}_y)$

Categorical Probability Models

$$\begin{array}{ccccccc} 1 & 2 & 3 & 4 & 5 & 6 \\ \vec{\alpha}(1) & \vec{\alpha}(2) & \vec{\alpha}(3) & \vec{\alpha}(4) & \vec{\alpha}(5) & \vec{\alpha}(6) \end{array}$$

- **Categorical:** a distribution over a single multi-category event

$$p(x) = \prod_{m=1}^M \vec{\alpha}(m)^{\vec{x}_{(m)}} \sum_m \vec{\alpha}(m) = 1 \quad \vec{x} \in \mathbb{B}^M : \sum_m \vec{x}(m) = 1$$

- Encode events as binary indicator vectors

$$\begin{array}{ccccccc} \vec{x}(1) & \vec{x}(2) & \vec{x}(3) & \vec{x}(4) & \vec{x}(5) & \vec{x}(6) \end{array}$$

- Related to the more general **multinomial** distribution

- Find α using Maximum Likelihood (with IID assumption):

$$\sum_{i=1}^N \log p(\vec{x}_i | \vec{\alpha}) = \sum_{i=1}^N \log \prod_{m=1}^M \vec{\alpha}(m)^{\vec{x}_i(m)} = \sum_{i=1}^N \sum_{m=1}^M \vec{x}_i(m) \log (\vec{\alpha}(m))$$

- Can't just take gradient over α , use sum= 1 constraint:

- Insert constraint using Lagrange multipliers

$$\frac{\partial}{\partial \alpha_q} \sum_{i=1}^N \sum_{m=1}^M \vec{x}_i(m) \log (\vec{\alpha}(m)) - \lambda \left(\sum_{m=1}^M \vec{\alpha}(m) - 1 \right) = 0$$

$$\sum_{i=1}^N \left[\vec{x}_i(q) \frac{1}{\vec{\alpha}(q)} \right] - \lambda = 0 \Rightarrow \vec{\alpha}(q) = \frac{1}{N} \sum_{i=1}^N \vec{x}_i(q)$$

Categorical Maximum Likelihood

- Taking the gradient with Lagrangian gives this formula for each q:

$$\vec{\alpha}(q) = \frac{1}{N} \sum_{i=1}^N \vec{x}_i(q)$$

- Recall the constraint: $\sum_m \vec{\alpha}(m) - 1 = 0$

- Plug in α 's solution: $\sum_m \frac{1}{N} \sum_{i=1}^N \vec{x}_i(m) - 1 = 0$

- Gives the lambda: $\lambda = \sum_m \sum_{i=1}^N \vec{x}_i(m)$

- Final answer: $\vec{\alpha}(q) = \frac{\sum_{i=1}^N \vec{x}_i(q)}{\sum_m \sum_{i=1}^N \vec{x}_i(m)} = \frac{N_q}{N}$

- Example: Rolling dice

$$\begin{array}{cccccc} x=1 & x=2 & x=3 & x=4 & x=5 & x=6 \\ \hline 0.1 & 0.1 & 0.1 & 0.1 & 0.1 & 0.5 \end{array}$$

Multinomial Probability Model

- The multinomial is a categorical over *counts* of events
Dice: 1,3,1,4,6,1,1 Word Dice: the, dog, jumped, the
- Say document i has $W_i=2000$ words, each an IID dice roll
 $p(doc_i) = p(\vec{x}_i^1, \vec{x}_i^2, \dots, \vec{x}_i^{W_i}) = \prod_{w=1}^{W_i} p(\vec{x}_i^w) \propto \prod_{w=1}^{W_i} \prod_{d=1}^D \bar{\alpha}(\vec{x}_i^w(d))$
- Get count of each time an event occurred
 $p(doc_i) \propto \prod_{w=1}^{W_i} \prod_{d=1}^D \bar{\alpha}(\vec{x}_i^w(d)) = \prod_{d=1}^D \bar{\alpha}(\vec{x}_i^d)^{\sum_{w=1}^{W_i} \vec{x}_i^w(d)} = \prod_{d=1}^D \bar{\alpha}(\vec{x}_i^d)^{\vec{X}_i(d)}$
- BUT: order shouldn't matter when "counting" so multiply by # of possible choosings. Choosing $X(1), \dots, X(D)$ from N

$$\binom{W_i}{\vec{X}_i(1), \dots, \vec{X}_i(D)} = \frac{W_i!}{\prod_{d=1}^D \vec{X}_i(d)!} = \frac{(\sum_{d=1}^D \vec{X}_i(d))!}{\prod_{d=1}^D \vec{X}_i(d)!}$$
- Multinomial: over discrete integer vectors X summing to W

$$p(\vec{X}_i) = \frac{W_i!}{\prod_{d=1}^D \vec{X}_i(d)!} \prod_{d=1}^D \bar{\alpha}(\vec{X}_i(d))^{vec(X_i(d))} \text{ s.t. } \sum_d \vec{X}_i(d) = 1, \vec{X}_i \in \mathbb{Z}_{+}^D, \sum_{d=1}^D \vec{X}_i(d) = W$$

Text Modeling via Multinomial

- Also known as the bag-of-words model
 - Each document is 50,000 dimensional vector
 - Each dimension is a word, set to # times word in doc
- | | X_1 | X_2 | X_3 | X_4 |
|---------------|-------|-------|-------|-------|
| Dim1: "the" | 9 | 3 | 1 | 0 |
| Dim2: "hello" | 0 | 5 | 3 | 0 |
| Dim3: "and" | 6 | 2 | 2 | 2 |
| Dim4: "happy" | 2 | 5 | 1 | 0 |
- Each document is a vector of multinomial counts

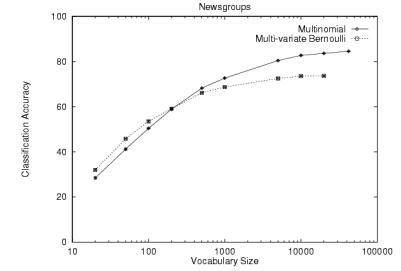
$$p(doc_i) = p(\vec{X}_i) = \frac{[\sum_{d=1}^D \vec{X}_i(d)]!}{\prod_{d=1}^D \vec{X}_i(d)!} \prod_{d=1}^D \bar{\alpha}(\vec{X}_i(d))^{vec(\vec{X}_i(d))} \quad \sum_d \vec{X}_i(d) = 1 \quad X \in \mathbb{Z}_{+}^D$$
 - Log-likelihood:

$$l(\vec{\alpha}) = \sum_{i=1}^N \log p(\vec{X}_i) = \sum_{i=1}^N \log \frac{[\sum_{d=1}^D \vec{X}_i(d)]!}{\prod_{d=1}^D \vec{X}_i(d)!} \prod_{d=1}^D \bar{\alpha}(\vec{X}_i(d))^{vec(\vec{X}_i(d))}$$

$$= \sum_{i=1}^N \sum_{d=1}^D \vec{X}_i(d) \log \bar{\alpha}(\vec{X}_i(d)) + const$$
 - Find α just like the multinomial maximum likelihood formula!

