University of Potsdam Dept. of Computer Science



Transfer Learning Adversarial Learning

Tobias Scheffer

Shopping for PhD Students, Postdocs

- Have you got a Master's Degree in CS or Mathematics?
- Or do you have a PhD in a field related to machine learning with a number of first-tier publications?
- Do you feel connected to statistics, programming, and data?
- Do you feel attracted to research in a highly relevant area and a growing research community?
- Do you feel like machine learning is a part of you that has been missing all along?
- Act on these feelings, give me your cv!

University of Potsdam Dept. of Computer Science



Transfer Learning

Tobias Scheffer

Prerequisites

- Statistics
 - Random variables, distributions
 - Bayes' equation
- Linear algebra
 - Vectors and matrices
 - Transposed, inverted matrices
 - Eigenvalues and eigenvectors
- Calculus
 - Derivative, partial derivative
 - Gradient

Respawning Points



- In places, the material may get technically somewhat involved.
- This icon markes points at which you can catch up, if you dropped out earlier.

Overview

- Recap: Supervised learning
 - Empirical inference, graphical models
 - Logistic regression, linear models
- Learning multiple distinct, related problems (transfer learning, domain adaptation)
 - Hierarchical Bayesian Inference
- Learning under covariate shift (differing marginal input distribution)
 - Importance sampling
 - Direct estimation of important weights
- Importance sampling for domain adaptation

Classification

- Input: instance $\mathbf{x} \in X$
 - X may be vector space.
 - Instance is a set of values for these attributes
- Output: class $y \in Y$; finite set Y.
 - Class y also called class label.

Classification: Example

- Input: instance $\mathbf{x} \in X$
 - X = Set of all combinations of a base set of substances



Classification: Learning Problem

• Input to learning problem: Data T_n .



• Training data:

$$T_n = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$$

Classification: Learning Problem

• Input to learning problem: Data T_n .



Training data:

 $T_n = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$

Output: Model



• For instance, $y_{\theta}(\mathbf{x}) = \begin{cases} & & \text{if } \phi(\mathbf{x})^{T} \mathbf{\theta} \ge 0 \\ & & \text{otherwise} \end{cases}$ Linear classifier with parameters $\mathbf{\theta}$.

- What is the modt likely class y given instance x and given the training data?
 - $y^* = \arg \max_{y} P(y | \mathbf{x}, \mathbf{X}, \mathbf{y})$
- Need assumptions about data generation process to solve.

Classification: Graphical Model

- Graphical model defines stochastic process
- Modeling assumptions on data generation process
- First, parameter vector $\boldsymbol{\theta}$ is drawn
- This $\boldsymbol{\theta}$ parameterizes training data $P(y_i | \mathbf{x}_i, \boldsymbol{\theta})$
- Marginal distribution p(x_i) is not part of the model; instances are treated as if constant: discriminative model!



Example

- Evolution guides physiological parameters of humans.
- Given these parameters and a combination of substances, nature rolles dice to determine whether an individual survives ingestion.



• Survival is governed by $p(y_i | x_i, \theta)$

Inference of the probability of y given x and training data:

•
$$P(y | \mathbf{x}, \mathbf{X}, \mathbf{y}) = \int p(y, \boldsymbol{\theta} | \mathbf{x}, \mathbf{X}, \mathbf{y}) d\boldsymbol{\theta}$$



Inference of the probability of y given x and training data:

• $P(y | \mathbf{x}, \mathbf{X}, \mathbf{y}) = \int p(y, \theta | \mathbf{x}, \mathbf{X}, \mathbf{y}) d\theta$ = $\int P(y | \mathbf{x}, \theta) p(\theta | \mathbf{X}, \mathbf{y}) d\theta$

Integration over space of all model parameters: Bayesian model averaging

Independence assumptions from graphical model

- No closed-form solution for classification
- Numerical integragtion over space of all model parameters generally infeasible.

Inference of the probability of y given x and training data:

•
$$P(y | \mathbf{x}, \mathbf{X}, \mathbf{y}) = \int p(y, \theta | \mathbf{x}, \mathbf{X}, \mathbf{y}) d\theta$$

= $\int P(y | \mathbf{x}, \theta) p(\theta | \mathbf{X}, \mathbf{y}) d\theta$
 $\approx P(y | \mathbf{x}, \theta_{MAP})$ with $\theta_{MAP} = \arg \max_{\theta} p(\theta | \mathbf{X}, \mathbf{y})$

 Approximation using only the single most likelyparameter vector: MAP model.

Inference of the probability of y given x and training data:

•
$$P(\textcircled{K}, \mathbf{y}) = \int p(\textcircled{K}, \mathbf{y}) d \textcircled{K}, \mathbf{y}) d \textcircled{K}$$

 $= \int P(\textcircled{K}, \mathbf{y}) d \textcircled{K}, \mathbf{y}) d \textcircled{K}$
 $\approx P(\textcircled{K}, \mathbf{y}) d \textcircled{K}, \mathbf{y}) d \textcircled{K}$
 $\approx P(\textcircled{K}, \mathbf{y}) d \textcircled{K}, \mathbf{y}) d \textcircled{K}$

 Approximation using only the single most likelyparameter vector: MAP model.

Model

• Inference of θ_{MAP} :

•
$$\boldsymbol{\theta}_{MAP} = \arg \max_{\boldsymbol{\theta}} p(\boldsymbol{\theta} \mid \mathbf{X}, \mathbf{y})$$

= $\arg \max_{\boldsymbol{\theta}} \frac{P(\boldsymbol{\theta}, \mathbf{X}, \mathbf{y})}{p(\mathbf{X}, \mathbf{y})}$
= $\arg \max_{\boldsymbol{\theta}} P(\mathbf{y} \mid \mathbf{X}, \boldsymbol{\theta}) p(\boldsymbol{\theta})$

θ
$\mathbf{x}_i \bullet_n$

• Inference of θ_{MAP} :

$$\boldsymbol{\theta}_{\text{MAP}} = \arg \max_{\boldsymbol{\theta}} p(\boldsymbol{\theta} \mid \mathbf{X}, \mathbf{y})$$

= $\arg \max_{\boldsymbol{\theta}} \frac{P(\boldsymbol{\theta}, \mathbf{X}, \mathbf{y})}{p(\mathbf{X}, \mathbf{y})}$
= $\arg \max_{\boldsymbol{\theta}} P(\mathbf{y} \mid \mathbf{X}, \boldsymbol{\theta}) p(\boldsymbol{\theta})$
= $\arg \max_{\boldsymbol{\theta}} \prod_{i=1}^{n} P(y_i \mid \mathbf{x}_i, \boldsymbol{\theta}) p(\boldsymbol{\theta})$







 $\Sigma \bullet$ Inference of θ_{MAP} : • $\boldsymbol{\theta}_{\text{MAP}} = \arg \max_{\boldsymbol{\theta}} p(\boldsymbol{\theta} | \mathbf{X}, \mathbf{y})$ θ $= \arg \max_{\boldsymbol{\theta}} \frac{P(\boldsymbol{\theta}, \mathbf{X}, \mathbf{y})}{p(\mathbf{X}, \mathbf{y})}$ $= \arg \max_{\boldsymbol{\theta}} P(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}) p(\boldsymbol{\theta})$ = $\arg \max_{\boldsymbol{\theta}} \prod_{i=1}^{n} P(y_i | \mathbf{x}_i, \boldsymbol{\theta}) N[0, \Sigma](\boldsymbol{\theta})$ $= \arg\min_{\boldsymbol{\theta}} \sum_{y'}^{n} \left(\log \sum_{y'} e^{\phi(\mathbf{x}_{i})^{\mathrm{T}} \boldsymbol{\theta}^{\phi,y'}} - \phi(\mathbf{x}_{i})^{\mathrm{T}} \boldsymbol{\theta}^{\phi,y_{i}} \right) + \frac{1}{2} \boldsymbol{\theta}^{\mathrm{T}} \boldsymbol{\Sigma}^{-1} \boldsymbol{\theta}$ Inference of most likely class: $\arg \max P(y_i | \mathbf{x}_i, \mathbf{\theta}) = \arg \max_{v} \phi(\mathbf{x}_i)^{\mathrm{T}} \mathbf{\theta}^{\phi, y_i} + b_{v}$

Linear model in features ϕ .



Inference of most likely class:

• arg max
$$P(\textcircled{o}_{x_i}, \boldsymbol{\theta}) = \arg \max_{y} \phi(\textcircled{o}_{y_i})^{\mathrm{T}} \boldsymbol{\theta}^{\phi, y_i} + b_y$$





Tobias Scheffer

Tobias Scheffer

Empirical Inference: reg. ERM



Tobias Scheffer

Empirical Inference: reg. ERM

Inference of θ^* : $\boldsymbol{\theta}^* = \arg \max_{\boldsymbol{\theta}} p(\boldsymbol{\theta} \mid \mathbf{X}, \mathbf{y})$ θ $= \arg \max_{\boldsymbol{\theta}} \frac{P(\boldsymbol{\theta}, \mathbf{X}, \mathbf{y})}{p(\mathbf{X}, \mathbf{y})}$ $= \arg \max_{\mathbf{\theta}} P(\mathbf{y} | \mathbf{X}, \mathbf{\theta}) p(\mathbf{\theta})$ $= \arg \max_{\boldsymbol{\theta}} \prod_{i=1}^{n} P(y_i \mid \mathbf{x}_i, \boldsymbol{\theta}) p(\boldsymbol{\theta})$ $= \arg\min_{\boldsymbol{\theta}} \sum_{i=1}^{n} \log P(y_i | \mathbf{x}_i, \boldsymbol{\theta}) + \log p(\boldsymbol{\theta})$ $= \arg\min_{\boldsymbol{\theta}} \sum_{i=1}^{n} \ell(y_i, y_{\boldsymbol{\theta}}(\mathbf{x}_i)) + \Omega(\boldsymbol{\theta})$ Regularized empirical risk minimization

Empirical Inference: reg. ERM



θ

Empirical Inference: reg. ERM

- Alternative rationalization: minimize risk = expected loss
 - $\mathbf{\theta}^{\otimes} = \arg\min_{\mathbf{\theta}} \int \ell(y_i, y_{\mathbf{\theta}}(\mathbf{x}_i)) p(\mathbf{x}_i, y_i)$
- Distribution $p(\mathbf{x}_i, y_i)$ not known
 \rightarrow approximateby sum over sample

$$\approx \arg\min_{\theta} \sum_{i=1}^{n} \ell(y_i, y_{\theta}(\mathbf{x}_i))$$

 Minimization problem is ill-posed; small change in data can lead to large change in parameters.
 Smooth by adding Tikhonov regularizer

$$\approx \arg\min_{\boldsymbol{\theta}} \sum_{i=1}^{n} \ell(y_i, y_{\boldsymbol{\theta}}(\mathbf{x}_i)) + \Omega(\boldsymbol{\theta})$$

Regularized empirical risk minimization

29

- Goal: Minimize function $L(\mathbf{\theta}) = \sum_{i=1}^{n} \ell(f_{\mathbf{\theta}}(\mathbf{x}_i), y_i) + \Omega(\mathbf{\theta})$ for given loss function and regularizer
- Numeric solutions:
 - Gradient descent
 - Cutting plane method
 - Interior point method

- Goal: Minimize function $L(\mathbf{\theta}) = \sum_{i=1}^{n} \ell(f_{\mathbf{\theta}}(\mathbf{x}_i), y_i) + \Omega(\mathbf{\theta})$ for given loss function and regularizer
- Gradient: vector of partial derivatives
- Ascent direction for function $L(\theta)$.



- Goal: Minimize function $L(\mathbf{\theta}) = \sum_{i=1}^{n} \ell(f_{\mathbf{\theta}}(\mathbf{x}_i), y_i) + \Omega(\mathbf{\theta})$ for given loss function and regularizer
- Gradient descent:
 - Iterative procedure.
 - In each step, move into direction of steepest descent
 - Descent direction given by negative gradient.



• Goal: Minimize function $L(\mathbf{\theta}) = \sum_{i=1}^{n} \ell(f_{\mathbf{\theta}}(\mathbf{x}_i), y_i) + \Omega(\mathbf{\theta})$ for given loss function and regularizer

Gradient descent:

```
RegERM(Data (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n))

Let \mathbf{\theta}^0 = \mathbf{0} and k = 0

DO

determine gradient \nabla L(\mathbf{\theta}^k)

determine step size \alpha^k

Let \mathbf{\theta}^{k+1} = \mathbf{\theta}^k - \alpha^k \nabla L(\mathbf{\theta}^k)

Let k = k + 1

WHILE \|\mathbf{\theta}^k - \mathbf{\theta}^{k-1}\| > \varepsilon

RETURN \mathbf{\theta}^k
```



- Goal: Minimize function $L(\mathbf{\theta}) = \sum_{i=1}^{n} \ell(f_{\mathbf{\theta}}(\mathbf{x}_i), y_i) + \Omega(\mathbf{\theta})$ for given loss function and regularizer
- Gradient descent:
 - $L(\theta)$ decreases in each step.
 - If L(θ) is convex, converges
 to global minimum
 - If loss function and regularizer are convex, L(θ) is convex.



Setting the Regularization Parameter

• Optimization criterion has a parameter:

 $\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} \sum_{i=1}^n \ell(y_i, y_{\boldsymbol{\theta}}(\mathbf{x}_i)) + \boldsymbol{\theta}^{\mathrm{T}} \Sigma^{-1} \boldsymbol{\theta}$ $\approx \arg\min_{\boldsymbol{\theta}} \sum_{i=1}^n \ell(y_i, y_{\boldsymbol{\theta}}(\mathbf{x}_i)) + \frac{1}{\lambda} \boldsymbol{\theta}^{\mathrm{T}} \boldsymbol{\theta}$ One parameter per element of $\boldsymbol{\theta}$, even if Σ is diagonal

 To set parameter, use grid search and n-fold cross validation (when training sample is large, one training-and-test split)

- Graphical model formulates assumptions that lead to
 - $\boldsymbol{\theta}_{\text{MAP}} = \arg\min_{\boldsymbol{\theta}} \sum_{i=1}^{n} \log P(y_i \mid \mathbf{x}_i, \boldsymbol{\theta}) + \log p(\boldsymbol{\theta})$
- θ_{MAP} is the most likely model given prior and given training data.
- Substituting log-likelihood for general loss function and log-prior for general regularizer leads to regularized empirical risk minimization

•
$$\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} \sum_{i=1}^n \ell(y_i, y_{\boldsymbol{\theta}}(\mathbf{x}_i)) + \Omega(\boldsymbol{\theta})$$

Minimum can be found by gradient descent.





Linear Models

- Hyperplane defined by normal vector and offset: $H_{\theta,b} = \{ \mathbf{x} \mid f_{\theta}(\mathbf{x}) = \phi(\mathbf{x})^{\mathrm{T}} \mathbf{\theta} + b = 0 \}$
- Class probability (two classes, logistic regression):
 P(y_i = +1 | x_i, θ) = 1 + e^{φ(x_i)^Tθ^φ+b}
 1 + e^{φ(x_i)^Tθ^φ+b}
- Decision function:

•
$$f_{\boldsymbol{\theta}}(\mathbf{x}_i) = \boldsymbol{\phi}(\mathbf{x}_i)^{\mathrm{T}} \boldsymbol{\theta}^{\phi} + b$$

Classifier:

•
$$y_{\theta}(\mathbf{x}_i) = \operatorname{sign}(f_{\theta}(\mathbf{x}_i))$$


Hyperplane defined by normal vector and offset:

$$H_{\boldsymbol{\theta},b} = \{ \mathbf{x} \mid f_{\boldsymbol{\theta}}(\mathbf{x}) = \boldsymbol{\phi}(\mathbf{x})^{\mathrm{T}} \boldsymbol{\theta} + b = 0 \}$$

- Class probability (two classes, logistic regression): $P(y_i = +1 | \mathbf{x}_i, \mathbf{\theta}) = \frac{1}{1 + e^{\phi(\mathbf{x}_i)^T \mathbf{\theta}^{\phi} + b}}$
- Decision function:

•
$$f_{\boldsymbol{\theta}}(\mathbf{x}_i) = \boldsymbol{\phi}(\mathbf{x}_i)^{\mathrm{T}} \boldsymbol{\theta}^{\phi} + b$$

Classifier:

•
$$y_{\theta}(\mathbf{x}_i) = \operatorname{sign}(f_{\theta}(\mathbf{x}_i))$$



 $\phi(\mathbf{X}_i) = \mathbf{X}_i$

Hyperplane defined by normal vector and offset:

$$H_{\boldsymbol{\theta},b} = \{ \mathbf{x} \mid f_{\boldsymbol{\theta}}(\mathbf{x}) = \boldsymbol{\phi}(\mathbf{x})^{\mathrm{T}} \boldsymbol{\theta} + b = 0 \}$$

- Class probability (two classes, logistic regression):
 P(y_i = +1 | x_i, θ) = 1 + e^{φ(x_i)^Tθ^φ+b}
- Decision function:

•
$$f_{\boldsymbol{\theta}}(\mathbf{x}_i) = \boldsymbol{\phi}(\mathbf{x}_i)^{\mathrm{T}} \boldsymbol{\theta}^{\phi} + b$$

Classifier:

•
$$y_{\theta}(\mathbf{x}_i) = \operatorname{sign}(f_{\theta}(\mathbf{x}_i))$$

 $f_{\theta}(\mathbf{x}) = 0$

Tobias Scheffer

 $\phi(\mathbf{X}_i) = \mathbf{X}_i$

Hyperplane defined by normal vector and offset:

$$H_{\boldsymbol{\theta},b} = \{ \mathbf{x} \mid f_{\boldsymbol{\theta}}(\mathbf{x}) = \boldsymbol{\phi}(\mathbf{x})^{\mathrm{T}} \boldsymbol{\theta} + b = 0 \}$$

- Class probability (two classes, logistic regression): x₂
 P(y_i = +1 | x_i, θ) = $\frac{1}{1 + e^{\phi(x_i)^T \theta^{\phi} + b}}$
- Decision function:

•
$$f_{\boldsymbol{\theta}}(\mathbf{x}_i) = \boldsymbol{\phi}(\mathbf{x}_i)^{\mathrm{T}} \boldsymbol{\theta}^{\phi} + b$$