Text Modeling Experiments

- For text modeling (McCallum & Nigam '98)
Bernoulli better for small vocabulary
Multinomial better for large vocabulary



Machine Learning

Topic 9

- Continuous Probability Models
- Gaussian Distribution
- Maximum Likelihood Gaussian
- Sampling from a Gaussian



Gaussian Distribution

- Recall 1-dimensional Gaussian with mean parameter μ translates Gaussian left & right

$$p(x | \mu) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(x - \mu)^2\right)$$

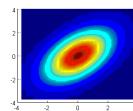
- Can also have variance parameter σ^2 widens or narrows the Gaussian

$$p(x | \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right)$$

Note: $\int_{x=-\infty}^{\infty} p(x) dx = 1$

Multivariate Gaussian

- Gaussian can extend to D-dimensions
- Gaussian mean parameter μ vector, it translates the bump
- Covariance matrix Σ stretches and rotates bump
- Mean is any real vector $\vec{x} \in \mathbb{R}^D$, $\vec{\mu} \in \mathbb{R}^D$, $\Sigma \in \mathbb{R}^{D \times D}$
- Max and expectation = μ
- Variance parameter is now Σ matrix
- Covariance matrix is positive definite
- Covariance matrix is symmetric
- Need matrix inverse (inv)
- Need matrix determinant (det)
- Need matrix trace operator (trace)



Max Likelihood Gaussian

- Have IID samples as vectors $i=1..N$: $\mathcal{X} = \{\vec{x}_1, \vec{x}_2, \dots, \vec{x}_N\}$

- How do we recover the mean and covariance parameters?

- Standard approach: Maximum Likelihood (IID)

- Maximize probability of data given model (likelihood)

$$\begin{aligned} p(\mathcal{X} | \theta) &= p(\vec{x}_1, \vec{x}_2, \dots, \vec{x}_N | \theta) \\ &= \prod_{i=1}^N p(\vec{x}_i | \vec{\mu}, \Sigma_i) \quad \text{independent Gaussian samples} \\ &= \prod_{i=1}^N p(\vec{x}_i | \vec{\mu}, \Sigma) \quad \text{identically distributed} \end{aligned}$$

- Instead, work with maximum of log-likelihood

$$\sum_{i=1}^N \log p(\vec{x}_i | \vec{\mu}, \Sigma) = \sum_{i=1}^N \log \frac{1}{(2\pi)^{D/2} \sqrt{|\Sigma|}} \exp\left(-\frac{1}{2}(\vec{x}_i - \vec{\mu})^T \Sigma^{-1} (\vec{x}_i - \vec{\mu})\right)$$

Max Likelihood Gaussian

$$\begin{aligned} \text{• Max over } \mu &\quad \frac{\partial}{\partial \mu} \left(\sum_{i=1}^N \log \frac{1}{(2\pi)^{D/2} \sqrt{|\Sigma|}} \exp\left(-\frac{1}{2}(\vec{x}_i - \vec{\mu})^T \Sigma^{-1} (\vec{x}_i - \vec{\mu})\right) \right) = 0 \\ &\quad \frac{\partial}{\partial \mu} \left(\sum_{i=1}^N -\frac{D}{2} \log 2\pi - \frac{1}{2} \log |\Sigma| - \frac{1}{2}(\vec{x}_i - \vec{\mu})^T \Sigma^{-1} (\vec{x}_i - \vec{\mu}) \right) = 0 \\ &\quad \boxed{\frac{\partial \vec{x}^T \vec{x}}{\partial \vec{x}} = 2\vec{x}^T} \quad \sum_{i=1}^N (\vec{x}_i - \vec{\mu})^T \Sigma^{-1} = \vec{0} \end{aligned}$$

see Jordan Ch. 12, get sample mean...

$$\vec{\mu} = \frac{1}{N} \sum_{i=1}^N \vec{x}_i$$

$$\text{• For } \Sigma \text{ need Trace operator: } \text{tr}(A) = \text{tr}(A^T) = \sum_{d=1}^D A(d, d)$$

and several properties:

$$\text{tr}(AB) = \text{tr}(BA)$$

$$\text{tr}(BAB^{-1}) = \text{tr}(A)$$

$$\text{tr}(\vec{x}\vec{x}^T A) = \text{tr}(\vec{x}^T A \vec{x}) = \vec{x}^T A \vec{x}$$

Continuous Probability Models

- Probabilities can have both discrete & continuous variables

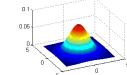
- We will discuss:

- 1) discrete probability tables

x=T	x=H
0.4	0.6

x=1	x=2	x=3	x=4	x=5	x=6
0.1	0.1	0.1	0.1	0.1	0.5

- 2) continuous probability distributions

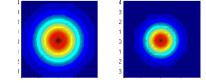


- Most popular continuous distribution = Gaussian

Multivariate Gaussian

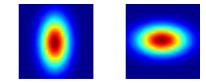
- Spherical:

$$\Sigma = \sigma^2 I = \begin{bmatrix} \sigma^2 & 0 & 0 \\ 0 & \sigma^2 & 0 \\ 0 & 0 & \sigma^2 \end{bmatrix}$$



- Diagonal Covariance:

$$\begin{aligned} \text{dimensions of } x \text{ are independent product of multiple 1d Gaussians} \\ p(\vec{x} | \vec{\mu}, \Sigma) = \prod_{d=1}^D \frac{1}{\sqrt{2\pi\sigma_d^2}} \exp\left(-\frac{(\vec{x}(d) - \vec{\mu}(d))^2}{2\sigma_d^2}\right) \\ \Sigma = \begin{bmatrix} \vec{\sigma}(1)^2 & 0 & 0 & 0 \\ 0 & \vec{\sigma}(2)^2 & 0 & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & 0 & \vec{\sigma}(D)^2 \end{bmatrix} \end{aligned}$$



Max Likelihood Gaussian

- Likelihood rewritten in trace notation:

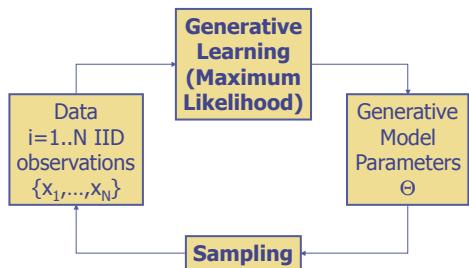
$$\begin{aligned} l &= \sum_{i=1}^N -\frac{D}{2} \log 2\pi - \frac{1}{2} \log |\Sigma| - \frac{1}{2}(\vec{x}_i - \vec{\mu})^T \Sigma^{-1} (\vec{x}_i - \vec{\mu}) \\ &= -\frac{ND}{2} \log 2\pi + \frac{N}{2} \log |\Sigma| - \frac{1}{2} \sum_{i=1}^N \text{tr}\left[(\vec{x}_i - \vec{\mu})^T \Sigma^{-1} (\vec{x}_i - \vec{\mu})\right] \\ &= -\frac{ND}{2} \log 2\pi + \frac{N}{2} \log |\Sigma| - \frac{1}{2} \sum_{i=1}^N \text{tr}\left[(\vec{x}_i - \vec{\mu})(\vec{x}_i - \vec{\mu})^T \Sigma^{-1}\right] \\ &= -\frac{ND}{2} \log 2\pi + \frac{N}{2} \log |A| - \frac{1}{2} \sum_{i=1}^N \text{tr}\left[(\vec{x}_i - \vec{\mu})(\vec{x}_i - \vec{\mu})^T A\right] \end{aligned}$$

- Max over $A = \Sigma^{-1}$

$$\begin{aligned} \text{use properties: } \frac{\partial l}{\partial A} &= -0 + \frac{N}{2} (A^{-1})^T - \frac{1}{2} \sum_{i=1}^N \left[(\vec{x}_i - \vec{\mu})(\vec{x}_i - \vec{\mu})^T \right]^T \\ &= \frac{N}{2} \Sigma - \frac{1}{2} \sum_{i=1}^N (\vec{x}_i - \vec{\mu})(\vec{x}_i - \vec{\mu})^T \\ \frac{\partial l}{\partial A} &= \frac{\partial \log |A|}{\partial A} = (A^{-1})^T \quad \frac{\partial \text{tr}(BA)}{\partial A} = B^T \\ &= \frac{\partial \log |A|}{\partial A} = (A^{-1})^T \quad \frac{\partial \text{tr}(BA)}{\partial A} = B^T \end{aligned}$$

$$\text{• Get sample covariance: } \frac{\partial l}{\partial A} = 0 \rightarrow \Sigma = \frac{1}{N} \sum_{i=1}^N (\vec{x}_i - \vec{\mu})(\vec{x}_i - \vec{\mu})^T$$

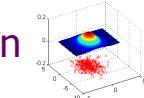
Sampling & Max Likelihood



Sampling from a Gaussian

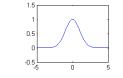
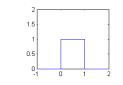
- Fit Gaussian to data, how is this Generative?
- Sampling! Generating discrete data easy:

0.73	0.1	0.17
------	-----	------
- Assume we can do uniform sampling:
i.e. rand between (0,1)
if $0.00 \leq \text{rand} < 0.73$ get A
if $0.73 \leq \text{rand} < 0.83$ get B
if $0.83 \leq \text{rand} < 1.00$ get C
- What are we doing?
Sum up the Probability Density Function (PDF) to get Cumulative Density Function (CDF)
- For 1d Gaussian, Integrate Probability Density Function get Cumulative Density Function Integral is like summing many discrete bars



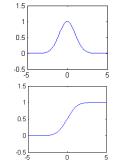
0.73	0.1	0.17
------	-----	------

0.73	0.83	1.00
------	------	------



Sampling from a Gaussian

- Integrate 1d Gaussian to get CDF:
 $p(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}x^2\right)$
 $F(x) = \int_{-\infty}^x p(t) dt = \frac{1}{2} \operatorname{erf}\left(\frac{1}{\sqrt{2}}x\right) + \frac{1}{2}$
- If sample from uniform, get: $u \sim \text{uniform}(0,1)$
- Compute mapping: $x = F^{-1}(u) = \sqrt{2}\operatorname{erfinv}(2u - 1)$
- This is a Gaussian sample: $x \sim N(x | 0, 1)$
- For D-dimensional Gaussian $N(z | 0, I)$ concatenate samples:
 $\vec{x} = [\vec{x}(1) \dots \vec{x}(D)]^T \sim p(\vec{x} | 0, I) = \prod_{d=1}^D \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} \vec{x}(d)^2\right)$
- For $N(z | \mu, \Sigma)$, add mean & multiply by root cov
 $\vec{z} = \Sigma^{1/2} \vec{x} + \vec{\mu} \sim p(\vec{z} | \vec{\mu}, \Sigma)$
- Example code: gendata.m



Machine Learning

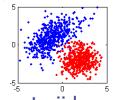
Topic 10

- Classification with Gaussians
- Regression with Gaussians
- Principal Components Analysis

Classification with Gaussians

- Have two classes, each with their own Gaussian: $\{(x_1, y_1), \dots, (x_N, y_N)\} \quad x \in R^D \quad y \in \{0,1\}$
- Given parameters $\theta = \{\alpha, \mu_0, \Sigma_0, \mu_1, \Sigma_1\}$ we can generate iid data from $p(x, y | \theta) = p(y | \theta) p(x | y, \theta)$ by:
 - flipping a coin to get y via Bernoulli $p(y | \theta) = \alpha^y (1 - \alpha)^{1-y}$
 - sampling an x from y 'th Gaussian $p(x | y, \theta) = N(x | \mu_y, \Sigma_y)$
- Or, recover parameters from data using maximum likelihood

$$\begin{aligned} l(\theta) &= \log p(data | \theta) = \sum_{i=1}^N \log p(x_i, y_i | \theta) \\ &= \sum_{i=1}^N \log p(y_i | \theta) + \sum_{i=1}^N \log p(x_i | y_i, \theta) \\ &= \sum_{i=1}^N \log p(y_i | \alpha) + \sum_{y_i \in 0} \log p(x_i | \mu_0, \Sigma_0) + \sum_{y_i \in 1} \log p(x_i | \mu_1, \Sigma_1) \end{aligned}$$