Classifier:

•
$$y_{\theta}(\mathbf{x}_i) = \operatorname{sign}(f_{\theta}(\mathbf{x}_i))$$

 $\phi(\mathbf{x}_i) = \begin{vmatrix} \mathbf{x}_i \\ \mathbf{x}_i \end{vmatrix}$ $p(\mathbf{x} \mid y = +1, \mathbf{\theta})$

 $p(\mathbf{x} \mid y = -1, \mathbf{\theta})$

Tobias Scheffer

 x_1

Hyperplane defined by normal vector and offset:

$$H_{\boldsymbol{\theta},b} = \{ \mathbf{x} \mid f_{\boldsymbol{\theta}}(\mathbf{x}) = \boldsymbol{\phi}(\mathbf{x})^{\mathrm{T}} \boldsymbol{\theta} + b = 0 \}$$

- Class probability (two classes, logistic regression):
 P(y_i = +1 | x_i, θ) = 1 + e^{φ(x_i)^Tθ^φ+b}
- Decision function:

•
$$f_{\boldsymbol{\theta}}(\mathbf{x}_i) = \boldsymbol{\phi}(\mathbf{x}_i)^{\mathrm{T}} \boldsymbol{\theta}^{\phi} + b$$

Classifier:

•
$$y_{\theta}(\mathbf{x}_i) = \operatorname{sign}(f_{\theta}(\mathbf{x}_i))$$

 $f_{\theta}(\mathbf{x}) = 0$



Linear Models – Multi-Class Case

Hyperplanes defined by normal vector and offset:

$$H_{\boldsymbol{\theta}^{\phi, y_i}, b_{y_i}} = \{ \mathbf{x} \mid f_{\theta}(\mathbf{x}, y_i) = \phi(\mathbf{x})^{\mathrm{T}} \boldsymbol{\theta}^{\phi, y_i} + b_{y_i} = 0 \}$$

- Class probability (logistic regression): • $P(y_i | \mathbf{x}_i, \mathbf{\theta}) = \frac{e^{\phi(\mathbf{x})^T \mathbf{\theta}^{\phi, y_i} + b_{y_i}}}{\sum_{y'} e^{\phi(\mathbf{x})^T \mathbf{\theta}^{\phi, y'} + b_{y'}}}$
- Decision function:

•
$$f_{\boldsymbol{\theta}}(\mathbf{x}_i, y_i) = \boldsymbol{\phi}(\mathbf{x}_i)^{\mathrm{T}} \boldsymbol{\theta}^{\boldsymbol{\phi}, y_i} + b_{y_i}$$

Classifier:

•
$$y_{\theta}(\mathbf{x}_i) = \arg \max_{y} f_{\theta}(\mathbf{x}_i, y)$$





Decision function of generalized linear classifiers:

•
$$f_{\boldsymbol{\theta}}(\mathbf{x}_i) = \boldsymbol{\phi}(\mathbf{x}_i)^{\mathrm{T}} \boldsymbol{\theta}^{\phi} + \boldsymbol{b}$$

Multi-class case:

•
$$f_{\boldsymbol{\theta}}(\mathbf{x}_i, y_i) = \boldsymbol{\phi}(\mathbf{x}_i)^{\mathrm{T}} \boldsymbol{\theta}^{\boldsymbol{\phi}, y_i} + b_{y_i}$$

- Linear case: $\phi(\mathbf{x}) = \mathbf{x}$
- General feature mapping leads to kernel machines.

Multiple Kids Experimenting with Drugs

- So far: only one distribution $p(y_i | x_i, \theta)$ for everyone.
- What about different metabolic rates? Differing genetic factors? Levels of tolerance?



Multiple Kids Experimenting with Drugs

 We could learn distinct models for each user *j*:

•
$$\mathbf{\theta}_{j}^{*} = \arg \max_{\mathbf{\theta}} p(\mathbf{\theta} | \mathbf{X}_{j}, \mathbf{y}_{j})$$

•
$$\boldsymbol{\theta}_{j}^{*} = \arg\min_{\boldsymbol{\theta}} \sum_{i=1}^{n_{j}} \ell(y_{ji}, y_{\boldsymbol{\theta}}(\mathbf{x}_{ji})) + \Omega(\boldsymbol{\theta})$$

 But this would treat individuals as independent and disregard much data (if many previous users died, maybe this does tell you something)



Tobias Scheffer

Forms of Transfer Learning

- Learning under covariate shift
 - Only marginal distributions $p(\mathbf{x}_i | \tau_{\text{train}})$ and $p(\mathbf{x}_i | \tau_{\text{test}})$ differ; conditional $P(y_i | \mathbf{x}_{i}, \boldsymbol{\theta})$ is constant.
 - Goal: minimize risk over test distribution.
- Multi-task learning, domain adaptation
 - Both, marginal distributions and conditional distributions P(y_i |x_i,θ_{train/test}) may differ



- Instead, model a common prior over individuals (tasks) $p(y_{ji} | x_{ji}, \theta_j)$
- Physiology of everyone has been produced by the same process of evolution.



- Instead, model a common prior over individuals (tasks) $p(y_{ji} | x_{ji}, \theta_j)$
- Physiology of everyone has been produced by the same process of evolution.
- $p(\boldsymbol{\mu} \mid \boldsymbol{\Sigma}') = N[0, \boldsymbol{\Sigma}'](\boldsymbol{\mu})$

$$p(\mathbf{\theta}_j | \mathbf{\mu}, \mathbf{\Sigma}) = N[\mathbf{\mu}, \mathbf{\Sigma}](\mathbf{\theta}_j)$$

• Substitution: $\boldsymbol{\theta} = \boldsymbol{\mu} + \boldsymbol{v}$ $\Rightarrow p(\boldsymbol{v}_j | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = N[0, \boldsymbol{\Sigma}](\boldsymbol{\theta}_j)$



- Instead, model a common prior over individuals (tasks) $p(y_{ji} | x_{ji}, \theta_j)$
- Physiology of everyone has been produced by the same process of evolution.
- $p(\boldsymbol{\mu} \mid \boldsymbol{\Sigma}') = N[0, \boldsymbol{\Sigma}'](\boldsymbol{\mu})$

$$p(\mathbf{\theta}_j | \mathbf{\mu}, \mathbf{\Sigma}) = N[\mathbf{\mu}, \mathbf{\Sigma}](\mathbf{\theta}_j)$$

• Substitution: $\boldsymbol{\theta} = \boldsymbol{\mu} + \boldsymbol{v}$ $\Rightarrow p(\boldsymbol{v}_j | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = N[0, \boldsymbol{\Sigma}](\boldsymbol{v}_j)$



• Inference of
$$\mathbf{\theta}^*$$
:
 $\mathbf{\theta}_{MAP} = \arg \max_{\mu+\nu} p(\mu+\nu | \mathbf{X}, \mathbf{y}, \mathbf{\Sigma}, \mathbf{\Sigma}')$

$$= \arg \max_{\mu+\nu} \prod_{j} P(\mathbf{y}_{j} | \mathbf{X}_{j}, \mu+\nu) \prod_{j} p(\mathbf{v}_{j} | \mathbf{\Sigma}) p(\mu | \mathbf{\Sigma}')$$

$$= \arg \min_{\mu+\nu} \left(\sum_{j=1}^{m} \sum_{i=1}^{n_{j}} \log P(y_{ji} | \mathbf{x}_{ji}, \mu+\nu) + \sum_{j} \log p(\mathbf{v}_{j} | \mathbf{\Sigma}) + \log p(\mu | \mathbf{\Sigma}') \right)$$

$$= \arg \min_{\mu+\nu} \left(\sum_{j=1}^{m} \sum_{i=1}^{n_{j}} \ell(y_{ji}, y_{\mu+\nu}(\mathbf{x}_{ji})) + \sum_{j} \Omega(\mathbf{v}_{j}) + \Omega(\mu) \right)$$
• For general loss and regularizers:
 $\mathbf{\theta}^* = \arg \min_{\mu+\nu} \sum_{j=1}^{m} \sum_{i=1}^{n_{j}} \ell(y_{ji}, y_{\mu+\nu}(\mathbf{x}_{ji})) + \sum_{j} \Omega(\mathbf{v}_{j}) + \Omega(\mu)$

- Parameters = sum of population-specific and individual parameters.
- Decision function for individual j:

•
$$f_{\boldsymbol{\theta},j}(\mathbf{x}) = (\mathbf{\mu} + \mathbf{v}_j)^{\mathrm{T}} \boldsymbol{\phi}(\mathbf{x})$$





- A simple trick for implementing hierarchical Bayes
- Decision function for individual *j*:

•
$$f_{\theta,j}(\mathbf{x}) = (\mathbf{\mu} + \mathbf{v}_j)^T \phi(\mathbf{x})$$

= $\mathbf{\mu}^T \phi(\mathbf{x}) + \mathbf{v}_1^T \phi(\mathbf{x}) \delta(j = 1)$
+ ... + $\mathbf{v}_m^T \phi(\mathbf{x}) \delta(j = m)$



- A simple trick for implementing hierarchical Bayes
- Decision function for individual *j*:

•
$$f_{\theta,j}(\mathbf{x}) = (\mathbf{\mu} + \mathbf{v}_j)^{\mathrm{T}} \phi(\mathbf{x})$$
$$= \mathbf{\mu}^{\mathrm{T}} \phi(\mathbf{x}) + \mathbf{v}_1^{\mathrm{T}} \phi(\mathbf{x}) \delta(j=1)$$
$$+ \dots + \mathbf{v}_m^{\mathrm{T}} \phi(\mathbf{x}) \delta(j=m)$$
$$= \begin{pmatrix} \mathbf{v}_1 \\ \dots \\ \mathbf{v}_m \\ \mathbf{\mu} \end{pmatrix}^{\mathrm{T}} \begin{pmatrix} \phi(\mathbf{x}) \delta(j=1) \\ \dots \\ \phi(\mathbf{x}) \delta(j=m) \\ \phi(\mathbf{x}) \end{pmatrix}$$



[Daume III et al., Frustratingly easy domain adaptation, 2010]

- A simple trick for implementing hierarchical Bayes
- Decision function for individual *j*:



- A simple trick for implementing hierarchical Bayes
- Decision function for individual j:

$$f_{\boldsymbol{\theta},j}(\mathbf{x}) = \begin{pmatrix} \mathbf{v}_1 \\ \cdots \\ \mathbf{v}_m \\ \boldsymbol{\mu} \end{pmatrix}^{\mathrm{T}} \begin{pmatrix} \phi(\mathbf{x})\delta(j=1) \\ \cdots \\ \phi(\mathbf{x})\delta(j=m) \\ \phi(\mathbf{x}) \\ \phi(\mathbf{x}) \\ \boldsymbol{\Phi}_j(\mathbf{x}) \end{pmatrix}$$

Learning problem:

•
$$\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} \sum_{j=1}^{m} \sum_{i=1}^{n_j} \ell(\boldsymbol{y}_{ji}, \boldsymbol{\theta}^{\mathrm{T}} \boldsymbol{\Phi}_j(\boldsymbol{x}_{ji})) + \Omega(\boldsymbol{\theta})$$



- A simple trick for implementing hierarchical Bayes
- Decision function for individual j:

$$f_{\boldsymbol{\theta},j}(\mathbf{x}) = \begin{pmatrix} \mathbf{v}_1 \\ \dots \\ \mathbf{v}_m \\ \boldsymbol{\mu} \end{pmatrix}^{\mathrm{T}} \begin{pmatrix} \phi(\mathbf{x})\delta(j=1) \\ \dots \\ \phi(\mathbf{x})\delta(j=m) \\ \phi(\mathbf{x}) \\ \boldsymbol{\phi}(\mathbf{x}) \\ \boldsymbol{\phi}_j(\mathbf{x}) \end{pmatrix}$$

• Learning problem:

•
$$\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} \sum_{j=1}^{m} \sum_{i=1}^{n_j} \ell(\boldsymbol{y}_{ji}, \boldsymbol{\theta}^{\mathrm{T}} \boldsymbol{\Phi}_j(\boldsymbol{x}_{ji})) + \Omega(\boldsymbol{\theta})$$

• Implement using any learning algorithm by constructing appropriate Φ



Hierarchical Bayesian Inference: Exercise

- Three kids, training set
 - { { (\mathbf{x}_{11}, y_{11}), (\mathbf{x}_{12}, y_{12}), (\mathbf{x}_{13}, y_{13}) }, { (\mathbf{x}_{21}, y_{21}), (\mathbf{x}_{22}, y_{22}) }, { (\mathbf{x}_{31}, y_{31}) }
- Construct the training set that maps this multi-task problem to a regular SVM
- Kid 3 takes substances x. Will he live? How do you use the resulting SVM for inference?
- Kid 4 takes substances x. Can we say anything about whether she will live?





- Graphical model: population-specific and individual factors determine model parameters.
- When $p(\boldsymbol{\mu} \mid \boldsymbol{\Sigma}') = N[0, \boldsymbol{\Sigma}'](\boldsymbol{\mu}), \ p(\mathbf{v}_j \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}) = N[0, \boldsymbol{\Sigma}](\mathbf{v}_j),$ and $\boldsymbol{\theta} = \boldsymbol{\mu} + \mathbf{v}$, then

 $\Phi_{\text{MAP}} = \arg \min_{\boldsymbol{\mu} + \boldsymbol{\nu}} \left(\sum_{j=1}^{m} \sum_{i=1}^{n_j} \log P(y_{ji} \mid \mathbf{x}_{ji}, \boldsymbol{\mu} + \boldsymbol{\nu}) + \sum_{j} \log p(\mathbf{v}_j \mid \boldsymbol{\Sigma}) + \log p(\boldsymbol{\mu} \mid \boldsymbol{\Sigma}') \right)$

For general loss and regularizers,

• $\boldsymbol{\theta}^* = \operatorname{arg\,min}_{\boldsymbol{\mu}+\boldsymbol{\nu}} \sum_{j=1}^m \sum_{i=1}^{n_j} \ell(y_{ji}, y_{\boldsymbol{\mu}+\boldsymbol{\nu}}(\mathbf{x}_{ji})) + \sum_j \Omega(\mathbf{v}_j) + \Omega(\boldsymbol{\mu})$

 Can be mapped to regular learning problem using mapping

•
$$f_{\theta,j}(\mathbf{x}) = \begin{pmatrix} \mathbf{v}_1 \\ \cdots \\ \mathbf{v}_m \\ \mathbf{\mu} \end{pmatrix}^{\mathrm{T}} \underbrace{\begin{pmatrix} \phi(\mathbf{x})\delta(j=1) \\ \cdots \\ \phi(\mathbf{x})\delta(j=m) \\ \phi(\mathbf{x}) \\ \phi(\mathbf{x}) \\ \Phi_j(\mathbf{x}) \end{pmatrix}}_{\boldsymbol{\Phi}_j(\mathbf{x})}$$





Importance Sampling for Covariate Shift

Goal: minimize risk on test distribution

• $R_{\text{test}} = \int \int \ell(y_i, f_{\theta}(\mathbf{x}_i)) P(y_i | \mathbf{x}_i) p(\mathbf{x}_i | \tau_{\text{test}}) d\mathbf{x}_i dy_i$

 Idea: Write as expected value over training distribution by using appropriate weights.

•
$$R_{\text{test}} = \int \int \frac{p(\mathbf{x}_i \mid \tau_{\text{test}})}{p(\mathbf{x}_i \mid \tau_{\text{train}})} \ell(y_i, f_{\theta}(\mathbf{x}_i)) P(y_i \mid \mathbf{x}_i) p(\mathbf{x}_i \mid \tau_{\text{train}}) d\mathbf{x}_i dy_i$$

Regularized empirical risk:

•
$$\hat{R}_{\text{test}} = \sum_{i=1}^{n_{\text{train}}} \frac{p(\mathbf{x}_i \mid \tau_{\text{test}})}{p(\mathbf{x}_i \mid \tau_{\text{train}})} \ell(y_i, f_{\theta}(\mathbf{x}_i)) + \Omega(\theta)$$

Sum over labeled
training sample



Integral over

test distribution

training

Integral over

distribution

Importance Sampling for Covariate Shift

Regularized empirical risk:

•
$$\hat{R}_{\text{test}} = \sum_{i=1}^{n_{\text{train}}} \frac{p(\mathbf{x}_i \mid \tau_{\text{test}})}{p(\mathbf{x}_i \mid \tau_{\text{train}})} \ell(y_i, f_{\theta}(\mathbf{x}_i)) + \Omega(\theta)$$

- Densities $p(\mathbf{x}_i | \tau_{\text{test}})$ and $p(\mathbf{x}_i | \tau_{\text{train}})$ unknown and highdimensional (impractical to estimate).
- Density ratio is just a number for each training data point (should be easier to estimate).
- Direct estimation of optimal reweighting factors.