Classification with Gaussians

- Max Likelihood can be done separately for the 3 terms

$$l = \sum_{i=1}^N \log p(y_i | \alpha) + \sum_{y_i \in 0} \log p(x_i | \mu_0, \Sigma_0) + \sum_{y_i \in 1} \log p(x_i | \mu_1, \Sigma_1)$$
- Count # of pos & neg examples (class prior): $\alpha = \frac{N_0}{N_0 + N_1}$
- Get mean & cov of negatives and mean & cov of positives:

$$\begin{aligned} \mu_0 &= \frac{1}{N_0} \sum_{y_i \in 0} x_i & \Sigma_0 &= \frac{1}{N_0} \sum_{y_i \in 0} (x_i - \mu_0)(x_i - \mu_0)^T \\ \mu_1 &= \frac{1}{N_1} \sum_{y_i \in 1} x_i & \Sigma_1 &= \frac{1}{N_1} \sum_{y_i \in 1} (x_i - \mu_1)(x_i - \mu_1)^T \end{aligned}$$
- Given (x, y) pair, can now compute likelihood $p(x, y)$
- To make classification, a bit of Decision Theory
- Without x , can compute prior guess for y $p(y)$
- Give me x , want y , I need posterior $p(y | x)$
- Bayes Optimal Decision: $\hat{y} = \arg \max_{y=\{0,1\}} p(y | x)$
- Optimal iff we have true probability

Posterior gives Logistic

- Bayes Optimal Decision: $\hat{y} = \arg \max_{y=\{0,1\}} p(y | x)$
- To get conditional:

$$p(y | x) = \frac{p(x, y)}{p(x)} = \frac{p(x, y)}{\sum_y p(x, y)} = \frac{p(x, y)}{p(x, y=0) + p(x, y=1)}$$
- Check which is greater: $p(y=0 | x) \geq ? \leq p(y=1 | x)$
- Or check if this is > 0.5

$$\begin{aligned} p(y=1 | x) &= \frac{p(x, y=1)}{p(x, y=0) + p(x, y=1)} \\ &= \frac{1}{\frac{p(x, y=0)}{p(x, y=1)} + 1} \\ &= \frac{1}{\exp(-\log \frac{p(x, y=1)}{p(x, y=0)}) + 1} \\ &= \text{sigmoid}(\log \frac{p(x, y=1)}{p(x, y=0)}) \end{aligned}$$
- Get logistic squashing function of log-ratio of probability models

Linear or Quadratic Decisions

- Example cases, plotting decision boundary when $= 0.5$

$$\begin{aligned} p(y=1 | x) &= \frac{p(x, y=1)}{p(x, y=0) + p(x, y=1)} \\ &= \frac{\alpha N(x | \mu_1, \Sigma_1)}{(1-\alpha) N(x | \mu_0, \Sigma_0) + \alpha N(x | \mu_1, \Sigma_1)} \end{aligned}$$
- If covariances are equal:
 - linear decision
 -
- If covariances are different:
 - quadratic decision
 -

Regression with Gaussians

- Have input and output, each Gaussian:

$$\{(x_1, y_1), \dots, (x_N, y_N)\} \quad x \in R^{D_x} \quad y \in R^{D_y}$$
- concatenate $z_i = \begin{bmatrix} x_i \\ y_i \end{bmatrix}$
- $p(z | \mu, \Sigma) = \frac{1}{(2\pi)^{D/2} \sqrt{|\Sigma|}} \exp\left(-\frac{1}{2}(z - \mu)^T \Sigma^{-1} (z - \mu)\right)$
- Maximum Likelihood is as usual for a multivariate Gaussian

$$\mu = \frac{1}{N} \sum_{i=1}^N z_i \quad \Sigma = \frac{1}{N} \sum_{i=1}^N (z_i - \mu)(z_i - \mu)^T$$
- Bayes optimal decision: $\hat{y} = \arg \max_{y \in \mathbb{R}} p(y | x)$
- Or we can use: $\hat{y} = E_{p(y|x)} \{y\}$
- Have joint, need conditional: $p(y | x) = \frac{p(x, y)}{p(x)} = \frac{p(x, y)}{\int_y p(x, y)}$

Gaussian Marginals/Conditionals

- Conditional & marginal from joint: $p(y | x) = \frac{p(x, y)}{p(x)} = \frac{p(x, y)}{\int_y p(x, y)}$
- Conditioning the Gaussian:

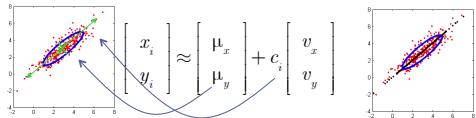
$$\begin{aligned} p(z | \mu, \Sigma) &= \frac{1}{(2\pi)^{D/2} \sqrt{|\Sigma|}} \exp\left\{-\frac{1}{2}(z - \mu)^T \Sigma^{-1} (z - \mu)\right\} \\ p(x, y) &= \frac{1}{(2\pi)^{D/2} \sqrt{|\Sigma|}} \exp\left\{-\frac{1}{2}\left[\begin{array}{c|c} x & y \end{array}\right] - \left[\begin{array}{c|c} \mu_x & \mu_y \end{array}\right]\right]^T \left[\begin{array}{cc} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{array}\right]^{-1} \left[\begin{array}{c|c} x & y \end{array}\right] - \left[\begin{array}{c|c} \mu_x & \mu_y \end{array}\right]\right\} \\ p(x) &= \frac{1}{(2\pi)^{D_x/2} \sqrt{\Sigma_{xx}}} \exp\left\{-\frac{1}{2}(x - \mu_x)^T \Sigma_{xx}^{-1} (x - \mu_x)\right\} \\ &= N(x | \mu_x, \Sigma_{xx}) \\ p(y | x) &= N(y | \mu_y + \Sigma_{yx} \Sigma_{xx}^{-1} (x - \mu_x), \Sigma_{yy} - \Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xy}) \end{aligned}$$
- Here argmax is expectation which is conditional mean: $\hat{y} = \mu_y + \Sigma_{yx} \Sigma_{xx}^{-1} (x - \mu_x)$

Principal Components Analysis

- Gaussians: for Classification, Regression... & Compression!
- Data can be constant in some directions, changes in others
- Use Gaussian to find directions of high/low variance

Principal Components Analysis

- Idea: instead of writing data in all its dimensions, only write it as mean + steps along one direction



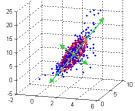
- More generally, keep a subset of dimensions C from D (i.e. 2 of 3)

$$\vec{x}_i \approx \vec{\mu} + \sum_{j=1}^C c_{ij} \vec{v}_j$$

Compression method: $\vec{x}_i \gg \vec{c}_i$

Optimal directions: along eigenvectors of covariance

Which directions to keep: highest eigenvalues (variances)