Regularized empirical risk:

•
$$\hat{R}_{\text{test}} = \sum_{i=1}^{n_{\text{train}}} \frac{p_{i}}{p_{i}} \ell(y_{i}, f_{\theta}(\mathbf{x}_{i})) + \Omega(\boldsymbol{\theta})$$



Density ratio can be rephrased:

$$\frac{p(\mathbf{x}_{i} \mid \tau_{\text{test}})}{p(\mathbf{x}_{i} \mid \tau_{\text{train}})} = \frac{p(\tau_{\text{train}})}{p(\tau_{\text{test}})} \frac{p(\mathbf{x}_{i} \mid \tau_{\text{test}})}{p(\mathbf{x}_{i} \mid \tau_{\text{train}})} \frac{p(\tau_{\text{test}})}{p(\tau_{\text{train}})}$$

$$= \frac{p(\tau_{\text{train}})}{p(\tau_{\text{test}})} \left(1 + \frac{p(\mathbf{x}_{i} \mid \tau_{\text{test}})}{p(\mathbf{x}_{i} \mid \tau_{\text{train}})} \frac{p(\tau_{\text{test}})}{p(\tau_{\text{train}})} - 1 \right)$$

$$= \frac{p(\tau_{\text{train}})}{p(\tau_{\text{test}})} \left(\frac{p(\tau_{\text{train}})p(\mathbf{x} \mid \tau_{\text{train}}) + p(\tau_{\text{test}})p(\mathbf{x} \mid \tau_{\text{test}})}{p(\tau_{\text{train}})p(\mathbf{x} \mid \tau_{\text{train}})} - 1 \right)$$

$$= \frac{p(\tau_{\text{train}})}{p(\tau_{\text{test}})} \left(\frac{p(\tau_{\text{train}} \mid \mathbf{x}) + p(\tau_{\text{test}} \mid \mathbf{x})}{p(\tau_{\text{train}})} - 1 \right)$$

$$= \frac{p(\tau_{\text{train}})}{p(\tau_{\text{test}})} \left(\frac{p(\tau_{\text{train}} \mid \mathbf{x}) - 1}{p(\tau_{\text{train}} \mid \mathbf{x})} - 1 \right)$$

[Bickel, Brückner, Scheffer, Discriminative learning for differing training and test distributions, ICML 2007]

Density ratio can be rephrased:



Density ratio can be estimated directly:

$$\frac{p(\mathbf{x}_i \mid \tau_{\text{test}})}{p(\mathbf{x}_i \mid \tau_{\text{train}})} = \frac{p(\tau_{\text{train}})}{p(\tau_{\text{test}})} \left(\frac{1}{p(\tau_{\text{train}} \mid \mathbf{x}_i)} - 1\right)$$

• Train logistic regression model $p(\tau_{\text{train}} | \mathbf{x}_i, \mathcal{G}) = \frac{1}{1 + e^{\phi(\mathbf{x}_i)^T \mathcal{G} + b}}$ that discriminates training from test data.



1. Train model $p(\tau_{\text{train}} | \mathbf{x}_i, \vartheta) = \frac{1}{1 + e^{\phi(\mathbf{x}_i)^T \vartheta}}$ using training data as positive class and unlabeled test data as negative class.



1. Train model $p(\tau_{\text{train}} | \mathbf{x}_i, \theta) = \frac{1}{1 + e^{\phi(\mathbf{x}_i)^T \theta}}$ using training data as positive class and unlabeled test data as negative class.

2. Infer Weights
$$\frac{p(\mathbf{x}_i \mid \tau_{\text{test}})}{p(\mathbf{x}_i \mid \tau_{\text{train}})} = \frac{p(\tau_{\text{train}})}{p(\tau_{\text{test}})} \left(\frac{1}{p(\tau_{\text{train}} \mid \mathbf{x}_i, \mathcal{G})} - 1\right)$$



1. Train model $p(\tau_{\text{train}} | \mathbf{x}_i, \theta) = \frac{1}{1 + e^{\phi(\mathbf{x}_i)^T \theta}}$ using training data as positive class and unlabeled test data as negative class.

2. Infer Weights
$$\frac{p(\mathbf{x}_i \mid \tau_{\text{test}})}{p(\mathbf{x}_i \mid \tau_{\text{train}})} = \frac{p(\tau_{\text{train}})}{p(\tau_{\text{test}})} \left(\frac{1}{p(\tau_{\text{train}} \mid \mathbf{x}_i, \mathcal{G})} - 1\right)$$

3. Train final classifier on weighted training data.



- The logistic regression of the first step also has a regularization parameter.
- Set parameter using grid search and n-fold cross validation.
- (Label information is whether instance is training or test instance.)
- Hence, possible to tune regularization parameter.



- Toxicity study conducted in country X has yielded labeled data.
- In country Y, physicians prefer different drugs, leading to a different distribution over drug prescriptions.
- How would you learn a toxicity model that works well for country Y?
- Does importance sampling make a difference at all? How does the decision whether or not a combination of substances is toxic depend on the distribution of substances administered?

- Goal: Model that minimizes risk on $p(\mathbf{x}_i | \tau_{\text{test}})$
- Training data governed by $p(\mathbf{x}_i | \tau_{\text{train}})$
- Optimal resampling weights:

•
$$\frac{p(\mathbf{x}_i \mid \tau_{\text{test}})}{p(\mathbf{x}_i \mid \tau_{\text{train}})} = \frac{p(\tau_{\text{train}})}{p(\tau_{\text{test}})} \left(\frac{1}{p(\tau_{\text{train}} \mid \mathbf{x}_i)} - 1\right)$$

- Train logistic regression model for $p(\tau_{train} | \mathbf{x}_i, \vartheta)$ using training data as positive data and unlabeled test data as negative class.
- Calculate weights (above formula) and train classifier on weighted training data.

KLIEP

- Alternative way of estimating density ratio.
- Define weights $w(\mathbf{x}_i)$ such that $w(\mathbf{x}_i) p(\mathbf{x}_i | \tau_{\text{train}}) = p(\mathbf{x}_i | \tau_{\text{test}})$
- Minimize KL-divergence between distributions of test data and weighted training data.

•
$$KL[p(\mathbf{x}_{i} | \tau_{\text{test}}) || w(\mathbf{x}_{i}) p(\mathbf{x}_{i} | \tau_{\text{train}})]$$

$$= \int p(\mathbf{x}_{i} | \tau_{\text{test}}) \log \frac{p(\mathbf{x}_{i} | \tau_{\text{test}})}{w(\mathbf{x}_{i}) p(\mathbf{x}_{i} | \tau_{\text{train}})} d\mathbf{x}_{i}$$

$$= \int p(\mathbf{x}_{i} | \tau_{\text{test}}) \log \frac{p(\mathbf{x}_{i} | \tau_{\text{test}})}{p(\mathbf{x}_{i} | \tau_{\text{train}})} d\mathbf{x}_{i} - \int p(\mathbf{x}_{i} | \tau_{\text{test}}) \log w(\mathbf{x}_{i}) d\mathbf{x}_{i}$$

$$= -\int p(\mathbf{x}_{i} | \tau_{\text{test}}) \log w(\mathbf{x}_{i}) d\mathbf{x}_{i} + const$$

$$\approx -\frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} \log w(\mathbf{x}_{i}) = -\frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} \log \phi(\mathbf{x}_{i})^{\mathrm{T}} \mathcal{G}$$

KLIEP

 Minimize KL-divergence between distributions of test data and weighted training data.

•
$$KL[p(\mathbf{x}_i | \tau_{\text{test}}) || w(\mathbf{x}_i) p(\mathbf{x}_i | \tau_{\text{train}})] \approx -\frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} \log \phi(\mathbf{x}_i)^{\mathrm{T}} \mathcal{G}$$

- Weighted training data $w(\mathbf{x}_i) p(\mathbf{x}_i | \tau_{train}) = p(\mathbf{x}_i | \tau_{test})$ have to be nortmalized, weights have to be positive.
- Optimization problem:

•
$$KL[p(\mathbf{x}_i | \tau_{\text{test}}) || w(\mathbf{x}_i) p(\mathbf{x}_i | \tau_{\text{train}})] \approx$$

 $\mathcal{G}^* = \arg \max_{\mathcal{G}} \sum_{i=1}^{n_{\text{test}}} \log \phi(\mathbf{x}_i)^T \mathcal{G}$
subject to $\sum_{i=1}^{n_{\text{train}}} \log \phi(\mathbf{x}_i)^T \mathcal{G} = n_{\text{train}}$
for all $j : \mathcal{G}_j \ge 0$

[Sugiyama, Suzuki, Nakajima, Kashima, H., von Bünau, Kawanabe. Direct importance estimation for covariate shift adaptation. Annals of the Institute of Statistical Mathematics, vol.60, no.4, 2008.]

Kernel Mean Matching

- Alternative way of estimating density ratio.
- Define weights $w(\mathbf{x}_i)$ such that $w(\mathbf{x}_i) p(\mathbf{x}_i | \tau_{\text{train}}) = p(\mathbf{x}_i | \tau_{\text{test}})$
- Expected mean feature mapping:

• $\mu_{\tau}(\phi(\mathbf{x})) = \int \phi(\mathbf{x}) p(\mathbf{x} \mid \tau) d\mathbf{x}$

- When $k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^{T} \phi(\mathbf{x}')$ is a universal kernel, then there is a one-to-one relationship between $\mu_{\tau}(\phi(\mathbf{x}))$ and $p(\mathbf{x} | \tau)$
- Idea: set weights such that

• $\mu_{\tau_{\text{test}}}(\phi(\mathbf{x})) = \mu_{\tau_{\text{train}}}(w(\mathbf{x})\phi(\mathbf{x}))$
Kernel Mean Matching

Idea: set weights such that

• $\mu_{\tau_{\text{test}}}(\phi(\mathbf{x})) = \mu_{\tau_{\text{train}}}(w(\mathbf{x})\phi(\mathbf{x}))$

- Weights have to be positive and normalized.
- Optimization problem:

•
$$w^* = \operatorname{arg\,min}_{w} \| \frac{1}{n_{\operatorname{train}}} \sum_{i=1}^{n_{\operatorname{train}}} w(\mathbf{x}_i) \phi(\mathbf{x}_i) - \frac{1}{n_{\operatorname{test}}} \sum_{i=1}^{n_{\operatorname{test}}} \phi(\mathbf{x}_i) \|$$

- Optimization problem is convex.
- Regularization: impose upper bound on individual weights and sum over weights.
- Problem: need labeled test data to tune regularizer.