Principal Components Analysis

- If we have eigenvectors, mean and coefficients:

$$\vec{x}_i \approx \vec{\mu} + \sum_{j=1}^C c_{ij} \vec{v}_j$$

- Get eigenvectors (use `eig()` in Matlab): $\Sigma = V \Lambda V^T$

$$\begin{bmatrix} \Sigma(1,1) & \Sigma(1,2) & \Sigma(1,3) \\ \Sigma(2,1) & \Sigma(2,2) & \Sigma(2,3) \\ \Sigma(3,1) & \Sigma(3,2) & \Sigma(3,3) \end{bmatrix} = \begin{bmatrix} [\vec{v}_1] & [\vec{v}_2] & [\vec{v}_3] \end{bmatrix} \begin{bmatrix} \lambda_1 & 0 & 0 \\ 0 & \lambda_2 & 0 \\ 0 & 0 & \lambda_3 \end{bmatrix} \begin{bmatrix} [\vec{v}_1] & [\vec{v}_2] & [\vec{v}_3] \end{bmatrix}^T$$

- Eigenvectors are orthonormal: $\vec{v}_i^T \vec{v}_j = \delta_{ij}$

- In coordinates of v , Gaussian is diagonal, $\text{cov} = \Lambda$

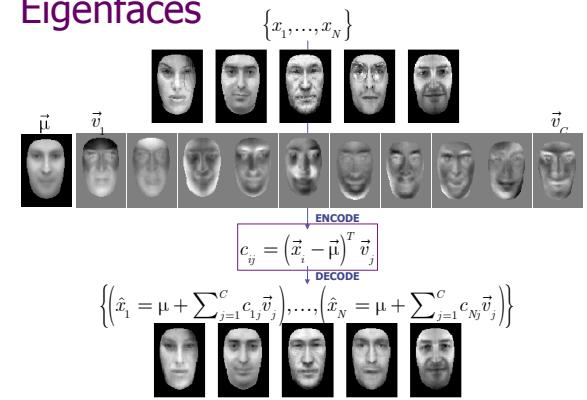
- All eigenvalues are non-negative $\lambda_i \geq 0$

- Higher eigenvalues are higher variance, use the top C ones

$$\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \lambda_4 \geq \dots$$

- To compute the coefficients: $c_{ij} = (\vec{x}_i - \vec{\mu})^T \vec{v}_j$

Eigenfaces



Machine Learning

Instructor: Tony Jebara

Bayesians & Frequentists

- Frequentists (Neymann/Pearson/Wald). An orthodox view that sampling is infinite and decision rules can be sharp.
- Bayesians (Bayes/Laplace/de Finetti). Unknown quantities are treated probabilistically and the state of the world can always be updated.



de Finetti: $p(\text{event})$ = price I would pay for a contract that pays 1\$ when event happens

- Likelihoodists (Fisher). Single sample inference based on maximizing the likelihood function and relying on the Birnbaum's Theorem. Bayesians – But they don't know it.

Bayesian Inference

- Bayes rule gives rise to maximum likelihood
- Assume we have a prior over models $p(\theta)$

$$p(\theta | x) = \frac{p(x | \theta) p(\theta)}{p(x)}$$

likelihood prior
posterior evidence

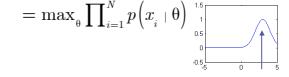
- How to pick $p(\theta)$?
Pick simpler θ is better
Pick form for mathematical convenience
- We have data (can assume IID): $\mathcal{X} = \{x_1, x_2, \dots, x_N\}$
- Want to get a model to compute: $p(x)$
- Want $p(x)$ given our data... How to proceed?

Topic 11

- Maximum Likelihood as Bayesian Inference
- Maximum A Posteriori
- Bayesian Gaussian Estimation

Why Maximum Likelihood?

- So far, assumed max (log) likelihood (IID or otherwise)
- Philosophical: Why? $\max_{\theta} L(\theta) = \max_{\theta} p(x_1, \dots, x_N | \theta) = \max_{\theta} \prod_{i=1}^N p(x_i | \theta)$
- Also, why ignore $p(\theta)$?
- Hint: Recall Bayes rule: $p(\theta | x) = \frac{p(x | \theta) p(\theta)}{p(x)}$
- Everyone agrees on probability theory: inference and use of probability models when we have computed $p(x)$
- But how get to $p(x)$ from data? Debate...
- Two schools of thought: Bayesians and Frequentists



Bayesians & Frequentists

- Frequentists:
 - Data are a repeatable random sample- there is a frequency
 - Underlying parameters remain constant during this repeatable process
 - Parameters are fixed
- Bayesians:
 - Data are observed from the realized sample.
 - Parameters are unknown and described probabilistically
 - Data are fixed

Bayesians & Frequentists

- Frequentists:** classical / objective view / no priors every statistician should compute same $p(x)$ so no priors can't have a $p(\text{event})$ if it never happened avoid $p(\theta)$, there is 1 true model, not distribution of them permitted: $p_\theta(x, y)$ forbidden: $p(x, y | \theta)$
Frequentist inference: estimate one best model θ use the **ML estimator** (unbiased & minimum variance)
do not depend on Bayes rule for learning
- Bayesians:** subjective view / priors are ok put a distribution or pdf on all variables in the problem even models & deterministic quantities (i.e. speed of light)
use a prior $p(\theta)$, on the model θ before seeing any data
Bayesian inference: use Bayes rule for learning, integrate over all model (θ) unknown variables

Bayesian Inference

$$\begin{aligned} \text{Want } p(x) \text{ given our data... } p(x | \mathcal{X}) &= p(x | x_1, x_2, \dots, x_N) \\ p(x | \mathcal{X}) &= \int_{\theta} p(x | \theta, \mathcal{X}) d\theta \end{aligned}$$

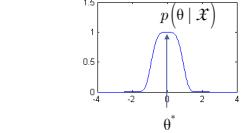
$$\begin{aligned} &= \int_{\theta} p(x | \theta, \mathcal{X}) \frac{p(\mathcal{X} | \theta) p(\theta)}{p(\mathcal{X})} d\theta \quad \text{Prior} \end{aligned}$$

$$\begin{aligned} &= \int_{\theta} p(x | \theta) \frac{\prod_{i=1}^N p(x_i | \theta) p(\theta)}{p(\mathcal{X})} d\theta \quad \text{Many models} \\ &\quad \text{Weight on each model} \end{aligned}$$

Bayesian Inference to MAP & ML

- The full Bayesian Inference integral can be mathematically tricky. Maximum likelihood is an approximation of it...

$$\begin{aligned} p(x | \mathcal{X}) &= \int_{\theta} p(x | \theta) \frac{\prod_{i=1}^N p(x_i | \theta) p(\theta)}{p(\mathcal{X})} d\theta \\ &\approx \int_{\theta} p(x | \theta) \delta(\theta - \theta^*) d\theta \\ \text{where } \theta^* &= \begin{cases} \arg \max_{\theta} \frac{\prod_{i=1}^N p(x_i | \theta) p(\theta)}{p(\mathcal{X})} & \text{MAP} \\ \arg \max_{\theta} \frac{\prod_{i=1}^N p(x_i | \theta) \text{uniform}(\theta)}{p(\mathcal{X})} & \text{ML} \end{cases} \end{aligned}$$



- Maximum A Posteriori (MAP) is like Maximum Likelihood (ML) with a prior $p(\theta)$ which lets us prefer some models over others

$$l_{MAP}(\theta) = l_{ML}(\theta) + \log p(\theta) = \sum_{i=1}^N \log p(x_i | \theta) + \log p(\theta)$$



Bayesian Inference Example

- For Gaussians, we CAN compute the integral (but hard!)

$$p(x | \mathcal{X}) = \int_{\theta} p(x | \theta) \frac{\prod_{i=1}^N p(x_i | \theta) p(\theta)}{p(\mathcal{X})} d\theta$$

$$\propto \int_{\theta} p(x | \theta) \prod_{i=1}^N p(x_i | \theta) p(\theta) d\theta$$

- Example:... assume 1d Gaussian & Gaussian prior on mean

$p(x | \theta) = \text{Gaussian}$

$p(\theta) = \text{Gaussian}$

$$p(x | \mathcal{X}) \propto \int_{\mu} \left(\frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x-\mu)^2} \right) \prod_{i=1}^N \left(\frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x_i-\mu)^2} \right) \left(\frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(\mu-\mu_0)^2} \right) d\mu$$

Bayesian Inference Example

- Solve integral over all Gaussian means with variance=1

$$p(x | \mathcal{X}) \propto \int_{\mu=-\infty}^{\mu=\infty} \left(\frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x-\mu)^2} \right) \prod_{i=1}^N \left(\frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x_i-\mu)^2} \right) \left(\frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(\mu-\mu_0)^2} \right) d\mu$$

$$\propto \int_{\mu=-\infty}^{\mu=\infty} \exp \left(-\frac{1}{2} (x - \mu)^2 - \sum_i \frac{1}{2} (x_i - \mu)^2 - \frac{1}{2} (\mu_0 - \mu)^2 \right) d\mu$$

$$\propto \int_{\mu=-\infty}^{\mu=\infty} \exp \left(-\frac{1}{2} [(N+2)\mu^2 - 2\mu(x + \mu_0 + \sum_i x_i) + x^2] \right) d\mu$$

$$\propto \int_{\mu=-\infty}^{\mu=\infty} \exp \left(-\frac{1}{2} [(N+2)\mu^2 - 2\mu(x + \mu_0 + \sum_i x_i) + x^2] + [\quad]^2 - [\quad]^2 \right) d\mu$$

$$\propto \exp \left(-\frac{1}{2} \left[\frac{-(x+\mu_0 + \sum_i x_i)^2}{N+2} + x^2 \right] \right)$$

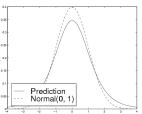
$$\tilde{\mu} = \frac{\mu_0 + \sum_i x_i}{N+1}$$

$$\tilde{\sigma}^2 = \frac{N+2}{N+1}$$

$$= N(x | \tilde{\mu}, \tilde{\sigma}^2)$$

- Can integrate over μ and Σ for multivariate Gaussian (Jordan ch. 4 and Minka Tutorial)

$$p(x | \mathcal{X}) = \frac{\Gamma((N+1)/2)}{\Gamma((N+1-d)/2)} \left| \frac{1}{(N+1)\pi} \Sigma^{-1} \right|^{d/2} \left(\frac{1}{N+1} (x - \tilde{\mu})^T \Sigma^{-1} (x - \tilde{\mu}) + 1 \right)^{-(N+1)/2}$$



Machine Learning

Instructor: Tony Jebara

Mixtures as Hidden Variables

- Consider a dataset with K subpopulations but don't know which subpopulation each point belongs to
- I.e. looking at height of adult people, we see K=2 subpopulations: males & females
- I.e. looking at weight and height of people we see K=2 subpopulations: males & females
- Because of the 'hidden' variable (y can be 1 or 2), these distributions are not Gaussians but **Mixture of Gaussians**

$$\begin{aligned} p(\vec{x}) &= \sum_y p(\vec{x}|y) = \sum_y p(y)p(\vec{x}|y) = \sum_y \pi_y N(\vec{x}|\vec{\mu}_y, \Sigma_y) \\ &= \sum_{y=1}^K \pi_y \frac{1}{(2\pi)^{D/2} \sqrt{\Sigma_y}} \exp\left(-\frac{1}{2}(\vec{x} - \vec{\mu}_y)^T \Sigma_y^{-1} (\vec{x} - \vec{\mu}_y)\right) \end{aligned}$$

K-Means Clustering

- Geometric, each point goes to closest Gaussian
 - Recompute the means by their assigned points
 - Essentially minimizing the following cost function:
- $$\min_{\vec{\mu}} \min_z J(\vec{\mu}_1, \dots, \vec{\mu}_K, \vec{z}_1, \dots, \vec{z}_N) = \sum_{n=1}^N \sum_{i=1}^K \vec{z}_n(i) \|\vec{x}_n - \vec{\mu}_i\|^2$$
- $$\vec{z}_n(i) = \begin{cases} 1 & \text{if } i = \arg \min_j \|\vec{x}_n - \vec{\mu}_j\|^2 \\ 0 & \text{otherwise} \end{cases} \quad \vec{\mu}_i = \frac{\sum_{n=1}^N \vec{z}_n(i) \vec{x}_n}{\sum_{n=1}^N \vec{z}_n(i)}$$

- Guaranteed to improve per iteration and converge
- Like **Coordinate Descent** (lock one var, maximize the other)
- A.k.a. **Axis-Parallel Optimization** or **Alternating Minimization**



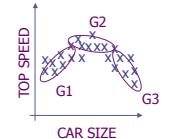
Topic 12

- Mixture Models and Hidden Variables
- Clustering
- K-Means
- Expectation Maximization

Mixtures for More Flexibility

- With mixtures (e.g. mixtures of Gaussians) we can handle more complicated (e.g. multi-bump, nonlinear) distributions.

subpopulations: G1=compact car
G2=mid-size car
G3=cadillac



- In fact, if we have enough Gaussians (maybe infinite) we can approximate any distribution...