[Huang, Smola, Gretton, Borgwardt, Schölkopf, Correcting sample selection bias by unlabeled data. NIPS 2007.]

Importance Sampling for Domain Adaptation

Minimize risk on test distribution:

•
$$R_{\text{test}} = \int \int \ell(y_i, f_{\theta}(\mathbf{x}_i)) P(y_i | \mathbf{x}_i) p(\mathbf{x}_i | \tau_{\text{test}}) d\mathbf{x}_i dy_i$$

- Rephrase as expected weighted training loss • $R_{\text{test}}(\boldsymbol{\theta}) = \int \int \frac{p(\mathbf{x}_i, y_i | \tau_{\text{test}}, \boldsymbol{\theta}_{\text{train}})}{p(\mathbf{x}_i, y_i | \tau_{\text{train}}, \boldsymbol{\theta}_{\text{train}})} \ell(y_i, f_{\boldsymbol{\theta}}(\mathbf{x}_i)) p(\mathbf{x}_i, y_i | \tau_{\text{train}}, \boldsymbol{\theta}_{\text{train}}) d\mathbf{x}_i dy_i$
- Estimate on training sample:

$$\hat{R}_{\text{test}}(\boldsymbol{\theta}) = \sum_{i=1}^{n_{\text{train}}} \frac{p(\mathbf{x}_i, y_i | \tau_{\text{test}}, \boldsymbol{\theta}_{\text{test}})}{\underbrace{p(\mathbf{x}_i, y_i | \tau_{\text{train}}, \boldsymbol{\theta}_{\text{train}})}_{W(\mathbf{x}_i, y_i)}} \ell(y_i, f_{\boldsymbol{\theta}}(\mathbf{x}_i)) + \Omega(\boldsymbol{\theta})$$



Importance Sampling for Domain Adaptation

Density ratio can be rephrased:



[Bickel, Sawade, Scheffer. Transfer Learning by Distribution Matching for Targeted Advertising. NIPS 2008.]

Direct Estimation of Importance Weights

1. Train model $p(\tau_{\text{train}}, \boldsymbol{\theta}_{\text{train}} | \mathbf{x}_i, y_i, \boldsymbol{\vartheta}) = \frac{1}{1 + e^{\phi(\mathbf{x}_i, y_i)^T \boldsymbol{\vartheta}}}$ using training data as positive class and labeled test data as negative class.

2. Infer Weights
$$\frac{p(\mathbf{x}_i \mid \tau_{\text{test}}, \boldsymbol{\theta}_{\text{test}})}{p(\mathbf{x}_i \mid \tau_{\text{train}}, \boldsymbol{\theta}_{\text{train}})} = \frac{p(\tau_{\text{train}}, \boldsymbol{\theta}_{\text{train}})}{p(\tau_{\text{test}}, \boldsymbol{\theta}_{\text{test}})} \left(\frac{1}{p(\tau_{\text{train}}, \boldsymbol{\theta}_{\text{train}} \mid \mathbf{x}_i, y_i, \boldsymbol{\beta})} - 1\right)$$

3. Train final classifier on weighted training data.



Estimation of Importance Weights



- Importance weights match weighted training data to distribution of test data.
- Learn logistic regression model $p(\tau_{train} | \mathbf{x}_i, \boldsymbol{\vartheta})$ and let

$$w(\mathbf{x}_i) = \frac{p(\mathbf{x}_i \mid \tau_{\text{test}})}{p(\mathbf{x}_i \mid \tau_{\text{train}})} = \frac{p(\tau_{\text{train}})}{p(\tau_{\text{test}})} \left(\frac{1}{p(\tau_{\text{train}} \mid \mathbf{x}_i, \mathcal{G})} - 1\right)$$

- Alternatively, minimize KL divergence between weighted training and test data
- Alternatively, minimize distance between mean of weighted training and test data
- Applies to domain adaptation too (some labeled data from target distribution required).

University of Potsdam Dept. of Computer Science



Adversarial Learning

Tobias Scheffer

Adersarial Learning



Adersarial Learning



Overview

- Learning with Invariances
 - Make result of learning process invariant to specific types of transformations Φ(x, y)
- Minimax probability machine
- Game Theory Basics
- Game-Theoretic Learning Models

Robust Learning Models

- Standard learning assumes training data drawn iid from distribution at application time.
- Optimistic, if adversary tampers with distribution.
- Robust learners minimize maximum risk for any distribution over a certain class of distributions.
- Underlying assumption: adversary can choose distribution within certain class, and will try to inflict greatest possible damage for learner.

Learning with Invariances

- Training set is drawn according to $p(\mathbf{x}_i, y_i | \tau_{\text{train}}, \boldsymbol{\theta}_{\text{train}})$
- Adversary gets to exercise transformation on distribution: τ_{test} , $\theta_{test} = \Phi(\tau_{test}, \theta_{test})$
- Test set is drawn according to $p(\mathbf{x}_i, y_i | \tau_{\text{test}}, \boldsymbol{\theta}_{\text{test}})$
- Equivalent to drawing according to $p(\mathbf{x}_i, y_i | \tau_{\text{train}}, \boldsymbol{\theta}_{\text{train}})$ and transforming instances: $\mathbf{x}', y' = \Phi(\mathbf{x}, y)$
- Goal: Minimize regularized empirical risk:

•
$$\hat{R} = \sum_{i=1}^{n_{\text{test}}} \ell(f_{\theta}(\mathbf{x}_i), y_i) + \Omega(\boldsymbol{\theta})$$

Heuristic SVM Learning with Invariances

Dual SVM decision function

•
$$f_{\boldsymbol{\beta}}(\mathbf{x}) = \sum_{i=1}^{n_{\text{train}}} \beta_i y_i \phi(\mathbf{x}_i)^T \phi(\mathbf{x})$$

Is determined by support vectors

• $S = {\mathbf{x}_i : \beta_i > 0}$

- Heuristic invariant learning algorithm:
 - Learn SVM on training data L
 - If set of transformations is finite: add all possible $\{\mathbf{x}', y' = \Phi(\mathbf{x}, y) : \mathbf{x}, y \in L \cap S, \Phi \in \Phi\}$ to training set
 - Otherwise: draw some $\Phi \in \Phi$ and $(\mathbf{x}, y) \in L \cap S$ and add $\mathbf{x}', y' = \Phi(\mathbf{x}, y)$ to *L*.
 - Learn from enhanced sample.

Reminder: SVM-Struct

- Hinge-loss for binary SVM:
 - $\ell(y, f_{\theta}(\mathbf{x})) = \max(0, 1 yf_{\theta}(\mathbf{x}))$
- Loss for multi-class SVM and SVM-Struct

 $\bullet \quad \ell(y, f_{\theta}(\mathbf{x}, \cdot)) = \max_{y' \neq y} (0, 1 - f_{\theta}(\mathbf{x}, y) + f_{\theta}(\mathbf{x}, y'))$

Corresponding constraint optimization problem

•
$$\min_{\boldsymbol{\theta},\boldsymbol{\xi}} \lambda \sum_{i=1}^{n} \boldsymbol{\xi}_{i} + \frac{1}{2} \sum_{i=1}^{k} \boldsymbol{\theta}^{i\mathrm{T}} \boldsymbol{\theta}^{i}$$

subject to the constraints:

$$\forall y \neq y_i : f_{\theta}(\mathbf{x}_i, y_i) \ge f_{\theta}(\mathbf{x}_i, y) + 1 - \xi_i$$

and $\xi_i \ge 0$.

SVM Learning with Invariances

Define loss



[Teo, Globerson, Roweis, Smola, Convex learning with invariances, NIPS 2007]

SVM Learning with Invariances

Define loss

 $\ell(y, \mathbf{x}, f_{\theta}) = \max_{\Phi \in \Phi, y' \neq y} (0, 1 - f_{\theta}(\Phi(\mathbf{x}), y) + f_{\theta}(\Phi(\mathbf{x}), y'))$

- Loss is a convex upper bound on zero-one loss.
- Constrained optimization problem • $\min_{\alpha \in \mathcal{X}} \lambda \sum_{i=1}^{n} \mathcal{E}_{i} + \frac{1}{2} \sum_{i=1}^{k} \mathbf{\theta}^{iT} \mathbf{\theta}^{i}$

$$\lim_{\boldsymbol{\theta},\boldsymbol{\xi}} \lambda \sum_{i=1}^{k} \boldsymbol{\zeta}_i + \frac{1}{2} \sum_{i=1}^{k} \boldsymbol{\theta}$$

subject to the constraints:

$$\forall y' \neq y_i \forall \Phi \in \mathbf{\Phi} : f_{\theta}(\Phi(\mathbf{x}_i), y_i) \ge f_{\theta}(\Phi(\mathbf{x}_i), y') + 1 - \xi_i$$

and $\xi_i \ge 0$.