Unlabeled data → Clustering

- Recall classification problem:

maximize the log-likelihood of data given models:

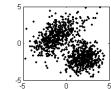
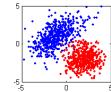
$$\begin{aligned} l &= \sum_{n=1}^N \log p(\vec{x}_n | y_n, \pi, \mu, \Sigma) \\ &= \sum_{n=1}^N \log \pi_{y_n} N(\vec{x}_n | \vec{\mu}_{y_n}, \Sigma_{y_n}) \end{aligned}$$

- If we don't know the class treat it as a hidden variable

maximize the log-likelihood with unlabeled data:

$$\begin{aligned} l &= \sum_{n=1}^N \log p(\vec{x}_n | \pi, \mu, \Sigma) = \sum_{n=1}^N \log \sum_{y=1}^K p(\vec{x}_n, y | \pi, \mu, \Sigma) \\ &= \sum_{n=1}^N \log (\pi_1 N(\vec{x}_n | \vec{\mu}_1, \Sigma_1) + \dots + \pi_K N(\vec{x}_n | \vec{\mu}_K, \Sigma_K)) \end{aligned}$$

- Instead of classification, we now have a **clustering** problem



K-Means Clustering

- K-means solves a Chicken-and-Egg problem:
If knew classes, we can get model (max likelihood!)
If knew the model, we can predict the classes (classifier!)
- Kmeans: guess a model, use it to classify the data, use classified data as labeled data to update the model, repeat.
- Assumes each point x has a discrete multinomial vector z

0) Input dataset $\{\vec{x}_1, \dots, \vec{x}_N\}$

- 1) Randomly initialize means $\vec{\mu}_1, \dots, \vec{\mu}_K$
- 2) Find closest mean for each point $\vec{z}_n(i) = \begin{cases} 1 & \text{if } i = \arg \min_j \|\vec{x}_n - \vec{\mu}_j\|^2 \\ 0 & \text{otherwise} \end{cases}$
- 3) Update means $\vec{\mu}_i = \sum_{n=1}^N \vec{z}_n(i) / \sum_{n=1}^N \vec{z}_n(i)$
- 4) If any z has changed go to 2

Expectation-Maximization (EM)

- EM is a soft/fuzzy version of K-Means (which does winner-takes-all, closest Gaussian Mean completely wins datapoint)

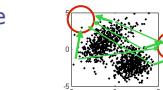
$$\vec{z}_n(i) = \begin{cases} 1 & \text{if } i = \arg \min_j \|\vec{x}_n - \vec{\mu}_j\|^2 \\ 0 & \text{otherwise} \end{cases} = \arg \max_j N(\vec{x}_n | \vec{\mu}_j, I) = \arg \max_j p(\vec{x}_n | \vec{\mu}_j)$$

- Instead, consider soft percentage assignment of datapoint

$$assign \propto \pi_j \frac{1}{(2\pi)^{D/2}} \exp\left(-\frac{1}{2} \|\vec{x}_n - \vec{\mu}_j\|^2\right)$$

- EM is 'less greedy' than K-Means uses $\tau_{n,i} = p(\vec{z} = \vec{\delta}_i | \vec{x}_n, \theta)$ as shared responsibility for \vec{x}_n

- Update for the means are then 'weighted' by responsibilities:



$$\begin{aligned} \tau_{n,1}, \dots, \tau_{n,K} &= \frac{\pi_1 N(\vec{x}_n | \vec{\mu}_1, \Sigma_1)}{\sum_j \pi_j N(\vec{x}_n | \vec{\mu}_j, \Sigma_j)} \\ \mu_i &= \frac{\sum_{n=1}^N \tau_{n,i} \vec{x}_n}{\sum_{n=1}^N \tau_{n,i}} \quad \pi_i = \frac{\sum_{n=1}^N \tau_{n,i}}{N} \\ \Sigma_i^{(t+1)} &= \frac{\sum_{n=1}^N \tau_{n,i}^{(t)} (\vec{x}_n - \vec{\mu}_i^{(t+1)}) (\vec{x}_n - \vec{\mu}_i^{(t+1)})^T}{\sum_{n=1}^N \tau_{n,i}^{(t)}} \end{aligned}$$

Expectation-Maximization

- EM uses expected value of $\vec{z}_n(i)$ rather than max $\tau_{n,i} = E\{\vec{z}_n(i) | \vec{x}_n\} = p(\vec{z}_n = \vec{\delta}_i | \vec{x}_n, \theta)$
- EM updates covariances, mixing proportions AND means...
- The algorithm for Gaussian mixtures:

EXPECTATION: $\tau_{n,i}^{(t)} = \frac{\pi_i N(\vec{x}_n | \vec{\mu}_i^{(t)}, \Sigma_i^{(t)})}{\sum_j \pi_j N(\vec{x}_n | \vec{\mu}_j^{(t)}, \Sigma_j^{(t)})}$

MAXIMIZATION: $\vec{\mu}_i^{(t+1)} = \frac{\sum_n \tau_{n,i}^{(t)} \vec{x}_n}{\sum_n \tau_{n,i}^{(t)}} \quad \pi_i^{(t+1)} = \frac{\sum_n \tau_{n,i}^{(t)}}{N}$
 $\Sigma_i^{(t+1)} = \frac{\sum_n \tau_{n,i}^{(t)} (\vec{x}_n - \vec{\mu}_i^{(t+1)}) (\vec{x}_n - \vec{\mu}_i^{(t+1)})^T}{\sum_n \tau_{n,i}^{(t)}}$

- DEMO... like an iterative divide-and-conquer algorithm
- But, divide&conquer is not a guarantee. Can we prove EM?

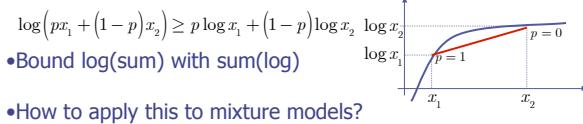
Machine Learning

Topic 13

- Expectation Maximization as Bound Maximization
- EM for Maximum A Posteriori

Jensen's Inequality

- An important general bound from Jensen (1906)
- For convex f: $f(E\{x\}) \leq E\{f(x)\}$
- For concave f: $f(E\{x\}) \geq E\{f(x)\}$
- Expectation in discrete case is sum weight by probability
- For convex f: $f\left(\sum_{i=1}^M p_i x_i\right) \leq \sum_{i=1}^M p_i f(x_i)$ when $\sum_{i=1}^M p_i = 1$, $p_i \geq 0$
- For concave f: $f\left(\sum_{i=1}^M p_i x_i\right) \geq \sum_{i=1}^M p_i f(x_i)$ when $\sum_{i=1}^M p_i = 1$, $p_i \geq 0$
- Example: $f(x) = \log(x)$ is concave and M=2



EM as Expected Log-Likelihood

- Incomplete Log-Likelihood $l(\theta) = \log p(\text{observed} | \theta) = \sum_{n=1}^N \log \sum_z p(x_n, z | \theta)$
- Complete Log-Likelihood $l^C(\theta) = \log p(\text{observed, hidden} | \theta) = \sum_{n=1}^N \log p(x_n, z_n | \theta)$
- We don't know the hidden variables z
- EM computes expected values of hidden z under current θ_t
- EM chooses Q to be the Expected Complete Log-Likelihood