SVM Learning with Invariances

Define loss

 $\ell(y, \mathbf{x}, f_{\theta}) = \max_{\Phi \in \Phi, y' \neq y} (0, 1 - f_{\theta}(\Phi(\mathbf{x}), y) + f_{\theta}(\Phi(\mathbf{x}), y'))$

- Constrained optimization problem
 - $\min_{\boldsymbol{\theta},\boldsymbol{\xi}} \lambda \sum_{i=1}^{n} \boldsymbol{\xi}_{i} + \frac{1}{2} \sum_{i=1}^{k} \boldsymbol{\theta}^{iT} \boldsymbol{\theta}^{i}$ subject to the constraints: $\forall i \forall y' \neq y_{i} \forall \boldsymbol{\Phi} \in \boldsymbol{\Phi} : f_{\boldsymbol{\theta}}(\boldsymbol{\Phi}(\mathbf{x}_{i}), y_{i}) \geq f_{\boldsymbol{\theta}}(\boldsymbol{\Phi}(\mathbf{x}_{i}), y') + 1 - \boldsymbol{\xi}_{i}$ and $\boldsymbol{\xi}_{i} \geq 0$.
- SVM-Struct working set algorithm:
 - Iterate over examples, find the one y' and transformation that violates margin the most:

 $y'^*, \Phi^* = \arg \max_{y', \Phi} 1 - f_{\theta}(\Phi(\mathbf{x}_i), y_i) + f_{\theta}(\Phi(\mathbf{x}_i), y')$

Add to working set, solve OP, reiterate.

Minimax Probability Machine

- Let $\mathbf{P}_{\mu,\Sigma}$ be the set of *all* distributions with mean μ and covariance Σ .
- Distribution of training data:
 - Positive class has mean μ_{+1} , covariance matrix Σ_{+1} .
 - Negative class has mean μ_{-1} , covariance matrix Σ_{-1} .
- At application time, adversary can choose *any* input distribution $p(\mathbf{x},+1) \in \mathbf{P}_{\mu_{+1},\Sigma_{+1}}, p(\mathbf{x},-1) \in \mathbf{P}_{\mu_{-1},\Sigma_{-1}}$
- What is the right optimization criterion for learner?
 - Based on zero-one loss

Minimax Probability Machine

- Let $\mathbf{P}_{\mu,\Sigma}$ be the set of *all* distributions with mean μ and covariance Σ .
- Distribution of training data:
 - Positive class has mean μ_{+1} , covariance matrix Σ_{+1} .
 - Negative class has mean μ_{-1} , covariance matrix Σ_{-1} .
- At application time, adversary can choose *any* input distribution $p(\mathbf{x},+1) \in \mathbf{P}_{\mu_{+1},\Sigma_{+1}}, p(\mathbf{x},-1) \in \mathbf{P}_{\mu_{-1},\Sigma_{-1}}$
- Optimization criterion for learner:

•
$$\max_{p(\mathbf{x},+1)\in \mathbf{P}_{\mu_{+1},\Sigma_{+1}}, p(\mathbf{x},-1)\in \mathbf{P}_{\mu_{-1},\Sigma_{-1}}} \sum_{y} \int \ell_{0/1}(f_{\theta}(\mathbf{x}), y) p(\mathbf{x}, y) d\mathbf{x}$$

[Lanckriet et al., Minimax Probability Machine, NIPS 2002]

Minimax Probability Machine

- Let $\mathbf{P}_{\mu,\Sigma}$ be the set of *all* distributions with mean μ and covariance Σ .
- Distribution of training data:
 - Positive class has mean μ_{+1} , covariance matrix Σ_{+1} .
 - Negative class has mean μ_{-1} , covariance matrix Σ_{-1} .
- At application time, adversary can choose *any* input distribution $p(\mathbf{x},+1) \in \mathbf{P}_{\mu_{+1},\Sigma_{+1}}, p(\mathbf{x},-1) \in \mathbf{P}_{\mu_{-1},\Sigma_{-1}}$
- Optimization criterion for learner:

• $\max_{p(\mathbf{x},+1)\in \mathbf{P}_{\mu_{+1},\Sigma_{+1}}, p(\mathbf{x},-1)\in \mathbf{P}_{\mu_{-1},\Sigma_{-1}}} \sum_{y} \int \ell_{0/1}(f_{\theta}(\mathbf{x}), y) p(\mathbf{x}, y) d\mathbf{x}$

Can be rewritten as convex optimization problem.

[Lanckriet et al., Minimax Probability Machine, NIPS 2002]

Invariant SVMs



- Family of methods for making SVM more robust against transformations of instances by adversary.
- Simplistic approach: add transformed support vectors to training set.
- More sophisticated: minimize, for all examples, maximal loss over transformation space.
- Minimax probability machine: minimize maximal loss over all input distributions with observed mean value and variance.

- Game defined in terms of
 - Players P_1, \dots, P_n
 - Action spaces A_1, \dots, A_n
 - Cost functions C_1, \dots, C_n with $C_i : A_1 \times \dots \times A_n \to \mathbb{R}$
 - Cost functions are interleaved optimization problems.
- Zero-sum game: two players, $C_1(A_1, A_2) + C_2(A_1, A_2) = 0$
 - One player's loss is other player's gain.
- Non-cooperative game:
 - Players cannot exchange messages

- Static game:
 - Players act simultaneously, not knowing their opponents' move.
 - Learner and adversary act simultaneously.
- Dynamic, extensive-form game:
 - Some ordering defined over moves, information about previous actions.
 - Stackelberg competition: learner committs first to model, then adversary transforms distribution.

- Complete information:
 - Each player know other players' cost functions.
- Incomplete information:
 - At least one player is uncertain about opponents' cost functions.
 - Bayesian players infer expected cost functions over their beliefs.

Game Theory: Zero Sum Games

- Security level for player P_1
 - $\bar{C}_1 = \min_{a_1 \in A_1} \max_{a_2 \in A_2} C_1(a_1, a_2)$
- Security level for player P₂

•
$$\bar{C}_2 = \min_{a_2 \in A_2} \max_{a_1 \in A_1} C_2(a_1, a_2) = -\max_{a_2 \in A_2} \min_{a_1 \in A_1} C_1(a_1, a_2)$$

Example

•
$$\bar{C}_1 = 2$$
, row 3

• $\bar{C}_2 = 1$, column 1

		P2	
P1	(4, -4)	(-1, 1)	(<mark>-2</mark> , 2)
	(-1, 1)	(- <mark>2</mark> , <mark>2</mark>)	<mark>(3</mark> , -3)
	(1, -1)	(<mark>2</mark> , -2)	(1, -1)

Game Theory: Zero Sum Games

- Minimax strategy for player P_1
 - $\bar{a}_1 = \operatorname*{argmin}_{a_1 \in A_1} \max_{a_2 \in A_2} C_1(a_1, a_2)$
- Minimax strategy for player P₂
 - $\bar{a}_2 = \underset{a_2 \in A_2}{\operatorname{argmin}} \max_{a_1 \in A_1} C_2(a_1, a_2)$
- Example
 - $\bar{a}_1 = 3, \bar{a}_2 = 1$
 - $C_1(\overline{a}_1, \overline{a}_2) = 1$
- It holds that
 - $\overline{C}_1 \ge C_1(\overline{a}_1, \overline{a}_2) \ge -\overline{C}_2$
 - No guaranteed stability

		P2	
P1	(4, -4)	(-1, 1)	(<mark>-2</mark> , 2)
	(-1, 1)	(- <mark>2</mark> , <mark>2</mark>)	<mark>(3</mark> , -3)
	(1, -1)	(<mark>2</mark> , -2)	(1, -1)

Game Theory: Non-Zero Sum Games

- Generalization of zero-sum games
 - $C_1 = -C_2$ no longer prerequisite
- Minimax strategy
 - Still guarantees least maximal costs
 - Not well motivated if opponent wants to minimize own costs and costs not directly antagonistic.

- Equilibrium (stable points)
 - (a_1^*, a_2^*) is Nash equilibrium if

*
$$a_1^* = \operatorname*{argmin}_{a_1 \in A_1} C_1(a_1, a_2^*)$$

★
$$a_2^* = \underset{a_2 \in A_2}{\operatorname{argmin}} C_2(a_1^*, a_2)$$





- Equilibrium (stable points)
 - (a_1^*, a_2^*) is Nash equilibrium if

*
$$a_1^* = \operatorname*{argmin}_{a_1 \in A_1} C_1(a_1, a_2^*)$$

★
$$a_2^* = \operatorname*{argmin}_{a_2 \in A_2} C_2(a_1^*, a_2)$$

 In a Nash equilibrium, each player reacts optimally to their opponent's move.

- Equilibrium (stable points)
 - (a_1^*, a_2^*) is Nash equilibrium if

*
$$a_1^* = \operatorname*{argmin}_{a_1 \in A_1} C_1(a_1, a_2^*)$$

★
$$a_2^* = \operatorname*{argmin}_{a_2 \in A_2} C_2(a_1^*, a_2)$$

- In a Nash equilibrium, each player reacts optimally to their opponent's move.
- Simple algorithm: imagine any move; infer opponent's optimal reaction; infer own optimal reaction; reiterate until fixed point is reached.
- Fixed point is Nach equilibrium.

- When does it make sense to acto according to Nash equilibrium?
- Minimax is overly pessimistic:
 - Adversaries do not want to maximize your costs but want to minimize their own costs.
- When adversary will act according to Nash equilibrium, using this equilibrium too is optimal.
 - Only makes sense to assume if equilibrium exists.
- When players act according to different equilibria, then outcome can be arbitrarily bad for all players.
 - Acting according to Nash equilibrium makes most sense if equilibrium is unique.

Game Theory: Extensive-Form Games

- So far static games: players act simultaneously
 - No information about adversary's action exploitable.
- Extensive form games: players act in defined order.
 - Game takes the for of a tree.
- Simplest form: Stackelberg competition
 - First player acts first
 - Then, adversary acts (knowing the first player's move).
 - Then, costs are inferred.

Game Theory: Extensive-Form Games

- Simplest form: Stackelberg competition
 - First player acts first
 - Then, adversary acts (knowing the first player's move).
- Adversary solves simple minimization problem:

•
$$a_2^* = \underset{a_2 \in A_2}{\operatorname{argmin}} \quad V_2(a_1^*, a_2).$$

First player has to account for adversary:

•
$$a_1^* = \underset{a_1 \in A_1}{\operatorname{argmin}} \max_{a_2 \in A_2} V_1(a_1, a_2)$$



- Static (simultaneous actions) vs. extensive-form.
- Non-cooperative (no message passing)
- Zero-sum (antagonistic costs) vs non-zero-sum
- Zero-sum games: minimax strategies minimizes worst-case (over adversarial action space) costs.
- Nash equilibrium: deviating unilaterally increases costs for either player.
- Non-zero-sum games:
 - Nash equilibrium optimal if adversary will act according to the same equilibrium.
 - Does the game have an equilibrium? Is it unique?
 Will the adversary be rational enough to infer it?

Game-Theoretic Learning Models



Game-Theoretic Learning Models

- Action of the data generator (adversary)
 - Theoretically: transform input distribution $p(\mathbf{x}, y) \rightarrow \overline{p}(\mathbf{x}, y)$
 - Empirically: transform data $D = \{(\mathbf{x}_i, y_i)\} \rightarrow D' = \{(\mathbf{x}'_i, y'_i)\}$
- Action of the learner:

• Choose model parameters $f_{\theta}(\mathbf{x}) = \theta^{\mathrm{T}} \phi(\mathbf{x})$

Game-Theoretic Learning Models

Empirical costs of the data generator (adversary)

•
$$C_{+1}(\boldsymbol{\theta}, D') = \frac{1}{n} \sum_{i=1}^{n} \ell_{+1}(\boldsymbol{\theta}^{\mathrm{T}} \boldsymbol{\phi}(\mathbf{x}'_{i}), y_{i}) + \rho_{+1} \Omega_{+1}(D, D')$$

Adversary's loss function Transformation costs

Empirical costs of the learner

•
$$C_{-1}(\boldsymbol{\theta}, D') = \frac{1}{n} \sum_{i=1}^{n} \ell_{-1}(\boldsymbol{\theta}^{\mathrm{T}} \boldsymbol{\phi}(\mathbf{x}'_{i}), y_{i}) + \rho_{-1} \Omega_{-1}(\boldsymbol{\theta})$$

Learner's loss function Regularizer
(adversary) D, D')

Game-Theoretic Learning Models

• Empirical costs of the data generator (adversary) • $C_{+1}(\mathbf{\theta}, D') = \frac{1}{n} \sum_{i=1}^{n} \ell_{+1}(\mathbf{\theta}^{\mathrm{T}} \phi(\mathbf{x}'_{i}), y_{i}) + \rho_{+1} \Omega_{+1}(D, D')$



• Empirical costs of the data generator (adversary) • $C_{+1}(\mathbf{\theta}, D') = \frac{1}{n} \sum_{i=1}^{n} \ell_{+1}(\mathbf{\theta}^{\mathrm{T}} \phi(\mathbf{x}'_{i}), y_{i}) + \rho_{+1} \Omega_{+1}(D, D')$

- Empirical costs of the learner • $C_{-1}(\boldsymbol{\theta}, D') = \frac{1}{n} \sum_{i=1}^{n} \ell_{-1}(\boldsymbol{\theta}^{\mathrm{T}} \boldsymbol{\phi}(\mathbf{x}'_{i}), y_{i}) + \rho_{-1} \Omega_{-1}(\boldsymbol{\theta})$
- Empirical costs cannot be minimized directly because opponent's action is not known.

- Worst-case solution for learner:
 - $\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} \max_{D'} C(\boldsymbol{\theta}, D')$
 - Iearning with invariances
 - Minimizes $C_{-1}(\theta, D')$ if D=D'.
 - Adversary will do nothing; iid assumption
- Unrealistically optimistic.

- Naive solution for learner:
 - $\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} C(\boldsymbol{\theta}, D) = \arg\min_{\boldsymbol{\theta}} \frac{1}{n} \sum_{i=1}^n \ell_{-1}(\boldsymbol{\theta}^T \boldsymbol{\phi}(\mathbf{x}_i), y_i) + \rho_{-1} \Omega_{-1}(\boldsymbol{\theta})$ • = learn from training data
- Minimizes C₋₁(θ, D') if adversary will try to maximize learner's costs (as opposed to minimizing own costs).
- Too pesstimic.

Nash Equilibrial Prediction Models



- Learner's action is optimal reaction to adversary's action which is optimal reaction...
- Fixed point of *optimal reaction* is a Nash equilibrium.
- Probably not quite realistic for most applications, but perhaps more realistic than previous solutions.

Nash-Equilibrial Prediction Models

- Theorem: Nash equilibrium for prediction games exist and is unique if
 - Both loss functions have identical curvature, and
 - Regularization parameters are sufficiently large

[Brückner, Kanzow, Scheffer. Static Prediction Games for Adversarial Learning Problems. JMLR, 13:2617-2654, 2012]

Tobias Scheffer

Nash-Equilibrial Prediction Models

Nach equilibrium meets:

•
$$C_{-1}(\boldsymbol{\theta}^*, D^*) = \min_{\mathbf{w}} C_{-1}(\boldsymbol{\theta}, D^*)$$

• $C_{+1}(\boldsymbol{\theta}^*, D^*) = \min_{D'} C_{+1}(\boldsymbol{\theta}^*, D')$

- Nikaido-Isoda function: summed improvement that players can enjoy by changing their action
 - $\vartheta(\boldsymbol{\theta}, D', \dot{\boldsymbol{\theta}}, \dot{D}') = C_{-1}(\boldsymbol{\theta}, D') C_{-1}(\dot{\boldsymbol{\theta}}, D') + C_{-1}(\dot{\boldsymbol{\theta}}, D') C_{-1}(\dot{\boldsymbol{\theta}}, \dot{D}')$
 - Has to be zero at Nash equilibrium
- Theorem: descent direction for Nikaido-Isoda function exists and can be calculated.
 - Leads to Nash-SVM that finds Nash equilibrial prediction model.

Exercise

Loss / cost functions of spammer and email service provider?

Regularizers for learner and spammer?

Possible transformations of spammer?

Exercise

Loss / cost functions of spammer and email service provider?

• $\ell_{learner}(y, y') = \ell_{0/1}(y, y') \begin{cases} 1 \text{ if } yy' < 0\\ 0 \text{ otherwise} \end{cases}$

•
$$\ell_{spammer}(y, y') = \ell_{0/1}(y, -1)$$

Regularizers for learner and spammer?

•
$$\Omega_{learner}(\boldsymbol{\theta}) = \boldsymbol{\theta}^{\mathrm{T}}\boldsymbol{\theta}$$

•
$$\Omega_{spammer}(D, \dot{D}) = |D - \dot{D}|_F^2 = \sum_{i} \sum_{j} (D_{ij} - \dot{D}_{ij})^2$$

- Possible transformations of spammer?
 - Can transform any training matrix into any other matrix.

Nash-Equilibrial Prediction Models



- Learner and adversary act simultaneously.
- Both have distinct, conflicting but not necessarily antagonistic cost functions
- Unique Nash equilibrium exist if curvature of loss function matches and if regularizers are sufficiently large
- Equilibrium point is optimal strategy for learner if adversary also acts according to it.
- Nash-SVM finds equilibrium point
- Empirically, more robust for spam filtering than SVM and SVM with inveriances.

Stackelberg Prediction Games

- Two-stage game:
 - Learner fixes parameters of predictive model
 - Adversary then transforms input distribution.
- Learner chooses optimal parameters under the assumption that adversary will minimize his costs.
- Bilevel optimization problem:

•
$$\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} \max_{D^*} C_{-1}(\boldsymbol{\theta}, D^*)$$

s.t. $C_{+1}(\boldsymbol{\theta}, D^*) = \min_{D'} C_{+1}(\boldsymbol{\theta}, D')$

Stackelberg Prediction Games

- Unique Stackelberg equilibrium always exists.
- Can be rephrased as constraint optimization problem (KKT conditions of lower-level OP).
- Locally optimal solution by SQP solver or interiorpoints method.

[Brückner, Scheffer. Stackelberg games for adversarial prediction problems. ACM KDD 2011.]



- Interests of learner and adversary modeled as loss functions.
- Transformation costs and prior on model parameteres modeled as regularizers.
- Interleaved optimization problems constitute a game between learner and adversary.
- Ragularizers and instance-specific costs make costs non-antagonistic.
 - Minimax solution too pessimistic
- Nash-SVM chooses predictive model according to Nash equilibrium.
- Stackelberg-SVM: learner acts first, then adversary.