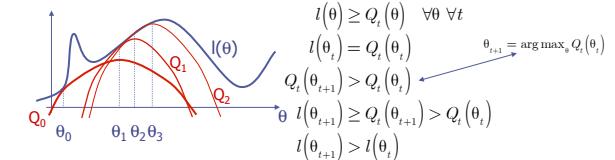
$$\begin{aligned} E\{l^C(\theta)\} &= \sum_{\text{hidden}} p(\text{hidden} | \text{observed}, \theta_t) l^C(\theta) \\ &= \sum_{z_1} \dots \sum_{z_N} p(z_1, \dots, z_N | x_1, \dots, x_N, \theta_t) l^C(\theta) \\ &= \sum_{z_1} \dots \sum_{z_N} \prod_n p(z_n | x_n, \theta_t) l^C(\theta) \\ &= \sum_{z_1} \dots \sum_{z_N} \prod_n p(z_n | x_n, \theta_t) \sum_n \log p(x_n, z_n | \theta) \\ &= \sum_n \sum_{z_n} p(z_n | x_n, \theta_t) \log p(x_n, z_n | \theta) \sum_{z_{\text{rest}}} \dots \sum_{z_N} \prod_{i \neq n} p(z_i | x_i, \theta_t) \\ &= \sum_n \sum_{z_n} p(z_n | x_n, \theta_t) \log p(x_n, z_n | \theta) = Q(\theta | \theta_t) \end{aligned}$$

Expectation-Maximization

$$\begin{aligned} l(\theta) &= \sum_{n=1}^N \log p(x_n | \theta) && \text{Original Log-Likelihood} \\ &= \sum_{n=1}^N \log \sum_z p(x_n, z | \theta) && \text{Has Hidden Variables (messy)} \\ &= \sum_{n=1}^N \log \sum_z p(x_n, z | \theta) \frac{p(z | x_n, \theta_t)}{p(z | x_n, \theta_t)} && \text{Multiply by 1 Ratio of hidden posterior density} \\ &= \sum_{n=1}^N \log \sum_z p(z | x_n, \theta_t) \frac{p(x_n, z | \theta)}{p(z | x_n, \theta_t)} && \text{Rearrange} \\ &\geq \sum_{n=1}^N \sum_z p(z | x_n, \theta_t) \log \frac{p(x_n, z | \theta)}{p(z | x_n, \theta_t)} && \text{Jensen log}(\sum_i p_i x_i) \\ &= \sum_{n=1}^N \sum_z p(z | x_n, \theta_t) \log p(x_n, z | \theta) \\ &\quad - \sum_{n=1}^N \sum_z p(z | x_n, \theta_t) \log p(z | x_n, \theta_t) \\ &= Q(\theta | \theta_t) - \text{const} && \text{New auxiliary function called Q (not messy)} \end{aligned}$$

EM as Bound Maximization

- Let's now show that EM indeed maximizes likelihood
- Bound Maximization: optimize a lower bound on $l(\theta)$
- Since log-likelihood $l(\theta)$ not concave, can't max it directly
- Consider an auxiliary function $Q(\theta)$ which is concave
- $Q(\theta)$ kisses $l(\theta)$ at a point and is less than it elsewhere



- Monotonically increases log-likelihood
- But how to find a bound and guarantee we max it?

EM as Bound Maximization

- Now have the following bound and maximize it:

$$\begin{aligned} l(\theta) &\geq Q(\theta | \theta_t) - \sum_{n=1}^N \sum_z p(z | x_n, \theta_t) \log p(z | x_n, \theta) \\ \theta^{t+1} &= \arg \max_{\theta} Q(\theta | \theta_t) = \arg \max_{\theta} \sum_{n=1}^N \sum_z p(z | x_n, \theta_t) \log p(x_n, z | \theta) \\ &= \arg \max_{\theta} \sum_{n=1}^N \sum_z \tau_{n,z} \log p(x_n, z | \theta) \end{aligned}$$

- $Q(\theta | \theta_t)$ is called Auxiliary Function... take derivatives of it
- This is easy for e-families... just weighted max likelihood!
- For example, Gaussian mixture:

$$\begin{aligned} \frac{\partial Q(\theta)}{\partial \vec{\mu}_k} &= \frac{\partial}{\partial \vec{\mu}_k} \sum_{n=1}^N \sum_k \tau_{n,k} \log \pi_k N(\vec{x}_n | \vec{\mu}_k, \Sigma_k) \\ 0 &= \sum_{n=1}^N \tau_{n,k} \frac{\partial}{\partial \vec{\mu}_k} \left(-\frac{1}{2} (\vec{x}_n - \vec{\mu}_k)^T \Sigma_k^{-1} (\vec{x}_n - \vec{\mu}_k) \right) \\ \vec{\mu}_k &= \frac{\sum_{n=1}^N \tau_{n,k} \vec{x}_n}{\sum_{n=1}^N \tau_{n,k}} \end{aligned}$$

... similarly get π_k and Σ_k

EM for Max A Posteriori

- We can also do MAP instead of ML with EM (stabilizes sol'n)
- $\log \text{posterior}(\theta) = \sum_{n=1}^N \log \sum_z p(x_n, z | \theta) + \log p(\theta)$
- Prior doesn't have log-sum
- The E-step remains the same: lower bound log-sum
- $\log \text{posterior}(\theta) = l(\theta) + \log p(\theta) \geq E\{l^C(\theta)\} + \text{const} + \log p(\theta)$
- The M-step becomes slightly different for each model
- For example, mixture of Gaussians with prior on covariance
- $\log \text{posterior}(\theta) = \sum_{n=1}^N \log \sum_k \pi_k N(\vec{x}_n | \vec{\mu}_k, \Sigma_k) + \log \prod_k p(\Sigma_k | S, \eta)$
- $\log \text{posterior}(\theta) = \sum_{n=1}^N \sum_k \tau_{n,k} \log N(\vec{x}_n | \vec{\mu}_k, \Sigma_k) + \sum_k \log p(\Sigma_k | S, \eta) + \text{const}$
- Updates on π and μ stay the same, only Σ is:

$$\Sigma_k \leftarrow \frac{1}{\sum_{n=1}^N \tau_{n,k} + \eta} \left[\sum_{n=1}^N \tau_{n,k} (\vec{x}_n - \vec{\mu}_k)(\vec{x}_n - \vec{\mu}_k)^T + \eta S \right]$$

Typically, we use the identity matrix I for S and a small eta.